

Investigating Drivers of Repeated Behaviors in Field Data

Thesis by
Anastasia Buyalskaya

In Partial Fulfillment of the Requirements for the
Degree of
Doctor of Philosophy

The logo for the California Institute of Technology (Caltech), featuring the word "Caltech" in a bold, orange, sans-serif font.

CALIFORNIA INSTITUTE OF TECHNOLOGY
Pasadena, California

2021
Defended May 14, 2021

© 2021

Anastasia Buyalskaya
ORCID: 0000-0002-1848-1661

All rights reserved

ACKNOWLEDGEMENTS

First and foremost, I am very grateful to my advisor Colin Camerer. I cannot think of a better person to guide me through the process of asking big scientific questions and answering them using a range of innovative methods. A lot of my growth as a researcher is a function of the knowledge and advice you have imparted. I have particularly appreciated how thoughtful you have been about giving me research opportunities which met my desire for breadth across topics and methods. Few social scientists can say that their PhD included running a field experiment with a vending machine, A/B tests with a dating app company, lab studies of curiosity using eye trackers, and helping out at a haunted house.

I am very grateful to my wonderful thesis committee. Thank you Marina Agranov for taking an interest in me and offering to run an experiment together all those years ago. You taught me how to think like an experimental economist, and every graduate student interested in experimental work would benefit from spending time with you. Thank you Matthew Shum for being consistently nice and supportive through what have sometimes been difficult times. I am delighted that one of the chapters in this thesis is one we wrote together, and I look back on our organized and collaborative research process as a role model of how co-authorship can work. Finally, I thank the world expert on human habit Wendy Wood, for joining my thesis committee and providing thoughtful feedback on our efforts to study habitual behavior in the field. I know you have plenty of USC students vying for your time and attention, and I appreciate you taking an interest in my work.

I want to thank the broader SDN faculty. John O'Doherty, thank you for teaching me enough about fMRI that I can speak about it intelligently without ever having used it. Antonio Rangel, thank you for being the biggest fan of my academic presentations and communication style. Dean Mobbs, thank you for sending me to Santa Monica my first term at Caltech to ask strangers if they wanted to swim with sharks - I learned very quickly that research does not have to be constrained to the lab. Ralph Adolphs, thank you for being a role model on how to create a sense of community in a lab and reminding us all that life is more than research.

It can be hard to do research alone, and I am grateful to a small but mighty lab for their ongoing support and motivation. Thank you "team elder" Devdeepa Bose for making sure we remembered to laugh and have lunch. Thank you Marcos

Gallo for helping me advance the “golden age” and encouraging us to spruce up our windowless office. Thank you Virginia Fedrigo for being endlessly helpful, especially when it came to installing and stocking a vending machine, which I know was not the most intellectually stimulating part of our days.

It has been an honor to work with fabulous co-authors in addition to those on my committee. Thank you to the PCS team - the prolific Katherine Milkman and Angela Duckworth, and the superhuman Hung Ho - I would love for all our paths to cross again. Thank you as well to Megan Hunter, with whom we started three different research projects and I have faith we will be able to actually finish one someday!

I am grateful to several granting agencies which made my research possible. Thank you to the Linde and Chen Institutes at Caltech. Thank you especially to the Sloan Foundation for generous support of our vending machine work (via grant G2018 11259 to PI Colin Camerer), as well as Peter Landry, who had the original idea of using a vending machine to run experiments studying habit formation.

Thank you to a number of extraordinary mentees who helped with various research projects: Jingjing Li, for discovering eye tracking alongside me and helping me run my first set of lab experiments. Anthony Kukavica, for spending a whole summer meticulously linking weather and gym data and investigating relevant empirical models. Tads Ciecierski-Holmes, for remaining very optimistically British while managing a difficult dataset and several rote research tasks. Maya Srikanth, for being a very talented computer scientist and enthusiastically helping me wrangle data on food purchases from our canteen. I am eager to see the great things that each of you will accomplish.

Most importantly, I thank my partner Alexandre Metz, without whose support I would have never started nor finished this program. Thank you for agreeing to move across the world with me just a few months after we moved in together in London. I’m looking forward to our next adventure.

ABSTRACT

This dissertation investigates the influences on frequently repeated human behaviors (e.g. eating, exercising, washing hands) using empirical tests on field data. While some of the phenomena discussed have been studied in lab settings (e.g., self-regulation failures, insensitivity to reward devaluation), these studies present some of the first tests of these behavioral phenomena in the field. This dissertation also assembles a number of methodologies which can be used to study individual-level field data, informed by an interdisciplinary perspective on social and decision science research.

The first chapter uses field data to study spillovers across behavioral domains, namely exercise and food choice. This work joins a small group of papers which document field evidence related to domain spillovers and failures of self-regulation. Most of the existing research on self-regulation has been conducted in controlled laboratory settings, where participants are either asked to imagine making hypothetical restrained choices or exert effort on a laboratory task as a proxy for making a restrained choice. As is the critique of many lab studies without direct field equivalents however, it is debatable whether the self-regulation behaviors observed in survey and laboratory settings necessarily generalize to the field. We fill this gap by looking at how natural (rather than incentivized) changes in exercise systematically affect food choice, thus empirically identifying spillovers across two behavioral domains in field data. We find that, even after controlling for individual fixed effects, there is a robust effect of morning exercise on the healthiness of a lunch choice. We complement the analysis of field data with surveys to better understand the mechanism driving this result.

The second chapter presents a novel methodology for identifying behaviors that are highly and predictably context-sensitive, and thus candidates for being habitual. While there is a large body of laboratory research documenting the mechanisms underlying well-developed habits in animals and humans, there is much less field research on how human habits naturally develop over time. Using two large datasets on gym attendance and handwashing behavior, we use machine learning to statistically classify when choices are predicted by an identifiable set of context variables. This technique generates a person-specific measure of behavioral predictability, which can then be used to study individual differences in predictability and speed of habit formation. This allows us to establish two important discoveries. First, the sets of context cues that are predictive of individual-level behavior are different for

different people. Specifically, while historical behavior is an important universal predictor, other context variables such as day of the week or month of the year have more heterogeneous effects. Second, contrary to common wisdom, there is no “magic number” for how long it takes to form a habit. Instead, the speed of habit formation appears to vary significantly, both between behavioral domains and between individuals within domains.

The third chapter uses a novel methodology to run a field experiment testing the effect of a price promotion on consumer behavior. The goal of this “pilot study” is to credibly dissociate predictions made by brand loyalty/habit formation from reference-dependence theories. A customizable vending machine serves as a “mini-retailer,” allowing for full control of price promotion details in an ecologically valid setting. The vending machine allows controlling for stockpiling behavior, an important concern for empirical work analyzing price promotions in the marketing literature. Analysis of the data collected from this pilot study suggests that price promotions increase the sales of both discounted and non-discounted items, as well as the total number of unique customers making purchases. Furthermore, in line with the loss leader hypothesis, more items are purchased during the sale period overall.

TABLE OF CONTENTS

Acknowledgements	iii
Abstract	v
Table of Contents	vii
List of Illustrations	ix
List of Tables	x
Chapter I: Introduction	1
1.1 The Benefits of Field Data as a Compliment to Lab and Survey Measures	2
1.2 Obtaining Lab-Level Control in a Field Experiment	6
1.3 Embracing an Interdisciplinary Approach	13
Chapter II: Identifying Self-Regulation Failures in the Field	20
Abstract	21
2.1 Introduction	22
2.2 Theoretical Framework	25
2.3 Data	27
2.4 Behavioral Results	31
2.5 Mechanism	39
2.6 Discussion	42
Appendix A - Robustness Checks	45
Appendix B - Additional Results	49
Appendix C - Additional Data Sources	50
Appendix D - Survey Format and Results	51
Chapter III: Predicting Context-Sensitivity of Behavior in Field Data	55
Abstract	56
3.1 Introduction	57
3.2 Study 1: Gym Attendance	59
3.3 Study 2: Hand Washing among Hospital Workers	67
3.4 Discussion	73
Appendix A - Literature Review	75
Appendix B - Dataset Descriptions	92
Appendix C - Analysis Details	94
Appendix D - Field Tests of Insensitivity to Reward Devaluation	99
Appendix E - Demographic Predictors of AUC	107
Chapter IV: Using a Vending Machine “Retailer” to Study Repeat Purchases in Consumer Behavior	115
Abstract	116
4.1 Introduction	117
4.2 Literature Review	121
4.3 Experimental Design	128

4.4 Theoretical Predictions	133
4.5 Data	137
4.6 Results	139
4.7 Discussion	145
Appendix A - Vending Machine Details	150
Appendix B - Pre-Testing Survey	151
Appendix C - Vending Machine Announcements	154
Appendix D - Binary Choice Model	154
Appendix E - Difference in Difference Regression	156
Appendix F - Substitute Pair Analysis	157
Bibliography	162

LIST OF ILLUSTRATIONS

<i>Number</i>	<i>Page</i>
3.1 Two examples of gym members: one with high and one with low predictability (measured by AUC).	62
3.2 Development of habit formation of gym attendance.	64
3.3 Development of habit formation of hospital hand washing.	72
3.4 Relationship between holdout AUC and outcome frequency	96
3.5 Reward devaluation sensitivity in gym attendance	105
3.6 Reward devaluation sensitivity in handwashing	106
4.1 The front of the customizable vending machine.	131
4.2 The inside of the customizable vending machine.	131
4.3 Illustrative willingness to pay distribution	135
4.4 Sales of Coke and Izzе through the study period	139

LIST OF TABLES

<i>Number</i>	<i>Page</i>
2.1 Lunch transactions summary by demographics	28
2.2 Ten least healthy products	29
2.3 Ten most healthy products	29
2.4 Summary statistics for key variables	30
2.5 Fixed effects regressions	33
2.6 Healthiness of lunch and post-lunch gym use	34
2.7 Results using sweet instead of unhealthy indicators	35
2.8 Truncated results using only non-sweet foods to define healthiness of lunch	35
2.9 Baseline regressions: using IV's for <i>earlygym</i> and <i>lategym</i>	36
2.10 Food choice differences between early and late gym-goers	38
2.11 Timing of lunch and gym use	39
2.12 Including timelag between morning gym use and lunch	40
2.13 Breakdown by gender	45
2.14 Breakdown by position type	46
2.15 Breakdown by midterm/final weeks	48
2.16 First stage regressions	49
3.1 Summary statistics on $n = 30,110$ gym goers	60
3.2 Context predictors of gym attendance	63
3.3 Demographic predictors of AUC	66
3.4 Summary statistics on $n = 3,124$ hospital caregivers' hand washing	68
3.5 Context predictors of hospital hand washing	70
3.6 Time to habit formation	98
3.7 Summary statistics on gym goers by quality of asymptotic fit	98
3.8 Summary statistics on hospital caregivers by quality of asymptotic fit	99
3.9 Summary statistics on weather data	101
3.10 Correlation matrix of continuous variables	108
4.1 Top ten products, by net sales	132
4.2 Theoretical predictions - Sale period	133
4.3 Theoretical predictions - Post-Sale period	135
4.4 Summary statistics on vending machine purchases	138
4.5 Weekly Transaction Summary by Treatment Period	140
4.6 Bundle-Level Summary by Treatment Period	141

4.7	Purchases During Sale and Post-Sale Periods	142
4.8	Purchases During Sale and Post-Sale Periods	143
4.9	Weekly Customer Summary by Treatment Period	144
4.10	Customer Composition During Sale and Post-Sale Periods	145
4.11	Pre-testing survey results	153
4.12	Binary Logit Regression Results	155
4.13	DID Regressions - Purchases Post-Sale	156
4.14	T-Tests of Mean Daily Purchases in Substitute Pairs	157

Chapter 1

INTRODUCTION

1.1 The Benefits of Field Data as a Compliment to Lab and Survey Measures

The pros and cons of data collected from laboratory experiments

Experiments are a vital part of the methodological toolbox of economics and other social science disciplines given they represent an important instrument of causal inference. The highly controlled setting of a laboratory, where most economics experiments take place, allows for designing a test of whether a variable of interest can *cause* an outcome of interest. The experimental method is sometimes uniquely positioned to answer the research question at hand. To quote Plott and Smith, 2008, “nature may never create a situation that clearly separates the predictions of competing models or may never create a situation that allows a clear view of the underlying principles at work.”

Laboratory experiments are also well designed for reproducibility - the ability to replicate a research finding several times. Reproducing a finding serves as a robustness test of the finding’s validity and is viewed as an important goal of social science research (Jasny et al., 2011, McNutt, 2014). Unlike field experiments, which can be very hard to replicate even if the experimental design and details are clearly documented, well-documented laboratory experiments - with clear directions on how the experiment was designed, what content was in the instructions and experimental task, and so forth - can be replicated by other researchers at any other place and time. Large-scale efforts to reproduce experiments have occurred in both psychology and economics (OSC, 2012, Camerer, Dreber, et al., 2016). Several high-profile experiments in economics (as determined by their publication in one of the “Top 5” journals in the field) have been put to the test, with 11 of the 18 experiments replicated having a significant effect in the same direction as the original study (Camerer, Dreber, et al., 2016).¹

Laboratory experiments are also particularly useful for understanding the mechanism(s) driving a specific behavior, which can be much harder to test or measure in the field (Kessler and Vesterlund, 2015). The high level of control inside a lab setting allows for abstracting away as much “extraneous” detail as possible in order to focus on the most important aspects of a decision task facing an economic agent. Lab experiments are therefore very useful for testing predictions made by economic theory, namely whether human behavior is consistent with the assumptions and predictions made by specific models of behavior and how incentives and institutional frameworks can influence the behavior of economic actors (Roth, 1986). They are

¹The effect size was also typically smaller than that in the original study.

also low-cost ways to “pilot” whether certain behavioral interventions may work in the direction expected, before a costlier policy is implemented in the field (e.g. Agranov and Buyalskaya, 2021).

These same levels of control have raised criticism around whether the findings from laboratory experiments would generalize to other populations and other settings (Levitt and List, 2007, Levitt and List, 2008). The concerns raised include that participant actions are not “anonymous” (subject behavior is clearly being watched and recorded), the context in which decisions are made is important and almost always hard to replicate in a lab setting (given it is not an ecologically valid setting), and stakes are low (lab payments tend to represent a fraction of real-world stakes for similar decision tasks). Some of these critiques are valid, particularly that the population studied may not represent the population of interest. Since laboratory experiments are largely done in universities, the subject pool tends to be university students, which research has shown is not representative of behavior across the broader population (Henrich et al., 2005).

A review of economics experiments which were specifically designed to test lab-field generalizability found that most of the laboratory findings *could* be generalized to their relevant field setting counterparts (Camerer, 2010). Certain behaviors in particular demonstrated a high lab-field correlation. For instance, generosity in a laboratory task tends to be highly correlated with pro-social behavior in field settings (Franzen and Pointner, 2013). However other social preference “games” (the main target of the Levitt and List, 2007 critique) may not replicate as well in real-world situations. A recent meta-analysis from Galizzi and Navarro-Martinez, 2019, found that some real-world behaviors were rather different in laboratory settings as compared to their field equivalents, with only 37.5% of lab-field regressions showing a statistically significant association between game and field behaviors, and the average lab-field correlation reported amongst all papers being 0.14.

While the debate on whether lab and field behavior is sufficiently correlated or not continues - with the true answer likely depending on behavioral domain and decision task - all agree that the laboratory environment is typically not an ecologically valid setting (Hogelst, Schoevers, and Rot, 2015). For example, Camerer and Mobbs, 2017, argue that there is evidence of differential brain activity in response to hypothetical and real-world choices in various domains, arguing that “imaginative realistic paradigms” are necessary in the lab in order to evoke behavior as close to real as possible. While threats to the external validity of laboratory experiments

may have validity, the same critique can sometimes be made of field data when field studies are focused on groups which are not representative of a broader population (e.g. inferring how all individuals make trade-offs with respect to money based on a study in a developing country with high rates of poverty).

The pros and cons of data collected using surveys

Economic research tends to focus almost exclusively on choice data of revealed preferences (Samuelson, 1948), with several economists actively voicing skepticism regarding the use of any non-choice measures (Faruk and Pesendorfer, 2008). The rationale for this view is that economic theories are for the most part not equipped to deal with predictions about any data other than revealed choice. Other social science disciplines however, namely psychology, have embraced other methods of data collection. Specifically, psychologists frequently use surveys to gather data from individuals on both their choices and preferences, as well as psychological or contextual features which may help elucidate the mechanism that is driving choice (Breakwell and Rose, 2006).

The benefits of survey measures are clear: similarly to a lab experiment, researchers have more control over the data collection process. In survey design, a researcher has the ability to choose the set of variables they wish to collect data on and then proceed to survey the population of interest in order to collect that data. The range of data they collect can also be, and often is, broader than choices - researchers seek to assess mood, emotional state, and desires - which can be hard if not impossible to observe in revealed preference data. Furthermore, unlike lab experiments which often rely on narrow subject pools, surveys are easier to administer to subsets of a population which are unlikely to be observed in a lab setting, including older individuals and professionals.

The downside is that most of these surveys require individuals to answer questions about their own behavior (i.e. self-report). As a case study, habit is still largely studied using self-report scales (Gardner, 2015), with one scale called the SRHI, or "Self-Report Habit Index," dominating the literature (Verplanken and Orbell, 2003). The popularity of this self-report survey is explained in part by its brevity (it is a 12-item scale) and efficacy, making it the modern standard in psychology research on habitual behavior. Since habits are highly automatic behaviors, one of the questions asks individuals to rate their agreement with the statement "I do this behavior without thinking" on a Likert scale.

This is problematic because accurately self-reporting automaticity assumes one has very good meta-cognition (thinking about their thinking). But if automatic behavior occurs with little to no awareness, people will have a hard time explaining it. Further, they may misattribute automatic behaviors to behaviors of volition instead (Adriaanse et al., 2014, Gillebaart and Adriaanse, 2017, Wood and R unger, 2016).

So while survey methods allow for the collection of a broader set of variables than empirical data might contain, there have been some concerns, including from inside the field, about the exclusive use of self-report methods to study the natural habit formation processes (Harrington, 2017, Rebar et al., 2018). In particular, it is unlikely that aspects of habitual behavior like automaticity and context-sensitivity can be accurately captured using self-report measures alone.

How field data can help us address gaps in lab and self-report measures

If revealed choice data in natural settings is the economics “gold standard,” making some economists skeptical about both self-report and laboratory choice data, then field data is the antidote. While it is not perfect (as discussed earlier, the threats to external validity may still hold, for example), it does tend to validate research findings in a way that the previous two data sources struggle to do at times. Even the most compelling economic ideas tend to receive more interest when they can account for field data better than status quo economic models do. For example, despite Tversky and Kahneman, 1979 being the most cited economics paper of all time, the tenets of Prospect Theory were made more compelling by Camerer, 2011, who described ten well-documented empirical regularities in naturally-occurring field data which proved to be anomalies for Expected Utility theory (the status quo decision theory model in economics), but could be easily explained using three elements of Prospect Theory (Tversky and Kahneman, 1979, Tversky and Kahneman, 1992). Field data has therefore been used to test classic economic models in “natural settings” across a range of economics sub-fields, including game theory (Walker and Wooders, 2001,  stling et al., 2011, Brown, Camerer, and Lovallo, 2013).

The main advantage of research using field data is the ability to test whether theoretically and/or experimentally informed research findings replicate in natural choice settings. Field data is by definition data on choices made by humans in their natural environments, such that their behavior in that environment is supposedly driven by the individual’s response to the context and incentives present, rather than anything unrelated to the task at hand (e.g. their beliefs about what the experimenter expects them to do, per Zizzo, 2010).

Chapter overview

In Chapter 2, we use field data to study spillovers across behavioral domains, namely exercise and food choice. This work joins a small group of papers which document field evidence related to domain spillovers and failures of self-regulation (Karmarkar and Bollinger, 2015, Dolan and Galizzi, 2014). Besides these studies, most of the existing research on self-regulation has been conducted in controlled laboratory settings, where participants are either asked to imagine making hypothetical restrained choices (Khan and Dhar, 2006) or exert effort on a laboratory task as a proxy for making a restrained choice (De Witt Huberts, Evers, and De Ridder, 2012). As is the critique of many lab studies without direct field equivalents, it is debatable whether the self-regulation behaviors observed in survey and laboratory settings necessarily generalize to the field.

The research in Chapter 2 fills this gap by analyzing choice spillovers and potential failures of self-regulation in observational data. The novelty of this work is looking at how natural (rather than incentivized) changes in exercise systematically affect food choice, thus empirically identifying spillovers across two behavioral domains in field data. We find that, even after controlling for individual fixed effects, there is a robust effect of morning exercise on the healthiness of a lunch choice. We complement analysis of the field data with survey data in order to better understand the mechanism driving this result, and the survey results further support our assumption that choices regarding gym time and lunch foods are exogenous.

1.2 Obtaining Lab-Level Control in a Field Experiment

Running experiments in the field

Field experiments are difficult and time-consuming to set up, and a lot can go wrong in the process. Karlan and Appel, 2016 describe five leading causes of failures in field research, which serve as a useful framework for what to avoid or look out for when setting up a field experiment. They are:

1. Inappropriate research setting - this includes considering both the place and timing of the field study, the feasibility of the intervention, and ensuring that the treatment is appropriate for the context at hand.
2. Technical design flaws - this includes following standard experimental protocols, like thoughtful survey design (surveys being “the most common method

for collecting data in the field"), ensuring sufficient sample size, and randomizing treatment and control groups correctly.

3. Partner organization challenges - for field experiments which require partnering with an external organization, ensuring the partner has sufficient buy-in and bandwidth to complete the experiment and see it through to fruition.
4. Survey and measurement execution problems - collecting data carefully, including removing bias from any survey procedure (e.g. not using human surveyors which may pose the survey questions in a leading way) and imperfect measurement tools (e.g. sensors which may not be 100% accurate).
5. Low participation rates - recognizing that individuals may be hesitant, if not outright skeptical, about your experiment and that participation may be lower than anticipated, including lower than a pilot which may have had the draw of "novelty" in attracting participants.

Setting up a field experiment on campus

In Chapter 4, we present a novel method for running field experiments on campus using a customizable vending machine. As with many field experiments, it took over a year to set it up, and what follows are some extracts from the journey of that process. Firstly, to address how, if any, of the common research obstacles affected this experiment:

1. Inappropriate research setting - we faced a number of restrictions with respect to possible locations for installing the vending machine such that it would meet campus criteria and still maximize engagement and use. In the end, the undergraduate dormitory laundry room where the vending machine is currently set up was accepted as appropriate given that our goal was to target undergraduates and ensure they have access to snacks at all times of the day (e.g. the machine was often used in the middle of the night, when other on-campus options were closed and it was unattractive to go outside the dormitory for alternatives).
2. Technical design flaws - we spent a lot of time ensuring that the vending machine would be built to meet our design specifications and understanding how and when the machine might break down. This included lots of manual tests, like ensuring we stocked the machine correctly such that products would

not get stuck while vending, experimenting with the temperature control to ensure drinks came out satisfactorily cool, and so forth. Survey methods were used as part of the pre-test, but were not the main source of data collection during the study.

3. Partner organization challenges - part of the incentive to buy a vending machine, as discussed in more detail in Chapter 4, was to own the entire process of running the experiment such that we would not be dependent on partner organizations. Having said that, there were a number of times when the vending machine required a software update and we had to rely on timely communication with the hardware manufacturer. The credit card machine also relies on a partner organization (Nayax) which ran smoothly for the most part, but at some point became the target of a cyber-security attack such that we received phishing emails from hackers posing as Nayax accounts (we did not respond).
4. Survey and measurement execution problems - we did not use survey methods as our main measurement process, instead relying on the objective data collection of the machine. There were no major issues here, in that the machine was on around the clock as expected without experiencing power outages, and the inventory always matched the physical stock, so we gained confidence that all of the sensors and hardware designed to collect data were working properly.

The vending machine journey

The process began with us researching and choosing the right vending machine provider. This involved some preliminary market research to identify companies which would be able to build and customize a vending machine such that we could collect the following data types (and/or machine functionality)²:

- Unique ID – either via the student’s ID card or a credit card, so that we could track individuals through time using panel data (we ended up using unique credit card numbers, as it was too difficult to integrate the student ID card reader into the machine).

²We also considered video monitoring, which would have enabled us to capture everything from response times (e.g. comparing time of arrival with times of purchase) to emotional responses (e.g. using emotion face recognition technology). The video monitor could have been installed on top of the vending machine hardware, but we did not end up including it in our design specification given it would have violated the anonymity of customers.

- Accurate time (and day) stamp of each purchase (this was built in to the functionality of the machine we purchased).
- Ability for us to change prices, items, and run promotions dynamically (this was built in to the functionality of the machine we purchased).

The market research then informed a set of deeper diligence calls with two vendors who met our initial criteria. We prepared questions to guide the discussion and to ensure we covered everything we believed would be important to know about the machines and companies. These calls allowed us to narrow in on a specific provider - Digital Media Vending ("DMV") - which offered smart vending machine solutions using touch screen technology. The DMV machines do not typically have sensors or functionality which records people, but they were open to installing additional hardware (i.e. a camera) and doing additional customizations to meet our design needs.

Once we decided to contract with DMV, the next step was to agree on the exact specifications we wanted for the machine, and the timeline for the hardware build. We went with DMV's "Option 4 Touchscreen Elevator Vending Machine" (the hardware is discussed in more detail in Chapter 4), and formalized our engagement with DMV by allowing them to issue invoices and begin the customized build of our machine. We agreed to a 3-month build, with an understanding that it might take longer (actual arrival of the machine was closer to 5 months after purchase, per Hofstadter's law).

Once the machine was ordered, we started the Caltech IRB process. Since the vending machine was purchased with the intention to be used for collecting data on human subjects, it was necessary to draft and submit a research proposal for review and approval by our IRB committee. Our IRB application explained that we were partnering with a vending machine company and Caltech Dining Services to install a dynamically programmable vending machine in a building on campus. The goal of our experiment would be to track consumer behavior and, eventually, seek to influence behavior by changing variables related to the choice environment (such as the order in which items are presented on the screen). We went through multiple revisions of the IRB application to address multiple rounds of questions and concerns, and eventually agreed to adapt our experimental design to fit their criteria (e.g. posting a sign on the machine letting customers know that their data was being recorded).

In parallel to completing the IRB, we had multiple conversations with Caltech Dining Services about the prospect of replacing one of their existing vending machines with our experimental machine, and seeing whether it might be possible to have the restocking of the machine managed by the campus dining staff. In the end, it was not possible to replace an existing machine (they are all outsourced to a vendor, and the vendor contracts are very long-term with costly break clauses) nor have it restocked by Dining Services. Instead, we had to find a new location in which we could install our machine and planned to handle the restocking (with physical restocking help from a research assistant in our lab, and ordering products help from my advisor's administrative assistant).

Given the novelty of this experiment, it also took some time to identify and obtain all necessary approvals for proceeding with the field experiment. In addition to the IRB Committee and the Student Dining Services, the following departments needed to be notified (and/or approval gotten from) in order for us to successfully install the vending machine on campus:

- Student Affairs: (1) Approval for installation of a vending machine on campus (including verification of any possible conflicts with existing vendors); (2) Final decision regarding the installation location of the machine.
- Undergraduate Housing: (1) Approval for installation of a vending machine in a laundry space shared by several of the undergraduate houses; (2) Obtaining card access to the laundry space so that we could regularly visit and stock the machine; (3) Directed communication with undergraduates living in the relevant houses.
- Maintenance Team: (1) Physical unpacking and installation of the machine (since the delivery service did not provide enough manpower).

Once the machine was installed, we had a number of tutorials (followed by several weeks of trial-and-error learning) to learn how to use the hardware and software of the vending machine. David Ashforth (a DMV employee and our main contact) visited the campus at the beginning of December 2019 to set up the machine with our broader lab present. This installation involved installing and starting up the Nayax payment terminal (which is embedded in the machine, but had shipped separately), guiding us through setting up the vending machine software, and testing

the dispensing of pre-purchased snacks. At the close of this setup, the machine was functioning and able to dispense all stocked foods.

Credit card purchasing was enabled via a Nayax terminal, which is embedded on the right-hand side of the machine and operates separately from the machine, having its own web portal interface. The terminal accepts cards (swipe, chip, or contactless) as well as NFC payment (Apple, Google, Samsung Pay). For internal testing purposes, we were issued a number of Nayax payment cards which we could use when running dispenser tests on the machine.

The front of the vending machine is a large touchscreen running on Android OS. The operator view has several functionalities that are currently working and necessary for the operation of the machine, including:

- **Inventory:** This allows for viewing what items are currently in stock in the machine, and updating that inventory following a restocking.
- **Complete Refill:** This functionality allows for a one-touch way to update the inventory on every snack to maximum capacity (if doing a complete restock of the machine, it is easier to do this than manually update the inventory of each item).
- **Product Sort:** This allows for the sorting of products that are offered by the machine, as shown on the front display. If the order is not manually set, the machine will randomly reshuffle the order in which products are presented at preset intervals.
- **Motor Testing:** A function which permits the testing of the vending machine motors in order to ensure the product will dispense properly when a customer buys it.

In addition to the hardware, we had to become familiar with the software which allowed us to control the vending machine remotely. This was an online vending tracker which was created by DMV (or rather, which DMV outsourced to their team in China - some of the tutorials were exclusively in Chinese so we had to leverage Chinese speakers in our lab to understand them). The software allowed us to remotely specify the products we wanted to be in the machine (we then had to stock those products and update inventory on the hardware of the machine). Each of these set-up pieces required experimentation and trial-and-error. For example, we

had to be careful to use high resolution images, but the images had to be less than 1 GB and a square shape, in order for the machine to process them. In the back-end, we also controlled the “capacity” of each row in the vending machine - for example, telling the machine that there could be up to 8 Cokes on a specific row. This would allow for shortcuts to restocking where we could indicate we made a “full restock” instead of increasing the counter from 1 to 8.

Finally, the software allowed us to track sales through time, and had a number of reporting features which provided live data of the machine stock, purchase behavior, and so forth, as well as raw data files, which were used for the analysis.

While it took some time to set up, the vending machine now opens up many exciting future directions for research. With a customizable vending machine that can be controlled remotely, it becomes possible to run lab-like experiments in the field. This may include testing the impact of prices on consumer sentiment and behavior, analyzing the role of attention in consumer choice, and teasing apart predictions made by closely related theoretical models which require a carefully controlled experimental design.

Chapter overview

In Chapter 4, we use this customizable vending machine to run a field experiment testing the effect of a price promotion on consumer behavior. The goal of this “pilot study” was to credibly dissociate predictions made by brand loyalty/habit formation from reference-dependence theories. The vending machine served as a “mini-retailer,” allowing us to control all details of a price promotion treatment in an ecologically valid setting, and collect granular panel data on purchases occurring at all hours of the day over the course of 10 weeks. The unique contribution of the vending machine was that, in contrast to past work, it allowed me to control for stockpiling behavior, which is an important concern for empirical work analyzing price promotions in the marketing literature. Unfortunately, we began collecting data in January 2020, and was forced to stop data collection in early March 2020 when the COVID-19 pandemic was spreading across the globe. While the data collected from this first study is more limited than the data we anticipate being able to draw conclusions from (i.e. running multiple studies over the course of the following year), we were still able to learn a few valuable takeaways from these early results.

In particular, analysis of the data collected for Chapter 4 suggests that price promo-

tions increase the sales of both discounted and non-discounted items, as well as the number of unique customers making purchases during the sale period. Furthermore, in line with the loss leader hypothesis, more items are purchased during the sale period overall. Finally, a habit perspective (as opposed to a reference-dependence perspective) appears to do a better job explaining consumer behavior following the sale given that purchases of discounted items increase during the post-sale period. These results provide support in favor of retailers running occasional price promotions if the strategic goal is to increase brand loyalty for a specific product. In contrast to the predictions made by the reference-dependence hypotheses, customers do not appear to experience loss aversion following a price promotion. Instead, they appear equally likely to purchase the same products which experienced a discount when they are back to their full price.

1.3 Embracing an Interdisciplinary Approach

Seeing the same behavior from more than one lens

Different social science disciplines may approach investigating the same human behavior with a different set of tools and assumptions. It is becoming easier and more important to stretch across disciplines when conducting social science research. As we write in Buyalskaya, Gallo, and Camerer, 2021, it is “the confluence of data, diverse teams, and difficult challenges which makes [now] a unique and exciting time for social scientists to tackle important research questions.” We argue that embracing this “golden age” requires developing a lingua franca, or a language which transcends disciplinary borders and is built on the concept of disciplines trading the best definitions, methods and theories with one another in a parsable way.

The two social science disciplines this work has pulled from the most are economics and psychology. Economics has been useful in demonstrating that even complex human decisions can be formalized using simple mathematical models with clearly identified assumptions and parameters that can be tweaked and tested using experimental or empirical methods. Psychology has demonstrated that some of the “extraneous” variables which these simple mathematical models seek to abstract away from - such as context or mood - may actually be incredibly important to determining the outcome of the decision task, and we sometimes need to rely on data sources outside of revealed choice data to measure these variables. This research has also been informed by an understanding of cognitive neuroscience, a field

which investigates the physical machinery that implements, to cite Marr, 1982, the algorithms we use to solve the computations in front of us.

Despite being incredibly important, interdisciplinary work is not always easy. Some of the challenges are worth mentioning briefly here (and a longer discussion of each of these can be found in Buyalskaya, Gallo, and Camerer, 2021):

- Silos between journals - a lot of modern-day journals still cater to the readership of a very specific discipline or even sub-discipline. This leads to scholars reading, and citing work largely in their own discipline, even if relevant (and high-quality) research has been done in adjacent disciplines.
- Career incentives - early in one's research career, scholars are encouraged to remain focused on making a notable contribution to one subject area, with broader research viewed as a distraction. In some fields, including economics, there are still authorship norms which create strong incentives against interdisciplinary work (i.e. the emphasis on solo-authored papers).
- Unifying frameworks - since different disciplines each bring their own set of theories and methods to a problem, there are occasionally multiple (and competing) explanations for the same data. The development of a true lingua franca which constrains this set to an agreed upon explanation takes time, effort and humility that scholars may not always have.

Despite these challenges, interdisciplinary work is incredibly rewarding and worth pursuing. As a case study of how an interdisciplinary perspective can lead researchers close to the scientific "truth," we will again take a look at the study of habit formation.

Case study: Habit

Habit research is of particular interest to social scientists because many of the behaviors we study are candidates for becoming habitual. The study of habit naturally crosses disciplinary boundaries, and the most promising understanding of it is likely to come from integrating evidence and methods across disciplines (Rebar et al., 2018; pg. 42). In particular, the disciplines of psychology, computational neuroscience, and economics have thought most deeply about the way to define and measure habits in human behavior.

In psychology, habit is defined as a behavior which is prompted automatically by contextual cues as a result of learned context-action associations (Wood and Neal, 2009). This definition combines two key attributes of habitual behavior which guide a lot of the psychology research: predictable context-sensitivity (do features of an environment cue the behavioral execution of the habit?) and automaticity (is the behavior executed with very little cognitive control?). The most common measures used in psychology to assess context-sensitivity and automaticity of behavior are self-report surveys (the benefits and cons of which were discussed in Section 1.1).

In computational neuroscience, habit is defined as cognitive activity which has moved from the region of the brain known as the associative striatum (DMS), which controls more goal-directed activities where habits “originate,” to the sensorimotor striatum (DLS), which is where behavior is controlled once it has become habitual. Functional MRI studies are used to identify which brain regions control the behavior in question. With respect to behavioral measures, an important test used to determine whether a behavior is habitual or not is a test of sensitivity to reward devaluation. The procedure originated in animal learning studies and involves devaluing a reward following an extensive period of training to see whether the highly-trained habit continues to be executed. This phenomenon has been termed insensitivity to reward devaluation, and is a behavioral hallmark of habitual processing.

Finally in economics, habit is typically defined (somewhat narrowly) as history-dependence. The theories of habit formation from economics are motivated by evidence of an empirical correlation between past and current consumption. These models therefore specify consumption utility as a function of actual immediate consumption, and make the prediction that current behavior depends on expectations of the future. People are assumed to be “self-aware” enough to quit a habit if they anticipate a future price hike, for example. With respect to measurement, economists almost exclusively use data on revealed choice (as discussed in Section 1.1) when testing the predictions from these models.

Chapter overview

In Chapter 3, we consider and incorporate findings from all of these disciplines into our research approach. Specifically, we adopt the definition of habit used in psychology, but seek to replace self-report measures (as would be of interest to economists), on which we are able to run a test of reward devaluation insensitivity (as would be of interest to computational neuroscientists).

We present a novel methodology for identifying behaviors that are highly and predictably context-sensitive, and thus candidates for being habitual.³ While there is a large body of laboratory research documenting the mechanisms underlying well-developed habits in animals and humans, there is much less field research on how human habits naturally develop over time (Verplanken, 2018, pg.7). Our investigation introduces a new approach to studying context-sensitive behavior in the wild: using two large datasets on gym attendance and handwashing behavior, we use machine learning to statistically classify when choices are predicted by an identifiable set of context variables. This machine learning technique generates a person-specific measure of behavioral predictability, which can then be used to study individual differences in predictability and speed of habit formation.

Chapter 3 establishes two important discoveries using the field data on gym attendance and handwashing behaviors. First, the sets of context cues that are predictive of individual-level behavior are different for different people. Specifically, while historical behavior is an important universal predictor, other context variables such as day of the week or month of the year have more heterogeneous effects. Second, contrary to common wisdom, there is no “magic number” for how long it takes to form a habit. Instead, the speed of habit formation appears to vary significantly, both between behavioral domains and between individuals within domains. This research was informed by an interdisciplinary perspective on habitual behavior, which sought to take the best evidence and methods from across multiple social science fields.

References

- Adriaanse, M. A. et al. (2014). “Effortless inhibition: Habit mediates the relation between self-control and unhealthy snack consumption”. In: *Frontiers in Psychology* 5, p. 444.
- Agranov, M. and A. Buyalskaya (2021). “Deterrence effects of enforcement schemes: An experimental study”. In: *Management Science*, forthcoming.
- Breakwell, G. M. and D. Rose (2006). “Preface to the Handbook”. In: *Research Methods in Psychology*. Ed. by G. M. Breakwell et al. Sage Publications, pp. 2–23.
- Brown, A., C. F. Camerer, and D. Lovallo (2013). “Estimating structural models of limited strategic thinking in the field: The case of missing movie critic reviews”. In: *Management Science* 59.3, pp. 733–747.

³Unfortunately, the field data do not enable us to measure automaticity, which means we cannot be certain whether the behaviors we observe are truly habits or simply a function of strong preferences.

- Buyalskaya, A., M. Gallo, and C. F. Camerer (2021). “The golden age of social science”. In: *Proceedings of the National Academy of Sciences* 118.5.
- Camerer, C. F. (2010). “The promise and success of lab-field generalizability in experimental economics: A critical reply to Levitt and List”. In: *The Methods of Modern Experimental Economics*. Ed. by G. Frechette and A. Schotter. Oxford University Press.
- (2011). “Prospect theory in the wild: Evidence from the field”. In: *Advances in Behavioral Economics*. Ed. by C. F. Camerer, G. Loewenstein, and M. Rabin. Princeton University Press, pp. 148–161.
- Camerer, C. F., A. Dreber, et al. (2016). “Evaluating replicability of laboratory experiments in economics”. In: *Science* 351.6280, pp. 1433–1436.
- Camerer, C. F. and D. Mobbs (2017). “Differences in behavior and brain activity during hypothetical and real choices”. In: *Trends in Cognitive Sciences* 21.1, pp. 46–56.
- De Witt Huberts, J., C. Evers, and D. De Ridder (2012). “License to sin: Self-licensing as a mechanism underlying hedonic consumption”. In: *European Journal of Social Psychology* 42.4, pp. 490–496.
- Dolan, P. and M. Galizzi (2014). “Because I’m worth it. A lab-field experiment on the spillover effects of incentives in health”. In: *LSE CEP Discussion Paper CEPDP 1286, London*.
- Faruk, G. and W. Pesendorfer (2008). “The case for mindless economics”. In: *The Foundations of Positive and Normative Economics*. Ed. by A. Caplin and A. Schotter. Oxford University Press, pp. 3–40.
- Franzen, A. and S. Pointner (2013). “The external validity of giving in the dictator game: A field experiment using the misdirected letter technique”. In: *Experimental Economics* 16, pp. 155–169.
- Galizzi, M. and D. Navarro-Martinez (2019). “On the external validity of social preference games: A systematic lab-field study”. In: *Management Science* 65.3, pp. 976–1002.
- Gardner, B. (2015). “A review and analysis of the use of ‘habit’ in understanding, predicting and influencing health-related behavior”. In: *Healthy Psychology Review* 9.3, pp. 277–295.
- Gillebaart, M. and M. A. Adriaanse (2017). “Self-control predicts exercise behavior by force of habit, a conceptual replication of Adriaanse et al. (2014)”. In: *Frontiers in Psychology* 8, p. 190.
- Harrington, N. (2017). “Commentary: Why it doesn’t pay to ask consumers about habitual behaviors”. In: *Journal of the Association for Consumer Research: The Habit-Driven Consumer* 2.3.

- Henrich, J. et al. (2005). ““Economic man” in cross-cultural perspective: Behavioral experiments in 15 small-scale societies”. In: *Behavioral and Brain Sciences* 28.6, pp. 795–815.
- Hogenelst, K., R. A. Schoevers, and M. Rot (2015). “Studying the neurobiology of human social interaction: Making the case for ecological validity”. In: *Social Neuroscience* 10.3, pp. 219–229.
- Jasny, B. R. et al. (2011). “Again, and Again, and Again . . .” In: *Science* 334.6060, pp. 1225–1225.
- Karlan, D. and J. Appel (2016). “Part 1: Leading causes of research failures”. In: *Failing in the Field*. Princeton University Press, pp. 17–70.
- Karmarkar, U. and B. Bollinger (2015). “BYOB: How bringing your own shopping bags leads to treating yourself and the environment”. In: *Journal of Marketing* 79.4, pp. 1–15.
- Kessler, J. and L. Vesterlund (2015). “The external validity of laboratory experiments: The misleading emphasis on quantitative effects”. In: *Handbook of Experimental Economic Methodology*. Ed. by G. Frechette and A. Schotter. Oxford University Press, pp. 391–406.
- Khan, U. and R. Dhar (2006). “Licensing effect in consumer choice”. In: *Journal of Marketing Research* 43.2, pp. 259–266.
- Levitt, S. and J. List (2007). “What do laboratory experiments measuring social preferences reveal about the real world?” In: *Journal of Economic Perspectives* 21.2, pp. 153–174.
- (2008). “Homo economicus Evolves”. In: *Science* 319.5865, pp. 909–910.
- Marr, D. (1982). *Vision: A computational investigation into the human representation and processing of visual information*. Cambridge, Mass: MIT Press.
- McNutt, M. (2014). “Reproducibility”. In: *Science* 343.229.
- OSC (2012). “An open, large-scale, collaborative effort to estimate the reproducibility of psychological science”. In: *Perspectives on Psychological Science* 7.6.
- Östling, R. et al. (2011). “Testing game theory in the field: Swedish LUPI lottery games”. In: *American Economic Journal: Microeconomics* 3.3, pp. 1–33.
- Plott, C. R. and V. L. Smith (2008). “Preface to the Handbook”. In: *Handbook of Experimental Economics Results*. Ed. by Charles R. Plott and Vernon L. Smith. Handbook of Experimental Economics Results. Elsevier, pp. 1–2.
- Rebar, A. et al. (2018). “The measurement of habit”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 31–49.
- Roth, A. (1986). “Laboratory experimentation in economics”. In: *Economics and Philosophy* 2, pp. 245–273.

- Samuelson, P. (1948). “Consumption theory in terms of revealed preference”. In: *Economica* 15, pp. 243–253.
- Tversky, A. and D. Kahneman (1979). “Prospect theory: An analysis of decision under risk”. In: *Econometrica* 47.2, pp. 263–291.
- (1992). “Advances in prospect theory: Cumulative representation of uncertainty”. In: *Journal of Risk and Uncertainty* 5, pp. 297–323.
- Verplanken, B. (2018). “Introduction”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 1–10.
- Verplanken, B. and S. Orbell (2003). “Reflections on past behavior: A self-report index of habit strength”. In: *Journal of Applied Social Psychology* 33.6, pp. 1313–1330.
- Walker, M. and J. Wooders (2001). “Minimax play at Wimbledon”. In: *The American Economic Review* 91.5, pp. 1521–1538.
- Wood, W. and D. Neal (2009). “The habitual consumer”. In: *Journal of Consumer Psychology* 19, pp. 579–592.
- Wood, W. and D. Rünger (2016). “Psychology of Habit”. In: *Annual review of psychology* 67, pp. 289–314.
- Zizzo, D.J. (2010). “Experimenter demand effects in economic experiments”. In: *Experimental Economics* 13, pp. 75–98.

Chapter 2

IDENTIFYING SELF-REGULATION FAILURES IN THE FIELD

ABSTRACT

Choice spillovers across behavioral domains are important for marketing, as they suggest that interventions in one choice domain may have far-reaching, and countervailing influences, in other choice domains. Using a novel dataset drawn from a university population, we study individual-level exercise behavior and food choices over the course of a year, and identify one such choice spillover. Specifically, we find that whether exercise is performed before lunch affects the healthiness of meal choice at lunch. Furthermore, the direction in which it affects meal choice is predictable: attending the gym in the morning is systematically correlated with a less healthy lunch choice. We hypothesize that the behavior may be driven by a self-licensing mechanism, as substantial time lags between gym use and meals cast doubt on a purely biological mechanism driving the results, and demonstrate that it does not appear to be modulated by stressful periods.

2.1 Introduction

A fundamental challenge faced by empirical researchers in consumer behavior is that most datasets offer only a narrow window into each consumer's life. Take a consumer named Joe. Depending on the data they have, a researcher may observe Joe's regular purchases of chocolate cake, or his daily jogging activity, but they will rarely observe both of these behaviors for the same Joe on the same day. The result is that researchers often analyze consumer behaviors in silos, and cannot pick up on possibly compensatory behaviors in other silos: other choice scenarios, or other locales, separated by space and time. One can measure how much chocolate cake Joe eats, and how often he jogs, but cannot determine whether Joe jogs more on exactly those days when he indulges in an extra slice of cake. Yet such spillovers across separate domains and locales are of key relevance for marketing, as they suggest that effective interventions in one choice domain may have countervailing influences in other, seemingly-unrelated domains. For instance, an advertising campaign which encourages Joe to exercise more may lead him to consume more desserts, thus undoing the benefits from a more rigorous exercise regimen.

The goal of this paper is to empirically investigate one such choice spillover from one common consumer choice domain - exercise - to another - food choice. We find field evidence of such a spillover, and specifically of a failure in self-regulation (or self-control), where a positive behavior in the exercise domain leads to a more indulgent behavior in the food domain. Our study exploits a unique administrative dataset from a university located on the west coast of the United States. This dataset combines the databases from the on-campus cafeterias and the fitness recreation center, thus allowing us to track food choices and exercise frequency for members of the university community (students, faculty, staff) over the course of one year.

This paper joins a small group of papers which document field evidence related to domain spillovers and failures of self-regulation. One such study by Karmarkar and Bollinger, 2015 found that bringing one's own shopping bag to a supermarket increases demand for less healthy items, albeit within only one choice locale (a supermarket). Another closely related study by Dolan and Galizzi, 2014 conducted a lab-field experiment which found that one-off high stakes financial incentives to increase exercise intensity (step counts) partially backfire as they significantly increase the number of calories consumed following the brisk exercise. The novelty of our paper is looking at how natural (rather than incentivized) changes in exercise systematically affect food choice, thus empirically identifying spillovers across two

behavioral domains (separate physical locations) in field (observational) data.

We also contribute to two active strands of research into consumers' exercise behaviors and food choices. In the context of gym use, DellaVigna and Malmendier, 2006 found significant inefficiencies in gym memberships, with a large proportion of gym goers apparently overspending on gym memberships given how infrequently they use the gym. In follow-up work, Acland and Levy, 2015 found a possible explanation for these results in terms of projection bias, which causes gym goers to systematically over-predict their future gym usage. Given exercise is typically seen in "net positive" terms, behavior change interventions have largely focused on increasing gym use. Several papers have evaluated the effectiveness of such interventions (e.g. Charness and Gneezy, 2009; Royer, Steher, and Sydnor, 2015). Other behavior change research agendas similarly focus on interventions addressed at individual behaviors: from increasing savings (Thaler and Benartzi, 2004) to decreasing smoking (Adda and Cornaglia, 2006, Volpp, 2009) and junk food consumption (Dubois, Griffith, and O'Connell, 2018).

A common feature of papers investigating consumer choices is that they focus on a single consumer behavior. Such a "partial-equilibrium" approach implicitly assumes that whatever happens in one realm (for example, exercise) does not affect other realms (for example, food choices). However, this idea is in stark contrast to evidence that individuals struggle with limited amounts of willpower and self-control (Duckworth, Milkman, and Laibson, 2018), and empirical evidence of *self-licensing* behavior, in which making a more restrained choice causes an individual to permit themselves to make more indulgent choices later (De Witt Huberts, Evers, and De Ridder, 2013). Taking our example of jogging and chocolate cake, the self-licensing hypothesis would predict that jogging would increase Joe's likelihood of later eating more chocolate cake: after having performed the restrained exercise of jogging, Joe indulges in the atypical piece of cake because he believes he has "earned" it, so to speak. Such links in consumer behavior across multiple choice domains requires a more holistic "general equilibrium" approach to human behavior.

Besides the two studies cited earlier, most of the existing research on self-regulation has been conducted in controlled laboratory settings, where participants are either asked to imagine making hypothetical restrained choices in the past (Khan and Dhar, 2006) or exert effort on a lab task as a proxy for making a restrained choice (De Witt Huberts, Evers, and De Ridder, 2012). As is the critique of many lab studies without direct field equivalents however, it is debatable whether the self-regulation

behaviors observed in survey and laboratory settings necessarily generalize to the field.

Our paper fills this gap by analyzing choice spillovers and potential failures of self-regulation in observational data. We find that, even after controlling for individual fixed effects, there is a robust effect of morning exercise on the healthiness of a lunch choice. Specifically, healthy lunch items such as salad are less likely, and unhealthy items (such as fries or energy drinks) more likely to be purchased at lunchtime following morning gym use. Across different slices of the university population, faculty are no less likely to exhibit this behavior than undergraduates. Women are also no less likely to exhibit this behavior than men. In addition, this result obtains even given the large gap (three hours on average) between morning gym use and lunch, which suggests that post-workout dehydration or calorie depletion cannot fully explain the result.

While we collect survey data which supports our assumption that gym time and lunch choices are exogenous, we also run robustness tests to acknowledge the possibility that they are not. Specifically, that there may be a common individual-level shock which affects both the timing of gym attendance as well as lunch healthiness. To accommodate this heterogeneity, we consider specifications in which gym time is allowed to be endogenous, and instrumented with day-of-week dummies. These day-of-week dummies capture exogenous variation in morning gym use arising from institutional features of undergraduate class scheduling at the university (e.g. students cannot go to the gym when they have classes). Our results remain robust and significant with this additional specification.

We complement these results with survey data meant to confirm the mechanism driving these decisions. When asked directly, only 24% of the population surveyed (which is a subset of the population on which we have empirical data) believe the reason for their unhealthy lunch choice is due to having fatigued themselves earlier in the day. In other words, even if some biological need is *influencing* the observed behavior, most people do not seem to be consciously attributing their unhealthy purchases to having depleted their energy resources earlier in the day.

These results have direct managerial relevance. Recently, management teams of companies have aimed to encourage “positive” practices amongst their employees, as a means of reducing health insurance costs or their carbon footprint. For instance, Google engages directly in the health of their employees by implementing initiatives aimed at healthy eating and redesigning food options within the company cafeterias

(Black, 2020). Our results imply that such initiatives may backfire if the spillover effects from other choice domains (e.g. exercise) onto the choice domain of eating are not taken into account.

2.2 Theoretical Framework

"It has been my experience that folks who have no vices have very few virtues."

- *Abraham Lincoln*¹

Choice spillover

A common feature of papers investigating consumer choices is that they focus on a single consumer behavior. This approach implicitly assumes that whatever happens in one realm has no effect on other realms. However, this idea is intuitively - and empirically - questionable. Individuals struggle with limited amounts of willpower and self-control (Duckworth, Milkman, and Laibson, 2018), where performing one good but effortful behavior depletes their ability to perform future good behaviors in other domains.

Dolan and Galizzi, 2015 go further to say that almost no behavior exists in a vacuum, finding that behavioral spillovers are present across choice domains when those behaviors are executed sequentially and somehow linked by an underlying motive. Given that exercise and nutrition are linked in their motive to promote a healthier lifestyle, it is reasonable to believe that decisions in one domain *will* influence decisions in the other.

Hypothesis 1: An individual's decision to go to the gym in the morning will influence their decision regarding food choice later in the day.

Self-licensing

Furthermore, the directional impact of the spillover may be predictable. In one scenario, the first behavior may lead to another behavior in the same direction - a sort of behavioral "momentum". For example, donating to charity in the morning may increase the likelihood of giving extra lunch money to one's child shortly after. In another scenario, the first behavior leads to a behavioral "reversal", with a spillover in the opposite direction.

¹As quoted in (Blaisdell, 2012).

Literature on self-licensing behavior has shown that making a more *restrained* choice causes an individual to permit himself to make more indulgent choices later on (De Witt Huberts, Evers, and De Ridder, 2014). Some research has looked at how consumers believe they have “earned the right to indulge.” In a series of studies, Kivetz and Simonson, 2002 found that when consumers put a lot of effort into earning their frequency program points, they were more likely to choose an indulgent luxury reward (and more likely to choose a frequency program which offers a luxury reward). Therefore, self-licensing theory predicts that morning exercise - a behavior which requires a lot of restraint and effort on the part of the individual - should lead to more indulgent behavior later on. In the context of lunch choice, this indulgent behavior would look like choosing a less healthy food at lunch following morning exercise.

Hypothesis 2: Going to the gym before lunch will increase an individual’s likelihood of choosing an unhealthy meal at lunch.

H1 and H2 are the primary hypotheses being tested in this paper. Looking at these two hypotheses together, it is important to acknowledge potential confounds. Specifically, that there may exist a third causal mechanism which explains both the decision to exercise and to eat less healthy foods. The most likely candidate is stress - which may lead to an increase in gym attendance (if students use exercise as a coping mechanism for stress) and has been shown to increase craving of unhealthy foods. Thus our empirical work also tests the following hypothesis:

Hypothesis 3: The choice of unhealthy food after morning exercise is not modulated by higher levels of stress, as proxied using exam periods.

Potential modulators

In addition to stress, there is some evidence that intense physical activity may lead to cognitive control fatigue and increased choice impulsivity (Blain et al., 2019). However, recent literature has disagreed on whether consuming more glucose following such physical exertion reinstates cognitive control, with some recent papers indicating that the behavior may be the result of a “false belief” more than of biological necessity (Job, Walton, et al., 2013). In order to understand whether the population we study holds these “false beliefs,” we run a survey on a subset of the

population to test this idea.

In the following set of results, we combine empirical analysis with survey data to demonstrate that there are spillovers from the domain of exercise to the domain of food choice. In Appendix A, we conduct a number of robustness checks to demonstrate that these results cannot be explained (do not appear to be modulated by) stress or other factors.

2.3 Data

This paper analyzes a comprehensive dataset on individuals who study and work at a university on the US west coast. In what follows, we describe the various components of this dataset.

Lunch purchases

We obtain daily panel data of lunch purchases spanning the year 2018 from the university's dining services. The campus has three dining locations: the main cafeteria, a deli (selling a smaller selection of salads and sandwiches), and a coffee shop (selling mainly hot and cold beverages and packaged, non-perishable items during our sample period). Our data includes all purchases made at any of these three locations from January 1 - December 31, 2018. While we have information for all meal times, we primarily restrict our attention to lunch purchases (defined as food purchases between 11am - 4pm, which is when lunch is available at the main cafeteria). We focus on lunch because many individuals do not regularly eat breakfast or dinner at campus facilities; in contrast, a large number of students and faculty consistently eat lunch on campus given the narrow range of outside options within short walking distance. The majority (61%) of unique purchases in our dataset are made at the main cafeteria.²

Each individual in the dataset is identified by a unique Customer ID which allows us to observe the same people through time. Individuals in the dataset are in one of five categories: Faculty, Undergraduate Students, Graduate Students, Staff, and Postdoctoral Scholars. The gender of each individual is also observed. Most of the food purchases in our dataset were made using the individual's university ID card, which is linked to a university account which is periodically loaded with funds drawn

²As a robustness check, we truncate our dataset to those individuals who purchased lunch at the cafeteria, the modal food venue in our dataset. The regression results are robust, indicating that the observed effect is not driven by individuals switching away from eating at a cafeteria with a wide range of choices to smaller cafes with more limited, and typically less healthy, items for sale.

directly from the individual’s bank account. A small fraction of food purchases were made using cash or credit card, for which the purchaser’s ID is not available; we removed these purchases from the dataset.³

The timestamp of each food purchase is observed, along with a detailed product description, and associated price. Most food items (burgers, pizza, hot entrees) are priced on an *a la carte* basis, but some key items (notably the salad bar and “build your own burrito” bar) are priced by weight. A small number of non-food purchases (such as OTC medicine, which can be bought in the cafes) were removed from the dataset. Our cleaned and trimmed down dataset is made up of 288,605 transactions made by 2,999 customers.

Table 2.1: Lunch transactions summary by demographics

Person Type	n	% Male
Faculty	17,969	88.4
Graduate	85,564	73.7
Postdoctoral Scholar	13,358	73.5
Staff	14,313	54.0
Undergraduate	157,401	58.2

Healthiness ratings for food items

Altogether, there are over 600 unique food products, which we aggregated into 41 product types. Our analysis requires a measure of the healthiness of each food item. We obtained these by directly canvassing undergraduates at the university and asking them to rate whether each of the 41 product categories is healthy, unhealthy, or neither.⁴

We then aggregate the responses obtained from students to obtain, for each product, the percentage of students who considered this item healthy, the percentage of students who considered this item unhealthy, and the percentage who choose neither. From there, we create a health score h_x of the form below, as well as a binary indicator of whether the product is sweet or not, s_x . Examples of items which result in being ranked as the least and most healthy are in Tables 2.2 and 2.3, respectively.

³Less than 1% of the dataset was made up of transactions made using cash or credit/debit card. As such, the vast majority of on-campus purchases are made using university ID cards.

⁴We used the university’s students rather than an outside population because a number of the lunch items require having experienced them in the cafeteria.

$$h_x = -1(\%_{Unhealthy}) + 1(\%_{Healthy})$$

Table 2.2: Ten least healthy products

Product (x)	$\%_{Unhealthy}$	$\%_{Healthy}$	h_x	s_x
Candies	91.23	1.75	-0.89	1
Fries	84.21	0.00	-0.84	0
Dessert	84.21	1.75	-0.84	1
Energy Drink	82.46	0.00	-0.82	1
Ice Cream	87.72	5.26	-0.82	1
Soda	84.21	3.51	-0.81	1
Processed Food	77.19	1.75	-0.75	0
Chips	68.42	1.75	-0.67	0
Baked Goods and Cookies	64.91	5.26	-0.60	1
Burger	56.14	12.28	-0.44	0

Table 2.3: Ten most healthy products

Product (x)	$\%_{Unhealthy}$	$\%_{Healthy}$	h_x	s_x
Water	0.00	96.49	0.96	0
Salad	0.00	92.98	0.93	0
Veggies	1.75	91.23	0.89	0
Fruit	5.26	91.23	0.86	1
Quinoa	1.75	84.21	0.82	0
Tea	1.75	80.70	0.79	0
Eggs	1.75	77.19	0.75	0
Yogurt	7.02	82.46	0.75	1
Deli	3.51	73.68	0.70	0
Soup	7.02	77.19	0.70	0

Gym usage data

We also obtain daily panel data for gym check-ins spanning the year 2018 from the university's recreation services. In order to use *any* of the athletic facilities or fitness offerings, individuals are required to sign in using their ID card prior to entering the gym. Hence, our data includes all gym check-ins made from January 1 - December 31, 2018. Each gym check-in has an associated Customer ID (same as those in the Lunch Purchases dataset) and an exact time stamp of the associated check-in.⁵

⁵While most individuals use the gym at most once per day, we do observe 2% of users who have multiple check-ins on the same day.

The university gym offers a comprehensive array of fitness offerings, including: treadmills, stair-climbers, weight lifting rooms, spinning bicycles and classes, a range of fitness classes from yoga to high intensity workouts, two long-range heated outdoor swimming pools with multiple lanes and a lifeguard always on duty, squash courts available for reservation, a basketball court, several tennis courts, and a competition size (400 meters) running track. Most of these facilities are available for use, and relatively un-congested at all times of day, and the gym is complimentary to all registered students and faculty on campus. Moreover, there are also two large locker rooms with shower facilities and towel service, so that even if some individuals venture off-campus to exercise (eg. running or cycling), it is likely that they will still check-in to the gym to use these facilities afterwards.

Table 2.4: Summary statistics for key variables

Variable name	<i>N</i>	mean	stdev	min	max
<i>Independent variables:</i>					
earlygym	221,507	0.0351	0.1840	0	1
lategym	221,507	0.1009	0.3012	0	1
neither	221,507	0.8640	0.3428	0	1
<i>Dependent variables:</i>					
weighthealth	221,507	0.1749	0.5394	-0.8948	0.9649
mealhealth	221,507	0.1571	0.5406	-0.8948	0.9649
mealprice	221,507	6.7693	3.8087	0.04	115.66
weightsweet	221,507	0.2204	0.3825	0	1
mealsweet	221,507	0.2490	0.3894	0	1

Combined dataset

The common Customer ID allows us to link the gym usage and food choice databases. For some of the empirical specifications, we restrict attention only to days in which a given individual both ate lunch as well as worked out at the gym (which makes up about 9.4% of the dataset). We also aggregate food purchases up to the meal level, using health and price information for the bundle of items consumed.

Table 2.4 contains summary statistics on the key variables used in our empirical specifications. Each observation is a day in which an individual eats lunch at the university's dining services; gym usage is registered on around 14% of these days, with 3.5% of them being morning gym sessions and the remainder (10%) afternoon/evening sessions.

Survey evidence

To complement our observational data, we also administered a small survey to the university population ($n = 146$).⁶ We asked a number of questions aimed at validating our dataset. Exactly 50% of our respondents report using the university gym at least once a week, with 40.4% reporting regularly using it more than once a week. Furthermore, 32% report “typically eating lunch at either the cafeteria, the deli or the coffee shop” on most weekdays while a significant 62% report “typically bringing lunch from home.” This confirms that we are observing frequent (daily to weekly) decisions made by a portion of the campus population, as opposed to infrequent decisions made by a select group. However, there will inevitably be some missing data which we cannot account for, such as the healthiness of lunches brought from home.

Second, we use the survey to confirm that decisions regarding exercise time were orthogonal to lunch choice. Our respondents tend to prefer going to the gym after lunch rather than before lunch, with an average of 74% visits reported to be after lunch. Among those who reported going in the morning, the most common rationales for going before lunch was having more time before than after lunch for exercise, a preference for going before lunch, and the group fitness schedule. Among those who reported going in the afternoon, the overwhelming rationale was having more time after lunch than before for exercise, followed by preferences for going after lunch and the fitness schedule. These responses suggest that, for the most part, decisions about when to go to the gym are independent of lunch choice (i.e. nobody responded that they go to the gym in the morning in order to eat a less healthy lunch). Nevertheless, in our empirical work below, we will also consider specifications in which we allow for simultaneity in the timing of gym use and lunch food choices.

Additional questions from the survey were included for the purpose of shedding light on our empirical results. We will discuss these later in the paper.

2.4 Behavioral Results

Estimating the effect of morning exercise on lunch choice

The key empirical exercise is to estimate the impact of going to the gym before or after lunch on peoples’ lunch choices. We use the following fixed effects regression specification:

⁶Appendix D contains additional details on the results of the survey, including the questionnaire itself.

$$y_{it} = \beta \cdot \text{earlygym}_{it} + \gamma \cdot \text{lategym}_{it} + \alpha_i + \epsilon_{it} \quad (2.1)$$

where i indexes each individual and t indexes each day. y_{it} represents one of three measures of lunch choice: the overall health of the meal (average of the health scores of each of the items), the total price for the meal, and a weighted health of the meal (weighted average by price of each item). earlygym_{it} is a binary indicator denoting whether individual i went to the gym before lunch on day t , and lategym_{it} is a binary indicator denoting whether individual i went to the gym after lunch on day t . Finally, α_i denotes individual fixed effects, which are allowed to be arbitrarily correlated with choice of gym time on any given day. We assume that our error term, ϵ_{it} , is uncorrelated with the regressors. These regressions collapse all purchases at the individual day level. Specifically, we calculate the number of unique customers⁷ which made a purchase in each day over the observed 10-week time period. We then regress the number of unique customers on indicators of whether it was a Sale period or Post-Sale period at the time of purchase.

The estimation results of Equation (2.1) for three measures of lunch choice are presented in Table 2.5. From the table, we can see that early gym visits are consistently associated with negative effects on lunch healthiness (using two different metrics) and overall lunch cost. This result is clearest in the top panel, where we consider only days during which a given individual both ate lunch at the university cafeteria *and* used the university gym. In other words, conditional on using the gym and eating lunch on campus on a given day, we find that going to the gym before lunch (as opposed to after) is associated with meal health which is 4.7-4.9% lower. Furthermore, meal cost is reduced by 40.6%, indicating that individuals are not necessarily eating *more* food but rather eating less healthy food. (We also concede that some healthy choices, like the salad bar or sushi, tend to be more expensive than average.) This presents our first evidence that individuals fail in their self-regulation, as a restrained behavior made earlier in the day (going to the gym) subsequently leads to a more indulgent behavior (eating an unhealthy lunch) later on.

When we add days on which an individual did not go to the gym at all (but still had lunch on campus), going to the gym before lunch begins to look similar to not going to the gym at all. As can be seen in the second panel, later gym visits are consistently correlated with healthier lunches and higher overall lunch cost. Lunch healthiness tends to be 2.7-2.8% higher for individuals who go to the gym in the afternoon.

⁷We treat each unique credit card number as a unique customer.

The results in Table 2.5 indicate that the choice of a morning or afternoon gym visit has a significant impact on lunch choice, even when individual heterogeneity has been controlled for. However, given the sequential aspect of the events we are studying (where lategym necessarily occurs after lunch), we hesitate to claim that the observed effect is causal, as such a claim would involve a timing assumption that individuals plan their schedule (including what time they will go to the gym) at the beginning of the day, before making their meal choice decisions.

Table 2.5: Fixed effects regressions

	DV: mealhealth		DV: mealprice		DV: weighthealth	
	Est	StdErr	Est	StdErr	Est	Stderr
earlygym	-0.0473	0.0121†	-0.4061	0.0877†	-0.0493	0.0121†
Constant	0.2084	0.0031†	6.9010	0.0226†	0.2222	0.0031†
Indiv FE	yes		yes		yes	
N	30,119		30,119		30,119	

	DV: mealhealth		DV: mealprice		DV: weighthealth	
	Est	StdErr	Est	StdErr	Est	Stderr
earlygym	-0.0135	0.0086	-0.1303	0.0583○	-0.0164	0.0088*
lategym	0.0270	0.0052†	0.0937	0.0336○	0.0279	0.0053○
Constant	0.1549	0.0006†	6.7645	0.0040†	0.1726	0.0007†
Indiv FE	yes		yes		yes	
N	221,507		221,507		221,507	

“DV” denotes dependent variable in regression. Top panel: The sub-sample of days with both (i) lunch eaten on campus and (ii) gym use; Bottom panel: All days with lunch eaten on campus, regardless of gym use. * / ○ / †: significant at 10 / 5 / 1%.

A follow-up question then is whether a healthy lunch affects post-lunch gym use. We answer this in Table 2.6, where we regress our two measures of meal health (one of which is a pure average and another of which is a weighted average of the meal component health score) on our binary measure of whether an individual went to the gym after lunch. We find that both measures of meal health are positively correlated with later gym use. In other words, the opposite effect – that is, a restrained (healthy) lunch choice leading to a less healthy behavior (skipping exercise) later in the day – does not occur. It would appear then that this self-regulation failure operates asymmetrically, with gym use affecting subsequent lunch choices but not vice versa. Furthermore, we recognize that our health score may be capturing *perceived* as opposed to objective healthiness. Indeed, even among expert nutritionists, there is

debate regarding whether foods high in fat are unhealthy or not. We would like to say something about objective healthiness. Since most nutritionists agree that sugar intake should be limited, we can test whether we observe early gym goers being more likely to have sweeter lunches.

In Table 2.7, using two measures of meal sweetness (a straight average of the sweet scores associated with all items purchased, and a weighted average of the sweet scores of lunch items, weighted by price), we find that sweeter items are purchased by individuals when they go to the gym in the morning. Among all lunch purchases, those with morning gym visits have 3.8-3.9% sweeter meals, and among those with both lunch purchases and a gym visit at some point in the day, those with morning gym visits have 6.4-6.8% sweeter lunches.

These results raise the possibility that our results may be driven by a purely biological need for higher glucose intake following exercise. To examine this, in Table 2.8, we repeat the baseline regression, but only including *non-sweet* foods. The results are much the same as before, with early gym use associated with less healthy lunches subsequently. This suggests that higher glucose needs after exercise cannot fully explain the estimated effect.

Table 2.6: Healthiness of lunch and post-lunch gym use

	DV: lategym		DV: lategym	
	Est	StdErr	Est	StdErr
mealhealth	0.0077	0.0016†	—	—
weighthealth	—	—	0.0079	0.0016†
Constant	0.1033	0.0002†	0.1032	0.0003†
Indiv FE	yes		yes	
N	213,731		213,731	

“DV” denotes dependent variable in regression. Using all days with lunch eaten on campus and no gym use recorded before lunch. */o/†: significant at 10/5/1%.

Table 2.7: Results using sweet instead of unhealthy indicators

	DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr
earlygym	0.0640	0.0098†	0.0675	0.0099†
Constant	0.2188	0.0025†	0.1931	0.0026†
Indiv FE	yes		yes	
N	30,119		30,119	

	DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr
earlygym	0.0375	0.0072†	0.0392	0.0073†
lategym	-0.0157	0.0036†	-0.0163	0.0036†
Constant	0.2492	0.0005†	0.2206	0.0005†
Indiv FE	yes		yes	
N	221,507		221,507	

“DV” denotes dependent variable in regression. Top panel: The sub-sample of days with both (i) lunch eaten on campus and (ii) gym use; Bottom panel: All days with lunch eaten on campus, regardless of gym use. */o/†: significant at 10/5/1%.

Table 2.8: Truncated results using only non-sweet foods to define healthiness of lunch

	DV: mealhealth		DV: weighthealth	
	Est	StdErr	Est	StdErr
earlygym	-0.0380	0.0085†	-0.0389	0.0086†
Constant	0.2990	0.0022†	0.3084	0.0022†
Indiv FE	yes		yes	
N	30,119		30,119	

	DV: mealhealth		DV: weighthealth	
	Est	StdErr	Est	StdErr
earlygym	-0.0111	0.0061*	-0.0124	0.0062o
lategym	0.0191	0.0036†	0.0201	0.0038†
Constant	0.2626	0.0004†	0.2722	0.0005†
Indiv FE	yes		yes	
N	221,507		221,507	

“DV” denotes dependent variable in regression. Baseline regression with all *sweet* foods re-coded as *neither healthy nor unhealthy*. Top panel: The sub-sample of days with both (i) lunch eaten on campus and (ii) gym use; Bottom panel: All days with lunch eaten on campus, regardless of gym use. */o/†: significant at 10/5/1%.

IV specification

In the baseline model, we assume that *earlygym* and *lategym* are exogenous, after controlling for individual-level fixed effects. However, we worry that there may be common individual-level shocks which affect both the timing of gym use, and the healthiness of food choices. For instance, an individual may make New Year’s resolutions for healthier living encompassing both a morning gym regimen as well as a healthier diet.

Table 2.9: Baseline regressions: using IV’s for *earlygym* and *lategym*

	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-0.5668	0.0922†	-0.6510	0.0945†	0.7063	0.0666†	0.7412	0.0688†
Constant	0.3425	0.0242†	0.3776	0.0248†	0.0530	0.0168†	-0.0191	0.0174†
Indiv FE	yes		yes		yes		yes	
N	30,119		30,119		30,119		30,119	
	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-1.9009	0.3821†	-2.4520	0.4284†	2.0263	0.3505†	2.3086	0.3681†
lategym	1.4985	0.1425†	1.7935	0.1629†	-2.0984	0.1509†	-2.2068	0.1603†
Constant	0.0727	0.0219†	0.0801	0.0248†	0.3895	0.0210†	0.3619	0.0223†
Indiv FE	yes		yes		yes		yes	
N	221,507		221,507		221,507		221,507	

2SLS using dummies for day-of-week instruments for *earlygym*, *lategym*. “DV” denotes dependent variable in regression. First-stage regression results are in Appendix B. Standard errors in these regressions are computed via bootstrap.

Top panel: The sub-sample of days with both (i) lunch eaten on campus and (ii) gym use; Bottom panel: All days with lunch eaten on campus, regardless of gym use. */†: significant at 10/5/1%.

To accommodate such heterogeneity, we next consider specifications in which *earlygym* and *lategym* are allowed to be endogenous, and instrumented with day-of-week dummies. Day-of-week dummies capture exogenous variation in morning gym use arising from institutional features of undergraduate class scheduling at the university. The majority of classes at the university are on either a Tuesday-Thursday or Monday-Wednesday-Friday weekly cadence. Moreover, on M-W-F, classes are scheduled at hourly intervals, while on T-Th, classes are scheduled at 90-minute intervals. As such we expect that undergraduate students (who make up the bulk of our sample) may have differing time availability for morning exercise on M-W-F vs. T-Th. Indeed, for all three academic quarters in 2018, M-W-F classes experienced the bulk of their enrollment (49-55%) *before lunch*, whereas T-Th classes had the

majority of their enrollment (54-66%) *after lunch*. The first-stage regressions are shown in Table 2.16 in Appendix B.

As we see from Table 2.9, the findings remain robust and significant in the IV specifications, with early gym-goers being more likely to have less healthy and more sweet lunches than late gym-goers.

Which food items drive the main results?

Next we explore the differences in specific food items chosen by early and late gymgoers. For each of the forty specific lunch items tabulated in our dataset, we ran a regression of *earlygym* on an indicator for whether that item was included in the lunch meal. For comparison with the baseline regressions, we also included individual fixed effects in these regressions, so that the estimated coefficient measures the deviation from the average propensity to consume each food item for each individual.

Table 2.10: Food choice differences between early and late gym-goers

	Food item	Health score	Sweet	Coefficient	Significance
1	baked goods	-0.5965	1	0.0038	
2	breakfast	0.2281	0	0.0008	
3	burger	-0.4386	0	-0.0032	
4	burrito	-0.1404	0	0.0021	†
5	candies	-0.8948	1	0.0073	
6	cereals	-0.0877	1	0	
7	chips	-0.6667	0	-0.0002	
8	coffee	0.0178	0	0.021	†
9	deli	0.7017	0	-0.0107	†
10	dessert	-0.8246	1	-0.0021	
11	eggs/omelette	0.7544	0	0.0004	
12	energy drink	-0.8246	1	0.0136	†
13	fries	-0.8421	0	0.0011	*
14	fruit	0.8597	1	0.0002	
15	grill sandwich	0.2632	0	-0.0144	†
16	icecream	-0.8246	1	0.0011	
17	juice	-0.0176	1	0.0077	†
18	meat	0.3158	0	-0.0136	†
19	milk	0.6316	0	0.0019	
20	mongolian	0.1404	0	-0.0046	o
21	noodles	0.0714	0	0.0003	
22	pasta	0.1579	0	0	
23	pizza	-0.4035	0	-0.009	†
24	popcorn	-0.193	0	0.0007	
25	processed food	-0.7544	0	0.0007	
26	protein snacks/bars	0.4386	1	0.0095	o
27	salad	0.9298	0	-0.0227	†
28	seafood	0.5893	0	-0.0231	†
29	smoothie	0.2105	1	0.0239	†
30	snack	0.0176	0	0.0011	
31	soda	-0.807	1	-0.0008	
32	soup	0.7017	0	-0.0016	
33	special	0.1579	0	-0.0013	
34	taco	0.0877	0	-0.0001	
35	tea	0.7895	0	0.0073	o
36	veggies	0.8948	0	-0.0009	
37	vitamin water	0.2457	1	0.0001	
38	water	0.9649	0	0.004	o
39	yogurt	0.7544	1	-0.0002	

“DV” denotes dependent variable in regression. In each row, we report the coefficient from a regression of *earlygym* on an indicator for each lunch item. Subsample of days with both (i) lunch eaten on campus and (ii) gym use. Regression includes individual fixed effects. */o/†: significant at 10/5/1%.

Results are summarized in Table 2.10, and appear to indicate that a handful of food items drive the observed behaviors. Among healthy foods, early gymgoers choose salad (#27) and deli sandwiches (#9) less frequently. At the same time, they choose energy drinks (#12) and, more marginally, fries (#13) more frequently.⁸ Interestingly, other unhealthy foods, such as soda (#31), baked goods (#1), and ice cream (#16), which are also likely high in sugar content, are no more likely to be chosen by early gym goers. Pizza (#23), another prominent unhealthy food choice, is actually less likely to be chosen by early gymgoers.

2.5 Mechanism

Biological need - time lag between exercise and lunch

Thus far, we have focused on a failure of self-regulation or self-control as the explanation for the main empirical results. Here we consider an alternative explanation for our results, namely that after physical exertion and exercise, there is a biological need for more calorie-dense or high-glucose foods, which would explain why early gym-goers may eat less healthy and sweeter foods at lunchtime.

Table 2.11: Timing of lunch and gym use

Timestamp of lunch hh:mm:ss					
	min	25%	50%	75%	max
All days	11:00:00	12:04:20	12:39:53	13:42:39	16:00:00
earlygym=1	11:00:11	12:06:25	12:53:19	14:06:03	15:59:52
lategym=1	11:00:01	12:03:36	12:38:24	13:32:50	15:59:55
Timestamp of gym check-in hh:mm:ss					
	min	25%	50%	75%	max
All days	06:00:00	12:56:00	16:42:00	19:02:00	22:14:00
earlygym=1	06:00:00	08:08:00	09:35:00	10:55:00	15:19:00
lategym=1	11:12:00	16:11:00	17:55:00	19:46:00	22:14:00
Time difference lunch-gym (hours)					
	mean	stdev	min	max	<i>N</i>
All days	-2.78	4.34	-10.98	9.87	30,119
earlygym=1	3.49	1.99	0.01	9.87	7,776
lategym=1	-4.97	2.35	-10.98	-0.09	22,343

We start with some additional analyses of the data. We examine the time lag

⁸Energy drinks include NOS Energy (which contains 41g of sugar), Monster Energy (27g of sugar), and Monster Lo-Carb (7g of sugar). For context, the American Heart Association (AHA) recommends no more than 37.5 grams of sugar for men and no more than 25 grams for women per day.

between early gym use and lunchtime, as summarized in Table 2.11. As the bottom panel in this table shows, the time difference between gym use and lunch is quite large: on average, early gym-goers eat lunch 3.49 hours after the beginning of their gym session in the morning (the median, not reported, is 3.1 hours). This is not unexpected; as the top panel in the table shows, the median check-in time for early gym goers is 9:35am, while the median lunchtime is 12:53pm, suggesting that many early gym-goers work out before their morning classes, and have lunch after those classes. Such a long lag between gym use and lunch casts doubt on the biological mechanism as an explanation of our results. If post-workout thirst or hunger were of prime importance, we imagine early gym goers would have satisfied these urges during this lag, long before they went to lunch.

Nevertheless, to test this, we introduce this *timelag* variable into our regressions. Regression results are in Table 2.12. In the regressions involving *mealhealth* and *weighthealth*, the negative coefficients on the interaction terms of *timelag* and *earlygym* show that the *longer* the time lag between exercise and eating lunch, the less healthy a lunch choice will be, which suggests that the gym workout does not immediately trigger the less healthy lunchtime choices. In the regressions involving *mealsweet* and *weightsweet*, this interaction is insignificant. Overall, we have little evidence that a real or perceived biological mechanism plays a key role in explaining our results.

Table 2.12: Including *timelag* between morning gym use and lunch

	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	0.0089	0.0231	0.0107	0.0232	0.0506	0.0197†	0.0483	0.0199○
early* <i>timelag</i>	-0.0178	0.0058†	-0.0190	0.0058○	0.0043	0.0046	0.0061	0.0046
constant	0.2099	0.0032†	0.2239	0.0032†	0.2184	0.0025†	0.1926	0.0025†
Indiv FE	yes		yes		yes		yes	
<i>N</i>	30,119		30,119		30,119		30,119	

“DV” denotes dependent variable in regression. *Timelag* is defined as the number of hours between morning gym use and lunch. These regressions utilize a subsample of days with both (i) lunch eaten on campus and (ii) gym use. */○/†: significant at 10/5/1%.

In confirmation of these empirical findings, the medical literature also offers no clear-cut evidence that the body requires more glucose after exercise (Hopkins et al., 2011). Specifically, we know that little to no carbohydrates are metabolized at lower intensity workouts, and most of what the average exercise leads to metabolizing

is fat stores (Sherman, 1995).

However, there is some psychological evidence that exercisers may use a *perceived* bodily need for glucose as an *excuse* to consume more sugary foods after exercise (Job, Walton, et al., 2013). Such behavior is consistent with the self-regulation literature, which states that people typically *choose* to allow self-regulation to fail, often due to feeling tired which makes the exertion of self-control less appealing. To quote Baumeister, Heatherton, and Tice, 1994: “Although it is very difficult to obtain decisive empirical data regarding the issue of acquiescence, we suspect that acquiescence is the norm, not the exception.” Seen in this light, this perceived biological explanation may be less a confound than a *channel* through which self-licensing operates.

Corroborating survey evidence on gym use and lunch choices

To assess this alternative explanation directly, we included questions in our survey directly pertaining to the link between respondents’ gym use and lunch choices. Specifically, we assessed whether gym-goers would attribute their lunch choices to a perceived biological necessity. When asked why they might choose to eat a less healthy lunch, only 24% reported that it would be due to having fatigued themselves, either mentally or physically, earlier in the day. The overwhelming majority (60.3%) reported simply craving something which happened to be unhealthy, and the second most popular rationale, at 30.8%, was due to feeling down or stressed (participants were allowed to choose multiple reasons). This is consistent with previous work which found that when people are asked why they overeat, they “often report doing so because they are sad, bored, or otherwise in a bad mood.” (Baumeister, Heatherton, and Tice, 1994). In other words, even when asked directly, only a minority reference prior physical exertion as the rationale for their unhealthy meal choice.

Furthermore, the literature on self-licensing emphasizes that the restrained choice must “increase one’s sense of positive self-concept” in order for licensing to operate. As a direct check for self-licensing, we asked survey-takers how they felt about themselves after exercising. An overwhelming 92.5% of respondents reported “feeling better about themselves” after having worked out at the gym, with none saying they felt worse about themselves and only 7.5% reporting they felt no better or worse. Among the undergraduate respondents, 100% reported “feeling better about themselves” after having worked out at the gym. This generates support that going to the gym is a behavior which creates a positive self-concept, a necessary pre-condition

for licensing.

The survey was sent to a general university-wide electronic mailing list, and we did not oversample groups with low response rates to get a fully representative sample. Instead, the responses organically generated a partially representative survey across person types with respondents composed of 8.2% undergraduates, 35.6% graduate students, 17.1% post-doctoral scholars, 36.3% staff, and 2.7% other. If we only focus on the responses of undergraduates (who make up the bulk of our lunch transactions at nearly 55%), the results are very similar. Only 25% reported that an unhealthy lunch choice would be due to having fatigued themselves earlier in the day. Again the majority (75%) reported simply craving something which happened to be unhealthy, and the second most popular rationale, now at 33%, was due to feeling down or stressed.

2.6 Discussion

A key motivation for this paper was to seek evidence of spillovers from a behavior in one domain (exercise) onto another domain (lunch choice). Various literatures have looked at these behaviors in silos, with most studies implicitly assuming that increasing a “positive” behavior like exercise is net positive for an individual. Few studies have considered possible (negative) spillovers into other behavioral domains. Inspired by a small group of closely related papers, we wanted to see whether we could find field evidence of self-regulation failures, where a positive behavior in a domain like exercise could lead to a more indulgent choice in an unrelated domain like lunch choice. The novelty of our paper was to look at how natural (rather than incentivized) changes in exercise systematically affect food choice, thus empirically identifying spillovers across two separate behavioral domains.

To do so, we obtained daily panel data from a university’s on-campus cafeterias and fitness recreation center. Using this dataset, we were able to test whether exercising (a positive behavior in absolute terms) prior to having lunch increases the likelihood of making an unhealthy lunch choice (a negative spillover). We find that, controlling for individual fixed effects, there is a robust effect of morning exercise on the healthiness of a lunch choice. Specifically, healthy lunch items such as salad are less likely, and unhealthy items (such as fries or energy drinks) more likely, to be purchased at lunchtime following morning gym use. Across different slices of the university population, faculty are no less likely to exhibit this behavior than undergraduates. Women are also no less likely to exhibit this behavior than men.

Interestingly, we do not observe self-licensing in the opposite direction: those who did not use the gym in the morning and eat a healthy lunch are actually *more* likely to use the gym in the afternoon. We posit that this may be explained by reference-dependent preferences in which exercise accumulates more “virtue points” than a healthy meal choice. In addition, this effect obtains even given the large gap (three hours on average) between morning gym use and exercise, which suggests that post-workout dehydration or calorie depletion cannot fully explain the result.

We conclude with two broader implications of our findings. First, our evidence confirms that self-licensing is an empirically relevant feature of behavior “in the wild” and well-intentioned policies which promote virtuous or healthy behavior may backfire without recognizing that the benefits from that behavior may be undone by compensating non-virtuous or unhealthy behaviors. Second, individual choices may be interlinked across distinct domains; hence, looking at behaviors through a partial equilibrium lens, in which we assume that whatever happens in one realm does not affect other realms, may yield an incomplete view of human behavior and policy effects. Our field evidence therefore also fits into a broader literature documenting what economists call “unintended consequences” - typically studied in the context of public policies which backfire, such as how mandated safety measures may lead to more risky behaviors which “offset” the benefits of the regulation (e.g., Peltzman, 1975). A potential biological rationale is that people may require a certain type of “homeostasis,” such that intervening in one domain inevitably leads to spillover effects in other domains. If this is indeed the case, it would imply that behavior change initiatives may need to be more holistic (and indeed, perhaps more paternalistic) in order to have a long-lasting effect of behavior.

These findings have important implications for marketing. First, our results offer a new key insight for the design and administration of marketing campaigns. Broadly speaking, any campaign promoting a product or service in one domain may backfire if it causes countervailing choices in another domain. For instance, the promotion of a more active, healthy lifestyle – such as Michelle Obama’s “Let’s Move” exercise campaign – may, according to our results, lead to more indulgent and unhealthy food choices: more sweet drinks, less salad. To that end, the proper design of such campaigns (often done in advance of their official release) would do well to include a wider range of behavioral measures and outcomes in order to account for spillovers into unexpected domains. More rigorous testing should lead to better forecasts regarding whether the new release is a net positive addition for a company.

At the same time, our field evidence in support of self-regulation failure points to an opportunity for advertisers to caution consumers against possible failures of self-control. On the other hand, advertisers can exploit the self-regulation failures by encouraging indulgent purchases at times when people may be feeling restrained. Indeed, several famous taglines already seem to do so (for example, L'Oréal's famous slogan, "Because you're worth it.").

Finally, this analysis also highlights the scientific benefits of richer comprehensive datasets, which enable us to observe multiple actions by the same individual not only across time, but also across different domains. There is value in seeking and constructing "wide" holistic datasets which generate new insights about how behaviors in different domains interact for the same individuals. Follow-up studies could look at various combinations of other behavioral domains to identify the relevant realms across which behavioral spillovers exist.

Appendix A - Robustness Checks

Demographic decomposition

The self-regulation literature asserts that older individuals have a less limited view of willpower than younger individuals (Job, Sieber, et al., 2018), and some hypothesize that there are gender differences with respect to beliefs about willpower as well. In order to assess whether these demographic variables may be moderating the observed effect, we first introduce a gender variable which we then interact with our earlygym and lategym regressors. The results of this estimation are in Table 2.13. Gender seems to be largely insignificant, except for a small effect on meal sweetness, with women choosing less sweet lunches than men conditional on having exercised in the morning (significant at $p < 0.10$). When the data is truncated to individuals on days when a gym visit and a lunch purchase was observed, there is no significant effect of gender on lunch health or sweetness.

Table 2.13: Breakdown by gender

	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-0.0388	0.0170○	-0.0426	0.0171○	0.0740	0.0131†	0.0773	0.0133†
early*Fem	-0.0222	0.0231	-0.0174	0.0231	-0.0257	0.0196	-0.0253	0.0198
Constant	0.2081	0.0032†	0.2220	0.0032†	0.2185	0.0026†	0.1927	0.0026†
Indiv FE	yes		yes		yes		yes	
N	30,119		30,119		30,119		30,119	
	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-0.0166	0.0114	-0.0202	0.0116*	0.0470	0.0094†	0.0489	0.0095†
lategym	0.0222	0.0061†	0.0229	0.0062†	-0.0147	0.0044†	-0.0147	0.0044†
early*Fem	0.0090	0.0172	0.0111	0.0175	-0.0267	0.0141*	-0.0274	0.0145*
late*Fem	0.0149	0.0117	0.0154	0.0118	-0.0030	0.0077	-0.0050	0.0077
Constant	0.1549	0.0006†	0.1727	0.0007†	0.2492	0.0005†	0.2206	0.0005†
Indiv FE	yes		yes		yes		yes	
N	221,507		221,507		221,507		221,507	

“DV” denotes dependent variable in regression. Top panel: subsample of days with both (i) lunch eaten on campus and (ii) gym use. Bottom panel: using all days with lunch eaten on campus. * / ○ / †: significant at 10 / 5 / 1%.

To test the hypothesis with respect to age, we use “type” as a proxy for average age. We do not observe the actual age of each individual, but we do know whether someone is a student or a faculty member. We can therefore include interaction terms for each position type to see whether the observed effect is stronger for

younger members of the population (undergraduates, typically 18-22 years old) as opposed to graduate students and post docs (typically early-20's to mid-30's) and older members of the population (faculty and staff, typically mid-30's to 70's). Our regression results are displayed in Table 2.14.

Contrary to the hypothesis that older individuals may exhibit less of an effect, we see no significant differences in the interaction between *earlygym* and *faculty*. However, staff members, who are also generally older than undergraduates, appear to have a smaller *earlygym* effect, choosing healthier and less sweet lunches than the average individual. Graduate students, who are typically older than undergraduates but younger than faculty, also display less of an *earlygym* effect, choosing on average healthier and less sweet lunches. Overall, then, there is no consistent trend here, and we cannot draw strong conclusions about age moderating the observed effect.

Table 2.14: Breakdown by position type

	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-0.0641	0.0148†	-0.0655	0.0147†	0.0777	0.0124†	0.0827	0.0124†
early*fac	-0.0486	0.0748	-0.0355	0.0781	0.0004	0.0655	-0.0102	0.0681
early*grd	0.0712	0.0282○	0.0667	0.0284○	-0.0473	0.0197○	-0.0500	0.0200†
early*pcd	0.0688	0.0831	0.0616	0.0855	-0.0614	0.0482	-0.0558	0.0497
early*stf	0.0705	0.0344○	0.0741	0.0354○	-0.0805	0.0459*	-0.1028	0.0491○
Constant	0.2093	0.0036†	0.2229	0.0037†	0.2189	0.0030†	0.1933	0.0031†
Indiv FE	yes		yes		yes		yes	
N	30,119		30,119		30,119		30,119	
	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-0.0154	0.0109	-0.0168	0.0111	0.0467	0.0094†	0.0483	0.0095†
lategym	0.0390	0.0073†	0.0404	0.0073†	-0.0230	0.0052†	-0.0253	0.0051†
early*fac	-0.0417	0.0317	-0.0365	0.0333	-0.0194	0.0249	-0.0259	0.0272
early*grd	0.0209	0.0208	0.0140	0.0212	-0.0256	0.0150*	-0.0233	0.0152
early*pcd	-0.0180	0.0447	-0.0244	0.0459	-0.0156	0.0254	-0.0140	0.0265
early*stf	0.0877	0.0431○	0.0802	0.0448*	-0.1074	0.0432○	-0.1101	0.0459*
late*fac	-0.0035	0.0205	0.0026	0.0215	-0.0080	0.0149	-0.0063	0.0161
late*grd	-0.0333	0.0117†	-0.0345	0.0119	0.0203	0.0080○	0.0233	0.0080†
late*pcd	-0.0157	0.0204	-0.0213	0.0204	0.0105	0.0146	0.0191	0.0129
late*stf	0.0598	0.0196†	0.0569	0.0221†	-0.0201	0.0107*	-0.0166	0.0120
Constant	0.1549	0.0006†	0.1727	0.0007†	0.2493	0.0005†	0.2207	0.0005†
Indiv FE	yes		yes		yes		yes	
N	221,507		221,507		221,507		221,507	

“DV” denotes dependent variable in regression. Top panel: subsample of days with both (i) lunch eaten on campus and (ii) gym use. Bottom panel: using all days with lunch eaten on campus. */○/†: significant at 10/5/1%.

Stress and exam periods

The self-regulation literature also contains some evidence that higher levels of stress are associated with greater rates of self-regulation failure. A university is a demanding academic institution which imposes substantial challenges and stress on its students. We may therefore be worried that the observed effects (unhealthy lunch choices) are being driven, or amplified, by stress. While we do not have direct measures of stress (such as individual cortisol measurements during lunch purchase), we can infer levels of stress by time-specific variables, like days associated with a heavy exam load. Therefore, we label each day of 2018 as having fallen (or not) on a midterm or finals week. These exam weeks are notoriously stressful periods, particularly for students.

Our regression results are displayed in Table 2.15. When considering all days during which lunch is eaten on campus, we see an effect of the exams on their own, with midterms increasing the likelihood of a sweet lunch purchase, and decreasing meal health. However, we observe no significant interaction effects between these exam weeks and our earlygym regressor. Further, when truncating the data to a subsample of those with both lunch purchases on campus and observed gym activity, the pure exam effects largely go away (although midterms now seem to be negatively correlated with meal health), and again no interaction effects persist. The data do not speak unisonally here, but there is little evidence that exam periods, which we believe are good proxies for stress, are either moderating or confounding our main results.

Table 2.15: Breakdown by midterm/final weeks

	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-0.0456	0.0122†	-0.0484	0.0121†	0.0633	0.0098†	0.0670	0.0099†
midterm	-0.0211	0.0123*	-0.0268	0.0125○	-0.0025	0.0083	-0.0033	0.0081
final	-0.0120	0.0148	-0.0075	0.0150	0.0061	0.0109	0.0068	0.0108
early*mid	0.0064	0.0230	0.0150	0.0232	-0.0058	0.0181	-0.0054	0.0182
early*fin	-0.0430	0.0336	-0.0390	0.0341	0.0224	0.0258	0.0173	0.0250
constant	0.2106	0.0033†	0.2247	0.0033†	0.2187	0.0026†	0.1930	0.0026†
Indiv FE	yes		yes		yes		yes	
N	30,119		30,119		30,119		30,119	
	DV: mealhealth		DV: weighthealth		DV: mealsweet		DV: weightsweet	
	Est	StdErr	Est	StdErr	Est	StdErr	Est	StdErr
earlygym	-0.0119	0.0089	-0.0152	0.0091*	0.0378	0.0072†	0.0396	0.0074†
lategym	0.0271	0.0054†	0.0284	0.0055†	-0.0145	0.0037†	-0.0153	0.0037†
midterm	-0.0059	0.0042	-0.0064	0.0042	0.0130	0.0031†	0.0126	0.0031†
final	-0.0247	0.0048†	-0.0216	0.0048*	0.0035	0.0035	0.0017	0.0035
early*mid	-0.0113	0.0200	-0.0078	0.0199	-0.0153	0.0161	-0.0165	0.0163
early*fin	-0.0265	0.0293	-0.0206	0.0298	0.0196	0.0232	0.0174	0.0223
late*mid	-0.0121	0.0130	-0.0169	0.0131	-0.0154	0.0089*	-0.0154	0.0086*
late*fin	0.0098	0.0154	0.0113	0.0156	0.0015	0.0110	0.0033	0.0110
constant	0.1569	0.0008†	0.1745	0.0008†	0.2480	0.0006†	0.2196	0.0006†
Indiv FE	yes		yes		yes		yes	
N	221,507		221,507		221,507		221,507	

“DV” denotes dependent variable in regression. Top panel: subsample of days with both (i) lunch eaten on campus and (ii) gym use. Bottom panel: using all days with lunch eaten on campus. */○/†: significant at 10/5/1%.

Appendix B - Additional Results

Table 2.16: First stage regressions

	DV: earlygym		DV: lategym	
	Est	StdErr	Est	StdErr
Sunday	-0.0092	0.0052	0.0027	0.0054
Monday	-0.0112	0.0050 ^o	0.0518	0.0052 [†]
Tuesday	-0.0161	0.0050 [†]	0.0518	0.0050 [†]
Wednesday	-0.0147	0.0051 [†]	0.0443	0.0047 [†]
Thursday	-0.0199	0.0050 [†]	0.0443	0.0048 [†]
Friday	-0.0222	0.0049 [†]	0.0362	0.0050 [†]
Constant	0.0512	0.0047 [†]	0.0573	0.0043 [†]
Indiv FE	yes		yes	
<i>F</i> -stat (6,2995)	9.85		24.10	
<i>N</i>	221,507		221,507	

“DV” denotes dependent variable in regression. These are the first stage regressions underlying the instrumental variable specifications in Table 2.9. ^{*}/^o/[†]: significant at 10/5/1%.

Appendix C - Additional Data Sources

1. Registrar Data: We obtained 2018 undergraduate class enrollment information from the office of the registrar (via a data request for research purposes). This data listed all classes offered in 2018, the days of the week and times during which they were taught, and the total number of enrolled students. This data provided additional context for interpreting the self-licensing effect by day of the week.
2. Academic Calendar: We obtained the 2017-2018 and 2018-2019 academic calendars, which are publicly available online. These calendars allowed us to isolate midterm and final exam weeks in order to test the hypothesis that the observed self-licensing effect may be driven by stress.

Appendix D - Survey Format and Results

1. "I am a... (Please select one)"
 - Undergraduate (8.22%)
 - Graduate Student (35.62%)
 - Post-doctoral Scholar (17.12%)
 - Faculty (0.00%)
 - Staff (36.30%)
 - Other (2.74%)

2. "I typically use the University Gym and/or Recreation Center... (Please select one)"
 - More than once a week (40.41%)
 - Once a week (9.59%)
 - 1-3 times a month (10.96%)
 - Once a month (7.53%)
 - Less than once a month (31.51%)

3. "Consider 10 times you have gone to the University Gym and/or Recreation Center. What percentage of these visits happened after lunch?"
 - Mean response: 74%

4. "If I go to the University Gym and/or Recreation Center before lunch it's because... (Please select one)"
 - I have more time before lunch than after that day (23.29%)
 - I prefer going before lunch (19.86%)
 - I never go before lunch (45.89%)
 - Other (10.96%)

5. "If I go to the University Gym and/or Recreation Center after lunch it's because... (Please select one)"
 - I have more time after lunch than before that day (45.89%)

- I prefer going after lunch (23.29%)
 - I never go after lunch (15.07%)
 - Other (15.75%)
6. "When I work out, I tend to feel... (Please select one)"
- Better about myself (92.47%)
 - Worse about myself (0.00%)
 - Neither better nor worse about myself (7.53%)
7. "When I work out, I typically spend [how many?] minutes at the University facilities."
- Mean response: 61 minutes
8. "When I work out at the gym, I typically burn about [how many?] calories during a session."
- Mean response: 434 calories
9. "I typically buy lunch at... (Please select one)"
- University Cafeteria (22.60%)
 - University Deli (4.79%)
 - University Cafe (4.11%)
 - None, I bring my own lunch (61.64%)
 - Other (6.85%)
10. "The main reason I typically buy lunch at the place chosen above is... (Please select one)"
- Convenience (it is closest to me) (22.60%)
 - Better food choice (I prefer the food options here) (19.86%)
 - Cost (better value) (41.10%)
 - Other (16.44%)

11. "On days that I choose an unhealthy lunch, it's likely because... (Select all that apply)"

- I'm craving something which happens to be unhealthy (60.27%)
- I fatigued myself, mentally and/or physically, earlier in the day (23.97%)
- I'm feeling down or stressed (30.82%)
- Other (22.60%)

References

- Acland, D. and M.R. Levy (2015). "Naiveté, projection bias, and habit formation in gym attendance". In: *Management Science* 61.1, pp. 146–160.
- Adda, J. and F. Cornaglia (2006). "Taxes, cigarette consumption, and smoking intensity". In: *American Economic Review* 96.4, pp. 1013–1028.
- Baumeister, R., T. Heatherton, and D. Tice (1994). *Losing control*. San Diego: Academic Press.
- Black, J. (2020). "How Google Got Its Employees to Eat Their Vegetables". In: *OneZero, a Medium publication*.
- Blain, B. et al. (2019). "Neuro-computational impact of physical training overload on economic decision-making". In: *Current Biology* 29.19, pp. 3289–3297.
- Blaisdell, B. (2012). *The Wit and Wisdom of Abraham Lincoln*. Dover Publications.
- Charness, G. and U. Gneezy (2009). "Incentives to exercise". In: *Econometrica* 77.3, pp. 909–931.
- De Witt Huberts, J., C. Evers, and D. De Ridder (2012). "License to sin: Self-licensing as a mechanism underlying hedonic consumption". In: *European Journal of Social Psychology* 42.4, pp. 490–496.
- (2013). "'Because I am worth it' A theoretical framework and empirical review of a justification-based account of self-regulation failure". In: *Personality and Social Psychology Review* 18.2, pp. 119–138.
- (2014). "Thinking before sinning: Reasoning processes in hedonic consumption". In: *Frontiers in Psychology* 5, p. 1268.
- DellaVigna, S. and U. Malmendier (2006). "Paying not to go to the gym". In: *American Economic Review* 96.3, pp. 694–719.
- Dolan, P. and M. Galizzi (2014). "Because I'm worth it. A lab-field experiment on the spillover effects of incentives in health". In: *LSE CEP Discussion Paper CEPDP 1286, London*.
- (2015). "Like ripples on a pond: Behavioral spillovers and their implications for research and policy". In: *J. Econ. Psychol* 47, pp. 1–16.

- Dubois, P., R. Griffith, and M. O'Connell (2018). "The effects of banning advertising in junk food markets". In: *Review of Economic Studies* 85.1, pp. 396–436.
- Duckworth, A., K. Milkman, and D. Laibson (2018). "Beyond willpower: Strategies for reducing failures of self-control". In: *Psychological Science in the Public Interest* 19.3, pp. 102–129.
- Hopkins, M. et al. (2011). "The relationship between substrate metabolism, exercise and appetite control: Does glycogen availability influence the motivation to eat, energy intake or food choice?" In: *Sports Medicine* 41.6, pp. 507–21.
- Job, V., V. Sieber, et al. (2018). "Age differences in implicit theories about willpower: Why older people endorse a nonlimited theory". In: *Psychology and Aging* 33.6, pp. 940–952.
- Job, V., G. M. Walton, et al. (2013). "Beliefs about willpower determine the impact of glucose on self-control". In: *PNAS* 110.37, pp. 14837–14842.
- Karmarkar, U. and B. Bollinger (2015). "BYOB: How bringing your own shopping bags leads to treating yourself and the environment". In: *Journal of Marketing* 79.4, pp. 1–15.
- Khan, U. and R. Dhar (2006). "Licensing effect in consumer choice". In: *Journal of Marketing Research* 43.2, pp. 259–266.
- Kivetz, R. and T. Simonson (2002). "Earning the right to indulge: Effort as a determinant of customer preferences toward frequency program rewards". In: *Journal of Marketing Research* 39.2, pp. 155–170.
- Peltzman, S. (1975). "The Effects of Automobile Safety Regulation". In: *Journal of Political Economy* 83.4, pp. 677–726.
- Royer, H., M. Steher, and J. Sydnor (2015). "Incentives, commitments, and habit formation in exercise: Evidence from a field experiment with workers at a Fortune-500 company". In: *American Economic Journal: Applied Economics* 7.3, pp. 51–84.
- Sherman, W. M. (1995). "Metabolism of sugars and physical performance". In: *The American Journal of Clinical Nutrition* 62, 228S–241S.
- Thaler, R. and S. Benartzi (2004). "Save more tomorrow (TM): Using behavioral economics to increase employee saving". In: *Journal of Political Economy* 112.S1, S164–S187.
- Volpp, K. G. et al. (2009). "A randomized, controlled trial of financial incentives for smoking cessation". In: *The New England Journal of Medicine* 360.7, pp. 699–709.

*Chapter 3***PREDICTING CONTEXT-SENSITIVITY OF BEHAVIOR IN
FIELD DATA**

ABSTRACT

In this chapter, we introduce a machine learning approach to characterizing individual trajectories of behavioral predictability. Predicting Context-Sensitivity (PCS) identifies a person-specific set of variables that maximizes the prediction of behavior over successive occasions. We apply PCS to two large, long-term field panel data sets tracking (1) hospital caregivers' hand-sanitizing and (2) gym members' gym attendance. We find that while past behavior is nearly universally predictive of future behavior, different subsets of context variables (e.g., day of the week, peer behavior) are predictive for different people. We also find that the time it takes for behavior to become predictable is domain-specific: we estimate that it takes 5 to 10 months to develop a predictable gym behavior and 2 to 3 weeks to develop predictable handwashing behavior. PCS reveals heterogeneity in behavioral trajectories outside of the laboratory, both across individuals and domains, and lays a new analytic foundation for studying personalized behavior change.

3.1 Introduction

Much of human behavior is habitual. Unlike choices that are consciously deliberated, habits constitute “a specific form of automaticity in which responses are directly cued by the contexts (e.g. locations, preceding actions) that consistently covaried with past performance” (Wood and Neal, 2009). By definition, when a behavior becomes habitual, the next time a familiar context is encountered, the choice previously made in that context is repeated automatically with high probability (Neal, Wood, and Quinn, 2006, Wood and Neal, 2007, Wood and Neal, 2009, Camerer, Landry, and Webb, 2021). Much habitual behavior therefore has two key hallmarks: predictable context-sensitivity (Ji and Wood, 2007, Danner, Vries, and Aarts, 2008) and automaticity (Orbell and Verplanken, 2010, Gardner, Abraham, et al., 2012). In this paper, we present a novel methodology for identifying behaviors that are highly and predictably context-sensitive, and thus candidates for being habitual. Unfortunately, the field data we use in this demonstration study do not enable us to measure automaticity.

There is a large body of laboratory research documenting the mechanisms underlying well-developed habits in animals and humans (see Appendix A). However, there is much less field research on how human habits naturally develop over time (Verplanken, 2018, pg. 7). To our knowledge, only three observational studies have studied habit formation over time in the wild, all of which relied upon research volunteers who completed daily self-report questionnaires. The earliest of these is Lally et al., 2010, in which 96 undergraduate volunteers were asked to carry out an eating, drinking, or exercise behavior daily in the same context for 12 weeks and to self-report habit strength daily. This study, as well as Kaushal and Rhodes, 2015 and Fournier et al., 2017, suggested that “habit typically develops asymptotically and idiosyncratically, potentially differing in rate across people, cues and behaviors” (Gardner and Lally, 2018, pg. 220).

Our investigation introduces a new approach to studying context-sensitive behavior in the wild: we use machine learning to statistically classify when choices are predicted by an identifiable set of context variables, and we identify which context cues tend to be the same or different across people. Specifically, Predicting Context-Sensitivity (PCS) uses a least absolute shrinkage and selection operator (LASSO) regression to identify variables that best predict behavior for each individual in a dataset. This machine learning technique generates a person-specific measure of behavioral predictability, which can then be used to study individual differences in

predictability and speed of habit formation.

While some psychologists have investigated a broader set of time, space, and social cues (Mazar and Wood, 2018), the identification of person-specific contextual variables that predict behavior represents a leap forward for the study of habits. In economic models of habit, the only context variable of interest has been prior behavior (Becker and Murphy, 1988, Dubé, Hitsch, and Rossi, 2010). This narrow focus on history-dependence is also evident “in applied social psychology, as a well as other areas such as health, social medicine, or education, and may have stalled progress in habit theory for quite some time” (Verplanken, 2018, pg. 3).

The datasets in this field study are well-suited to the application of PCS. Importantly, our panel data extends over a full year for gym members and healthcare workers, respectively, including many observations of each individual’s behavior. For each of the 30,110 people analyzed in our gym attendance dataset, we have, on average, over 1,525 daily observations (over four years of attendance) per individual, and for each of the 3,124 people analyzed in the handwashing dataset, we have on average over 3,000 observations (over 100 hospital shifts) per individual. In both samples, we study objective data on behavior, rather than self-report questionnaires, thereby avoiding possible errors of memory and meta-cognition.

Traditionally, the measurement of habit has been constrained to self-report scales (Gardner, 2015). However if people do not have much awareness about how strongly their behaviors are cued by context, they may unwillingly tend to misattribute habits to volition instead (Adriaanse et al., 2014, Gillebaart and Adriaanse, 2017, Wood and Rüniger, 2016). This potential for misattribution has led to concerns about the exclusive use of self-report methods to study natural habit formation processes (Harrington, 2017, Rebar et al., 2018).

PCS establishes two important discoveries in our field data. First, the sets of context cues that are predictive of individual-level behavior are different for different people. Specifically, while historical behavior is an important universal predictor, other context variables such as day of the week or month of the year have more heterogeneous effects. Second, contrary to common wisdom, there is no “magic number” for how long it takes to form a habit. Instead, the speed of habit formation appears to vary significantly, both between behavioral domains and between individuals within domains.

3.2 Study 1: Gym Attendance

Human Subjects Protections

Before initiating this project, the California Institute of Technology Institutional Review Board and the University of Pennsylvania Institutional Review Board reviewed and approved this study. Because this study involved an analysis of de-identified, archival data, a waiver of informed consent was approved by the Institutional Review Boards per Federal Regulation HHS CFR 45.46.117(c) (2).

Data

We partnered with 24 Hour Fitness, a large North American gym chain, to obtain check-in and background data for 60,277 regular gym users across 560 gyms who consented to share their information with researchers when they signed up to be in a fitness program. Our gym attendance dataset spans fourteen years, ranging from 2006 to 2019 and it includes over 12 million data points, each corresponding to one gym check-in. Each data point is accompanied by a timestamp when the relevant gym visit occurred, the gym location where it occurred, and other relevant information about the gym (such as its number of amenities and wi-fi availability). We further infer several other attributes, such as the day of the week and individual-level variables such as the time since gym membership creation. In total, our set of unique candidate context cue variables that may relate to predictable gym attendance includes month of the year, day of the week, time lag (the number of days which have elapsed since the last visit), attendance rate in the past week, the number of consecutive days of gym attendance (streak), and the number of consecutive days of gym attendance for the day of the week in question (day-of-week streak). A full list of the variables analyzed and a longer description of the data can be found in Appendix B.

Our analytic sample is a subset of gym goers selected based on two criteria. First we exclude anyone without a valid active gym contract, and this removes 1,083 members from the sample. Second, we only include participants with at least a year of data (removing 28,878 members) and enough attendance for the LASSO model to run (removing 206 additional members who had more than a year of data but attended fewer than 6 times total, implying that there is one cross-validation fold with no attendance so the LASSO model could not be computed). This leaves $N = 30,110$ gym goers who are the main focus of our analysis.

Table 3.1 provides general summary statistics for the final analytic dataset. Gym

members are 62% female and have a median age of 34 years. The average individual in this dataset goes to the gym every 4-6 days. The median number of days an individual is observed (or “has an opportunity to go to the gym”) is 1,525 days, or just over four years.

Table 3.1: Summary statistics on $n = 30,110$ gym goers

	Mean	SD	Q1	Median	Q3
Age	36.76	12.35	27.00	34.00	45.00
Female	0.62	–	–	–	–
Avg. daily attendance	0.19	0.16	0.07	0.14	0.27
Number of days observed	2,020	1,453	658	1,525	3,655
Avg. days between gym visits	15.77	29.74	3.69	6.89	15.22

Analytic Approach

For each individual in our dataset, we first train a LASSO model to predict the likelihood of gym attendance (a binary outcome variable) day-by-day. We use LASSO because it “zeroes out” variables that have low predictive power and might be false positives, which results in a more compact model. This allows us to train the model on the wide dataset and obtain a smaller set of individual-specific context cue variables and their respective coefficients, which are predictive of a specific individual going to the gym. We use five-fold cross validation, training the model on 85% of the full time series data for each individual (see Appendix C). We use the remaining 15% of the data as our “test” set, allowing us to see how good our LASSO model is at predicting an individual’s attendance on days the model did not observe. This gives us an out-of-sample test-set predictability measure, called the area under the curve (AUC), for each gym member. The AUC measure obtained for each individual serves as an objective measure of context-sensitive predictability (i.e., how habitual an individual’s behavior might be) and allows us to avoid errors induced by self-report measures, which rely on an individual’s memories of the context cues present when they last executed a habitual behavior.

Next, for each individual, we determine when, if ever, behavior becomes habituated over time. Following Lally et al., 2010, for each individual i , we attempt to identify $D_i(t)$, an exponential function of form $a - b^{-ct}$ describing daily-level habit strength as a function of time. Likewise, following Lally et al., 2010, we define the inferred time to habit formation as the time it takes for $D_i(t)$ to reach 95% of its asymptote (or 90% and 98% to test robustness).

To infer daily habit strength $D_i(t)$, we proceed as follows. For each individual, we compute the AUC values for a series of LASSO models using data from longer and longer periods $[0, t]$ for increasing values of t ; these AUC's are denoted $A_i(t)$ (for person i and window ending with time t). The procedure starts with each person's first two weeks of observed gym attendance, then expands to include each person's first four weeks, each person's first six weeks, and so forth. As illustrated in Figure 3.2a, the shape of this function $A_i(t)$ is estimated using the corresponding sequence of AUCs. $A_i(t)$ represents the integration of the instantaneous habit strength function over time. Assuming that habits strengthen over time, then as we add more fortnights of data to the LASSO models, the series $A_i(t)$ is blending more pre-habit formation AUCs with post-habit formation AUCs. Finally, we infer $D_i(t)$ as the derivative of $tA_i(t)$ so that $A_i(t)$ is the average value of the function D_i over the interval $[0, t]$ (see Appendix C); this function is plotted in green in Figure 3.2a.

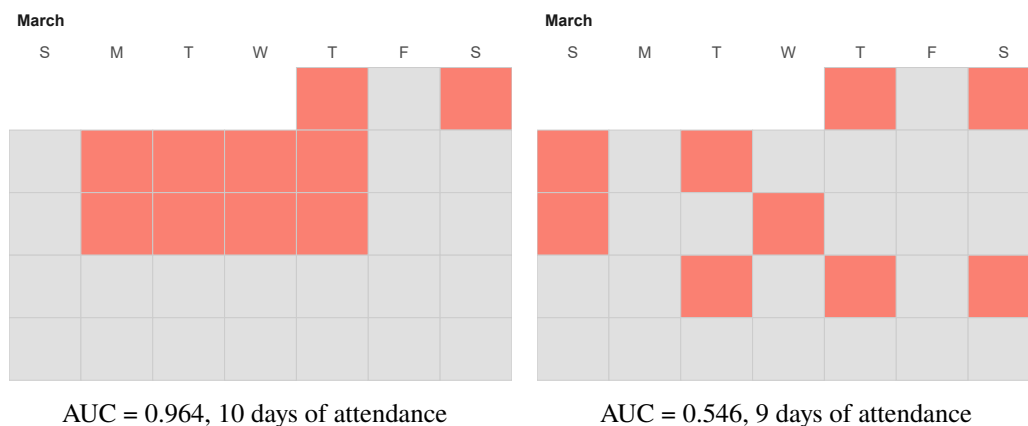
Results

In the LASSO training datasets for $N = 30,110$ gym members, the mean individual-level AUC is 0.806 (median of 0.811, interquartile range 0.750-0.868), where 0.5 is random and 1.0 is perfectly predictable. This indicates that the LASSO models tend to do a good job fitting the gym goers' gym attendance behavior. On the test datasets, these measures are slightly lower (as is expected), with a mean individual-level AUC of 0.768 (median of 0.778, and interquartile range 0.702-0.845).

Figure 3.1 illustrates two gym members' monthly attendance calendars from March 2018. While the two members go about equally often, one is highly predictable (AUC=0.946) and the other is not at all predictable (AUC=0.546).

Figure 3.1: Two examples of gym members: one with high and one with low predictability (measured by AUC).

One month (March 2018) of attendance history for two individuals (a) and (b) who have comparable attendance rates throughout their time series data but very different AUC values. The red squares indicate days that the individuals went to the gym, and the grey squares indicate days that they did not. This example illustrates the distinction between frequency of attendance (which is similar for both) and predictability (which is very different).



As shown in Table 3.2, the most important predictor of gym attendance across individuals is how much time has passed since the previous gym visit. The nature of this predictor is also very homogenous across individuals: for 76% of gym goers in our sample, the longer it has been since they last visited the gym, the less likely they are to go on a given day. Furthermore, there is substantial heterogeneity in the nature of calendar effects as predictors: most months of the year are nonzero for about half the gym goers, but are balanced between positive and negative month-specific effects. The exceptions are December and January, which are negative and positive, respectively (consistent with new year “fresh starts,” per Dai, Milkman, and Riis, 2014). This captures what are likely different “types” of gym goers - for example, those who quit the gym in favor of outdoor running in the summer and those who only hit the gym in the summer to maintain fitness. Finally, the days of the week were important predictors for more than half of the gym goers in the data (illustrated by the examples in Figure 1), but vary in homogeneity. Monday is the most homogeneous, positively predicting attendance for 57% of the sample, whereas Saturday is the most heterogeneous, positively predicting attendance for 36% of the sample and negatively predicting attendance for another 29%.

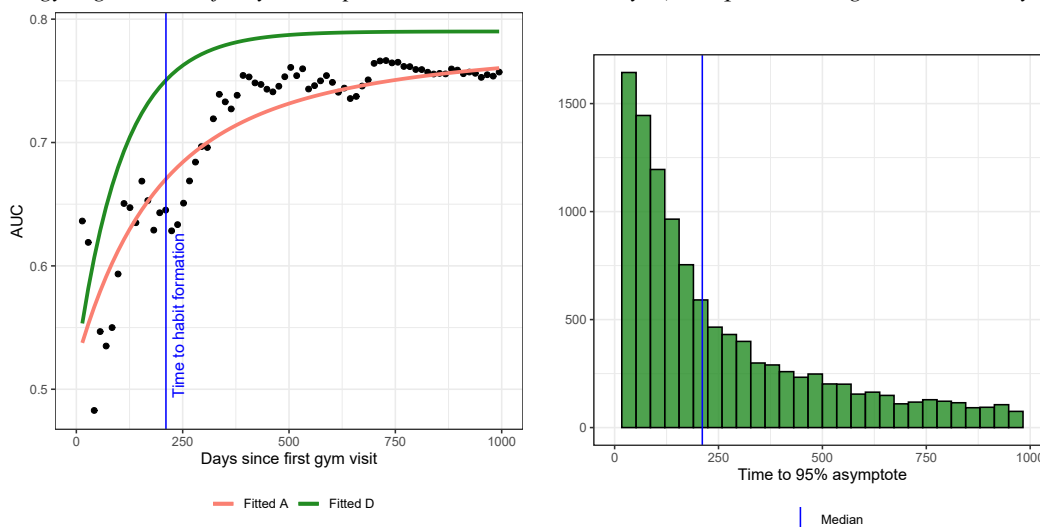
Table 3.2: Context predictors of gym attendance

	Variable				% zero	% positive	% negative	Homogeneity
	importance	Q1	Median	Q3				index (% pos. - % neg.)
Time lag	1.25	-1.40	-0.34	-0.02	22	2	76	74
(Time lag) ²	0.92	0.00	0.00	0.86	57	39	3	36
Monday	0.36	0.00	0.11	0.50	32	57	11	46
Tuesday	0.35	0.00	0.10	0.49	33	56	11	45
Wednesday	0.34	0.00	0.06	0.46	35	54	12	42
Attendance last 7 days	0.34	0.09	0.29	0.47	9	82	8	74
Thursday	0.31	0.00	0.00	0.40	37	49	14	35
Friday	0.28	0.00	0.00	0.27	36	39	24	15
Day-of-week streak	0.23	0.00	0.11	0.30	25	69	7	62
Streak	0.22	0.00	0.00	0.14	36	40	24	16
Saturday	0.22	-0.04	0.00	0.15	35	36	29	7
(Streak) ²	0.15	-0.13	0.00	0.00	46	13	42	29
(Day-of-week streak) ²	0.13	-0.16	0.00	0.00	48	9	43	34
December	0.11	-0.05	0.00	0.00	47	16	38	22
January	0.10	0.00	0.00	0.05	45	39	16	23
July	0.09	0.00	0.00	0.01	48	27	26	1
August	0.09	0.00	0.00	0.00	48	27	25	2
September	0.09	-0.01	0.00	0.00	49	25	27	2
October	0.09	-0.01	0.00	0.00	49	22	29	7
November	0.09	-0.02	0.00	0.00	49	21	31	10
February	0.08	0.00	0.00	0.02	48	31	21	10
March	0.08	0.00	0.00	0.01	49	29	22	7
April	0.08	0.00	0.00	0.01	49	29	22	7
May	0.08	0.00	0.00	0.00	49	26	25	1
June	0.08	0.00	0.00	0.01	48	29	23	6

Similarly to Lally et al. (who found 48% of subjects had a good model fit), we were able to fit the exponential curve to (and infer time to habit formation from) the sequences of AUCs and obtain a good fit ($R^2 > 0.5$) for 45% of individuals in our gym data (see Appendix C). Figure 3.2 contains a summary of results from fitting the exponential curves to gym goers' AUC sequences. The median estimated time to reach the 95% asymptote across all gym goers well fit by the exponential model is 211 days, or about 7 months. Model fit was not related to average frequency of gym attendance nor age of gym member. However, the well-fit sample is more female (0.64 vs 0.61, $p < 0.001$, Cohen's $d=0.062$) and this sample's base rate of attendance is slightly lower (0.18 vs 0.20, $p < 0.001$, Cohen's $d=0.125$).¹

Figure 3.2: Development of habit formation of gym attendance.

(a) An example of one individual's gym attendance behavior, where habit formation is modelled as an asymptotically increasing sequence of AUCs over time (orange curve $A(t)$) and the instantaneous strength of habit formation is the derivative of $tA(t)$ with respect to time (green curve $D(t)$). The time to habit formation is determined by when the individual reaches 95% of the asymptote and is marked by the blue line. (b) A summary of results from fitting the exponential curves to gym goers' AUC sequences. The median estimated time to reach the 95% asymptote across all gym goers well fit by the exponential model is 211 days (interquartile range is 85-597 days).



Additional Analyses of Reward Sensitivity and Individual Differences

We tested for insensitivity to reward change, a hallmark of strong habits that is seen in animal research, but is difficult to find reliably in humans (e.g. Wit et al., 2018, Pool et al., 2021) and has not been tested in natural field data, albeit experimentally

¹We don't have a strong prior for why this should be the case, but demographic analysis which follows found that men tend to have less predictable gym attendance behavior and may therefore be harder to fit.

replicated in field settings (Neal, Wood, Wu, et al., 2011). This analysis used the fact that estimated habit formation times divide an individual’s behavior into “pre-” and “post-” habit formation periods and that weather shocks could change the rewards for gym attendance. The maintained hypothesis is that unusual weather changes the reward value of going to the gym, similarly to how changing reward likelihood or value works in more controlled lab settings. We test whether unexpectedly better or worse weather (compared to the weather on recent days) leads to different behavior in the “pre-habit” period compared to that individual’s behavior in the “post-habit” period.² It does not (see Appendix D). This tentative null result is consistent with evidence that reward insensitivity is not as solid a hallmark of human habit as it seems to be in simpler animal learning.

Finally, we take advantage of the size and diversity of the gym goer sample to explore whether demographic and SES characteristics are correlated with predictability. Current research on demographic differences in predictability is limited, but a working paper which used an ethnographic approach (collecting self-report diary data) to study the effect of age on the proportion of behavior which is habitual did not find a direct correlation (Quinn and Wood, 2021). However, the authors did find that age is indirectly correlated with other lifestyle factors which *are* associated with a greater proportion of habitual behavior, such as employment (which increases the amount of habitual behavior) and living with others (which decreases the proportion).

To explore demographic predictors in our data, we link our individual-level AUC predictability measures with Census data using each individual’s home zip code and self-reported age and gender. We regress these demographic characteristics on the AUC of a truncated sample of 27,663 gym goers (removing those 2,447 people from our sample for whom we did not receive age or gender information from the gym, or whose zip code did not have data available). Regression results, which can be found in Table 3.3, confirm that demographic attributes are indeed predictors of AUC (see Appendix E), although most of the effects are small in magnitude. Specifically, older individuals living in more rural (low population density) areas where a large fraction of married couples have children have higher AUCs. Younger individuals living in more urban (high population density) areas have lower AUCs. This particular analysis was pre-registered on AsPredicted.org (59014), with the hypothesis that

²This analysis was done on a smaller subset of 7,355 individuals for whom we could match their gym’s zip code with weather data from NOAA (see Appendix D) and maintained their gym membership for at least 12 months in a row.

there exist systematic categorical differences in the degree of predictability and the speed of habit formation in different sub-groups of individuals.

Table 3.3: Demographic predictors of AUC

Multiple regression results summarizing the predictive power that various demographic variables have on individual-level AUC. Average household income, population density (sq. mi.), and the fraction of married and single households with children under 18, were calculated using ZCTA Census data and the gym goer's zip code. Gender and age information came from the gym chain, based on gym goer self-report. The rightmost column reports the t-statistic and effect size (Cohen's d) for each variable.

<i>Dependent variable: AUC</i>		
	Coefficient	Statistics
log(avg. household income)	-0.004** (0.002)	$t = -2.263$ $d = -0.027$
log(population density)	-0.005*** (0.001)	$t = -6.786$ $d = -0.082$
Fraction married with kids	0.019** (0.008)	$t = 2.328$ $d = 0.028$
Fraction single with kids	-0.022*** (0.007)	$t = -3.265$ $d = -0.039$
Age	0.001*** (0.0001)	$t = 9.409$ $d = 0.113$
Male	-0.008*** (0.001)	$t = -6.15$ $d = -0.074$
Constant	0.844*** (0.025)	
Observations	27,659	
R ²	0.007	
Adjusted R ²	0.007	
F Statistic	31.963*** (df = 6; 27652)	

Note: Standard errors in parentheses.

*p<0.1; **p<0.05; ***p<0.01

3.3 Study 2: Hand Washing among Hospital Workers

Human Subjects Protections

Prior to initiating this project, the California Institute of Technology Institutional Review Board and the University of Pennsylvania Institutional Review Board reviewed and approved this study. Because this study analyzed de-identified archival data, a waiver of informed consent was approved by the Institutional Review Boards per Federal Regulation HHS CFR 45.46.117(c) (2).

Data

We obtained hand-hygiene data from Proventix, a company that uses RFID technology to monitor whether individual healthcare providers wash their hands at every opportunity to do so throughout their hospital shifts. The initial dataset tracks 5,246 healthcare workers across 30 different hospitals. The dataset spans about a year, with over 40 million data points, each corresponding to whether an individual caregiver did or did not wash their hands in the face of an opportunity to do so (defined as a point in time when a caregiver entered or exited a patient's room with a Proventix sanitizer present).

Each data point is accompanied by a timestamp, as well as the room and hospital location where the opportunity to wash arose. We further infer several other attributes about each opportunity to wash, such as the day of the week when it arose and whether the healthcare worker in question had complied with handwashing guidelines (washed their hands) in this room previously. Our unique candidate context cue variables include the time of day, time spent working, previous room and shift compliance, and indicators for whether the hospital worker was entering or exiting the room. A full list of the variables used and a longer description of the data can be found in Appendix B.

We treat the introduction of Proventix's RFID surveillance technology as a "shock" that disrupted behavior (as documented by Staats et al., 2017), while acknowledging that handwashing may have habituated in some caregivers before we could observe them. We use two criteria to identify our final analytic sample. First, we remove any hospital workers who had fewer than 30 shifts (a month) of data (removing 2,115 hospital workers). Second, we remove anyone without enough hand washing compliance for the LASSO model to run.³ This gives us a subset of 3,124 hospital workers, on whom we run our analysis.

³This removes 7 hospital workers with over 30 shifts, but not enough hand washing observations to have variability in all 5 cross-validation folds and the holdout data.

Table 3.4 provides summary statistics about the workers in our analytic sample. The mean compliance with handwashing is 0.45 per opportunity. An average of 116 shifts are recorded per healthcare worker, and there are an average of 26 episodes (or visits to patient rooms, each with two opportunities to wash - one upon entry and one upon exit) per shift. We observe an average of 3,016 episodes per worker.

Table 3.4: Summary statistics on $n = 3,124$ hospital caregivers' hand washing

	Mean	SD	Q1	Median	Q3
Hand sanitizing compliance	0.45	0.23	0.26	0.43	0.63
Total number of shifts	116	77	56	98	153
Number of rooms visited	37	33	20	29	41
Avg. episode length (mins)	5.66	2.61	3.94	5.13	6.78
Avg. number of episodes per shift	25.72	16.49	13.95	24.2	34.54
Avg. shift length (mins)	511.91	213.75	408.38	581.2	645.74
Avg time between episodes (mins)	22.42	11.5	13.95	20.12	29.00
Avg time off between shifts (hours)	91.95	57.91	60.06	72.61	102.91

Analytic Approach

We use the same machine learning approach as described in Study 1, training a LASSO model to obtain person-specific sets of coefficients and predictability measurements (AUCs). We obtained an $R^2 > 0.5$ for 33% of individuals. Model fit was not related to the rate of hand washing compliance. However the well-fit group tended to have a slightly higher number of shifts (116 vs 115, $p < 0.001$, Cohen's $d=0.010$) and less time off between shifts (89.6 vs 93.1 hours, $p < 0.001$, Cohen's $d=0.062$) (see Appendix C for further discussion). These differences are highly significant given our statistical power but small in magnitude.

As in Study 1, we inferred habit formation time only for those individuals we could fit in the data, assuming others never became habituated to hand sanitizing. We adopted the same approach of fitting what we call pre- and post-habit data using lengthening pre-habit time intervals. Rather than adding two weeks of attendance data at a time to the pre-habit interval, we added two additional shifts at a time to the growing window blending pre- and post-habit data (up to, on average, 51 episodes).

Results

The LASSO model does a satisfactory job fitting hospital caregivers' hand washing behavior. In the training dataset, the mean (median) individual-level AUC is 0.788 (0.783), and the interquartile range is 0.742-0.828. In the test dataset these measures

are again lower as would be expected, with a mean (median) individual-level AUC of 0.781 (0.776), and interquartile range of 0.732-0.825. While our LASSO models have slightly less predictive power in this domain (compared to gym attendance), they still outperform random chance at predicting hospital caregivers' hand washing behavior.

As in Study 1, the AUC measure - which can be used in any behavioral domain - is produced for each individual, and it once again serves as an objective measure of context-sensitive predictability. Furthermore, PCS is again able to narrow down the set of context variables that are the most important predictors of hand washing at the aggregate level (see Table 3.5).

The most important context variable is handwashing compliance during their last shift. Surprisingly, a room entry indicator is negative for 77%.

Time of day intervals were not selected by the LASSO model as predictive of most people's hand washing behavior. However, consistent with previous research (Dai et al, 2015), the amount of time since the start of a caregiver's shift was a negative predictor of hand washing for 42% of the caregivers. The most important and homogenous predictors were a hospital worker's handwashing compliance during their last shift (a positive predictor for 100% of the hospital workers), room entry (which is a negative predictor for 77% of hospital workers, indicating most are more likely to wash their hands upon exiting, rather than entering, a room), and the room compliance of others (a positive predictor for 66% of hospital workers). The most heterogeneous predictors are room frequency (the rate at which a specific room is visited by the hospital worker, compared to other rooms) and time off work (that is, time off between shifts), both of which were equally likely to be positive or negative predictors of hand washing when predictive at all.

Table 3.5: Context predictors of hospital hand washing

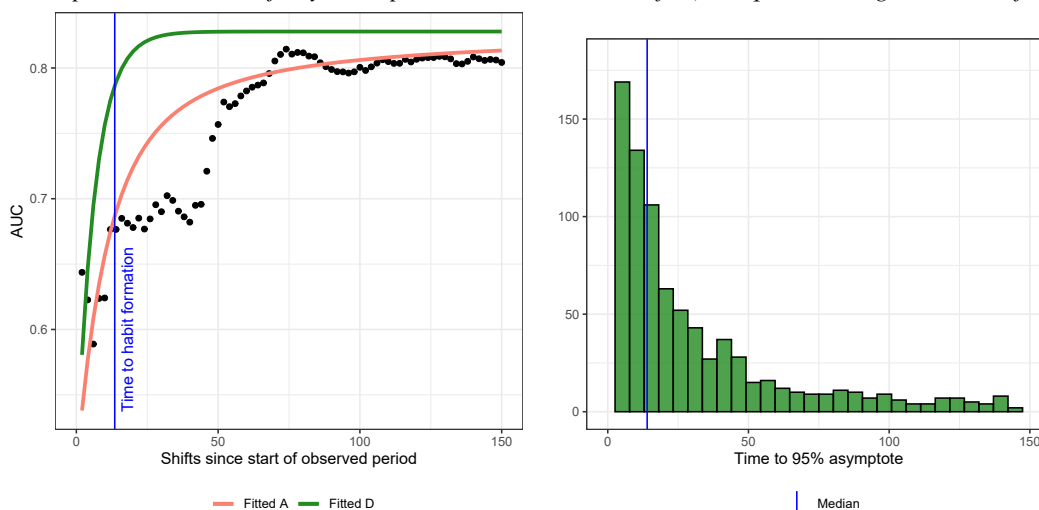
	Variable				% zero	% positive	% negative	Homogeneity index (% pos. - % neg.)
	importance	Q1	Median	Q3				
Compliance last shift	0.77	0.66	0.70	0.92	0	100	0	100
Entry indicator	0.35	-0.33	-0.28	-0.04	18	5	77	72
Compliance last opp.×Entry indicator	0.13	0.00	0.00	0.21	49	47	4	43
Compliance last opp.×Time since last opp.	0.12	0.00	0.00	0.00	54	1	45	44
Compliance within a room	0.12	0.00	0.01	0.14	33	51	16	35
Time since last opp.	0.09	0.00	0.00	0.00	61	24	15	9
(Time since last opp.) ²	0.08	0.00	0.00	0.00	74	7	18	11
Room compliance of others	0.08	0.04	0.05	0.12	32	66	2	64
Time at work	0.08	0.00	0.00	0.00	54	4	42	38
Compliance last opp.×(Time since last opp.) ²	0.07	0.00	0.00	0.00	74	20	5	15
Prev. room compliance	0.07	0.03	0.04	0.11	32	65	2	63
Compliance last opp.	0.05	0.00	0.00	0.07	47	45	7	38
Time at work×6am-12pm	0.05	0.00	0.00	0.00	78	10	12	2
Time since last compliance	0.05	0.00	0.00	0.00	64	9	27	18
Time at work×12pm-6pm	0.04	0.00	0.00	0.00	73	10	17	7
(Time since last compliance) ²	0.03	0.00	0.00	0.00	75	17	8	9
12am-6am	0.03	0.00	0.00	0.00	68	22	10	12
Frequency of patient encounter	0.03	0.00	0.00	0.01	58	31	12	19
Time at work×Patient encounter	0.03	0.00	0.00	0.00	64	8	28	20
Days since start	0.02	0.00	0.00	0.00	83	9	8	1
6am-12pm	0.02	0.00	0.00	0.00	80	7	13	6
12pm-6pm	0.02	0.00	0.00	0.00	77	12	11	1
Room frequency	0.02	0.00	0.00	0.00	63	19	19	0
Time at work×6pm-12am	0.02	0.00	0.00	0.00	82	10	7	3
(Time off) ²	0.01	0.00	0.00	0.00	84	8	8	0
October	0.01	0.00	0.00	0.00	81	10	9	1
November	0.01	0.00	0.00	0.00	82	10	8	2
December	0.01	0.00	0.00	0.00	81	10	9	1
March	0.01	0.00	0.00	0.00	82	9	10	1
April	0.01	0.00	0.00	0.00	80	10	11	1
May	0.01	0.00	0.00	0.00	80	9	10	1
June	0.01	0.00	0.00	0.00	80	10	10	0

July	0.01	0.00	0.00	0.00	79	11	10	1
August	0.01	0.00	0.00	0.00	78	11	11	0
September	0.01	0.00	0.00	0.00	82	9	9	0
Day-of-week frequency	0.01	0.00	0.00	0.00	77	13	10	3
Rooms visited in shift	0.01	0.00	0.00	0.00	83	8	9	1
6pm-12am	0.01	0.00	0.00	0.00	84	7	8	1
Prev. day-of-week compliance	0.01	0.00	0.00	0.00	79	8	13	5
Prev. unit compliance	0.01	0.00	0.00	0.00	78	10	13	3
Streak	0.01	0.00	0.00	0.00	78	8	15	7
Time off	0.01	0.00	0.00	0.00	85	8	7	1
Unit frequency	0.01	0.00	0.00	0.00	72	20	8	12
February	0.00	0.00	0.00	0.00	82	8	9	1

As in Study 1, we fit an exponential model to individuals' AUC sequences and plot these increasingly wider ranges of AUC values for individuals who were well fit by the model. This allows us to analyze the development of predictability over time (Figure 3.3), which in turn serves as a proxy for the speed of habit formation, as defined by reaching 95% of their asymptote of predictability. We see that the median habit is formed on the order of weeks, unlike gym attendance where the median time to habit formation was on the order of months.

Figure 3.3: Development of habit formation of hospital hand washing.

(a) An example of one individual's hand washing behavior, where habit formation is again modelled as an asymptotically increasing sequence of AUCs over time (orange curve $A(t)$) and the instantaneous strength of habit formation is the derivative of $tA(t)$ with respect to time (green curve $D(t)$). The time to habit formation is determined by when the individual reaches 95% of the asymptote and is marked by the blue line. (b) A summary of results from fitting the exponential curves to hospital workers' AUC sequences. The median estimated time to reach the 95% asymptote across all hospital workers well fit by the exponential model is 14 shifts (interquartile range is 5-37 shifts).



Additional Analysis of Reward Sensitivity

As in Study 1, we use individual-level habit formation time estimates to test whether post-habit behavior is less sensitive to a reward change. The hypothesized reward change is the last opportunity a caregiver has to wash their hands in the final room visit for their shift. The hypothesis is that they are less likely to wash their hands because it is less important to do so, for hygiene, when they are leaving. However, there is again no statistically significant effect (see Appendix D).

3.4 Discussion

Despite the obvious policy relevance of healthy (and unhealthy) habits, there is a notable absence of large-scale, observational field studies on the formation of habits over extended periods of time. We introduce a machine learning approach called Predicting Context-Sensitivity (PCS) to fill this gap. PCS is a machine learning method that identifies the set of contextual variables that best predict behavior for each individual in a dataset.

In two large longitudinal datasets, PCS confirms that past behavior robustly predicts future behavior. Specifically, hand sanitizing compliance during a caregiver's previous shift is the most important predictor variable of hand washing during the current shift, and time since last gym visit is the best predictor of today's gym attendance. Other relatively homogenous effects include the greater likelihood of going to the gym on weekdays versus weekends, and people's tendency to keep up a streak of gym visits. In the hospital, caregivers are more likely to wash their hands when exiting rather than entering a room. Compliance of other medical staff in a room is another reliable predictor of hand washing. But LASSO also reveals how context cues can differ across individuals. As shown in Table 3.2, the day of the week is highly predictive of gym attendance for over 60% of participants in our sample, but the sign of the effects are often different. Likewise, as shown in Table 3.2, time of day for hand washing has heterogeneous effects. Uncovering the importance and universality of these context cues extends the empirical literature on what contextual cues are associated with habitual behavior.

We also find heterogeneity in the speed of habit formation across different behavioral domains. Contradicting the popular belief in a "magic number" of days in which most habits form, we find that developing a habit of handwashing takes on the order of days to weeks, whereas developing a habit of gym attendance takes on the order of months. One possible explanation is that relative to handwashing, gym attendance is a less frequent and more effortful and intentional behavior. Handwashing among hospital caregivers is also likely to be a stronger social norm, with greater public accountability, than the personal decision to go to the gym. Handwashing is also more likely to involve chained sensorimotor action sequences which are more automatic (Balleine and Dezfouli, 2019). Another explanation may be that the context cues we have in our handwashing data are very granular, in some ways more so than those captured in our gym attendance data. It is therefore possible that if we had more (and more specific) predictors for gym attendance, we would see

predictability plateau at a faster rate. Alternatively, if gym attendance was broken out into its “micro behaviors”, such as going to a specific exercise machine, we would likely see predictability occur quicker. Applying PCS to additional datasets will help us more accurately attribute whether this difference in the speed of habit formation is a function of different behavioral attributes or a function of our data.

PCS opens up the possibility of identifying when and for whom personalized nudges could add the greatest value. For instance, PCS analysis could be used to determine the day of the week or time of day when an individual may especially benefit from intervention - not just when it comes to forming habits of handwashing and physical exercise, but also for medication adherence, healthy eating, and mores. We hope this machine learning innovation inaugurates a new era in the study and personalization of behavior change interventions.

Appendix A - Literature Review

Overview

Since habit naturally crosses disciplinary boundaries, the most promising understanding of it is likely to come from integrating evidence and methods across disciplines (Rebar et al., 2018; pg. 42). That is our approach. The purpose of the following section is to highlight key papers from the major disciplines we take evidence and methods from. Specifically, this section summarizes how habit is studied in psychology, computational neuroscience, economics, and political science. We first present a summary table comparing how these literatures have addressed the different hallmarks of habitual behavior. In bold are what we view as being the “gold standard” of measurement. Some of the attributes - such as time to habit formation - do not have a gold standard yet. We follow the table with more detailed reviews of each field’s major contributions.

	Psychology	Computational Neuroscience	Economics; Political Science	Our Paper
Analysis tools	Lab experiments, observational studies	Lab experiments	Econometrics, field experiments	Machine learning
Data types	Self-report surveys, behavior (subject to attrition)	Behavior (small samples), brain activity	Behavior	Behavior (with no attrition ¹)
Length of habit formation	The best estimate is about <u>two months</u> , with a wide range which appears to depend on the specific behavior. ²	The focus is on simple motor habits which can be created within <u>1-2 hours</u> in a lab setting.	A common assumption is that <u>one month</u> should be sufficient to form a habit. ³	Our model allows us to <u>try multiple durations, as well as individual-level parameters.</u>
Frequency	<u>Frequency is estimated based on self-report measures</u> ⁴ (occasionally complimented by objective behavioral data).	High repetition, or “ <u>overtraining</u> ” is used to <u>make a behavior become habitual</u> in the lab.	“Habit-forming” incentives are typically designed with the <u>goal of increasing the frequency of a behavior.</u>	We <u>separate frequency from predictability</u> , a binary measure that can be compared across behavioral domains ⁵ .
Automaticity	<u>Automaticity is measured using self-report</u> ⁶ , along with several newer but less frequently used implicit measures.	Variety of measurements are used in humans, from localization of brain activity to response times.	<u>Not measured</u> ; requires measuring something other than revealed choice behavioral data.	We do not have the necessary data types to measure automaticity.
Reward Devaluation	One study ⁷ found reward devaluation insensitivity in humans <u>when they were in a habit-cueing context.</u>	Has been found repeatedly in animal studies, but has proven difficult to replicate in humans ⁸ .	<u>Measured indirectly</u> ⁹ when <u>incentives are removed</u> from an intervention (ex. the monetary reward for going to the gym is removed).	We measure insensitivity to reward devaluation by looking at <u>individual model coefficients in response to ecologically relevant devaluations.</u>
Context-Sensitivity	Many studies look at how context-stability (measured using self-report) influences the execution of habits.	This is typically <u>not studied</u> given that most human studies are done in fMRI or lab settings.	The idea of state-dependence is well-established, but the <u>only ‘state’ considered is past consumption history.</u>	Our wide dataset of variables allows us to <u>objectively identify which contexts cue the relevant behavior for each individual.</u>
Individual Differences	The majority of studies model <u>between-person variation based on aggregate habit scores</u> ¹⁰ .	Oftentimes a “learning parameter” is customized to individual-level data.	<u>Not studied</u> ; studies tend to focus on highly aggregated data. Some is at the level of the household.	Our AUC measure is person-specific, allowing investigation of <u>how individual differences correlate with habit sensitivity.</u>

¹ We believe this is a unique feature of our paper given most studies focus exclusively on the existence (vs. absence) of a behavior. As noted by Mullan and Novorodovskaya (2018; pg. 86), “There is a difference between ‘a habit of doing’ versus ‘a habit of not doing’ (De Vries et al., 2011; Knussen & Yule, 2008). For example, those who are not used to wearing a seatbelt do not just have no habit of wearing a seatbelt, they also have a habit of not wearing a seatbelt. However, very few studies have actively explored this.”

² This estimate comes from Lally et al. (2010), who compiled panel observational data using self-report questionnaires and fit individual-level habituation models to estimate that it took anywhere from a few weeks to over half a year to reach the 95% asymptote of behavior.

³ While the origin of this standard is hard to pin down, it may originate with one of the first papers in economics to empirically test habit formation (Charness and Gneezy, 2009). While they do not explicitly justify the choice of one month, they do state (italics our own) “Our results indicate that it may be possible to encourage the formation of good habits by offering monetary compensation for a *sufficient number of occurrences*, as doing so appears to move some people past the “threshold” needed to engage in an activity.” Later studies continued to use one month as the standard, despite acknowledging that oftentimes this was not sufficient for habits to form (“An alternative explanation is that some subjects would have experienced an increase in postintervention attendance if the intervention period had been longer.” – Acland and Levy, 2015).

⁴ Ouellette and Wood’s (1998) BFCS (“Behavior Frequency x Context Stability”) measure, which is a self-report index covarying past behavior frequency with measured contextual variables, is based on the assumption that behaviors which are repeated frequently in familiar contexts are more likely to become habitual.

⁵ Different behaviors have different “natural” frequencies, making predictability a better universal measure of habit than frequency alone.

⁶ The most commonly used self-report measure is the SRHI (“Self-Report Habit Index”). A subscale has been designed specifically to capture automaticity and is called SRBAI (“Self-Report Behavioral Automaticity Index”), including questions like “I do this behavior without thinking.”

⁷ The study, Neal et al. (2011), found that individuals continue to eat stale popcorn beyond satiation if they are in a specific context which cues the habitual behavior (specifically, watching a film or eating with their dominant hand).

⁸ The only paper to find evidence of reward devaluation insensitivity with humans in the lab (Tricomi et al. 2009) has not been replicated across a number of studies which attempted to do so (de Wit et al. 2018).

⁹ Although economists end up measuring this indirectly, they do not use the language of “reward devaluation insensitivity” in the way that neuroscientists do. Furthermore, they often find *sensitivity* to the devaluation, with behavior decreasing following a devaluation.

¹⁰ As noted in Gardner (2015), “Theory has also been inadequately tested at the individual level. Most (80; 98%) studies have exclusively modelled between-person variation in habit, based on aggregates of individuals’ habit scores. Yet, habitual action is inherently idiosyncratic, based on personally acquired behavioural responses to personally meaningful cues. Within-person effects cannot be reliably interpreted from aggregations of processes that differ between people (Jaccard, 2012; Molenaar, 2004).”

Psychology

Psychologists define habit as a behavior which is prompted automatically by contextual cues as a result of learned context-action associations (Wood and Neal, 2009). This definition combines two key attributes of habitual behavior which guide a lot of the psychology research: predictable context-sensitivity and automaticity.

Some habit researchers make further distinctions about what should be considered a habit. For example, Gardner, 2015 argues that the initiation and performance of a behavior are distinct. He classifies behavior into one of three types: habitually initiated but consciously performed (his example: riding a bike to work every morning), consciously initiated but habitually performed (his example: exercising at the gym), or habitually initiated and habitually performed (eating a snack in the afternoon). This is a sensible distinction, but without measures of automaticity, we cannot apply it to our data. It is simply a reminder that repeated behaviors that we call habits need not be unconscious or automatic.

Context-Sensitivity

The focus on context-sensitivity came from evidence that habits arise when context-stable behavioral repetition creates a “transfer” from (internal) goals to (external) associations with environmental cues (Ji and Wood, 2007). In the language of animal learning and instrumental conditioning, an S-R-O relation in which an association is developed between a stimulus (S), the response (R) it elicits, and a reward outcome (O), becomes habitized as an S-R relation.

Habits are not “innate” to the behavioral repertoire in the way reflexes are (e.g. one is not born with, but must develop, the habit of tooth-brushing, unlike the reflex of being startled by something unexpected, which is present in newborns at birth). Instead, most habits begin as goal-directed behaviors. Eating solid food using a fork, for example, begins as a very deliberate goal-directed behavior in small children (one which requires a lot of motor and cognitive control in the beginning). It may take months or even years for eating to become an automatic motor sequences with little need for cognitive control, such that a habit can form. In adults, who have cheaper cognitive control and can eat “mindlessly,” the behavior of eating is ripe for developing associations with context and reward independent of nutritional goals. Specifically, the “trigger” to eat is often transferred to context elements of the environment which reliably co-occur with the behavior. For example, people who snack frequently in a stable context are no longer driven by an internal motivation to eat, but rather by an environmental cue (Danner, Vries, and Aarts, 2008).

The range of possible context cues is usually idiosyncratic, because they are likely to vary by the type of behavior and by individual. The context cues most often studied in psychology tend to be physical time, space, or social cues which are easily measurable (such as the location in which a behavior occurs, the day of the week, the time of day, whether other people were present when the behavior was executed, etc.) (Mazar and Wood, 2018). However, one can imagine less easily measurable contextual cues, such as a specific mood, sensory input, or a memory, as triggering a habit. These may be harder to measure objectively, for example relying on individuals' recollection of a memory or ability to verbally describe a feeling, but are still important. In clinical studies for example, stress and visual cues that induce craving states are often measured given their importance for behavior (Fox et al., 2007; Sinha, 2009; Ferguson and Shiffman, 2009). Psychology and applied psychology (e.g., health behavior research) are the most focused on, and seek to measure, context-sensitivity.

Automaticity

The other attribute that psychologists seek to measure to determine whether a behavior is truly a habit is automaticity (Gardner, Abraham, et al., 2012; Orbell and Verplanken, 2010). A behavior is considered automatic if it is “brought to mind by cognitive processes largely outside of conscious awareness” (Mazar and Wood, 2018; pg. 14).

An early start on this definition came from Bargh, 1994, who presented four criteria of automatic behavior. The first is awareness of the cognitive process which gives rise to the behavior. The second is intentionality – or control – over the initiation of the cognitive process. The third is efficiency – automatic processing requires fewer mental resources. And the fourth is control – the ability to stop or alter the cognitive process after it has begun. Even if there was an easy way to measure all four of these, Bargh noted that not all of the criteria need to be met in order for a behavior to be considered automatic. In fact, a behavior which only meets two or three may still be automatic, confusing the definition even further still. More recent theoretical models of automaticity have maintained the view that it is a multidimensional construct, continuing to emphasize the unintentionality, uncontrollability, and unconscious execution of behavior (Moors and De Houwer, 2006).

Animal learning studies also illustrate how simple theories of automaticity and habit are often hard to evaluate conclusively. Garr and Delamater, 2019 trained rats on a two-lever-press paradigm for 20 days or 60 days, then tested for automaticity

and sensitivity to reward devaluation. The more extensively trained rats performed the rewarding lever presses more often and more quickly (by these measures, their behavior became more automatic). But both groups exhibited similar insensitivity to reward devaluation and a difference in apparent goal-directed control of the two different levers.

Measurement

Next, we will examine the most common measures used in psychology to assess context-sensitivity and automaticity of behavior. Most of these measures are surveys which require individuals to self-report answers to questions about their own behavior. In a meta-analysis looking at 136 empirical studies which applied ideas from the habit literature to health behaviors over the years 1998-2013, Gardner, 2015 found that self-report scales are still the main methods used to measure habitual behavior. Two scales dominate the literature.

The first scale – relied on by 88% of the studies in Gardner’s meta-analysis – is the SRHI, or “Self-Report Habit Index” (Verplanken and Orbell, 2003). Its popularity stems from the fact that the questionnaire is short (a 12-item scale), direct, and has become the standard in psychology. One of the questions asks the subject to rate their agreement with the following statement on a Likert scale: “I do this behavior without thinking.” A subscale of SRHI was designed specifically to capture automaticity and is called SRBAI (“Self-Report Behavioral Automaticity Index”).

The challenge is that accurately self-reporting automaticity relies on good meta-cognition (our thinking about our thinking). But what if automatic behavior occurs with no awareness? If people do not have much awareness about how strongly their behaviors are cued by context, they may unwittingly misattribute habits to volition instead (Adriaanse et al., 2014, Gillebaart and Adriaanse, 2017, Wood and Runger, 2016).

Another popular measure – used by 12% of the studies in Gardner’s review – was Ouellette and Wood’s (1998) BFCS (“Behavior Frequency x Context Stability”) measure. This is a self-report index co-varying past behavior frequency with context stability. This measure is based on the assumptions outlined earlier that behaviors which are repeated frequently in familiar contexts are more likely to become habitual (Wood and Neal, 2009). The questions aim to assess both directly, phrased as “how often do you do this behavior?” and “when you do this behavior, how often is this cue present?”

This questionnaire also relies on high levels of accurate recall and metacognition. But how good is human recall for frequently performed behaviors? Take, for example, the behavior of checking one's phone. Using a smartphone app which calculated true frequency of phone use, Wilcockson, Ellis, and Shaw, 2018 were able to track the phone behavior of 27 participants over the course of 14 days. They found that there was no correlation between true phone-checking behavior and a self-report measure called the Mobile Phone Problem Use Scale ("MPPUS") which is a 27-item questionnaire that includes items such as "I can never spend enough time on my mobile phone." This anecdote points to another fault with self-report measures: they are inherently retrospective, relying heavily on hindsight. But memory degrades quickly – with the details of a morning becoming foggy as one enters their afternoon – meaning the timescale at which these questionnaires are administered is crucial.⁴

A more systematic review comes from Parry et al., 2021, who ran a meta-analysis of 47 studies to measure the link between logged and self-reported digital media use. To evaluate the association between self-reported and logged media use, 66 effect sizes from 44 studies were considered ($n = 52,007$) and correlations were calculated with robust variance estimation (RVE). Their analysis concluded that self-reported media use has a positive but medium-magnitude relationship with logged (objective) measurements ($r = 0.38$, 95% $CI = 0.33$ to 0.42 , $p < 0.001$). Furthermore, problematic media use showed an even smaller association with usage logs ($r = 0.25$, 95% $CI = 0.20$ to 0.29 , $p < 0.001$). These studies, along with other critiques (Harrington, 2017) point out the issues with self-reporting habits.

Besides these two most common scales, two other measures were used in just one study each. The EHS ("Exercise Habit Survey"), used in one study, is similar to BFCS. The other measure was an association test, designed to measure cue-behavior associations underpinning habitual behaviors (an implicit association test).

So while psychologists have identified two important elements of habitual behavior - context-sensitivity and automaticity - there have been some concerns about how good their current measurement tools are as proxies for true habitual behavior (Rebar et al., 2018). In particular, it is unlikely that automaticity and context-sensitivity can be accurately captured using self-report measures alone.

What behaviors can become habitual?

⁴A modern technique which the smartphone makes available is real-time experience sampling where people are prompted to discuss situational cues and whether they are executing a habit.

Psychologists study habits across a range of behavioral domains. Popular domains of study include activities which are done frequently: eating, exercising, and hygiene behavior. However there is some debate around how complex a behavior can be before it can no longer be considered a candidate for becoming habitual. This is in part due to research which has demonstrated that simpler actions like drinking water tend to become habitual more quickly than complex actions like exercise routines (Lally et al., 2010). The idea is also evident in animal learning, in which chained motor sequences are slower to habitize (Graybiel, 1998).

Focusing on the two behaviors covered in this paper, hand-washing seems to be ripe for becoming habitual because it involves a short motor sequences. Potthoff et al., 2018 (p.248), suggest that hand-washing habits “minimize[e] cognitive resources required for a given behavior to ensure that it can be performed with a maximum of patients and/or for when such resources are especially needed”.

Whether exercise can become habitual is more debatable (Rhodes and Rebar, 2018). Physical activity, particularly travelling to a gym for exercise, is different from other familiar habitual behaviors. Two differences worth noting are that it is a multi-step behavior, not a simple motor action, and that it takes a long time to perform.

However, the type of exercise which is done inside a gym is often a relatively straightforward motor action. Running on a treadmill, rowing, lifting weights – while requiring “control” and “awareness” and hence not meeting the definition of automaticity – are simple enough that many gym goers are able to multi-task while doing them – as is obvious by watching gym-goers listening on their headphones, holding a conversation, reading or watching TV while they exercise. Secondly, the other attribute of habitual behavior, context-sensitivity, is likely present for gym goers. Location, other people, time of day, or biological states (for very regular exercisers) are likely candidates for cuing the decision to attend the gym.

Speed of Habit Formation

Given the learning process of behaviors going from goal-directed to habitual through repetition in a context-stable state, some researchers have been interested in how long it takes for a habit to form. However, answering this question using traditional psychology tools is difficult because it requires a significant amount of data collection (obtaining regular SRHI responses over many days, as an example). This requires researcher time and persistent longitudinal engagement by subjects. Hence, only a handful of studies have been done to answer this question (Lally et al., 2010,

Kaushal and Rhodes, 2015, Fournier et al., 2017).

A seminal study is Lally et al., 2010. The researchers collected SRHI measures for 82 subjects daily over the course of 12 weeks for an eating, drinking, or physical activity behavior chosen by the subject. Lally et al., 2010 then fit a curve to each individual's self-report scores through time in order to measure the time it took them to reach 95% of the asymptote (their definition of when something became a habit). They were able to fit the model for 62 individuals and obtain a good fit for 39 out of those 62, finding that “performing the behaviour more consistently was associated with better model fit.” Their results showed that the median time to habit formation was 66 days, with a range of 18 to 254 days to habit formation depending in part on the complexity of the behavior (e.g. the relatively simple act of drinking a glass of water was quicker to habitize than a more complex physical activity).

Another study looked at the development of exercise habits by asking new gym members to complete surveys over the course of 12 weeks (Kaushal and Rhodes, 2015). They found that exercising at least four times per week for 6 weeks was the minimum requirement to establish an exercise habit, based on the time at which behavior appeared to reach an asymptote (i.e. not change significantly after that time period). The most recent observational study focused on the effect of circadian cortisol (modulated by time of day) on the development of a simple physical habit. Fournier et al., 2017 tracked 42 French students for 90 days as they did a stretching exercise behavior. Some students were assigned to do it in the morning (when cortisol levels are high) and some in the evening (when cortisol is low). The SRBAI was collected daily, and the speed of habit formation process was then modelled using learning curves by fitting a four-parameter logistic curve to SRBAI responses. The curve-fitting process was successful, converging for each participant (in contrast to the power function following Lally et al., 2010, which the researchers also tried, finding that only 48% had a moderate fit as defined by $R^2 > 0.70$). Their results showed that the morning group achieved automaticity at an earlier time point (106 days) than the evening group (154 days), concluding that time of day influences the speed of habit formation.

Of these three quantitative studies, all showed that “habit typically develops asymptotically and idiosyncratically, potentially differing in rate across people, cues and behaviors” (Gardner and Lally, 2018; pg. 220).

Computational Neuroscience

What does habitual behavior look like in brain activity? This has been the driving question for much research in computational neuroscience. This research tends to focus on the neural basis of the two types of cognitive processing mentioned in the last section: “goal-directed” behavior, a more deliberate cognitive functioning, and habitual behavior. The existence of these respective decision making systems is now well-accepted and commonly modeled theoretically as model-free (MF) and model-based (MB) decision-making (Gläscher et al., 2010; Dayan and Berridge, 2014; Daw et al., 2011). MF learning transitions to habit learning with extensive experience.

When a new habit is being learned, inputs to the midbrain dopamine system drive dopaminergic neural activity which encodes reward prediction errors (RPEs). These RPEs serve as learning signals. Learning an accurate prediction of a stable reward results in smaller and smaller reward prediction errors over time. These signals are thought to modulate synaptic plasticity in the striatum which in turn serves as the “gate-keeper for tentative motor plan representations” (Pauli et al. 2018). The striatum can be further segmented into two distinct areas: the dorsolateral striatum (DMS) and the dorsomedial striatum (DLS).

Instrumental behaviors which respond to reward values may start out as goal-directed actions largely controlled by the associative striatum (DMS), which controls more goal-directed activity, when they are first being learned. But under certain conditions and with enough repetition, these behaviors may become habitual and no longer contingent on reward. Then cognitive control shifts to the sensorimotor striatum (DLS), which controls more stimulus-driven behaviors (Yin and Knowlton, 2006; Knowlton and Patterson, 2016). Functional MRI studies which are used to localize brain activity during decision making have confirmed that habitual processing tends to occur in the “sensorimotor loop,” which connects the basal ganglia with the sensorimotor cortices and parts of the midbrain (Tricomi, Balleine, and O’Doherty, 2009, Yin and Knowlton, 2006). Brain scans have therefore been used to confirm that the brain has two independent sources of action control which govern behavior, and to help determine whether a behavior is habitual or goal-directed (Balleine and O’Doherty, 2010).

So what are the conditions necessary for a behavior to move from being goal-directed to being habitual? The animal literature suggests that habit formation requires a behavior to be repeated many times – a process known as “overtraining”

(Tricomi, Balleine, and O’Doherty, 2009). This process creates an association between the stimulus, the behavior, and the reward outcome (a form of instrumental learning) such that the behavior begins to depend on reward reliability rather than reward optimization (Dickinson, Nicholas, and Adams, 1983; Lee, Shimojo, and O’Doherty, 2014). Once a habit has been established, even if reward value changes such that it is no longer optimal to execute the behavior, the subject may continue to do so if such a response has reliably produced a reward following previous behavioral executions. A number of conditions have been found to speed up this process of shifting behavior from goal-directed towards habitual. The most notable one, which has been reproduced in a number of settings, is learning under stress - lab studies have found that inducing stress (in animals, including humans) leads to quicker formation and reliance on habitual behavior (Schwabe & Wolf 2009).

Habitual behavior that is automatic is accompanied by measurable psychological and biological features, including faster response times, limited attention during choice (Knowlton and Patterson, 2016), and degraded declarative memory (explaining the basis for choice when asked, see Seger and Spiering, 2011).⁵ These attributes can be studied using a range of measurement tools, some of which are more portable outside of a laboratory setting, including eye-tracking methods to measure attention.

One important test used to determine whether a behavior is habitual or not is a test of sensitivity to reward devaluation. The procedure originated in animal learning studies, with Adams and Dickinson, 1981, who studied how lever pressing in rats could become habitual. When they analyzed habit, they described it as a behavior which becomes so automatic that even devaluation of the reward value of an outcome will have little effect on the execution of the habitual behavior. Specifically, they found that poisoning a food pellet after a rodent has developed a highly-trained habit of lever-pressing for the pellets did not deter the rodent from continuing to press the lever. This phenomenon has been termed insensitivity to reward devaluation, and is a behavioral hallmark of habitual processing.

Insensitivity to reward devaluation has been established in humans. Tricomi, Balleine, and O’Doherty, 2009 trained participants to learn that responses to two different fractal images were associated with two different snack rewards. After

⁵Studying two patients with large MTL lesions, Bayley, Franscino, and Squire, 2005 found neurotypical-level performance in an overtrained discrimination task with no declarative memory or conscious awareness. Thus, lesion patients could perform the task automatically. However, performance was completely degraded to random on a minor task variant. The two patients also learned the task about as quickly as four monkeys did.

overtraining (choosing their preferred fractal many times in short succession), they were given one of the snacks to eat to satiety, which presumably devalued it. Subjects who had food devalued this way continued to choose the fractal associated with the devalued foods, indicating habit. This is evidence of human insensitivity to reward change similar to the animal experiments.

However, other researchers have not been able to replicate these findings (Wit et al., 2018). This raises the question of whether an experimental paradigm using rodents can be easily transferred to human behavior. Another concern which has been raised about the reward devaluation paradigm is that it implies that behavior which is not goal-directed is necessarily habitual. For example, the goal-independent behavior may not be context-sensitive (Mazar and Wood, 2018; pg.23). However, there remains an interest in replicating this effect with humans with different paradigms and training protocols.

One of the best studies showing insensitivity to reward devaluation in humans is a psychology study. While it does not have neuroscience data, it is included here because it is a clear illustration of this reward devaluation test. Neal, Wood, Wu, et al., 2011, found that people were more likely to overeat stale ("devalued") popcorn in a context which cued habitual behavior of eating popcorn (e.g. watching a movie in a cinema) but not when they were in an unfamiliar popcorn-eating context (e.g. watching a movie in a meeting room, or eating the popcorn with their non-dominant hand) which did not cue the habitual behavior. The effect captures a two-way interaction (cinema vs. meeting room or dominant vs. non-dominant hand and whether the popcorn received was stale or fresh) and is evident only among individuals classified as "high habit" (vs. medium or low habit) per self-reports on a 7-point scale used to assess habit strength for eating popcorn in movie theaters. The same study found that for low or medium habit individuals, or high habit individuals in novel contexts, like eating popcorn in a meeting room, behavior remained sensitive to reward value and decreased in frequency when the popcorn was stale (devalued).

Economics

Economic theories and empirical tests have generally used the term "habit" in one way: To describe history-dependent "adjacent complementarity" of goods or services. ⁶ The theories are motivated by strong evidence of empirical correlation

⁶Another form of habit is the idea that the discount factor depends on consumption (Shi and Epstein, 1993). They appeal to an intuitive concept of "habits of thrift" or luxurious spending hypothesized by Fisher, 1930 (pg. 337-338) (with no evidence) which link more income to less

between past and current consumption. These models therefore specify consumption utility as a function of actual immediate consumption *relative to a reference point* or ‘consumption habit’ (see Duesenberry, 1949, Ryder and Heal, 1973, Deaton, 1992).

This approach was never empirically microfounded in psychology or neuroscience but it is mentioned prominently in the earliest studies creating a foundation for intertemporal choice. Koopmans, 1960 wrote: “One cannot claim a high degree of realism for [consumption insensitivity], because there is no clear reason why complementarity of goods could not extend over more than one time period.”

In conventional microeconomic consumer theory, “complements” are pairs of goods X and Y which increase each other’s marginal utilities when consumed together—that is, the marginal utility of X is greater if you have more Y. Familiar examples of complements include hot dogs and hot dog buns, hammers and nails, and computer hardware and software. Koopmans’s point is that complementarity could extend to the same good consumed in adjacent periods (called “adjacent complementarity”). Rather than treating hot dogs and hot dog buns as complements, yesterday’s hot dogs and today’s hot dog consumption are considered as possible complements.

In one macro-finance specification, (see Campbell and Cochrane, 1999) the crucial variables are current consumption C_t and habit X_t . Utility depends on past aggregate consumptions $C_{t-1}, C_{t-2} \dots$ through another equation. In that specification $U_t = \frac{(C_t - X_t)^{1-\gamma} - 1}{1-\gamma}$ (and X_t is related to previous consumption levels in a complicated way).

Such preference assumptions were used in macroeconomics and finance to explain facts which are puzzling in specifications (see Sundaresan, 1989, Constantinides, 1990). Campbell and Cochrane, 1999 motivate their specification with the following hypothesis⁷: “repetition of a stimulus diminishes the perception of the stimulus and responses to it” (pg. 208). This is indeed a property of sensory systems which are adaptive. However, these types of “repetition suppression” are very short-run (e.g., seconds to minutes or days). Whether the same kind of history-dependent adaptation works for, say, quarterly consumption by a household is an open question.

Rozen, 2010 derives a set of axioms relating the functional form of habitual history-sensitivity to underlying principles that are mathematically equivalent. The func-

tion. This concept is theoretically interesting but appears to be empirically counterfactual, as much evidence suggests higher income is associated with more patience, rather than less patience.

⁷The phenomenon they are describing is similar to reward prediction learning or, in perception, is called “repetition suppression” (e.g., Gonsalves et al., 2005). It would be useful to explore even a highly speculative link between these psychological foundations and the hypothesized micro-foundation for macroeconomics further..

tional representation of utility is:

$$U_h(c) = \sum_{t=0}^{\infty} \delta^t u \left(c_t - \sum_{k=1}^{\infty} \lambda_k h_k^{(t)} \right)$$

where h_k is the habit consumption history for k periods in the past and λ_k is a decay factor which weighs more distant consumption history less.

A bolder extension of adjacent complementarity is called “rational addiction (RA)” (Becker and Murphy, 1988). In this approach, current utilities depend on consumption history, due to adjacent complementarity, much as in the Rozen, 2010 formalization. But it is also coupled with self-awareness of the history-dependent structure and planning about the future. In this model, “rationally addicted” people understand that if they consume more X today, they will value X tomorrow more highly.

The key prediction of the RA model is that current consumption will depend on current prices and will *also* depend on expected future prices. For example, once they hear that a large cigarette tax increase will take place soon, rationally-addicted smokers might quit a habit abruptly - *before* the increase occurs. They’ll quit right away because they prefer, today, to be an ex-smoker at time T when the tax goes up; otherwise, continuing to smoke at T will be too expensive.

Both the macro-finance and RA specifications are natural in economics because the primitives in economic analyses are stable preferences, Bayesian beliefs, and budget constraint. Habit can then enter into the theory in one of those three ways. The default approach is to define habit as current preference depending on past consumption.

Conventional economic theory with these ingredients does not have learning, RPE, reward reliability in it. There is also no implicit cost of mental effort. And there is no attempt to relate the history-dependent model to adaptive functionality or to neural implementation.

Most economic empirical studies using the RA approach treat the fact that history-dependent consumption could be present in a wide range of goods and activities as a provocative prediction. “People can be addicted not only to harmful goods like cigarettes, alcohol, and illegal drugs, but also to activities that may seem to be physically harmless, such as sports participation, shopping, listening to music, watching television, working, etc.” (Shen and Giles, 2006). The RA approach does

make the non-obvious prediction that current behavior depends on expectations of the future, in sharp contrast to the neuroeconomic habit model which is not forward-looking.

There are many studies of RA. There are two limits in these previous empirics: (1) Most of the early empirical evidence uses very coarse time scales (e.g., quarterly tax receipts to measure state-by-state cigarette consumption); and (2) estimates of the expected future price component are not very good. Expected future prices are usually proxied by past prices, and these proxies may not be independent of current consumption. Even very sophisticated tests on coarse quarterly data have very limited power to test whether there is actually forward-planned RA.

Auld and Grootendorst, 2004 demonstrate the kinds of biases that can lead to results consistent with RA even when the basic data-generating process has no actual adjacent complementarity mechanism. The central test of RA is whether current consumption is increasing in (expected) future consumption. Simulations show that when the consumption time series is highly auto correlated (as is typical), even if there is no history-sensitivity, the RA prediction can spuriously appear to hold. However, other diagnostic features of these tests (such as inferred discount factors reasonably close to 1) can also fail in both artificial and actual data sets.

An illustrative example of how history-sensitivity is used in empirical practice is Crawford, 2010. He derived a tractable way to test whether optimal consumption with habit can be rationalized nonparametrically, in the sense that one can find some set of inferred utilities, satisfying simple restrictions like GARP and extended to allow adjacent complementarity, which fits a data set on consumption. The logic of this exercise is that if no set of inferred utilities can “rationalize” the data, then the specification of stable utilities with adjacent complementarity is incorrect.

Crawford applied the method to data on quarterly smoking expenditures for 3,134 Spanish households. The best-fitting habit lag is two quarters. Most households’ (91%) data can be rationalized using two lags (compared to only 24% with one lag), but the power of the two-lag test is not very high (only 20% of random-generated data would fail the test for optimization).

History-sensitivity is seen again and again in many types of data: It is established in internet use (Kwon et al., 2016) and employment (Heckman, 1981). In marketing it is attributed to inertia or brand loyalty (Kuehn, 1962, Keane, 1997, Dubé, Hitsch,

and Rossi, 2010)⁸.

The boldest predictions of the RA theory seem to be just flat wrong. In theory, rational addicts should take advantage of volume discounts on addictive goods, because they will optimally self-ration the goods over time. There is no direct evidence of this pattern (e.g., alcoholics buying in bulk and self-rationing), although it could be that rational addicts are liquidity-constrained. Wertenbroch, 1998 found in lab and field data that “vice” goods, such as cigarettes, are often purchased in smaller quantities, have higher quantity discounts, and have lower price elasticities than similar virtue goods, regardless of liquidity-constraint. There is also substantial evidence that restricting hours at which addictive goods are sold (typically alcohol) reduces consumption (Middleton et al., 2010). This is inconsistent with rational forward-looking optimization by addicts, who should plan their shopping around reduced hours.

For the purposes of this paper, we also note that the economic RA model does not connect with what is known from psychology and cognitive neuroscience. The latter is loosely constrained by the philosophy that a good understanding of a behavior should have an explanation for adaptive functionality, algorithmic specificity, and neural implementation.

Laibson, 2001 introduced an economic model of a specialized idea of context-sensitivity, from clinical psychology and neuroscience, to explain cue-sensitivity of addiction. In the model, the presence of a state-dependent cue actually changes utility. If a cue value is x^i , and consumption activity is a^i (=0 or 1), then the period-specific utility is assumed to be $u(a^i, x^i) = u(a^i - \lambda x^i) + (1 - a^i)\eta$ where $(1 - a^i)\eta$ is the expected utility of the next-best activity if the target activity is not done.

This is a simple economic translation of the evidence about biological addiction from opponent processes to maintain homeostatis, but it is not a biologically plausible general model for everyday habits. An implication of the Laibson specification is that mere presence of the cue creates negative utility (through unpleasant craving) if the good is not consumed. In the PCS view, the presence of a cue is typically not pleasant or unpleasant; it just predicts behavior through a neural autopilot mechanism driven by reward reliability, rather than via unpleasant craving which addicts “self-medicate” to avoid.

⁸There is some evidence of what Belk, 1975 calls “situations” (the same as our cues or states) influencing choices but it has not been an active area of research.

Bernheim and Rangel, 2004) create a more general model tailor-made to understand addictive habits. Preferences are influenced both by a numerical state, which catalogs consumption history, how frequently states trigger an involuntary “hot” craving state, and some other features. Their model is not as much a specific theory, as it is a modelling language to describe different kinds of addiction patterns and invite empirical estimation.

The Laibson homeostatic cues model and Bernheim-Rangel M-states model are two examples of state-sensitivity of preferences which go beyond the history-sensitivity in so much empirical work. In their models, the relevant state, on which preferences depend, is a cue or history variable. The idea is that what people subjectively value could depend on an environmental or contextual state (e.g. Karni, 2008). Nothing is new or surprising about that—umbrella preference goes up when it’s raining. Historically, however, economists were reluctant to allow too broad a range of state-sensitivity of preferences for fear—probably legitimately—that doing so would lead to an erosion of falsifiability. Common examples in which state-sensitivity is central are examples like health, in which health quality (a physical state) clearly influences subjective value of leisure or work.

Political Science

Political scientists have studied habit in the domain of voting. While it is conceivable that a different mechanism leads to repeat voting behavior (a la Volpp and Loewenstein, 2020), within the range of behaviors studied by political scientists, voting is the most likely candidate to become habit forming given it might be cued by context variables.

Voting is interesting for our purposes because it is very infrequent— particularly compared to hand-washing or gym attendance, and to other activities studied in empirical applied psychology. It is similarly far from the animal learning-based concept of motor habit formation and insensitivity to reward change from hundreds of rapid trials in short time spans, on the time scale of hours or days. So can voting be habitual?

The answer seems to be yes, in the simple sense that voting exhibits context-sensitivity. Researchers have mostly focused on how a disruption to total voting (“turnout”) in one election affects subsequent turnout. The disruptions that are diagnostic are exogenous “natural experiments” which suggests possible causality, as if an experimental treatment changed voting for some people but not similar

others. If skipping voting one time breaks one's "taste for voting" – reducing the likelihood of voting in future elections – then voting is considered habitual, in the history-dependent sense. This concept of habit follows directly from the economics formulation of adjacent complementarity; relying on the assumption that more past voting behavior predicts more future behavior, as has indeed been empirically documented (Brody and Sniderman, 1977).

These studies are of three types:

1. Observational studies seek to isolate the impact of an "as-if random" inducement to vote in one year, on voting turnout in subsequent election years (Franklin and Hobolt, 2011; Denny and Doyle, 2009).
2. Experimental studies apply a truly random assignment to inducement to vote and test whether it increases future voting (Green and Gerber, 2002; Bedolla and Michelson, 2012)
3. "Quasi-experimental" causal identification studies use regression discontinuity designs which take advantage of strict voting eligibility requirements – e.g. to test whether two similar people born days apart (Meredith, 2009) vote more in the future, if one got a lucky chance to vote before while the other person did not.

A challenge, as pointed out by Coppock and Green, 2016, is that these designs often suffer from weak identification of short-run and long-run effects. For example, if an inducement to focus on the treated election is focused on encouraging people to "do their civic duty," this effect of social pressure may endure into the next election, independent of habit formation. Similarly, the early inducement may lead to increased interest in politics, which then causes the later turnout.

More recent work has acknowledged that behavior alone is not enough to label an action as habitual, citing the psychology literature on automaticity and context-sensitivity as inspiration for creating a self-report voting habit index akin to the SRHI (Cravens, 2020). Cravens argues that the "cost" of voting (Downs, 1957) will be lower when voting becomes habitual.

Other papers have looked at the consistency of environmental context voting behavior by looking at voting rates following a change in home address or voting location address. This approach is a special case of our general focus on PCS except

for a narrow range of context variables and a long time between behaviors (and unfortunately, also a change in cost).

For example, Brady and McNulty, 2011 found that the consolidation of voting precincts in Los Angeles county decreased overall turnout substantially (which was partially, but not fully, offset with an increase in absentee votes). This change is consistent with the hypothesis that removal of the environmental cue of the physical precinct deterred some individuals from voting. Aldrich, Montgomery, and Wood, 2011 found that both self-reported previous voting and not moving (situational consistency) were associated with voting. Research into other contextual cues, like time of day, which may be predictive of voting behavior has been more limited (Cravens, 2020).

Appendix B - Dataset Descriptions

The purpose of this section is to provide additional detail on the two main datasets used in this paper, along with a full list of the context variables which were used to train the LASSO models.

Hand Washing Data

Hand-hygiene data came from Proventix, a company which uses RFID technology to monitor whether a healthcare provider sanitized their hands during a hospital shift. The initial dataset tracks 5,246 hospital healthcare workers across 30 different hospitals. The dataset spans about a year, with over 40 million data points, each corresponding to whether an individual did or did not wash their hands. Each data point has a timestamp, room, and hospital location.

We further infer several other attributes, such as time of day and individual-level variables such as whether the healthcare worker complied (washed their hands) in this room previously. A full list of the variables that are used follows.

Gym Attendance Data

We obtain check-in data from a North American gym chain, containing information for 60,277 regular gym users across 560 gyms. The data spans fourteen years, from 2006 to 2019. There were initially over 12 million data points, each corresponding to one gym check-in. Each data point is accompanied by a timestamp, gym location, and other information about the gym (such as the number of amenities and wi-fi availability, which we do not use in this analysis).

We further infer several other attributes, such as the day of the week and individual-level variables such as the time since gym membership creation. A full list of the variables that are used follows.

Description of Context Variables in Hand Washing Data

- **Time at work:** minutes elapsed since the start of a person's shift.
- **Rooms visited in shift:** number of rooms the caregiver had visited previously during the shift.
- **Compliance last opportunity:** an indicator variable of whether the caregiver washed her hands at the last opportunity.
- **Time since last opportunity (mins):** minutes elapsed since the last opportunity.
- **Time since last compliance (mins):** minutes elapsed since the last compliance.
- **Frequency of patient encounter:** percentage of time in patient rooms as a fraction of time worked. At any moment in the shift, this is defined as $\frac{\text{cumulative time spent in patient room}}{\text{cumulative time elapsed in shift}}$.
- **Entry indicator (0-1):** an indicator of whether the opportunity to wash is an entry (1) into a room (as opposed to an exit (0) from a room).
- **Previous unit compliance:** average compliance (%) across previous shifts in the current hospital unit.
- **Unit frequency:** % of previous shifts in the current hospital unit.
- **Previous day-of-week compliance:** average compliance (%) across previous shifts in the current day of week.
- **Day-of-week frequency:** % of previous shifts in current weekday (compared to other weekdays).
- **Previous room compliance:** average compliance (%) across previous shifts in the current room.
- **Room frequency:** % of time spent working in current room (compared to other rooms in the same hospital).

- **Room compliance of others:** average compliance rate (%) of other caregivers in the current room.
- **Compliance last shift:** compliance rate in the last shift before the current one.
- **Days since start:** number of days worked since the observed start date.
- **Time off:** hours elapsed between end of the last shift and the current shift.
- **Streak:** number of consecutive shifts with less than 36 hours apart.
- **Hour-slot fixed effects:** time of day is divided into four categories: 12am-6am, 6am-12pm, 12pm-6pm, and 6pm-12am.
- **Compliance within a room:** an indicator of whether the caregiver washed her hands in this room in the current opportunity (e.g. if she washed upon entry, this variable value for the exit opportunity is equal to 1).
- **Month of the year.**

Description of Context Variables in Gym Attendance Data

- **Streak:** number of consecutive days with gym visits prior to the current day.
- **Day-of-week streak:** number of consecutive corresponding day-of-the-week gym visits prior to the current day.
- **Time lag:** number of days since the last gym visit.
- **Attendance last 7 days:** number of gym visits during the last 7 days.
- **Month of the year.**
- **Day of the week.**

Appendix C - Analysis Details

The purpose of this section is to provide additional detail on our analysis methodology. Specifically, we provide a formal description of our LASSO models and include a discussion of the model output (predictability) vs. a traditional measure of habit (frequency). We then provide a formal description of the exponential model used to fit the behavioral data to identify speed of habit formation, and discuss model.

Individual LASSO Regressions

We apply LASSO logistic regressions at the individual level. LASSO is ideal for our purpose because it can improve out-of-sample predictive accuracy by reducing variance without significantly increasing bias, while also effective at feature selection by shrinking insignificant variables towards 0.

For each individual, we select about 15% of their time series data as a holdout (“test”) set on which we will assess the performance of the model. For the remaining (“training”) data, we train the model based on the following logit specification:

$$\mathbb{P}(Y_t = 1) = \frac{\exp(\beta_0 + \mathbf{S}_t\beta_1)}{1 + \exp(\beta_0 + \mathbf{S}_t\beta_1)},$$

where t indexes time, Y_t is the binary outcome variable indicating whether a habit was executed at time t , and \mathbf{S}_t is a vector of state variables. With the LASSO penalty, the problem becomes minimizing the following loss function:

$$L(\beta \mid \lambda) = -\log \left[\prod_{Y_t=1} \mathbb{P}(Y_t = 1) \prod_{Y_t=0} (1 - \mathbb{P}(Y_t = 1)) \right] + \lambda \|\beta_1\|_1.$$

As is standard with machine learning applications, we use 5-fold cross validation to pick the optimal λ . The holdout set and the folds used in cross-validation are selected such that the proportions of observations with $Y_t = 1$ in each of them are the same.

AUC vs Frequency

To demonstrate the difference between AUC and behavioral frequency, we plot in Figure 1 the relationship between holdout AUC and frequency of behavioral execution for each individual in the two datasets. We see that there is no clear relationship between the two - specifically, increased frequency is not necessarily correlated with increased predictability (AUC), highlighting the importance of the latter as a novel measure of habit.

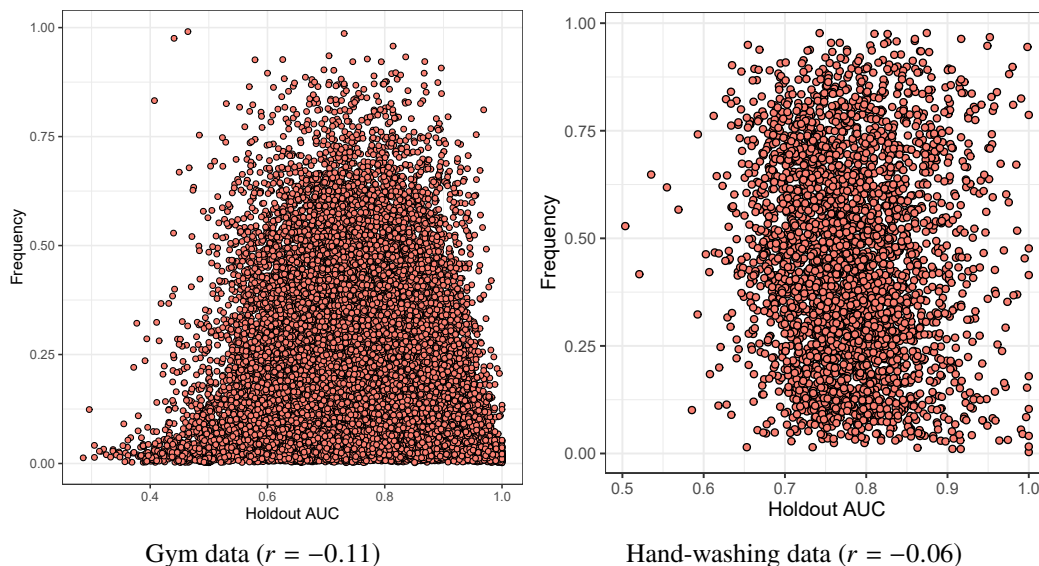


Figure 3.4: Relationship between holdout AUC and outcome frequency

Speed of Habit Formation

As we discussed in the literature review, there is not much good empirical evidence about the speed of habit formation. Most studies have typically relied on self-report measures and automaticity scales (Lally et al., 2010, Ersche et al., 2017). However, given that habitual behavior might, as part of its essence, be accompanied by degraded memory, self-report measures are not ideal. We seek to avoid errors induced by self-report by taking advantage of our granular observational data and using the predictability measurement AUC as a proxy of habit scores. In other words, we hypothesize that habit formation is manifested by an increasing sequence of AUCs over time. For each individual in the gym data, we looked at the AUC values obtained by fitting the LASSO model separately for a growing window of weeks. The sequence of LASSOs uses weeks $0 - 2, 0 - 4, \dots, 0 - W_i$, where W_i is the total number of weeks of observed data for individual i . This process generates a sequence of $\left\lfloor \frac{W_i}{2} \right\rfloor$ AUC values. For the hand-hygiene data, we considered shifts $0 - 5, 0 - 6, \dots, 0 - S_i$, where S_i is the total number of shifts observed in the data for individual i .

The maintained hypothesis is that habit formation is manifested as increasing predictability of behavior based on context variables, where predictability is measured by AUC.

Denote $A_i(t)$ as individual i 's AUC corresponding to period $0 - t$. Intuitively, $A_i(t)$ represents the *average* of accumulated habit strength over the first t periods. We

then compute instantaneous strength of habit formation at time t , denoted by $D_i(t)$, by relying on the following definition

$$A_i(t) = \frac{1}{t} \int_0^t D_i(s) ds.$$

Following Lally et al., 2010, we assume an asymptotic curve of the form $D_i(t) = a_i - b_i e^{-c_i t}$ for $D_i(t)$. The parameters a_i , $a_i - b_i$, and c_i represent the asymptote, the starting value $D_i(0)$, and the speed of adjustment (a higher value of c_i represents faster adjustment). We define the time to habit formation $T_i^* = -\ln(a_i/20b_i)/c_i$ as the time it takes for $D_i(t)$ to reach 95% of its estimated asymptote a_i . It follows from calculus that $A_i(t) = a_i - \frac{b_i[1 - \exp(-c_i t)]}{c_i t}$. We use nonlinear least squares to fit the empirical $A_i(t)$ to each individual i 's AUC sequence and obtain the estimates $\hat{a}_i, \hat{b}_i, \hat{c}_i$.

Model Fit

As mentioned in the main text of our paper, more than half of the individuals in our datasets were not well-fit by the exponential curve. This is a typical finding in previous studies, including Lally et al., 2010, who found 48% of subjects had a good model fit. In a related paper, Wood, Quinn, and Kashy, 2002 estimated that only a third to a half of human behavior becomes habitual, based on self-report automaticity scoring for a range of behaviors. Below we include additional information comparing those who were and were not well-fit by our model for gym attendance and handwashing behaviors.

With respect to gym attendance, 55% of individuals in the gym data were not well-fit by the model. It includes 2% for whom the model could not fit at all (meaning there were no convergent values produced for a, b, c by the optimization package). A significant portion of gym goers, 26%, had a linear fit, meaning they do not become predictable over time. Additionally, 4% each had model parameter $a < 0.7$ (meaning the asymptotic AUC is so low that the individual does not seem to be habitual) or $a > 1.5$ (which is not realistic, given that the asymptotic AUC must be less than 1). For 15% of gym goers, the model was simply a weak fit, as defined by $R^2 < 0.5$. Finally, we exclude 5% who had $R^2 > 0.5$ but had an extreme value of T^* (outside the 5-95% range).

Table 3.6: Time to habit formation

This table reports the estimated time to habit formation of the sample for whom we were able to obtain a good fit. We defined time to habit formation as reaching T^* as the time it takes to reach 95% of the asymptote. As a sensitivity check, we will also report the results when T^* is defined as the time it takes to reach 90% and 98% of the asymptote, respectively.

		Q1	Median	Q3
Gym data ($N = 13,449$)	Days to reach 90% asymptote	51	127	364
	Days to reach 95% asymptote	85	211	597
	Days to reach 98% asymptote	126	321	901
Hand-hygiene data ($N = 1,025$)	Shifts to reach 90% asymptote	3	8	21
	Shifts to reach 95% asymptote	5	14	37
	Shifts to reach 98% asymptote	7	22	58

Table 3.7 provides separate summary statistics for the gym goers who were and were not well fit by the exponential model. As seen in the table, there are no significant differences with respect to the age of the two groups. However, the well-fit sample is more female (0.64 vs 0.61, $t=5.584$, $p < 0.001$) and the base rate of attendance is slightly lower (0.18 vs 0.20, $t=-10.549$, $p < 0.001$). Furthermore, the average time between visits is slightly lower (14.37 vs 16.91, $t=-7.701$, $p < 0.001$) and the number of days observed is higher (2,127 vs 1,934, $t=11.468$, $p < 0.001$) for the well-fit sample.

While the difference in attendance rates does not have an intuitive explanation (it is unclear why individuals who go more often would be less well fit by the model), the most important difference is the number of days observed, which is lower for those not well fit by the model, as would be expected (more data generally increases the chance a model can be fit successfully).

Table 3.7: Summary statistics on gym goers by quality of asymptotic fit

	Well fit ($N = 13,449$)					Not well fit ($N = 16,661$)					t -statistic	p -value
	Mean	SD	Q25	Median	Q75	Mean	SD	Q25	Median	Q75		
Age	36.68	12.14	27.00	35.00	45.00	36.82	12.53	27.00	34.00	45.00	-1.006	0.315
Female	0.64	0.48	0.00	1.00	1.00	0.61	0.49	0.00	1.00	1.00	5.584	<0.001
Avg. daily attendance	0.18	0.15	0.06	0.14	0.25	0.20	0.17	0.07	0.16	0.29	-10.549	<0.001
Number of days observed	2,127	1,453	718	1,698	3,780	1,934	1,446	623	1,377	3,511	11.468	<0.001
Avg. days between gym visits	14.37	22.07	3.96	7.40	15.51	16.91	34.67	3.48	6.44	14.97	-7.701	<0.001

With respect to handwashing behavior, 67% of individuals in the hospital dataset were not well-fit by the exponential model. This includes 2% for whom the model could not fit at all (meaning the nonlinear fitting package could not produce values for a , b , c). Similarly to the gym attendees, a significant portion of the hospital workers

(37%) were better fit by a linear model than an exponential one. Additionally, 2% had a low asymptotic model parameter $a < 0.7$ and 12% had $a > 1.5$. For 12% of hospital workers, the model was simply a bad fit, as defined by $R^2 < 0.5$. Finally, as with the gym sample, we drop 5% who had $R^2 > 0.5$ but had an extreme value of T^* (outside the 5-95% range).

Table 3.8 provides separate summary statistics for the two groups of hospital workers whose handwashing behavior was, and was not, well fit by the exponential model. As seen in the table, there are no significant differences between the two groups with respect to the number of rooms visited (p -value = 0.949), the average episode length (p -value = 0.538) or the average time between episodes (p -value = 0.128).

However, t -tests reveal that there are statistically significant differences: Those well fit by the model are more likely to comply (.46 vs .44, $t=2.430$, $p < 0.001$), have a greater total number of shifts (116 vs 115, $t=6.292$, $p < 0.001$), large average number of episodes per shift (27 vs 25, $t=3.145$, $p = 0.002$), longer average shift lengths (525 vs 506 minutes, $t=3.216$, $p = 0.001$), and shorter periods of time off between shifts (90 vs 93 hours, $t=-1.523$, $p < 0.001$).

It is noteworthy that those who wash their hands more often (have higher rates of compliance) are better fit by the model. Furthermore, while not all of the differences between the two groups have a theoretically-informed explanation, the number of shifts observed is lower for those not well fit by the model as would be expected (since more data generally increases the chance a model can be fit successfully).

Table 3.8: Summary statistics on hospital caregivers by quality of asymptotic fit

	Well fit ($N = 1,025$)					Not well fit ($N = 2,099$)					t -statistic	p -value
	Mean	SD	Q25	Median	Q75	Mean	SD	Q25	Median	Q75		
Hand sanitizing compliance	0.46	0.23	0.26	0.45	0.64	0.44	0.23	0.26	0.43	0.62	2.430	0.015
Total number of shifts	116.07	72.13	60	101	148	115.34	78.85	53	96	156	6.292	<0.001
Number of visited rooms	36.79	35.44	20	29	39	36.44	31.17	20	29	41	0.064	0.949
Avg. episode length (mins)	5.6	2.36	3.93	5.09	6.69	5.69	2.73	3.95	5.14	6.81	0.616	0.538
Avg. number of episodes per shift	26.88	16.25	15.72	25.18	35.53	25.15	16.58	13.14	23.65	34.1	3.145	0.002
Avg. shift length (mins)	524.66	179.61	431.19	583.86	645.86	505.7	228.31	391.47	579.82	645.66	3.216	0.001
Avg. time between episodes (mins)	21.77	11.34	13.42	19.68	27.83	22.74	11.56	14.24	20.32	29.78	-1.523	0.128
Avg. time off between shifts (hours)	89.59	53.88	59.8	72.6	102.63	93.09	59.76	60.18	72.63	103.09	-5.810	<0.001

Appendix D - Field Tests of Insensitivity to Reward Devaluation

The purpose of this section is to describe our approach for running a test of reward devaluation insensitivity in our field data.

Weather Data

To test the impact of exogenous reward change on gym behavior in our truncated sample, we use unusual weather as an event which is plausibly random (i.e., not dependent on what gym goers did in the past) and may change the subjective reward value of going to the gym. We first map the ZIP codes of each gym to a latitude-longitude coordinate using data from 2013 collected by the US Census Bureau. As the average land area of a United States ZIP code is approximately 85 square miles, by modelling each ZIP code as a circle, we find that the average radius of each ZIP code is approximately 5 miles. It is assumed that daily weather is similar within such a radius. The National Ocean and Atmospheric Administration (NOAA) of the Department of Commerce provides a detailed list of the weather stations across the country and their respective coordinates.

For each date and gym combination that weather data is made available, we obtain the highest temperature, lowest temperature, average temperature, precipitation, and snowfall. Table 3.9 provides summary statistics on the main weather attributes and shows how these statistics differ on only those days when individuals attended the gym. We note that this restriction causes the means and standard deviations for both temperature attributes to increase slightly, and vice versa for precipitation and snowfall.

Not all weather stations provide measurements for all of the aforementioned attributes. Hence, to obtain recordings for each date and gym and each weather attribute, we searched through the list of nearby weather stations once for each attribute, in order from closest to farthest, until a station with a measurement of that attribute was found.

The mean distances to each weather station remained relatively small (from a minimum of 0.01 miles to a maximum of 13.67, means ranging between 3.05 for rain to 6.89 for average temperature). The slightly higher mean distance for average temperature measurements is due to the NOAA's classification of the other four as "core" elements of weather measurement, hence there are more stations that are equipped to regularly record them.

Table 3.9: Summary statistics on weather data

Statistic	TMAX (F)	TMIN (F)	PRCP (mm)	SNOW (mm)
Mean	72.11	51.50	18.19	0.64
Median	73	53	0	0
St. Dev.	15.24	14.05	81.13	9.85
Mean (Att = 1)	72.40	51.63	17.15	0.60
Median (Att = 1)	73	53	0	0
St. Dev. (Att = 1)	15.41	14.44	9.25	68.55
% Nonzero Days	N/A	N/A	21.2	1.3

In the case of gym attendance, the best exogenous reward variable is weather shocks. We aim to classify days into different weather categories based on snowfall, precipitation, and temperature fluctuations. Obviously, the relationship between temperature and perception of weather varies greatly with geographic locations, seasons and individual tolerances. For instance, while a 60°F spring day in Illinois is felt as warm and pleasurable, a winter day with similar temperature in California would feel cold.

Therefore, to link temperature and weather quality, we focus on the change in temperature relative to an “expected” level of temperature instead of raw measurements. Specifically, for each individual, we look at the distribution of average temperature over a year in the corresponding area. Temperatures between the 25th to 75th quantiles are considered normal. For each day in the individual’s time series, we say that it has adverse temperature fluctuation if (i) its temperature falls outside the normal zone, (ii) the average temperature of the previous 3 days is normal, and (iii) the change in temperature compared to the previous 3-day average is at least half the standard deviation of the average temperature distribution.

Similarly, a positive temperature fluctuation is observed when there is a change in temperature of at least half standard deviation in the opposite directions. A day is considered to have “unexpectedly bad” weather if it has either snowfall, persistent moderate rain to shower (defined as at least 5mm of rain per hour), or an adverse temperature fluctuation.

Conversely, a day with “unexpectedly good” weather has a positive temperature fluctuation and no precipitation or snow. The relation between unexpected weather could have either positive or negative effects on gym attendance. Bad weather can raise the cost of getting to the gym (reducing attendance) or lower the opportunity cost of being inside rather than exercising outside (increasing attendance). There-

fore, the analysis examines the absolute value of the coefficients associated with weather shocks.

Approach

A hallmark of strong habits used in animal learning, psychology, and neuroscience is insensitivity to reward devaluation.⁹ In animal studies, food rewards are devalued in various ways. One is by pausing a learning sequence and allowing the animal to freely eat food which is laced with an unpleasant but harmless toxin. The animals quickly develop taste-aversion toward the food. In early human experiments, people eat or drink freely until they are satiated; if they are truly satiated, more food or drink has zero value. In our data, we have no experimental control over reward value. Therefore, we can only hypothesize how exogenous changes might change the value of the subjective reward of gym attendance, and impact behavior differently in pre-habit and post-habit periods.

To examine the relationship between habit formation and sensitivity to weather shocks, we repeated the individual LASSO regressions described above with the addition of a set of weather dummy variables (unexpectedly good or unexpectedly bad), interacted with indicator variables for pre- and post-habit formation periods. We denote these interaction terms by $\hat{\beta}_{\text{pre}}$ and $\hat{\beta}_{\text{post}}$, respectively. We used the individual-specific estimated time to 95% asymptote T_i^* , described in the previous section, as a cutoff for pre- and post-habit formation.

Figure 3.5 shows the distribution of the LASSO coefficients of pre-habit and post-habit indicators interacted with each of the weather category variable. For both good and bad weather categories, the cumulative distribution curves of their interactions with post-habit term lie strictly above the cumulative distribution curves of their interactions with pre-habit term ($p < 0.001$ for all the one-tailed Kolmogorov-Smirnov tests), implying that these coefficients have smaller magnitudes and are shrunk to 0 more frequently by the LASSO algorithm than their pre-habit counterparts.

In the case of hand washing, it is slightly harder to find an exogenous shock which would reliably change the value of washing from one episode to the next. The ideal candidate would be rapid response events which occur at various times throughout the day. Any member of the hospital who is on the rapid response team would receive a notification about an urgent case requiring them to rush to the patient(s).

⁹There are also studies of insensitivity to reward contingency— e.g., the probability of reward contingent on a behavior such as a lever press is lowered, but the animal keeps responding at the baseline rate of presses per unit time.

Such situations are often a matter of life or death and require immediate attention, therefore potentially affecting the comparative value of hand washing behavior in the moment. Unfortunately, these events are only reported by individual hospitals and not coded in the dataset we have access to.

Instead, we use an indicator for whether a worker is exiting their last episode of the day as a proxy for reward devaluation. As with unexpected weather on gym behavior, the relation between exiting last episode and hand-washing can have either positive or negative effects. If the key driver for hand-washing is to not spread infection from one patient/episode to the next, then the last episode decreases the value of hand-washing. However, if the key driver for hand-washing is to keep oneself clean and free of infection after work, then the last episode increases the value of hand-washing. Hence, similar to the gym attendance, our analysis examines the absolute value of the sign of coefficients associated with the last episode. Specifically, we repeated the individual LASSO regressions with the addition of a last episode exit indicator variable interacted with indicator variables for pre- and post-habit formation periods. Figure 3.6 shows the distribution of the LASSO coefficients of pre-habit and post-habit indicators interacted with the last episode variable. We observe a pattern largely similar to the gym attendance case: the cumulative distribution of post-habit terms lies above that of pre-habit terms (again, $p < .001$ for all the one-tailed Kolmogorov-Smirnov tests).

A potential problem with the pre- and post-habit comparison is that the sample sizes are usually imbalanced; the pre-habit samples are usually much smaller. It could be that for various reasons, even if there is no change in sensitivity to unusual weather pre- vs. post-habit, that the imbalance of sample sizes creates a spurious difference.¹⁰

We therefore conduct a “placebo test” to examine whether the distributional difference in pre-habit and post-habit terms are true evidence of insensitivity to reward change. In a placebo test, a variable that is known to have no effect (or highly suspected to have no effect) is used. The original test statistic is the difference in coefficients for the unusual weather, when interacted pre- and post-habit. (Recall that the post-habit $\hat{\beta}_{\text{post}}$ are smaller in magnitude and have more zeros.) In the

¹⁰Recall that LASSO is a penalized regression. In small samples, even if a variable has no effect the unpenalized logit will generate $\hat{\beta}$ coefficients that are variable around zero. LASSO will shrink many of them to zero if they are too large. However, the net distribution of LASSO-penalized coefficients is not known to us. So it is hard to know a priori whether there will be a spurious effect without conducting the placebo test as we did.

placebo test, we create a artificial binary random variable instead of the unusual weather variable, which was hypothesized to have less effect post-habit. The artificial variable was generated, for each individual, such that it has the same frequency and is expected to be uncorrelated with the real reward revaluation variable, for both pre-habit and post-habit periods. For example, if the frequency of unusual good weather was 14% in the pre-habit period for person with id #432, then the placebo variable had 14% values of 1 in the pre-habit phase for that person.

Denote $D_i(\text{weather}) = |\hat{\beta}_{\text{pre}}| - |\hat{\beta}_{\text{post}}|$ as the difference in absolute values between the interactions of reward revaluation weather variables with pre-habit and post-habit terms for person i . $D_i(\text{placebo})$ is defined similarly for the artificial placebo random variable. If there is a sample-size bias generating a spuriously positive $D_i(\text{weather})$ effect, then $D_i(\text{placebo})$ will tend to be positive too. However, even controlling for this bias, it could be that there is a genuine change in the direction of $|\hat{\beta}_{\text{post}}| < |\hat{\beta}_{\text{pre}}|$.

The empirical question is whether there is such a placebo bias, and whether the measured weather sensitivity is greater than the placebo bias or not. Figures 3.5 and 3.6 below show that there is a bias. Comparing the within-person coefficient differences (pre minus post), we find that the average difference $D_i(\text{weather})$ is not significantly different from $D_i(\text{placebo})$ (gym: $t = -0.008$, $p = 0.993$, 95% CI = $[-0.024, 0.024]$; hand-washing: $t = 0.527$, $p = 0.598$, 95% CI = $[-0.031, 0.054]$). These confidence intervals exclude effects larger than around .04 in magnitude. Thus, there is no net evidence of much less sensitivity to hypothesized reward changes post-habit (compared to pre-habit)– i.e., we can be confident that any such effect, if it exists, is not greater than .04 in magnitude. One challenge is that in animal and human learning experiments, reward devaluation is carefully controlled– for example, by feeding animals food rewards to satiation. But here the reward devaluation is only crudely hypothesized based on unusual weather and the reduced value of having clean hands when leaving the hospital. Better reward change measures might detect an insensitivity, after strong habit formation, similar to the effects seen in experiments with animals (and some humans, e.g. Pool et al., 2021)

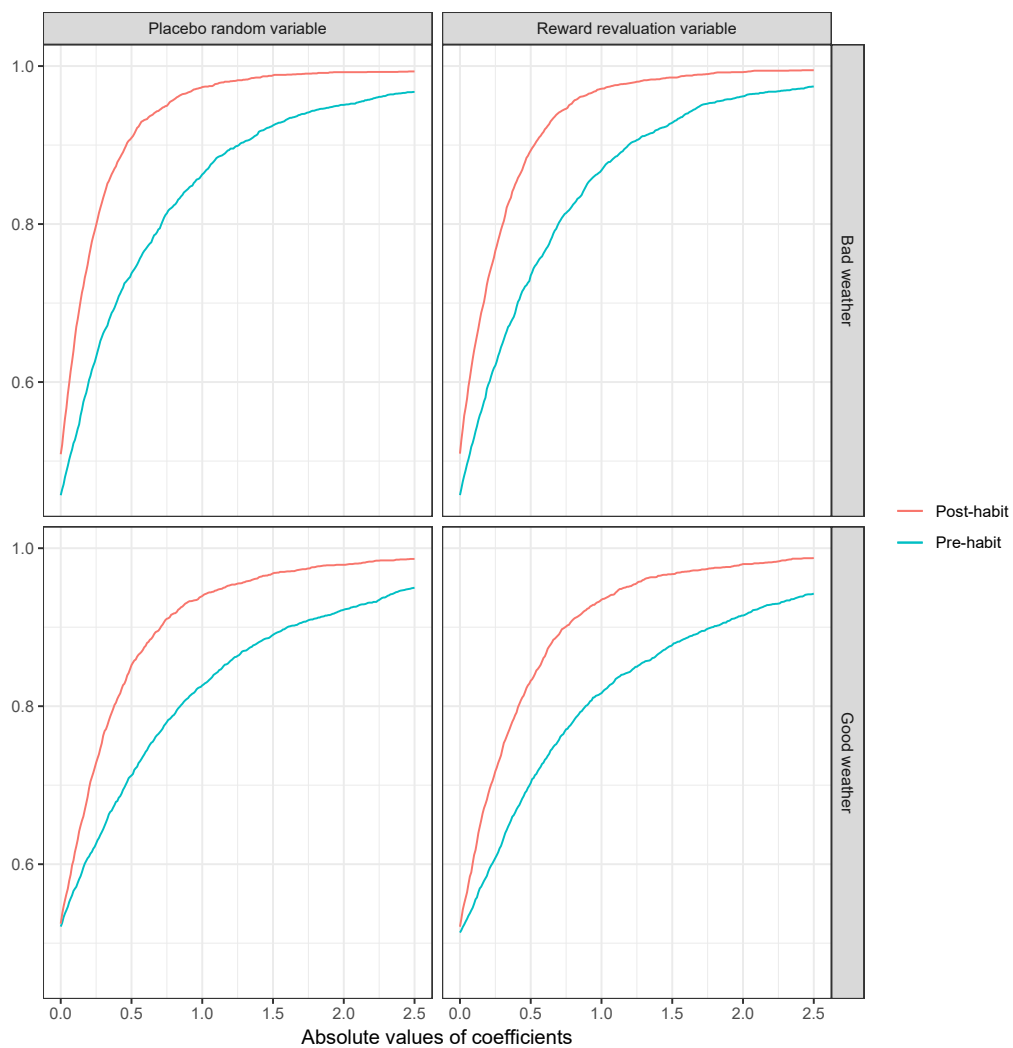


Figure 3.5: Reward devaluation sensitivity in gym attendance

This figure shows the empirical cumulative distributions in absolute values of (i) (right side) the interaction terms between pre-habit/post-habit and extreme weather indicators and (ii) (left side) the interaction terms between pre-habit/post-habit and placebo random variable.

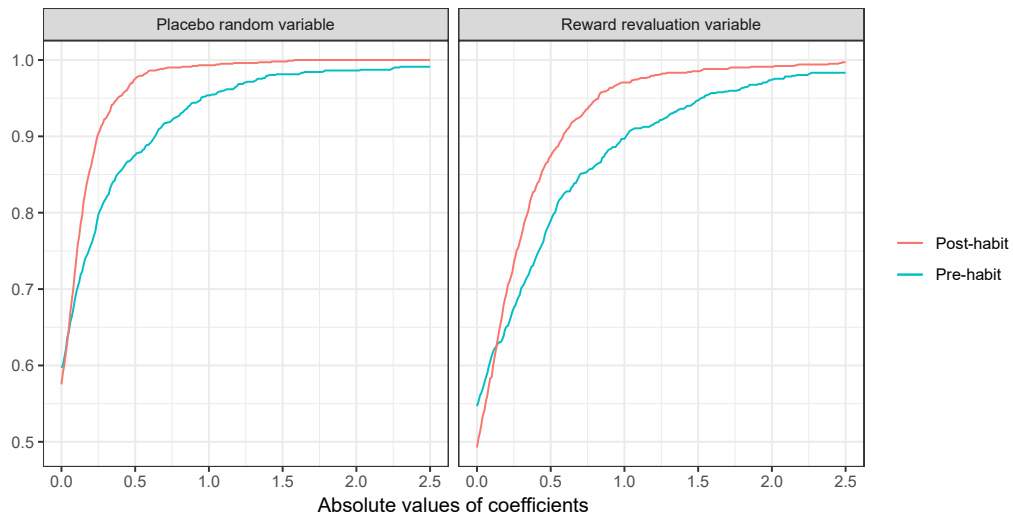


Figure 3.6: Reward devaluation sensitivity in handwashing

This figure shows the empirical cumulative distributions in absolute values of (i) (right side) the interaction terms between pre-habit/post-habit and last episode exit indicators and (ii) (left side) the interaction terms between pre-habit/post-habit and placebo random variable.

Appendix E - Demographic Predictors of AUC

Motivation

Given the rich individual-level data we work with, which includes a home zip code associated with each gym goer, it is possible to look for systematic categorical differences in the degree of predictability across sub-groups of gym goers. In order to run this analysis, we link the individual-level AUC data from gym goers with Census information from the year 2019. The Census data was purchased online from Income by Zip Code.¹¹ Unfortunately, the data on hospital workers does not come with zip code information, so we are unable to use this technique to analyze demographic differences with respect to the predictability of handwashing behavior.

The census variables discussed below, along with demographic data captured by the gym chain at time of registration (gender and age), allow us to estimate the demographic and SES profile of each individual gym goer and investigate demographic differences in gym attendance predictability.

Variable List

1. **Income:** As a proxy for individual income, we use the average household income of the individual's ZCTA.¹²
2. **Rural/Urban:** As a proxy for how rural or urban an individual's environment, we use a continuous measure of population density for the individual's ZCTA.
3. **Children:** As a proxy for an individual's likelihood of having children, we compute the fraction of married and single households in the gym goer's ZCTA who have children (under the age of 18).
4. **Age:** We have age data on the gym goers in our sample because they were required to report this at the time of gym membership registration. In addition, we calculate relative age by taking the difference between the median age in an individual's ZCTA (median age comes directly from the Census dataset) and their self-reported age.

¹¹The same data can be purchased at the following link under "Income by Zip Code List + Demographics (All US)." <https://www.incomebyzipcode.com/median-income-by-zip-code-list#pricing>

¹²ZCTAs, or ZIP Code Tabulation Areas, are generalized areal representations of United States Postal Service (USPS) ZIP Code service areas. While the latter is a trademark of the U.S. Postal Service, the former is a trademark of the U.S. Census Bureau. <https://www.census.gov/programs-surveys/geography/guidance/geo-areas/zctas.html>

5. **Gender:** We have gender data on the gym goers in our sample because they were required to report this at the time of gym membership registration.

Correlation Matrix

We analyze the correlation matrix of variables in our data which take on continuous values to see if there are significant and/or surprising correlations between demographic variables, as well as with individual-level AUC and base rates of attendance.

Significant correlations which are worth noting include the expected positive correlation between auc.test and auc.train ($\rho=0.661$) and the negative correlation between auc.train and base rate of attendance ($\rho=-0.237$), which underscores the difference between frequency of attendance and predictability. Also notable are the positive correlations between income and median age of neighborhood ($\rho=0.593$) as well as propensity to be married with kids ($\rho=0.373$) - which intuitively make sense given individuals tend to accumulate more income as they get older, and financial security gives people the ability to financially support children. Population density is negatively correlated with the median age ($\rho=-0.204$) (younger people are more likely to live in urban areas) and with having children if one is married or single ($\rho=-0.138$, $\rho=-0.229$, respectively).

Table 3.10: Correlation matrix of continuous variables

	Attendance rate	Age	Holdout AUC	Training AUC	Income	Pop. density	Median age	Married w/ kids	Single w/ kids
Base rate	1	0.002	0.073	-0.237	0.015	-0.010	-0.007	0.006	-0.004
Age	0.002	1	-0.007	0.073	0.062	-0.032	0.109	0.0005	0.076
Holdout AUC	0.073	-0.007	1	0.661	-0.002	-0.011	0.018	0.007	-0.004
Training AUC	-0.237	0.073	0.661	1	-0.020	-0.020	-0.010	0.011	-0.008
Income	0.015	0.062	-0.002	-0.020	1	0.027	0.593	0.373	0.063
Pop. density	-0.010	-0.032	-0.011	-0.020	0.027	1	-0.204	-0.138	-0.229
Median age	-0.007	0.109	0.018	-0.010	0.593	-0.204	1	-0.505	-0.258
Married w/ kids	0.006	0.0005	0.007	0.011	0.373	-0.138	-0.505	1	0.151
Single w/ kids	-0.004	0.076	-0.004	-0.008	0.063	-0.229	-0.258	0.151	1

References

- Adams, C. D. and A. Dickinson (1981). "Instrumental responding following reinforcer devaluation". In: *Quarterly Journal of Experimental Psychology* 33B, pp. 109–121.
- Adriaanse, M. A. et al. (2014). "Effortless inhibition: Habit mediates the relation between self-control and unhealthy snack consumption". In: *Frontiers in Psychology* 5, p. 444.
- Aldrich, J. H., J. M. Montgomery, and W. Wood (2011). "Turnout as a Habit". In: *Political Behavior* 33.4, pp. 535–563.

- Auld, M. C. and P. Grootendorst (2004). “An empirical analysis of milk addiction”. In: *Journal of Health Economics* 23.6, pp. 1117–1133.
- Balleine, B. and A. Dezfouli (2019). “Hierarchical Action Control: Adaptive Collaboration Between Actions and Habits”. In: *Frontiers in Psychology* 10.1, p. 2735.
- Balleine, B. and J. O’Doherty (2010). “Human and rodent homologies in action control: Corticostriatal determinants of goal-directed and habitual action.” In: *Neuropsychopharmacology : Official Publication of the American College of Neuropsychopharmacology* 35.1, pp. 48–69.
- Bargh, J. (1994). “The four horsemen of automaticity: Awareness, intention, efficiency, and control in social cognition”. In: *Handbook of Social Cognition: Basic Processes; Applications*. Ed. by R. S. Wyer and T.K. Srull. Lawrence Erlbaum Associates, Inc., pp. 1–40.
- Bayley, P., J. Franscino, and L. R. Squire (2005). “Robust habit learning in the absence of awareness and independent of the medial temporal lobe”. In: *Nature* 436.7050, pp. 550–553.
- Becker, G. S. and K. M. Murphy (1988). “A Theory of Rational Addiction”. In: *Journal of Political Economy* 96.4, pp. 675–700.
- Bedolla, L. G. and M. R. Michelson (2012). *Mobilizing Inclusion: Transforming the Electorate Through Get-Out-the-Vote Campaigns*. Yale University Press.
- Belk, R. W. (1975). “Situational variables and consumer behavior”. In: *Journal of Consumer Research* 2.3, pp. 157–164.
- Bernheim, B.D. and A. Rangel (2004). “Addiction and cue-triggered decision processes”. In: *American Economic Review* 94, pp. 1558–1590.
- Brady, H. E. and J. E. McNulty (2011). “Turning out to vote: The costs of finding and getting to the polling place”. In: *The American Political Science Review* 105.1, pp. 115–134.
- Brody, R. A. and P. M. Sniderman (1977). “From life space to polling place: The relevance of personal concerns for voting behavior”. In: *British Journal of Political Science* 7.3, pp. 337–360.
- Camerer, C. F., P. Landry, and R. Webb (2021). “The neuroeconomics of habit”. In: *The State of Mind in Economics*. Ed. by A. Kirman and M. Teschi.
- Campbell, J. Y. and J. H. Cochrane (1999). “By force of habit: A consumption-based explanation of aggregate stock market behavior”. In: *Journal of Political Economy* 107.2, pp. 205–251.
- Constantinides, G. M. (1990). “Habit formation: A resolution of the equity premium puzzle”. In: *Journal of Political Economy* 98.3, pp. 519–543.
- Coppock, A. and D. Green (2016). “Is voting habit forming? New evidence from experiments and regression discontinuities”. In: *American Journal of Political Science* 60.4, pp. 1044–1062.

- Cravens, M. (2020). "Measuring the strength of voter turnout habits". In: *Electoral Studies* 64, pp. 102–117.
- Crawford, I. (2010). "Habits revealed". In: *The Review of Economic Studies* 77.4, pp. 1382–1402.
- Dai, H., K.L. Milkman, and J. Riis (2014). "The fresh start effect: Temporal landmarks motivate aspirational behavior". In: *Management Science* 60.10, pp. 2563–2582.
- Danner, U., N. de Vries, and H. Aarts (2008). "Habit vs. intention in the prediction of future behaviour: The role of frequency, context stability and mental accessibility of past behaviour". In: *British Journal of Social Psychology* 47.2, pp. 245–265.
- Daw, N. et al. (2011). "Model-Based Influences on Humans' Choices and Striatal Prediction Errors". In: *Neuron* 69.6, pp. 1204–1215.
- Dayan, P. and K. Berridge (2014). "Model-based and Model-free pavlovian reward learning: Revaluation, revision and revelation". In: *Cognitive Affective Behavioral Neuroscience* 14.2, pp. 473–492.
- Deaton, A. (1992). *Understanding Consumption*. Oxford University Press.
- Denny, K. and O. Doyle (2009). "Does voting history matter? Analysing persistence in turnout". In: *American Journal of Political Science* 53.1, pp. 17–35.
- Dickinson, A., D. J. Nicholas, and C. D. Adams (1983). "The effect of the instrumental training contingency on susceptibility to reinforcer devaluation". In: *The Quarterly Journal of Experimental Psychology* 35.1, pp. 35–51.
- Downs, A. (1957). *An Economic Theory of Democracy*. New York: Harper & Row.
- Dubé, J.-P., G. J. Hitsch, and P. E. Rossi (2010). "State dependence and alternative explanations for consumer inertia". In: *The RAND Journal of Economics* 41.3, pp. 417–445.
- Duesenberry, J. S. (1949). *Income, Saving and the Theory of Consumption Behavior*. Cambridge, Mass.: Harvard University Press.
- Ersche, K. et al. (2017). "Creature of Habit: A self-report measure of habitual routines and automatic tendencies in everyday life". In: *Personality and Individual Differences* 116, pp. 73–85.
- Ferguson, S. G. and S. Shiffman (2009). "The relevance and treatment of cue-induced cravings in tobacco dependence". In: *Journal of Substance Abuse Treatment* 36.3, pp. 235–243.
- Fisher, I. (1930). *The theory of interest : As determined by impatience to spend income and opportunity to invest it*. New York: Macmillan Co.
- Fournier, M. et al. (2017). "Effects of circadian cortisol on the development of a health habit". In: *Health Psychology* 36.11, pp. 1059–1064.

- Fox, H. C. et al. (2007). "Enhanced sensitivity to stress and drug/alcohol craving in abstinent cocaine-dependent individuals compared to social drinkers". In: *Neuropsychopharmacology* 33.4, pp. 796–805.
- Franklin, M. N. and S. B. Hobolt (2011). "The legacy of lethargy: How elections to the European Parliament depress turnout". In: *Electoral Studies* 30.2, pp. 67–76.
- Gardner, B. (2015). "A review and analysis of the use of 'habit' in understanding, predicting and influencing health-related behavior". In: *Healthy Psychology Review* 9.3, pp. 277–295.
- Gardner, B., C. Abraham, et al. (2012). "Towards parsimony in habit measurement: Testing the convergent and predictive validity of an automaticity subscale of the Self-Report Habit Index". In: *International Journal of Behavioral Nutrition and Physical Activity* 9.102.
- Gardner, B. and P. Lally (2018). "Modelling Habit Formation and Its Determinants". In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 207–229.
- Garr, E. and A. Delamater (2019). "Exploring the relationship between actions, habits, and automaticity in an action sequence task". In: *Learning and Memory* 26, pp. 128–132.
- Gillebaart, M. and M. A. Adriaanse (2017). "Self-control predicts exercise behavior by force of habit, a conceptual replication of Adriaanse et al. (2014)". In: *Frontiers in Psychology* 8, p. 190.
- Gläscher, J. et al. (2010). "States versus rewards: Dissociable neural prediction error signals underlying model-based and model-free reinforcement learning". In: *Neuron* 66.4, pp. 585–595.
- Gonsalves, B. D. et al. (2005). "Memory strength and repetition suppression: Multimodal imaging of medial temporal cortical contributions to recognition". In: *Neuron* 47.5, pp. 751–761.
- Graybiel, A. M. (1998). "The basal ganglia and chunking of action repertoires". In: *Neurobiology of Learning and Memory* 70.1-2, pp. 119–136.
- Green, D. and A. S. Gerber (2002). "The downstream benefits of experimentation". In: *Political Analysis* 10.4, pp. 394–402.
- Harrington, N. (2017). "Commentary: Why it doesn't pay to ask consumers about habitual behaviors". In: *Journal of the Association for Consumer Research: The Habit-Driven Consumer* 2.3.
- Heckman, J. (1981). "Heterogeneity and State Dependence". In: *Studies in Labor Markets*. Ed. by Sherwin Rosen. National Bureau of Economic Research, Inc., pp. 91–140.
- Ji, M. F. and W. Wood (2007). "Purchase and consumption Habits: Not necessarily what you intend". In: *Journal of Consumer Psychology* 17.4, pp. 261–276.

- Karni, E. (2008). "State-dependent preferences". In: *The New Palgrave Dictionary of Economics*. Ed. by Palgrave Macmillan. Palgrave Macmillan, London.
- Kaushal, N. and R. Rhodes (2015). "Exercise habit formation in new gym members: a longitudinal study". In: *Journal of Behavioral Medicine* 38, pp. 652–663.
- Keane, M. P. (1997). "Modeling heterogeneity and state dependence in consumer choice behavior". In: *Review of Economics and Statistics* 15.3, pp. 310–327.
- Knowlton, B. and T. K. Patterson (2016). "Habit formation and the striatum". In: *Behavioral Neuroscience of Learning and Memory. Current Topics in Behavioral Neurosciences, vol 37*. Ed. by R.E. Clark and S. Martin. Springer.
- Koopmans, T. C. (1960). "Stationary ordinal utility and impatience". In: *Econometrica* 28.2, pp. 287–309.
- Kuehn, A. (1962). "Consumer brand choice as a learning process". In: *Journal of Advertising Research* 2, pp. 10–17.
- Kwon, H. E. et al. (2016). "Excessive dependence on mobile social apps: A rational addiction perspective". In: *Information Systems Research* 27.
- Laibson, D. (2001). "A Cue-Theory of Consumption". In: *Quarterly Journal of Economics* 116.1, pp. 81–119.
- Lally, P. et al. (2010). "How are habits formed: Modelling habit formation in the real world". In: *European Journal of Social Psychology* 40.6, pp. 998–1009.
- Lee, S., S. Shimojo, and J. O'Doherty (2014). "Neural computations underlying arbitration between model-based and model-free learning". In: *Neuron* 81.3, pp. 687–699.
- Mazar, A. and W. Wood (2018). "Defining habit in psychology". In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 13–29.
- Meredith, M. (2009). "Persistence in political participation". In: *Quarterly Journal of Political Science* 4.3, pp. 187–209.
- Middleton, J. C. et al. (2010). "Effectiveness of policies maintaining or restricting days of alcohol sales on excessive alcohol consumption and related harms". In: *American journal of Preventive Medicine* 39.6, pp. 575–589.
- Moors, A. and J. De Houwer (2006). "Automaticity: A theoretical and conceptual analysis". In: *Psychological Bulletin* 132.2, pp. 297–326.
- Neal, D., W. Wood, and J. Quinn (2006). "Habits—A repeat performance". In: *Current Directions in Psychological Science* 15.4, pp. 198–202.
- Neal, D., W. Wood, M. Wu, et al. (2011). "The Pull of the past: When do habits persist despite conflict with motives?" In: *Personality and Social Psychology Bulletin* 37.11, pp. 1428–1437.

- Orbell, S. and B. Verplanken (2010). “The automatic component of habit in health behavior: Habit as cue-contingent automaticity”. In: *Health Psychology* 29.4, pp. 374–383.
- Parry, D. et al. (2021). “A systematic review and meta-analysis of discrepancies between logged and self-reported digital media use”. In: *Nature Human Behaviour* (forthcoming).
- Pool, E. et al. (2021). “Determining the effects of training duration on the behavioral expression of habitual control in humans: A multi-laboratory investigation”. In: *PsyArXiv*: <https://psyarxiv.com/z756h>.
- Potthoff, S. et al. (2018). “Creating and breaking habit in healthcare professional behaviours to improve healthcare and health”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 247–265.
- Quinn, J. and W. Wood (2021). “Habits Across the Lifespan”. In: *Working Paper*.
- Rebar, A. et al. (2018). “The measurement of habit”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 31–49.
- Rhodes, R. and A. Rebar (2018). “Physical activity habit: Complexities and controversies”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 91–109.
- Rozen, K. (2010). “Foundations of intrinsic habit formation”. In: *Econometrica* 78.4, pp. 1341–1373.
- Ryder, H. E. and G. M. Heal (1973). “Optimal growth with intertemporally dependent preferences”. In: *The Review of Economic Studies* 40.1, pp. 1–31.
- Seger, C. A. and B. J. Spiering (2011). “A critical review of habit learning and the basal ganglia”. In: *Frontiers in Systems Neuroscience* 5.66.
- Shen, K. and D. E. Giles (2006). “Rational exuberance at the mall: Addiction to carrying a credit card balance”. In: *Applied Economics* 38.5, pp. 587–592.
- Shi, S. and L. Epstein (1993). “Habits and time preference”. In: *International Economic Review* 34.1, pp. 61–84.
- Sinha, R. (2009). “Modeling stress and drug craving in the laboratory: Implications for addiction treatment development”. In: *Addiction Biology* 14.1, pp. 84–98.
- Staats, B. R. et al. (2017). “Motivating process compliance through individual electronic monitoring: An empirical examination of hand hygiene in healthcare”. In: *Management Science* 63.5, pp. 1563–1585.
- Sundaresan, S. M. (1989). “Intertemporally dependent preferences and the volatility of consumption and wealth”. In: *The Review of Financial Studies* 2.1, pp. 73–89.
- Tricomi, E., B. Balleine, and J. O’Doherty (2009). “A specific role for posterior dorsolateral striatum in human habit learning”. In: *The European journal of neuroscience* 29.11, pp. 2225–2232.

- Verplanken, B. (2018). "Introduction". In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 1–10.
- Verplanken, B. and S. Orbell (2003). "Reflections on past behavior: A self-report index of habit strength". In: *Journal of Applied Social Psychology* 33.6, pp. 1313–1330.
- Volpp, K. G. and G. Loewenstein (2020). "What is a habit? Diverse mechanisms that can produce sustained behavior change". In: *Organizational Behavior and Human Decision Processes* 161, pp. 36–38.
- Wertenbroch, K. (1998). "Consumption self-control by rationing purchase quantities of virtue and vice". In: *Marketing Science* 17.4, pp. 317–337.
- Wilcockson, Thomas D., D. A. Ellis, and H. Shaw (2018). "Determining typical smartphone usage: What data do we need?" In: *Cyberpsychology, Behavior, and Social Networking* 21.6, pp. 395–398.
- Wit, S. de et al. (2018). "Shifting the balance between goals and habits: Five failures in experimental habit induction". In: *Journal of Experimental Psychology* 147.7, pp. 1043–1065.
- Wood, W. and D. Neal (2007). "A new look at habits and the habit-goal interface". In: *Psychological Review* 114.4, pp. 843–863.
- (2009). "The habitual consumer". In: *Journal of Consumer Psychology* 19, pp. 579–592.
- Wood, W., J. Quinn, and D. A. Kashy (2002). "Habits in everyday life: Thought, emotion, and action". In: *Journal of Personality and Social Psychology* 83.6, pp. 1281–1297.
- Wood, W. and D. Runger (2016). "Psychology of Habit". In: *Annual review of psychology* 67, pp. 289–314.
- Yin, H. H. and B. Knowlton (2006). "The role of the basal ganglia in habit formation". In: *Nature Reviews Neuroscience* 7, pp. 464–476.

*Chapter 4***USING A VENDING MACHINE "RETAILER" TO STUDY
REPEAT PURCHASES IN CONSUMER BEHAVIOR**

ABSTRACT

Installing a customizable vending machine on a university campus, we run a field experiment to understand how consumers respond to a price promotion in order to credibly dissociate predictions made by brand loyalty/habit formation from reference-dependence theories. The vending machine serves as a “mini-retailer,” allowing for the control of all price promotion treatment details in an ecologically valid setting, and collecting granular panel data on purchases occurring at all hours of the day over the course of 10 weeks. The vending machine can also be programmed to control for stockpiling behavior, which is an important concern for empirical work analyzing price promotions in the marketing literature. Analysis of the data collected during this study suggests that price promotions increase the sales of both discounted and non-discounted items, as well as the number of unique customers making purchases. Furthermore, in line with the loss leader hypothesis, more items are purchased during the sale period overall (although the difference is not statistically significant).

4.1 Introduction

Price promotions are of interest to marketing researchers given they are a ubiquitous strategy across real-world retailers. It was estimated that 28% of consumer goods in developed countries were bought on promotion in a recent market analysis (Eales, 2016). The report found that food tends to be promoted even more often than non-food in the U.S. market, with 37.4% of food items being on promotion at time of sale.¹ No sector is spared - price promotions are used on everything from technology to automobiles.

One practical reason for the use of price promotions is to liquidate current stock. This is particularly applicable to non-durable goods, like food items which are reaching their expiration date and risk spoiling (e.g. in the U.S., it is common to see "30% off turkey" immediately after Thanksgiving). Likewise, time-dependent goods such as a past season's clothing may need to be cleared in order to make space for next season's fashion (e.g. one might find "50% off swimwear" in several stores at the end of August).

Another reason to use price promotions has nothing to do with liquidating stock, but is rather oriented to achieve a specific outcome in purchasing behavior. Retailers might lower prices in hopes that existing consumers purchase more goods, thereby increasing sales volume in a denser time period. Retailers might also lower prices in order to acquire new consumers, or "gain market share." These include promotions akin to "25% off your first purchase," which are clearly designed to attract new consumers away from full-priced alternatives. This seems like a reasonable strategy. As the price drops, those with lower willingness to pay (i.e. "bargain shoppers" with low valuations) are now willing to purchase the item assuming the discounted price is lower than their WTP. Hence a sale may attract new customers who are now willing to pay for the reduced product. A lower price also increases the attractiveness of one item versus an equivalent substitute good, potentially strengthening brand loyalty as existing customers get more "value for money." Assuming existing customers do not adjust their reference points (discussed shortly), there is no reason to believe that purchasing more items at the sale price will harm or reduce brand loyalty once products are back to their original prices.

However findings from behavioral economics may overturn the benefits from these strategies. There is some evidence, for example, that consumers *do* update their

¹Compared to 25% of non-food items that same year.

reference points dynamically with changes in price.² This means that as a price reverts back to its non-discounted equivalent, consumers may suffer loss aversion. This would imply that putting items on sale actually risks making the firm's offering *less* attractive as compared to reasonable substitutes from competitors (Ray, Shum, and Camerer, 2015, Mela, Gupta, and Lehmann, 1997).

To make it more confusing, industry "intuition" regarding the benefit of price promotion varies. Some firms never put items on sale, presumably so that customers don't devalue their products (e.g. LVMH cuts up any unsold bags at the end of their season before putting them in the trash, to avoid the possibility of them being found and resold at a lower price). Other firms have frequent (usually seasonal) sales, presumably to remind and attract customers of their brand and steal away market share from competitors. Research suggests there may even be vertical differentiation in how firms think about price promotions (Marom and Seidmann, 2011, Bar-Isaac, Caruana, and Cuñat, 2012).

Many papers in the marketing literature have studied the impact of price promotions on consumer behavior. Several focus specifically on how details of the promotion itself can lead to different behavioral responses. For example, DelVecchio, H. S. Krishnan, and Smith, 2007 investigated how the framing of a price promotion in terms of percent or dollar terms leads to different outcomes. Meanwhile, Osborne, 2018 showed that the depth (i.e. magnitude) of a discount is more effective than the frequency with which discounts occur when it comes to increasing quantity of sale items sold and overall revenue - with increasing the depth of a discount while *decreasing* frequency being most effective.

Other studies have focused on the broader implication which price promotions have on subsequent consumer behavior. However most of these studies are interested in choice dynamics during or immediately after the promotion itself (Anderson and Simester, 2004, Hendel and Nevo, 2003, Chandon, Wansink, and Laurent, 2000, Mela, Gupta, and Lehmann, 1997). And many studies focus on how price promotions relate to perceptions of brand equity (Valette-Florence, Guizani, and Merunka, 2011, Ailawadi, Lehmann, and Neslin, 2003), rather than how they influence future choices.

Results from the small number of papers which *have* investigated the long-term impact (over one month following a price promotion) on consumer behavior have demonstrated that sales can lead to stockpiling behavior. Stockpiling describes

²We define reference points in the Literature Review on Reference-Dependence.

the phenomenon of consumers purchasing more goods at a lower price during the promotion than they would if those goods were being sold at full-price, with the goal of consuming them later (when it would be costlier to purchase them at their regular price). To quote DelVecchio, Henard, and Freling, 2006, "stockpiling leads to lower aggregate or per consumer sales for a brand following a sales promotion by taking consumers out of the market due to greater on-promotion purchase quantities." Stockpiling has been modelled and estimated in consumer data (Ching and Osborne, 2020, Sun, 2005, Erdem, Imai, and Keane, 2003, Pesendorfer, 2002, Helsen and Schmittlein, 1992). Stockpiling behavior can be a confound in empirical analysis because it makes it difficult to tease apart, for example, predictions made by a model of reference-dependence from predictions made by a model of habit formation.

Furthermore, research has looked at how stockpiling behavior can affect consumer expectations (Ailawadi, Gedenk, et al., 2007, Pauwels, Hanssens, and Siddarth, 2002, Gupta, 1988).³ For example, some older research indicates that households are more likely to wait for the next promotion before making any subsequent purchase decisions if they have had a lot of exposure to price promotions previously (Mela, Jedidi, and Bowman, 1998), and that consumers may become less likely to purchase frequently discounted items unless they are on sale (Helsen and Schmittlein, 1992). In other words, price promotions can make consumers more price sensitive and hurt long-term retailer profits.⁴

Stockpiling is therefore an important concern and one which is often accepted as a potential confound in marketing research that relies on empirical data from retailers who cannot stop or control stockpiling behavior. One advantage of our methodology is the ability to control for stockpiling behavior at the time of data collection, rather than seeking to control for it in the data analysis (Seiler, 2013, Chan, Narasimhan, and Q. Zhang, 2008, Hendel and Nevo, 2006, Erdem, Imai, and Keane, 2003).⁵ An

³While most of the research focuses on the negative impact of stockpiling behavior on consumer expectations, Ailawadi, Gedenk, et al., 2007 does consider two potential benefits - specifically, increased category consumption and preemptive brand switches (i.e. purchasing more of a discounted brand today means the consumer will be consuming the same brand tomorrow instead of switching to a competitor brand).

⁴Perhaps one disadvantage of our experimental setup is that we cannot account for these "anticipation" effects, since there is no way for the subjects in our study to anticipate the start or the end of discounts. This is in contrast to real markets, where seasonality and other cues can serve to increase the probability that a price promotion can be accurately anticipated.

⁵As discussed further in Ching and Osborne, 2020, current empirical literature on stockpiling employs the discrete choice dynamic programming framework from Rust, 1994 in which a dynamic structural model is estimated on the consumer purchasing data and then used to make counterfactual predictions.

additional advantage is the ability to see individual-level purchase data through time (as opposed to household-level panel data such as the Nielsen Retail Scanner Data).

The focus of this paper is understanding what happens to consumer behavior during the price promotion as well as several weeks to one month following.⁶ This objective is in line with recent interest in the marketing literature on the longer term impact of price promotions (and other marketing activities) (see Hanseens, 2018 for a compendium). In this study, we ask the question of whether price promotions can lead to the creation of brand loyalty.⁷ We run a field experiment using a vending machine intended to credibly dissociate brand loyalty/habit formation from reference-dependence behaviors. We use a tightly controlled price promotion treatment and collect granular panel data on purchases over the course of 10 weeks (something which is difficult, if not impossible, to do using lab studies).

A field experiment is particularly valuable in this scenario because it is very difficult to accurately answer this research question using empirical data alone. Commonly used consumer panel data, like the Nielsen Retail Scanner Data, does not have indicators for price promotions.⁸ Furthermore, in the absence of partnering with a retailer who is happy to cede all decision making about a price promotion to a researcher, researchers have no control over when price promotions are put in place, what size discounts are, or what is happening to close substitutes. The vending machine serves as a “mini-retailer” allowing us to control all details of a price promotion treatment in an ecologically valid setting. Importantly, it also allows us to control for stockpiling behavior, since we can program the machine to limit the number of purchases made at any one time (i.e. more than one of each good cannot be purchased at time of sale).⁹

In this paper, we first make three standard economics predictions about the behavior we should expect to see during the sale period. We then show that, as predicted, the price promotion leads to an increase in purchases of both discounted and non-

⁶We have 4 weeks of consumer behavior data following the end of the price promotion, which qualifies us to analyze “medium-term” effects. Unfortunately, the COVID-19 pandemic created an abrupt end to this and additional studies, so truly long-term effects (several months after the end of the price promotion) were not analyzed in this study.

⁷As we discuss shortly, we will use habit formation and brand loyalty interchangeably given they make the same *behavioral* predictions about the choice data.

⁸There are ways to impute price promotions in the data by looking at intermittent low prices, but this is an imperfect process which introduces various levels of measurement error.

⁹Stockpiling is not impossible - an individual could stand there and vend multiple items one after the other - but we assume that some social pressure of others being nearby, and the cost of time waiting for subsequent things to vend, is enough discouragement for serious stockpiling behavior.

discounted items during the sale period. Also as predicted, the total number of unique customers increases during the sale period.

Second, we compare two sets of predictions which are made by the brand loyalty/habit formation and neuro-autopilot theories (directionally the same as one another) versus reference-dependence (directionally opposite to those above). As predicted by brand loyalty/habit formation models, the purchases of discounted items appear to remain at elevated levels compared to the pre-sale period, although this is not statistically significant in our regression analysis. Finally, the total number of customers remains relatively unchanged during the post-sale period as compared to pre-sale behavior, which is not in line with predictions from either model. Further work needs to be done to fully disassociate brand loyalty, habit formation and neuro autopilot theories, which lead to similar behaviors via different underlying mechanisms, which we discuss towards the end of the paper.

Finally, the paper offers a methodological contribution by presenting a new way to run field experiments in marketing. By controlling a vending machine remotely as we do here, it is possible to carefully collect data on human behavior in the field using a tightly controlled experimental framework while reducing or entirely avoiding experimenter demand effects (Zizzo, 2010). In a lab setting, both the physical presence of an experimenter (typically a graduate student, post doc, or faculty member) who holds some position of authority over or compared to the subject (typically an undergraduate student), as well as the expectation that payment for the experimental task will be a function of the subjects' behavior, may lead to subjects' behaving differently in order to "please" the experimenter. In our setting, these risks are largely removed since the experimenter is not present at time of data collection, subjects receive no payments or incentives to purchase from the machine, and subjects are anonymous. Therefore, this "best of both worlds" method offers a number of opportunities for exploring promising future directions, such as measuring the role of attention and peer effects on choice.

4.2 Literature Review

In the following section, we provide an overview of three theories of consumer behavior from the marketing and economics literature which are most relevant to my research question. They are brand loyalty (known as habit formation in economics), neuro autopilot, and reference-dependence.

Brand Loyalty/Habit Formation

Brand loyalty, a subset of broader customer loyalty, has been defined as the strength of the relationship between a consumer's attitude to a specific brand and their subsequent repeat patronage (Dick and Basu, 1994). Repeat consumption has been studied in the marketing literature for many decades (for a recent review, see O'Brien, 2021), and brands themselves have become of immense interest to marketing scholars given their impact on consumer behavior (Keller and Lehmann, 2006). Loyalty, as the word suggests, implies that a customer feels a sense of allegiance with the brand such that if price or other attributes are equivalent between products, the brand itself would determine the customer's choice (Fishbach, Ratner, and Y. Zhang, 2011). Developing brand loyalty guides the marketing activities of many firms which aim to develop, maintain and enhance customers' loyalty towards its products and services.

The marketing literature is not overly prescriptive about the mechanisms which give rise to brand loyalty (Oliver, 1999). It may arise from accessibility, wherein consumers are more likely to develop associations with one brand due to its physical or psychological proximity (see Wisker, Kadirov, and Nizar, 2020 for a recent example of cultural accessibility). It may also arise from affective associations, like an emotional connection with a specific brand that has grown to represent something about the consumer's personal or collective identity (Coelho, Rita, and Santos, 2018, Yeh, Wang, and Yieh, 2016, He, Li, and Harris, 2012) which can be conceptualized as "brand love" (Batra, Ahuvia, and Bagozzi, 2012). There may also be more practical considerations like switching costs - monetary costs which are incurred when a consumer switches from one product to another, or psychological switching costs in the form of inertia (Dubé, Hitsch, and Rossi, 2009). Regardless of the mechanism, most consumers are aware of their brand loyalty - they have strong preferences for specific brands and would be able to claim as such in a brand association survey (Romaniuk and Nenycz-Thiel, 2013).

While we do not measure brand associations in our data, a model of habit formation from economics is similar to brand loyalty in the sense that a researcher would look to empirically observe repeat purchases in retail behavior. The traditional economic model of habit formation focuses solely on history-dependence and is motivated by evidence of an empirical correlation between past and current consumption. Specifically, this traditional model is the theory of rational addiction proposed by Becker and Murphy, 1988. They propose that current utilities depend on consumption

history and that highly habituated (or “addicted”) people, by consuming more of a product today will increase their value of that product tomorrow. Moreover these “rationally addicted” people are self-aware, in that they understand that consuming more today will result in them valuing the same product more (and hence spending more to purchase) tomorrow.

It is unsurprising that economics defines habit in such a way, given that economists tend to abstract everything to a function of preferences, beliefs and constraints. However, an understanding of psychology and neuroscience allows us to appreciate that habits have an adaptive function in releasing expensive cognitive capacity for highly predictable decisions. For example, a more psychologically-informed model of habits would predicts that a choice is made quickly (“automatically”), even if the reward is changed slightly, whereas brand loyalty would predict that a choice would be made slowly (a consumer is still doing a utility maximization problem) and slower still if the reward value has changed.

However the terms brand loyalty and habit formation will be used interchangeably in this chapter because the focus of study is the *behavior* observed, and not the underlying mechanism. To quote Tam, Wood, and Ji, 2009, “repeated patronage can reflect strong habits and be cued by stable features of purchase and consumption contexts or it can reflect brand loyalty and be influenced by strongly held, favorable brand evaluations that direct re-purchase and consumption intentions.” Since this experiment does not provide a way of disassociating the mechanism underlying the repeat purchases we observe, we focus on the behavioral predictions made by brand loyalty and habit formation theories together, in which we assume that more purchases lead to stronger habits/loyalty.

Neuro Autopilot

Introduced by Camerer, Landry, and Webb, 2021, neuro autopilot is a computational model which provides a psychologically-informed alternative to an economist’s definition of habit (which is essentially history-dependence, as discussed in the previous section). The motivation for the neuro autopilot model comes from research in computational neuroscience which demonstrated that the human brain operates in two decision making modes: a slow model-based mode which engages in classic utility maximization behavior and a fast habit mode which performs a simpler set of computations.

The habit mode is a valuable default if context-stability is high because utility

maximization requires costly cognitive effort. If one takes a break while working on their thesis chapter to get a drink from the cafe across the street, it is easier to choose the same reliable iced tea they had yesterday instead of standing and weighing the utility of each of the twelve different drink options – time and cognitive effort which is now being dedicated to drink choice instead of the thesis chapter. In other words, habit mode allows putting the brain into an automatic gear in order to conserve mental resources. But how does the human brain know when to use the model-based and when to use the habit system to make decisions?

Scientifically, there is evidence in favor of a neural arbitration mechanism. Lee, Shimojo, and O’Doherty, 2014 provide evidence for a system which allocates control over behavior across the model-based and habit systems based on the reliability of their predictions. The authors find evidence that the inferior lateral prefrontal and frontopolar cortex encode reliability signals and the output of a comparison between those signals, implicating these brain regions in the arbitration process. They go on to explain that this arbitration likely works through modulating the habit system in particular, when the arbitrator deems it should be “overwritten” in order for the model-based system to drive behavior. In other words, the habit system is used as the default so long as its predictions are reliable, and when it is no longer reliable, the model-based system is asked to come in.

This concept is modelled mathematically in the neuro autopilot theory. Camerer, Landry, and Webb, 2021 present a system of equations which mimic how the brain makes decisions when in habit mode. Specifically, if a consumer is choosing between two products a and b at time t , instead of choosing the product whose utility is highest (as modeled in economics using preference-based choice), the habitual consumer would simply repeat their previous choice without doing a utility calculation if their prediction about the value of the good has high reward reliability.

Say I am a consumer who has repeatedly chosen product a when again faced with the option between a and b at time t . Assuming I consumed a at time $t-1$, to form my prediction about a I would calculate a time-weighted average of my past utilities from consuming A as follows:

$$r_t(a) = (1 - \lambda_r)r_{t-1}(a) + \lambda_r u_{t-1}(a)$$

In the above, λ_r represents a learning rate parameter which captures how much weight I am placing on the present experienced utility versus the past (the closer λ_r ,

is to 1, the more weight is placed on present subjective value and the quicker my choices respond to value changes). Written in another way, this is essentially my prediction made at time $t-1$ plus the learning rate multiplied by the reward prediction error or RPE which captures how spot-on (or not) my prediction about value was. This is very similar to theories from reinforcement learning, which seek to model learning speeds in other situations.

In their theoretical framework, the customer also tracks the reliability of the options in her choice set by calculating what the authors dub a “doubt stock.” A low value of doubt stock means an option is reliable. So if I have low doubt stock in a , it means my past predictions $r(a)$ have been close to my experienced utility. On the contrary, if my doubt stock is high, it indicates that my historical predictions were not very close to my experienced utility, and that a is in fact not very reliable. Another learning rate parameter λ_d determines how much weight the consumer places on the most recent doubt stock in their calculation.

This doubt stock therefore evolves according to:

$$d_t(a) = (1 - \lambda_d)d_{t-1}(a) + |r_t(a) - u_t(a)|$$

This concept of doubt stock captures reward reliability in an elegant way. Essentially, if I continue to choose option a and the subjective value $u_t(a)$ does not change, my RPE will approach zero (in absolute value) and the doubt stock will shrink towards zero. But if the subjective value suddenly changes, I will experience a large RPE, which will increase my doubt stock (reducing my “trust” in the reward reliability of A). If the doubt stock remains below some threshold, I would remain in habit mode, exploiting the previous choice (A) without consideration of other available options.

But if enough large RPEs push the doubt stock to be outside some threshold, habit mode is “interrupted” by the arbitration system mentioned earlier, which now differs decision making to the model-based system. The model-based system then chooses the utility-maximizing option, which may be equivalent to or different from the previously habitual choice.

One analogy in the context of marketing is that when a customer is new to a store or set of products, they need to learn about all the options available (“sampling” akin to pulling different levers in a multi-arm bandit task). Once the customer has gained enough information to solve their utility maximization problem, they choose a specific option, their preferred a amongst the set. Assuming attributes of the set

remain static, the consumer continues to choose a so long as it is reliable (the utility, calculated from e.g. price and quality, remains highly predictable). If something changes such that option a is no longer reliable – e.g. the quality or price shifts considerably – the consumer is “jolted” back into the state they were in when they were brand new to the choice set, and is again forced to rely on their preference-based decision making to guide their choices between the available options.

The implication for behavior following a price promotion is that the neuro autopilot model would predict that there are some consumers for whom the doubt stock is too high at the time of their first purchase following a price promotion, such that they do not purchase a second time in the post-sale period.

This model is therefore a true “neuroeconomic” approach in that it takes the best of economics (a formal mathematical structure, which makes clear predictions about behavior) and the best of neuroscience (a deep understanding of the psychological and neuroscientific implementation of habit). To quote, “neuroeconomics seeks to establish what algorithms of economic choice achieve high-level function goals, and are actually implemented by neural circuitry and other biological forces” (Camerer, Landry, and Webb, 2021). These comments echo David Marr’s framework of Three Levels of Analysis, which guides a lot of computational neuroscience research. The neuro autopilot model of habit provides an answer to all three of Marr’s core questions: the “why” question (why is the brain functionally organized this way to run this specific computation?), the “what” question (what is the algorithm being used to solve this computation?) and the “how” question (how is this algorithm physically implemented in the brain?). (Marr, 1982).

Reference-Dependence

The idea of reference-dependence comes originally from perception research, and was introduced into economic theory with the publication of Prospect Theory (Tversky and Kahneman, 1979). This proposed alternative to expected utility theory posits that individuals do not make decisions about gains and losses in absolute terms but rather with respect to a specific reference point. Prospect Theory has four key tenets, of which reference dependence (the idea that utility is determined with respect to a reference point r) is only one. The other three are loss aversion (losses have coefficient λ , assumed to be greater than 1, which make them evaluated to be larger than equivalent gains), diminishing marginal sensitivity (the utility function is concave for gains and convex for losses), and a probability weighting func-

tion (which overweights low probability events and underweights high probability events). Formally, utility is frequently represented with the following specification:

$$\begin{aligned} &(x - r)^a \text{ if } x > r \\ &\lambda(x - r)^a \text{ if } x < r \end{aligned}$$

While these other three aspects of Prospect Theory have clear parameters, Kahneman and Tversky left the definition of reference points intentionally vague. So while we define them as data points which anchor the decision maker and influence their choice, there is disagreement around how these reference points are determined.

For example, hypothetical reference points may be backward looking (such as the “status quo” which may be, e.g., zero dollars at the start of an experiment) or forward looking (a goal, or a desired outcome). Aspirational reference points have some predictable properties. For example, marathon runners tend to run faster to meet a “round number” time (such as a 4:00 hour marathon, instead of a 4:05 marathon) with this reference point influencing male runner behavior more than female runner behavior (Allen and Dechow, 2020). Kőszegi and Rabin, 2007 argue that expectations are crucial for the formation of reference points, however this has been weakly supported in the literature. Abeler et al., 2011 found that a high probabilistic expectation about the wage received at the end of an experimental task led subjects to work longer and earn more money, but this result has been weakly replicated (Camerer, Dreber, et al., 2016 found a positive but insignificant effect size). In financial markets, reasonable reference points can be a relevant benchmark, expected returns, the risk-less rate, and so forth. In other words, the implementation of reference points is highly complex and reference points may even be determined jointly with information encoding (biologically, this would be to allow allocation of scarce attention in a useful way).

In a consumer setting, reasonable reference points might be willingness to pay for an item, WTP (which may be different from willingness to accept, WTA, per Kahneman, Knetsch, and Thaler, 1990) or the existing price of an item. If the existing price is used as the reference point, a suddenly lower price would make the consumer feel like they are purchasing “at a gain.” Likewise, a suddenly higher price would make them feel like they are purchasing “at a loss,” which should evoke a bigger response given the loss aversion parameter λ (K. and W., 1995, Putler, 1992). Such perceptions about the “expensiveness of a price” following a lowered reference

point are difficult to study in market data, but have been explored experimentally (Slonim and Garbarino, 1999). Reference points may also be a more complicated function incorporating the initially observed price as well as the latest price viewed (Baucells, Weber, and Welfens, 2011).

The formation of reference points has been studied in the context of how price promotions affect consumer choices as well (see Mazumdar, Raj, and Sinha, 2005 for a review). Some of this research has discovered that price promotions are “not created equal.” For example, DelVecchio, H. S. Krishnan, and Smith, 2007 found that framing a promotion in terms of percentage discount (versus absolute dollar amount) has less of an effect on price expectation (does not “reset the reference price” as dramatically) following the promotion itself.

In our analysis, we use the status quo as our reference point, and for ease of analysis, we assume that consumers update their reference points automatically after a new price is visible. In other words, immediately after a 50% sale is introduced, the new price enables an existing customer to buy more of the product they purchased previously for the same amount. Furthermore for customers whose WTP now exceeds the price, the price promotion moves them into the “domain of gains” when they are purchasing.

If indeed “increased discounting on previous purchase occasions results in lower reference prices on the current purchase occasion” (Mela, Jedidi, and Bowman, 1998), then we hypothesize that the consumer’s reference point is updated to be the sale price when they make a purchase during the price promotion. Once items are put back to pre-sale levels, purchasing now feels like being in the “domain of losses” since the new price is higher than their updated reference point of the sales price.

4.3 Experimental Design

In the following section, we walk through the experimental design used in this study. We look at the physical set-up of the vending machine used to collect data on consumer purchases. We then discuss the subject pool and additional details of the study design.

Vending Machine Physical Set-Up

We use a customizable vending machine to run our experiment. The machine was built by Digital Media Vending in California to meet our specifications¹⁰. We

¹⁰More detailed specifications regarding the hardware of the machine can be found in Appendix A.

installed the vending machine in a common area used by about 330 undergraduates on Caltech's campus. It was turned on and stocked for four weeks prior to the start of the experiment so that students could get comfortable with the introduction of the new machine in their common area. From the start of the experiment, it collected data non-stop ("24/7") over the course of 10 weeks, the length of a full academic term at Caltech, meaning individuals could purchase snacks and drinks at any time they liked during this period. It was kept fully stocked with 25 unique products (a combination of snacks and drinks), which were chosen to represent a range of healthy and unhealthy, as well as cheap and expensive snack items. The choice of products was informed by ratings from Caltech undergraduates collected using a pre-testing survey¹¹ conducted prior to the start of the experiment.

Several days prior to the launch of the experiment, an email was sent to all students who have access to that common area to let them know about the arrival of the machine,¹² so that all subjects would learn about the vending machine at the same time, as opposed to sequentially and dependent on when they next visited the common area. Finally, to comply with IRB protocol, we permanently attached the following notice to the machine letting people know that it was being used for research purposes:

"This machine is being used for research purposes, including analysis of your purchase data. All information is de-identified. If you have any questions, email caltech.vending@gmail.com. Minors may not purchase from this machine."

Subject Pool

All of our subjects were members of the Caltech community. We did not collect any identifiable information about them beyond a unique ID, which we obtain from the unique first and last four digits of their credit card number and use to link their purchases over time. Given the vending machine's location, we assume that the majority of users were Caltech undergraduates (typically between 18-22 years old). In addition, a number of staff members who regularly use the space may have also purchased from the machine. No attempt was made to collect individually identifying characteristics in order to protect the privacy of subjects and to maintain external validity in the experience of making a vending machine purchase.

¹¹Extracts of the pre-testing survey and aggregate results can be found in Appendix B.

¹²Exact wording of the email sent can be found in the Appendix C.

Programmable Interface

The vending machine uses an Android touchscreen interface, which needed to be programmed to the front-end desired. For each item that was displayed and sold in the machine, we uploaded a single photo of the item and assigned it an associated price.¹³ These were the only two pieces of information that participants could see when they encountered the machine, as can be seen from Figure 4.1.

Inside, the machine was stocked with the full range of items advertised (see Figure 4.2). Upon stocking the items, we updated the inventory amounts, which automatically updated on the software used to track the machine remotely. This software allowed me to track inventory live, such that we would know whether an item was close to being out of stock and could re-stock it in a short period of time. We were able to keep the machine well stocked and did our best to avoid stockouts given their ability to influence consumer perception and therefore our results (Anderson, Fitzsimons, and Simester, 2006). However, given a lot of purchasing activity occurred overnight and on the weekend, we did experience four instances of a product being out of stock during the experimental period, all of which were re-stocked within 24 hours. While it is possible that these stock-outs influenced results, we believe it is unlikely given we received no emails or complaints when items were temporarily unavailable.

If an item went out of stock or was unable to vend due to a technical malfunction, the photo of the item would be greyed out with a sign indicating that the item was "Sold Out." On the back-end, we would be able to see whether a certain rack needed attending to. Finally, since the vending machine uses weight-sensitive elevator technology, the machine never charged customers for items that it was not able to vend (avoiding any "the machine ate my dollar" issues).

Study Design

During weeks 1-4 of the experiment (hereafter referred to as the "Pre-Sale" period), we tracked "natural" consumer behavior to see how individuals use the machine with the basic set-up and prices. During this period, individuals could freely purchase the 25 products at their normal prices (all around a 1.4-2.2x retail mark-up). This allowed me to establish a "baseline" of purchasing behavior.

¹³We used cost-based pricing, first calculating the per-unit cost of the snack being sold and then identifying a price above this which was a multiple of \$0.50 and targeted a retail mark-up of 1-2, such that items put on the 50% discount would be sold below cost. See the Discussion for a more in-depth analysis on pricing.

Figure 4.1: The front of the customizable vending machine.

At the top of the machine, customers (mainly students) were shown a generic message prompting them to scroll through the snacks and drinks available for purchase, along with an email address to contact with any questions (or special requests, before the experiment started). The touchscreen then displayed a photo of all the items inside, along with a price. Customers could scroll through all of the items and add things to their cart, prior to checking out and purchasing everything at once.

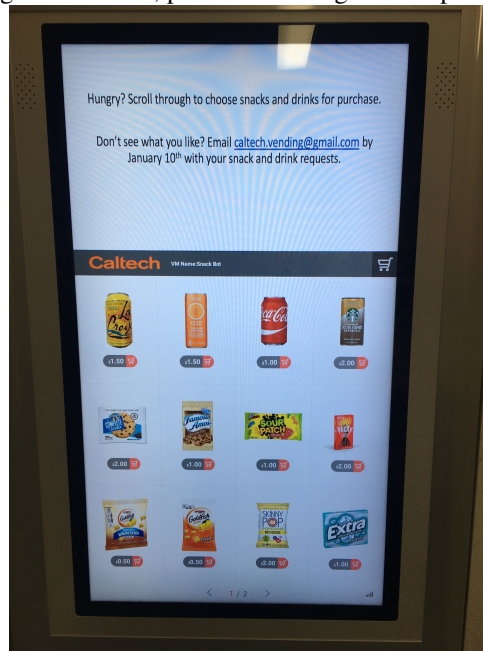


Figure 4.2: The inside of the customizable vending machine.

Each rack was stocked with a different snack item. A sensor, using the weight on the rack, was used to track real-time inventory, such that the machine could be restocked when items were running low and before they experienced a stockout.



Weeks 5-6 of the experiment (hereafter referred to as the "Sale" or "Treatment" period) involved observing consumer behavior during a price promotion which reduced the prices of 10 of the 25 products by 50%¹⁴. Each of the sale items had at least one close substitute product which remained at full price during the treatment.¹⁵ Table 4.1 lists the top 10 most purchased products at the aggregate level, half of which were products that remained full price throughout the treatment and half of which were discounted by 50% during the treatment period. Balance tests were performed to confirm that there were no significant differences across the products with respect to per-unit cost, original price, and retail mark-up confirming that the products were successfully randomly assigned to discounted and non-discounted treatments.

Finally, during weeks 7-10 of the experiment (hereafter the "Post-Sale" period), we tracked consumer behavior following the treatment. The interface of the machine was programmed to be just as it was during the pre-treatment period, with prices and photos identical to the pre-treatment phase.

Table 4.1: Top ten products, by net sales

Product Name	Pre/Post Price	Sale Price	Total Sales	Substitute
Coke	\$1.00	\$0.50	199	Izze
Sour Patch	\$1.00	\$0.50	44	Fruit Snacks
Complete Cookie	\$2.00	\$1.00	41	Famous Amos
Starbucks	\$2.00	\$1.00	36	Monster
Goldfish	\$0.50	\$0.25	30	Lays
Izze	\$1.50	N/A	90	Coke
Fruit Snacks	\$1.00	N/A	77	Sour Patch
Pringles	\$1.50	N/A	47	SkinnyPop
Famous Amos	\$1.00	N/A	33	Complete Cookie
Lays	\$0.50	N/A	32	Goldfish

Price and sales data on the top selling products during the 10-week study period. The top panel contains products which experienced a sale during the treatment period, and the bottom panel contains products which did not experience a sale during the treatment period. Examples of substitute products which were available at the time of purchase are provided.

¹⁴A copy of the sale banner used to advertise the price promotion can be found in Appendix C.

¹⁵Substitutes were chosen based on product features, such that each substitute pair of products shared the same form (liquid/solid), general taste (sweet/salty) and texture (soft/hard). Please see the Discussion for a longer description of how substitute pairs can be chosen in future studies.

4.4 Theoretical Predictions

In the following section, we present several theoretical predictions made by the models of consumer behavior reviewed earlier. We first present a set of standard economic predictions about behavior during the Sale period, which are the same across all the theories considered. We then present a set of predictions made about the Post-Sale period, which differ depending on the model of consumer behavior used.

Sale period

We make three general predictions about what we can expect to observe in the data during the Sale based on standard economic theory. Table 4.2 presents a summary of these predictions for the Sale period.

Table 4.2: Theoretical predictions - Sale period

Behavior	Prediction
Purchases of discounted items	increase (H1)
Purchases of non-discounted items	increase (H2)
Total unique customers	increase (H3)

First of all, we would expect the price promotion treatment to increase overall purchases of discounted items during the Sale period. This follows directly from the foundational law of demand in that a lower price increases the demand for a product. This prediction will serve as a “sense check” that the vending machine is indeed functioning like a “mini-retailer” and ensure that our Sale treatment worked to stimulate consumer behavior as intended.

Hypothesis 1: The price promotion will increase purchases of discounted items during the Sale period.

Furthermore, we predict that more non-discounted items will be purchased during the Sale period as well. This prediction follows from the loss leader literature (In and Wright, 2014, Chen and Rey, 2012), which states that retailers often sell one or multiple items at a steep discount (often below cost, as we do in our experiment) as a way of encouraging consumers to purchase additional items in the store¹⁶. We

¹⁶Another way we will test this loss leader hypothesis is to see whether the number of items

would expect this loss leader hypothesis to be relevant in the context of the vending machine, which serves as a “mini-retailer,” such that customers may be similarly lured to make a purchase due to a specific discounted item but find themselves adding additional items to their purchase cart than they would if everything was at full price. While this behavior may in part be due to an income effect (i.e. consumers having more money to spend since some of the items are now cheaper), it has also been shown to be a successful strategy employed by retailers for increasing consumer traffic and overall consumer spend (Hosken and Reiffen, 2004, Chevalier, Anil, and Rossi, 2003).

Hypothesis 2: The price promotion will increase purchases of non-discounted items during the Sale period.

We make one additional prediction about the total number of customers during the Sale period. Lowering the price of some of the available products should result in more unique customers “entering” the market and making purchases. This is based on the assumption that willingness to pay (WTP) for any individual product mirrors a standard normal distribution across consumers. This implies that for some subset of consumers, their WTP will be lower than the pre-sale price but higher than the sale price (see Figure 4.3 for a simple illustration). Furthermore, for customers who were already purchasing the product in the pre-sale period (hence had a WTP greater than the pre-sale price), the newly-discounted price becomes even more attractive than the pre-sale price and gives the consumer the perception of gaining greater value for their purchase (Monroe and R. Krishnan, 1985). The ability to test this hypothesis around the total number of unique customers is a unique feature of my dataset since we have unique identifiers for each customer¹⁷.

Hypothesis 3: The total number of unique customers will increase during the Sale period.

Post-Sale period

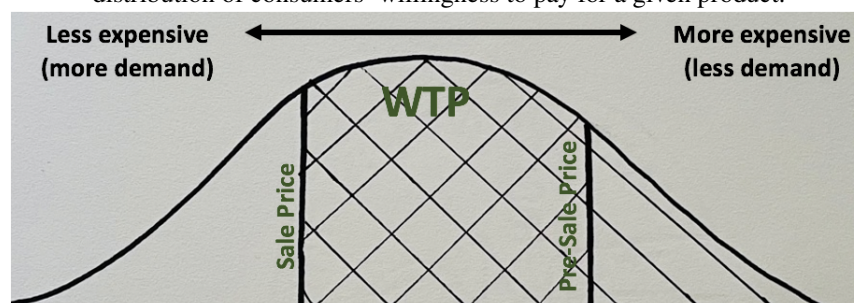
Table 4.3 presents a summary of the predictions made by theories of consumer behavior during the Post-Sale period. We first review predictions made from a Habit

purchased per individual is significantly higher during the Sale period than the Pre-Sale period, as would be anticipated.

¹⁷Previous studies who have looked at the effects of price promotion on consumer behavior without having unique consumer identifiers focus on overall volume of purchases, but cannot say whether that volume is coming from a significantly larger set of unique customers.

Figure 4.3: Illustrative willingness to pay distribution

An illustrative example of the interaction between price and demand, based on a normal distribution of consumers' willingness to pay for a given product.



Perspective (with predictions being directionally similar across Brand Loyalty/Habit Formation and Neuro Autopilot models). We then review predictions made from a Reference-Dependence perspective, based on a Prospect Theory model of consumer behavior.

Table 4.3: Theoretical predictions - Post-Sale period

Behavior	Habit	Reference-Dependence
Purchases of discounted items (compared to pre-sale)	increase (H4a)	decrease (H4b)
Total unique customers	increase (H5a)	decrease (H5b)

Predictions from a Habit Perspective

A person with history dependent preferences develops a “habit” following repeat consumption of a good whereby consuming more in the present will lead to more consumption in the future. Once the price promotion ends (the “Post-Sale” period begins), a history-dependent consumer will continue to consume a good despite the price increase.

The first prediction which is made by the economics specification of habit formation is that we should see an increase in the number of previously discounted items sold. This follows directly from the standard economic prediction that more discounted items will be sold during the sale period, and that a greater number of customers will be making purchases. This means that according to the model of rational addiction, the value of purchasing discounted items in the future will increase.¹⁸

¹⁸The neuro autopilot model makes the same directional prediction that post-sale purchases of discounted items should increase as compared to pre-sale. However, in contrast to the rational addiction formulation, which is independent of sale price and RPE (and based solely on the presence

Hypothesis 4a: Purchases of discounted items will *increase* during the Post-Sale period as compared to Pre-Sale behavior.

We can make a second prediction about the total number of customers, which we anticipate will increase following the sale period. This follows from our prediction that we anticipate previous habits to be strengthened (existing customers maintain or increase purchasing during the sale period) and new habits to be formed (customers who did not purchase previously but whose WTP is now higher than the sale price will start purchasing during the sale).

Hypothesis 5a: Total number of unique customers will *increase* during the Post-Sale period as compared to Pre-Sale behavior.

Post-Sale period: Predictions from a Reference-Dependence Perspective

A person with reference-dependent preferences makes consumption decisions with respect to a reference point, where anything below the reference point is viewed as a loss and anything above it is viewed as a gain. We consider the most recently observed price as the consumer's natural reference point. Hence, changes in price will be evaluated with respect to previous prices (Mazumdar, Raj, and Sinha, 2005, K. and W., 1995, Lattin and Bucklin, 1989).

The first prediction which is made by a model of reference-dependence is that since the Sale period will lead to more purchases of discounted items (as follows from the standard economic predictions), this will result in a decrease in discounted items sold in the Post-Sale period. Specifically, the Sale period will lead to customers updating their reference point to the (lower) discounted price for those items which are on sale (while remaining unchanged for the items which were not discounted). Following the Sale, when the discounted items are back to their "full price," customers will view

of previous choice), the neuro autopilot model predicts that this change in behavior will differ in magnitude depending on RPE. Specifically, behavior will respond more to larger RPE (discounts) and less to smaller RPE (discounts). Unfortunately, our experiment was not designed to tease apart predictions from the economics habit model and the neuro autopilot model, as future studies will be designed to do (see Discussion). For example, a more ideal design would vary price discounts, such that instead of all items being 50% off, some would experience a very small discount and others a large one, ensuring that there is variance in RPE depending on the product purchased. With the ability to measure additional empirical markers (such as speed of choice), we would be able to test whether a larger reward prediction error is more likely to move the consumer into a goal-directed mode, therefore spending more time looking at substitutes and potentially choosing an alternative other than the previously chosen or more discounted item (see Camerer et al. (2021) for two of four scenarios in which a price promotion may lead a habitual consumer to violate laws of demand).

purchasing these items as a “loss” because the price is higher than their reference point. In contrast, for the products which never experienced a discount, customers have the same reference point and do not experience a feeling of purchasing at a loss¹⁹.

Hypothesis 4b: Purchases of discounted items will *decrease* during the Post-Sale period as compared to Pre-Sale behavior.

Additionally, we can make a prediction about the total number of customers, which follows directly from the previous hypothesis in this section. The number of unique customers is likely to decrease in the Post-Sale period as compared to Pre-Sale. This prediction is based on the assumption that more consumers will experience loss aversion in the Post-Sale period in response to discounted items (as described above) and therefore have a lower likelihood of purchasing as compared to the Pre-Sale period.

Hypothesis 5b: The total number of unique customers will *decrease* during the post-sale period as compared to pre-sale behavior.

4.5 Data

In this section, we take a look at the data collected for empirical analysis. We first analyze high-level summary statistics on purchasing behavior, and then dig into a specific example of how purchasing behavior of products in a substitute product pair evolved through the study period.

Summary statistics

During the field experiment, 119 unique individuals made at least one purchase from the vending machine, and a total of 848 item-level transactions were observed. Table 4.4 provides a summary breakdown of purchase data over the time period observed. It is worth noting several outlier values - namely, that sales were largely concentrated in a few products (23% of sales were for Coke) and at least one individual purchased much more frequently than the average customer. Given the already small sample size we work with, we do not exclude these outlier values in our analysis; however, future studies with larger sample sizes could look to run analyses on a dataset which is winsorized at, e.g. 90%, in order to check that the result holds even following the

¹⁹These predictions rely on us observing customer-specific data, which we capture in the empirical analysis using individual-level fixed effects.

removal of extreme values.

Table 4.4: Summary statistics on vending machine purchases

Metric	Minimum	25%	Median	75%	Maximum	Mean
Total purchases of an individual product	2	13	23	36	199	33.9
Total purchases per individual*	1	1	2	5.5	129	6.3
Total unique products purchased per individual*	1	1	2	4	14	2.9

*As identified using a unique credit card number. A single purchase may include one or more products. It is of course possible that one person made purchases for others on their own card (which is our hypothesis for the outlier value). It is also possible that people have multiple credit cards which they used for purchases at the machine (which could explain what appears to be a large number of one-time purchases). Surveying customers (which has been recently approved by our IRB committee) would be helpful to test these assumptions.

Products were randomly assigned to experience a price promotion or not by using a random number generator.²⁰

A motivating example of two substitutes

Prior to looking at empirical results, it is worth considering a motivating example of two frequently-purchased substitute products in our data. Specifically, the two most popular items purchased throughout the 10-week study period were Coke and Izze, both cold fizzy drinks that are viewed as close substitutes. There was no difference in sales of these two items prior to the price promotion period (33 Cokes vs. 35 Izzes were sold). During the Sale period however, Coke was discounted while Izze was not, significantly increasing purchases of Coke (66 Cokes vs. 20 Izzes sold, a significant difference at $p = 0.002$ using a one-sided t-test of mean daily purchases). After the Sale, Coke continued to outsell compared to Izze (100 Cokes vs. 35 Izzes sold, significant at $p = 0.001$ using a one-sided t-test). By the end of the study, there were over twice as many purchases of Coke as of Izze (199 Cokes vs. 90 Izzes). The shift from a seeming indifference between these two substitutes prior to the price promotion to a clear preference for Coke can be seen visually in Figure 4.4.²¹ Unfortunately, this substitute pair is not a representative result (similar tests of other substitute pairs can be found in Appendix F).

²⁰Specifically, products and their substitutes were listed in the same row under columns A and B, respectively, and a random number $x < 0.5$ led to A being discounted whereas a random number $x \geq 0.5$ led to B being discounted.

²¹In Appendix D, we estimate a binary choice model of this decision process between two substitutes.

Figure 4.4: Sales of Coke and Izze through the study period



4.6 Results

In the following section, we present results testing the five hypotheses laid out earlier in the paper. We first show that the standard economics predictions about behavior during the Sale period were all born out in the data. We then take a look at the Post-Sale data.

Sale Increases Purchases of Discounted and Non-Discounted Items

Consistent with H1 and H2, there is an increase in discounted and non-discounted items during the Sale. The Sale period saw an increase in purchases as compared to the Pre-Sale period, across both discounted and non-discounted items, as can be seen in Table 4.5.

Table 4.5: Weekly Transaction Summary by Treatment Period

	Mean	SD	Q1	Median	Q3
Pre-Sale Discounted	27.75	4.92	25.50	26.00	28.25
Pre-Sale Non-Discounted	38.00	11.83	30.50	34.00	41.50
Pre-Sale Total Weekly Purchases	65.75	11.64	56.50	64.50	73.75
Sale Discounted	65.50	14.85	60.25	65.50	70.75
Sale Non-Discounted	54.00	7.07	51.50	54.00	56.50
Sale Total Weekly Purchases	119.50	21.92	111.80	119.50	127.20
Post-Sale Discounted	40.25	16.21	32.00	41.50	49.75
Post-Sale Non-Discounted	46.25	22.98	43.50	56.50	59.25
Post-Sale Total Weekly Purchases	86.50	38.00	75.50	98.5	109.50

Consistent with the loss leader hypothesis, consumers purchase more items per visit during (and after) the Sale. Aggregating individual-level purchases up to bundles, we compare the number of items purchased prior to the sale period with the number of items purchased during and after the sale. The number of items purchased during Sale is higher than Pre-Sale and this remains elevated during the Post-Sale period, increasing overall Post-Sale revenue, as seen in Table 4.6. While the composition of individual bundles between product categories does not shift dramatically (the percentage of the bundle which represents discounted items is slightly higher during the Sale, but comes back to Pre-Sale levels following), the bundles become and remain slightly larger Post-Sale, meaning customers are spending more per bundle on average (\$2.20 vs \$1.95 Pre-Sale). However, none of the differences were statistically significant using a Mann–Whitney U test.²²

²²The Mann–Whitney U test is a nonparametric test of the null hypothesis that, for randomly selected values from two experimental periods, the probability of one value being greater than the other value is equal to the probability of the latter value being greater than the former. Using this test, we find no significant differences between the number of items per bundle nor the percent of the bundle which is discounted items for each pair of experimental periods. Furthermore, we do find a significant difference between the dollar amount spent per bundle between Pre-Sale and Sale ($p < 0.01\%$) and Sale and Post-Sale ($p < 0.01\%$), but not between Pre-Sale and Post-Sale periods.

Table 4.6: Bundle-Level Summary by Treatment Period

	Mean	SD	Q1	Median	Q3
Pre-Sale items per bundle	1.49	0.93	1.00	1.00	2.00
Pre-Sale dollar amount spent per bundle	1.91	1.26	1.00	1.50	2.13
Pre-Sale percent of bundle which is discounted items	0.41	0.46	0.00	0.00	1.00
Sale items per bundle	1.93	1.69	1.00	1.00	2.00
Sale dollar amount spent per bundle	1.67	1.51	0.75	1.00	2.00
Sale percent of bundle which is discounted items	0.48	0.47	0.00	0.50	1.00
Post-Sale items per bundle	1.82	1.73	1.00	1.00	2.00
Post-Sale dollar amount spent per bundle	2.20	1.99	1.00	1.50	2.00
Post-Sale percent of bundle which is discounted items	0.39	0.44	0.00	0.00	1.00

Post-Sale Purchases of Discounted Items Increase Compared to Pre-Sale

Consistent with H4a, purchases of discounted items *increase* during the Post-Sale period. As can be seen in Table 4.5, the Post-Sale period saw a 48% increase in purchases of previously discounted items as compared to the Pre-Sale period, from a customer base which was comparable to baseline.

With respect to purchases of non-discounted substitute items, the sales here increase during the Post-Sale period, but decrease on a relative basis compared to discounted items. As can be seen in Table 4.5, the Post-Sale period saw an increase in purchases of items that were not discounted during the Sale period, albeit at a lower rate (a 24% increase).

We complement all of these descriptive statistics with regression results estimating model 4.1. In our purchase regression specification, y_t is the total number of purchases made at time t . We regress $duringsale_t$ and $postsale_t$, which are binary indicators of whether the purchase was made during the Sale or Post-Sale, respectively, and $week_t$, which is the number of the week of the experiment (1-10) and is intended to capture any time trends influencing the results. ϵ_t is our error term, which we assume is uncorrelated with the regressors. These regressions collapse all purchases at the individual day level. Specifically, we calculate the number of purchases made in each day over the observed 10-week time period. We then regress the number of purchases on indicators of whether it was a Sale period or Post-Sale period at the time of purchase.

$$y_t = \alpha + \beta_1 duringsale_t + \beta_2 postsale_t + \gamma(week_t)^2 + \epsilon_t \quad (4.1)$$

The regression results can be found in 4.7. We can see from this regression that purchases significantly increase during the Sale period, as compared to both Pre- and Post-Sale. When both time periods enter into our regression, Post-Sale purchases remain elevated compared to Pre-Sale as well.

Table 4.7: Purchases During Sale and Post-Sale Periods

	<i>Dependent variable:</i>		
	Total Purchases		
	(1)	(2)	(3)
DuringSale	4.614*** (1.524)		7.115*** (1.999)
PostSale		-1.716 (2.920)	6.775* (3.595)
week ²	0.043** (0.021)	0.061 (0.047)	-0.048 (0.053)
Constant	8.374*** (1.054)	9.403*** (1.100)	8.576*** (1.039)
Observations	68	68	68
R ²	0.159	0.045	0.203
Adjusted R ²	0.133	0.016	0.166
F Statistic	6.136*** (df = 2; 65)	1.543 (df = 2; 65)	5.436*** (df = 3; 64)

Note: *p<0.1; **p<0.05; ***p<0.01, two-tailed test of hypothesis

We run this regression including interaction terms for discounted items ("SaleItem" in our regression) to see whether the Sale had a disproportionate impact on discounted items versus non-discounted items. We do not find such an effect, and the results from this regression can be found in Table 4.8.²³

Total Customers Increases with Sale and Remains Unchanged Post-Sale

Consistent with H3, the total number of customers *increases* during the Sale period. As can be seen in Table 4.9, the Sale period saw a 42% increase in total unique customers as compared to Pre-Sale. The number of unique customers purchasing discounted items went up even more significantly, at a 73% increase. This finding

²³Difference-in-difference regressions, another way to measure the impact of the Sale treatment on future purchases of discounted items, also found no significant effect and can be found in Appendix E.

Table 4.8: Purchases During Sale and Post-Sale Periods

	<i>Dependent variable:</i>		
	Total Purchases		
	(1)	(2)	(3)
DuringSale	1.952* (1.012)		2.298* (1.188)
PostSale		-1.925 (1.476)	1.042 (1.716)
SaleItem	-1.192* (0.675)	-0.682 (0.770)	-1.241 (0.909)
week ²	0.021** (0.010)	0.046** (0.021)	0.007 (0.023)
DuringSale:SaleItem	1.592 (1.427)		1.641 (1.558)
PostSale:SaleItem		-0.449 (1.314)	0.111 (1.367)
Constant	4.958*** (0.604)	5.161*** (0.645)	5.035*** (0.683)
Observations	134	134	134
R ²	0.139	0.053	0.143
Adjusted R ²	0.113	0.024	0.102
F Statistic	5.217*** (df = 4; 129)	1.819 (df = 4; 129)	3.520*** (df = 6; 127)

Note: *p<0.1; **p<0.05; ***p<0.01, two-tailed test of hypothesis

is in line with our standard economic prediction based on the law of supply and demand.

Table 4.9: Weekly Customer Summary by Treatment Period

	Mean	SD	Q1	Median	Q3
Pre-Sale Discounted	15.00	3.65	12.50	15.00	17.50
Pre-Sale Non-Discounted	19.00	2.94	16.75	19.00	21.25
Pre-Sale Total Weekly Purchases	34.00	4.24	32.25	33.00	34.75
Sale Discounted	21.00	2.83	20.00	21.00	22.00
Sale Non-Discounted	20.00	1.41	19.50	20.00	20.50
Sale Total Weekly Purchases	41.00	4.24	39.50	41.00	42.50
Post-Sale Discounted	14.75	5.32	13.75	16.50	17.50
Post-Sale Non-Discounted	15.75	6.18	13.75	17.50	19.50
Post-Sale Total Weekly Purchases	30.50	11.09	29.75	35.00	35.75

Inconclusively with H5a or H5b, the number of unique customers remains *unchanged* in the Post-Sale period, as can be seen in Table 4.9. This finding is in contrast to both the prediction made by the brand loyalty/habit formation models (that total number of customers would increase post-sale) as well as the prediction made by the reference-dependence model (that total number of customers would decrease post-sale).

We complement these descriptive statistics with regression results estimating model 4.1. In this specification, y_t is now the total number of unique customers making a purchase at time t . As before, $during_{sale}_t$ and $post_{sale}_t$ are binary indicators of whether the purchase was made during the Sale or Post-Sale, respectively, and $week_t$ is the number of the week of the experiment (1-10) and is intended to capture any time trends influencing the results. ϵ_t is our error term, which we assume is uncorrelated with the regressors.

As observed in the Descriptive Analyses earlier, we find evidence in support of H3 - that total number of unique customers is significantly higher during the Sale period ($\beta = 1.292$, $p < 0.1$). We further find that the total number of customers remains unchanged in the Post-Sale period (does not drop or increase significantly, compared to the Pre-Sale period, in contrast to our two hypotheses).

Table 4.10: Customer Composition During Sale and Post-Sale Periods

	<i>Dependent variable:</i>		
	Unique Customers		
	(1)	(2)	(3)
DuringSale	1.292* (0.667)		2.767*** (0.855)
PostSale		0.695 (1.231)	3.998** (1.538)
week ²	-0.002 (0.009)	-0.014 (0.020)	-0.056** (0.023)
Constant	5.301*** (0.461)	5.741*** (0.464)	5.420*** (0.444)
Observations	68	68	68
R ²	0.057	0.007	0.147
Adjusted R ²	0.028	-0.023	0.107
F Statistic	1.962 (df = 2; 65)	0.244 (df = 2; 65)	3.677** (df = 3; 64)

Note: *p<0.1; **p<0.05; ***p<0.01, two-tailed test of hypothesis

4.7 Discussion

In this paper, I ran a field experiment using a novel methodology in order to answer the question of what medium-term effects price promotions have on consumer behavior. I began data collection in January 2020, and had to stop in early March when the COVID-19 pandemic was spreading across the globe. While this single study is more limited than the data I anticipated being able to draw conclusions from (i.e. running multiple studies over the course of the following year), we can still learn a few valuable takeaways from these early results.

Importantly, we see that this novel methodology – a customizable vending machine which collects data on consumer behaviors 24/7 – successfully serves as a “mini-retailer” for students. Consumer behavior during this pilot study was in line with what we would predicted it should be based on laws of supply and demand. In line with these general predictions, running a price promotion on the vending machine increased customer engagement by increasing the number of unique individuals making purchases in the machine (as determined by credit card data which gives us a unique identifier for each consumer) as well as increasing overall purchases. In

line with the loss leader literature in marketing, which states that retailers may sell a small number of items below cost in order to “lure” customers in to the store and purchase other items in addition to the heavily discounted one(s), we find that the sale of both discounted and non-discounted items increases during a price promotion (although this increase is not significant with our sample size).

Additionally, we set the foundation for credibly dissociating predictions made by the brand loyalty/habit formation literature and the reference-dependence literature (which make opposite sets of predictions about the impact that a price promotion should have on consumer behavior after the sale). We predict that given more data we would see that, in line with the habit formation and neuro autopilot theories, a price promotion leads to a subsequent increase in the purchases of discounted items following the price promotion, indicating that consumers who discovered or purchased more of the discounted items during the sale period may have become habitized to those items. Such a result would provide support in favor of retailers running occasional price promotions if the strategic goal is to increase brand loyalty to a specific firm or product (and they continue to be profitable during the promotion itself).

What Could be Improved

This study was intended to serve as a pilot, giving us the opportunity to learn how to operate and collect data using the vending machine. For this reason, we choose a relatively simple experimental design: implementing a price promotion of a standard magnitude across a subset of products for a two-week time period in order to test which one of several classic economic and marketing models does a better job explaining the consumer behavior observed.

As with any study, particularly one which uses a completely novel methodology, a lot could have been done better and should be noted as it can be improved upon in future research projects.

First of all, and perhaps more importantly, we know more about the purchase ecosystem than is fully used in this study. For example, we have the exact time stamps of purchases such that we can create a measure of “interaction likelihood” - students who purchase at the same time at least twice are likely friends or co-habitants (as a reminder, the vending machine is located in a dormitory). This means that we can test (and/or control) for peer effects which may lead to information spillover. For example, if one of the friends purchases something during the sale and tells her

friend about this purchase, this would theoretically influence the friend's reference point for that item, regardless of purchase. These spillovers are not unlikely - students (and friends, broadly) are known to eagerly share news of a good deal with one another - and hence it would be interesting to test the influence they might have on consumer behavior in different circumstances. If we are able to link credit card data with unique student identifiers, we could obtain further information about student profiles to use as controls in our analyses.

Secondly, the pricing decisions could have been made more scientifically. As mentioned in the main text, we used cost-based pricing to decide how much to charge for each product, first calculating the per-unit cost of the snack being sold and then choosing a price above this cost which met two criteria. The two criteria we used were that the price (1) was a multiple of \$0.50 (this was used to reduce cognitive load for purchasers and make total costs, if they were purchasing multiple items, easy to calculate) and (2) targeted a retail mark-up between 1-2. The combination of these two criteria meant that some items were sold just above cost (e.g. a snack which cost \$0.39 might be priced at \$0.50) and others were sold at a significant premium (e.g. a snack which cost \$0.52 would be priced at \$1.00). This is of course the "right" strategy for a real-world retailer who intends to make a profit while reducing cognitive load for prospective consumers. However, in retrospect we should have probably kept a more controlled experimental design where all items were sold at the same uniform mark-up, e.g. two times cost, even if it resulted in "unusual" prices, e.g. \$0.78 and \$1.04.

Finally, a third area of improvement is how the substitute products were chosen. In this study, I used a survey method to collect data on product attributes such as health and taste perceptions and then choose products informed by those survey results. In the future, I would collect a range of objective product attributes (e.g. calories, sugar content) and do clustering analysis on these products. This may lead to a more "scientific" set of substitute pairs (i.e. the idea being that similar products end up in the same cluster and can then be treated as substitutes, with one being randomly assigned to the treatment and other(s) to the control). This clustering work would allow me to further justify my choice of substitutes and strengthen the empirical case that differences in their purchase levels are indeed driven by the price promotion rather than any innate differences between the products. The clustering analysis would of course rely on having enough underlying variables which describe the relevant products.

Future Directions

I hope that this paper can serve as a proof of concept for a novel methodology, as well as motivation for more research in the area of consumer behavior. With a customizable vending machine that can be controlled remotely, it becomes possible to run lab-like experiments in the field. In other words, by having more control over other attributes which tend to be exogenous considerations in empirical data (such as stockpiling, individual identifiers, price promotion details), we can test the predictions made by specific decision theories of consumer behavior (just as one might do in an experimental lab setting) in the field under ecologically valid conditions.

Having shown that this is possible, this novel technique opens up new research paths with respect to better understanding the impact that attention, social preferences, and other variables have on consumer choices.

As an example, manipulating the order in which items are presented on the screen can allow for testing the role that attention plays in consumer choice. Specifically, a maximum of 12 items can be displayed on the vending machine touchscreen at any given time before consumers are forced to “swipe” to see additional items available. The experimenter controls which items are on which screen, and where they are positioned. A utility maximization model of consumer choice is agnostic to something like where items are placed on a screen - specifically, this should not matter for a utility-maximizing agent who would presumably consider all items and choose the optimal one. However, more psychologically-informed models of consumer choice may recognize that the probability of choosing items on the first screen is higher (“weighted more”) because of the costly cognitive and physical effort of scrolling through multiple screens before choosing an item.

Furthermore, the vending machine can be used to understand whether brand loyalty (or a “self-aware history dependence,” as modelled by economic theories of habit formation) or neuro autopilot are better models of human behavior. For example, imagine a consumer which regularly purchases an item, say Coke, which is always in a static position on the screen, say top left. They arrive at the machine on a regular basis, look to the top left corner, click on the Coke icon, and submit their payment in order for the item to dispense. Now let’s imagine changing the order in which items are displayed, such that the top left corner now has a new soda - Dr. Pepper - and Coke is hidden in the bottom right corner where the consumer never looks. Brand loyalty theory predicts that this location change should have no effect on behavior -

the consumer would do a visual search to find their Coke and purchase the Coke as usual. However, a model of neuro autopilot would predict that the consumer should experience a reward prediction error, or RPE, which “jolts” them out of habit mode, such that they may not decide to choose Coke. If they do choose Coke, the neuro autopilot model also predicts that it will take longer (a previous “fast” decision made in habit mode becomes a “slow” utility maximization decision).

The vending machine could also be used to run more studies testing the impact of prices on behavior. If the rational addiction theory from economics is an accurate model of behavior, then consumer expectations about future prices should influence current consumption. This is something which can be tested in our setting by manipulating expectations regarding future prices (i.e. advertising that a price change is upcoming) and testing whether these expectations change their purchasing behavior. In addition to and related to the question of prices, the machine can be used to study dynamic framing effects, as well as stock-outs. For example, if a consumer develops a predictable buying habit for a specific item, does removing that item from the choice set (an “artificial stock out”) lead to a new habit being formed? These experiments are straight-forward to set up using this novel methodology and present a clear way to test hypotheses from different models.

Appendix A - Vending Machine Details


Below is a summary of the model of vending machine purchased (“Solution Option 4”) which was taken from the Digital Media Vending sales presentation. As can be seen in the summary, this was a high-tech vending machine with a touchscreen interface, LED lighting and infrared sensors to track product movement.

www.digitalmediavending.com

SOLUTION OPTION 4

10 Column – up to 100 Selections

- 26" Touchscreen refrigerated elevator vending machine
- Conveyor belts (pre-configured for your products)
- Up to 10 Shelves, up to 10 products per shelf giving a maximum of 100 product selections
- Storage space available in touchscreen cabinet
- LED Lighting
- 20 infrared sensors track product movement for guaranteed product delivery
- Automatic delivery bin with locking mechanism
- Tray dimensions – 30" wide, 21" deep
- H73" x W53" x D38" (optional 31.5" depth)
- Weight 840lbs
- Unit pricing includes vending machine and screen & standard UI
 - 1 = \$13,495
 - 5 = \$12,795
 - 10 = \$11,995



On the following page is an extract from our contract with Digital Media Vending, which describes the specifications of the machine we eventually ordered and purchased. Custom aspects of our build included temperature control (refrigeration), the incorporation of a cashless device (using monitoring service Nayax), cloud-based vending management subscription, and adjustable product sequencing (the standard machine displays products alphabetically, and we requested the ability to set product order and randomize products).



SCHEDULE 1 (Exhibit B)

Floor standing touchscreen vending machine to dispense food & drinks (customizable trays)

- 50" Touchscreen elevator vending machine - your products do not fall
- Payter Cashless payment terminal accepts Credit, Debit and Apple Pay, Google Wallet etc (payter.com)
- Sensors track product movement for guaranteed product delivery
- Automatic delivery door with locking mechanism
- Up To 6 Shelves 30" wide and 13 ¾" deep
- Up to 60 selections, dependent on product sizes. For example, if all products are the size of a standard cigarette pack size, you would have 60 selections. If some products are much larger, then the number of selections reduces accordingly.
- Machine dimensions H76" x W41 ½" x D32"
- Weight 728lbs
- 110VAC power input
- 12 to 24VDC internal operating power
- Fully adjustable product trays with conveyor and elevator dispense system
- Vending management cloud-based system for remote audit and management

Appendix B - Pre-Testing Survey

Survey structure

On the following page are extracts of the pre-test ratings survey sent out to undergraduate students. These extracts display the wording of questions presented for each snack item displayed (only one snack item, Cheetos, is included here). Students were asked to rate their willingness to purchase the snack item at the given price (not always the price eventually used for the experiment), the taste of the item, and the healthiness of the item, each on a 7-point scale from "not at all" to "very".

Vending Machine Survey

We're interested in your opinions of foods that could be stocked in a vending machine on Caltech campus.

*Required

Oven Baked Flamin' Hot Cheetos



How willing would you be to purchase this product from a vending machine at \$0.75?

	1	2	3	4	5	6	7	
Not at all willing	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Very willing

How would you rank the taste of this item? *

	1	2	3	4	5	6	7	
Not at all tasty	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Very tasty

How healthy would you rate this item to be? *

	1	2	3	4	5	6	7	
Not at all healthy	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Very healthy

Survey results

The following page presents the aggregated ratings for each of the products in the pre-testing survey. Not all of these products were eventually purchased for the machine, and additional products not included in the pre-testing survey were eventually added to the machine as well.

Table 4.11: Pre-testing survey results

Product Name	Price (\$)	Purchase Likelihood	Health Rating	Taste Rating
Blue Diamond Almonds	0.50	3.6	5.6	4.0
Oven Baked Flamin' Hot Cheetos	0.75	3.8	2.3	4.4
Bare Apple Chips	0.75	4.3	5.6	4.5
Pirate's Booty, Aged White Cheddar	0.25	4.9	2.8	4.9
Baked Organic Crunchy Pea Snack	1.00	2.9	5.2	3.6
Grandma's Mini Sandwich Cremes	0.25	4.7	2.2	4.9
Stacy's Pita Chips, Cinnamon Sugar	0.50	4.4	3.2	4.6
Monster Rehab	1.25	2.9	1.6	3.1
Stacy's Pita Chips, Parmesan Garlic and Herb	1.25	3.9	3.9	5.0
Sour Patch Kids	0.50	4.8	1.8	5.6
Justin's Peanut Butter	0.50	3.5	4.6	4.5
Garden Veggie Straws, Zesty Ranch	0.75	3.9	4.4	4.5
BBQ Pop Chips	0.75	3.9	3.3	4.1
Pocky, Chocolate	2.00	3.6	2.8	5.9
Lipton Diet White Tea Raspberry	0.25	3.8	4.0	3.7
Famous Amos, Chocolate Chip	0.25	5.0	2.1	5.2
Extra Gum, Polar Ice	0.50	4.0	4.0	4.5
think! Creamy Peanut Butter High Protein Bar	1.50	3.4	4.9	4.1
RX Bar, Chocolate Sea Salt	1.50	3.8	5.0	4.3
Lenny & Larry's Complete Cookie, Chocolate Chip	1.25	3.4	3.2	4.1
Calpico	2.00	3.2	3.6	4.3
Premier Nutrition High Protein Shake, Vanilla	1.75	3.0	4.5	3.4
Brothers Strawberry Fruit Crisps	0.75	4.5	5.1	4.9
Gatorade, Strawberry Watermelon	0.50	4.7	4.1	4.8
Skinny Pop, White Cheddar	0.25	5.1	4.2	4.9
Flamin' Hot Cheetos	0.25	4.5	2.3	4.6
Cheezit	0.25	5.3	3.2	5.2
Organic Seaweed Snack	0.50	4.6	5.1	4.9
Ito En Green Tea	1.50	3.9	5.1	4.6
Goldfish Cheddar	0.25	5.2	3.5	5.3
Swedish Fish	0.75	4.1	2.1	5.0
Whole Grain Goldfish, Cheddar	0.25	4.4	3.9	4.4
Hi-Chew, Mango	1.50	3.7	2.8	5.1
Milk Chocolate M&M's	1.00	4.0	2.0	5.3
Black Forest Organic Gummy Bears	0.25	4.6	2.9	4.9
Starbucks Double Shot, Espresso Salted Caramel	1.50	3.8	2.9	4.2
Teriyaki Beef Jerky	0.75	3.8	3.7	4.2
Corn Nuts, Ranch	0.75	3.2	3.4	3.5
LaCroix Berry	0.75	3.3	4.1	3.4
Starbucks Double Shot, Espresso + Cream Light	1.50	3.8	3.1	3.9
Overall	0.83	4.0	3.6	4.5

Appendix C - Vending Machine Announcements

Vending machine introduction email

The following email was sent out to students following the installation of the vending machine. This email, announcing that the vending machine was now running, came from the Director of Student Housing. The caltech.vending@gmail.com address referenced was created and managed by the author.

Dear students,

The vending machine in the south houses laundry room is up and running. It is a touchscreen machine which offers snacks and drinks, and accepts only cashless payments. Feel free to email caltech.vending@gmail.com by January 10th if you would like to see any of your favorite snacks or drinks added to the machine.

Sincerely, [Director of Student Housing]

Sale banner

Below is the sale banner which was created and displayed at the top of the machine during the price promotion period.



Discounts of 50% off on
your favorite items!

Scroll through to choose snacks and drinks for purchase

Appendix D - Binary Choice Model

We can estimate a restricted model of the choice process of deciding between two substitute goods. Specifically in this case, conditional on having decided to

consume a sweet fizzy drink, how likely is the consumer will choose Coke (over Izze) depending on the consumer's historical behavior? We estimate the following fixed effects logit model:

$$P(y_{it} = 1) = \beta_0 + \beta_1(X(c)_{it}) + \beta_2(X(z)_{it}) + \alpha_i + \epsilon_{it} \quad (4.2)$$

Where $y_{it} = 1$ captures consumer i purchasing Coke at time t . X_{it} is a rolling variable which captures the number of times a drink was previously purchased by this consumer, where $X(c)$ is the number of times a Coke was previously purchased, and $X(z)$ is the number of times an Izze was previously purchased. α_i captures individual-level fixed effects.

In order to run this model, we "fill in" the data to include days when an individual did not purchase items as well as days when they purchased something. We then calculate the variables above for each person-date combination in the dataset. From there, we estimate the model, and the results of Equation (4.2) are presented in Table 4.12.

As can be seen from this table, the estimation presents additional evidence in favor of Brand Loyalty. The greater the number of Cokes an individual has consumed in the past, the more likely they are to choose Coke over Izze ($\beta = 0.242$, $p < 0.01$).

Table 4.12: Binary Logit Regression Results

	<i>Dependent variable:</i>	
	$y_{it} = 1$	
	(1)	(2)
$X(c)_{it}$	0.045*** (0.005)	0.242*** (0.044)
$X(z)_{it}$		-0.416*** (0.094)
Constant	-4.935*** (0.133)	-4.993*** (0.137)
Observations	7,973	7,973
Log Likelihood	-360.184	-351.937

Note: *p<0.1; **p<0.05; ***p<0.01

Appendix E - Difference in Difference Regression

Table 4.13: DID Regressions - Purchases Post-Sale

<i>Dependent variable:</i>	
Total Purchases	
SaleItem	-0.682 (0.781)
PostSale	0.633 (0.943)
SaleItem:PostSale	-0.449 (1.334)
Constant	5.932*** (0.553)
Observations	134
R ²	0.017
Adjusted R ²	-0.006
F Statistic	0.744 (df = 3; 130)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Appendix F - Substitute Pair Analysis

We used t-tests to compare the number of daily purchases between two goods for the top-selling substitute pairs across the three treatment periods. Table 4.14 shows the results of two-sided tests for the Pre-Sale period, and one-sided tests for Sale and Post-Sale periods.

The first substitute pair, Coke and Izzie were discussed in the Data section as a motivating example of how the sale period led to a strong preference for Coke over Izzie. However not all of the substitute pairs were similarly matched. For example, customers preferred Fruit Snacks to Sour Patch candy in all three periods. While the difference during the Pre-Sale and Sale periods was not statistically significant, the Sale and Post-Sale periods saw a clear preference for Fruit Snacks emerge despite Sour Patch being discounted in the previous period ($p = 0.012^{**}$ and $p = 0.001^{***}$, respectively). In other pairs the discounted item was preferred - specifically, Starbucks and Pringles were the preferred option before the Sale treatment began. This implies that we need to do a better job randomizing products in future studies, such that customers are truly indifferent between substitutes before one is discounted.

In other cases, substitutes were well-matched (in that customers seemed truly indifferent in the Pre-Sale period) but the price promotion did not appear to lead to habit formation. For example, customers were indifferent between Complete Cookie and Famous Amos cookies Pre-Sale, and while daily sales of Complete Cookie did increase when it was discounted, the difference was not statistically significant during the Sale nor the Post-Sale periods.

Table 4.14: T-Tests of Mean Daily Purchases in Substitute Pairs

Item pair X (discounted) and Y	Pre-Sale (X=Y)	Sale (X>Y)	Post-Sale (X>Y)
Coke and Izzie	$p = 0.876$	$p = 0.002^{***}$	$p = 0.001^{***}$
Sour Patch and Fruit Snacks	$p = 0.115$	$p = 0.988$	$p = 0.999$
Complete Cookie and Famous Amos	$p = 0.788$	$p = 0.140$	$p = 0.419$
Starbucks and Monster	$p = 0.018^{**}$	$p = 0.155$	$p = 0.042^{**}$
Goldfish and Lays	$p = 0.132$	$p = 0.758$	$p = 0.972$
Pringles and SkinnyPop	$p = 0.002^{***}$	$p = 0.006^{***}$	$p = 0.001^{***}$

References

- Abeler, J. et al. (2011). “Reference points and effort provision”. In: *American Economic Review* 101.2, pp. 470–92.
- Ailawadi, K. L., K. Gedenk, et al. (2007). “Decomposition of the sales impact of promotion-induced stockpiling”. In: *Journal of Marketing Research* 44.3, pp. 450–467.
- Ailawadi, K. L., D. Lehmann, and Scott A. Neslin (2003). “Revenue premium as an outcome measure of brand equity”. In: *Journal of Marketing* 67.4, pp. 1–17.
- Allen, E. and P. Dechow (2020). “Gender differences in the strategies used for task completion: An analysis of marathon runners”. In: *Available at SSRN: <https://ssrn.com/abstract=3536645>*.
- Anderson, E. T., G. J. Fitzsimons, and D. Simester (2006). “Measuring and mitigating the costs of stockouts”. In: *Management Science* 52.11, pp. 1751–1763.
- Anderson, E. T. and D. Simester (2004). “Long-run effects of promotion depth on new versus established customers: Three field studies”. In: *Marketing Science* 23.1, pp. 4–20.
- Bar-Isaac, H., G. Caruana, and V. Cuñat (2012). “Search, design, and market structure”. In: *American Economic Review* 102.2, pp. 1140–60.
- Batra, R., A. Ahuvia, and R. P. Bagozzi (2012). “Brand love”. In: *Journal of Marketing* 76.2, pp. 1–16.
- Baucells, M., M. Weber, and F. Welfens (2011). “Reference-point formation and updating”. In: *Management Science* 57.3, pp. 506–519.
- Becker, G. S. and K. M. Murphy (1988). “A Theory of Rational Addiction”. In: *Journal of Political Economy* 96.4, pp. 675–700.
- Camerer, C. F., A. Dreber, et al. (2016). “Evaluating replicability of laboratory experiments in economics”. In: *Science* 351.6280, pp. 1433–1436.
- Camerer, C. F., P. Landry, and R. Webb (2021). “The neuroeconomics of habit”. In: *The State of Mind in Economics*. Ed. by A. Kirman and M. Teschi.
- Chan, T., C. Narasimhan, and Q. Zhang (2008). “Decomposing purchase elasticity with a dynamic structural model of flexible consumption”. In: *Journal of Marketing Research* 45.4, pp. 487–498.
- Chandon, P., B. Wansink, and G. Laurent (2000). “A benefit congruency framework of sales promotion effectiveness”. In: *Journal of Marketing* 64.4, pp. 65–81.
- Chen, Z. and P. Rey (2012). “Loss Leading as an Exploitative Practice”. In: *American Economic Review* 102.7, pp. 3462–82.
- Chevalier, J. A., K. K. Anil, and P. E. Rossi (2003). “Why don’t prices rise during periods of peak demand? Evidence from scanner data”. In: *American Economic Review* 93.1, pp. 15–37.

- Ching, A. T. and M. Osborne (2020). “Identification and estimation of forward-looking behavior: The case of consumer stockpiling”. In: *Marketing Science* 39.4, pp. 669–848.
- Coelho, P. S., P. Rita, and Z. Raposo Santos (2018). “On the relationship between consumer-brand identification, brand community, and brand loyalty”. In: *Journal of Retailing and Consumer Services* 43, pp. 101–110.
- DelVecchio, D., D. H. Henard, and T. H. Freling (2006). “The effect of sales promotion on post-promotion brand preference: A meta-analysis”. In: *Journal of Retailing* 82.3, pp. 203–213.
- DelVecchio, D., H. S. Krishnan, and D. C. Smith (2007). “Cents or percent? The effects of promotion framing on price expectations and choice”. In: *Journal of Marketing* 71.3, pp. 158–170.
- Dick, A. and K. Basu (1994). “Customer loyalty: Toward an integrated conceptual framework”. In: *JAMS* 22, pp. 99–113.
- Dubé, J.-P., G. J. Hitsch, and P. E. Rossi (2009). “Do switching costs make markets less competitive?” In: *Journal of Marketing Research* 46.4, pp. 435–445.
- Eales, T. (2016). *Price and Promotion in Western Economies*. Tech. rep. IRI.
- Erdem, T., S. Imai, and M. P. Keane (2003). “Brand and quantity choice dynamics under price uncertainty”. In: *Quantitative Marketing and Economics* 1, pp. 5–64.
- Fishbach, A., R. K. Ratner, and Y. Zhang (2011). “Inherently loyal or easily bored? Nonconscious activation of consistency versus variety-seeking behavior.” In: *Journal of Consumer Psychology* 21, pp. 38–48.
- Gupta, S. (1988). “Impact of sales promotion on when, what, and how much to buy”. In: *Journal of Marketing Research* 25, pp. 342–355.
- Hanseens, D.M. (2018). *Long-Term Impact of Marketing*. World Scientific.
- He, H., Y. Li, and L. Harris (2012). “Social identity perspective on brand loyalty”. In: *Journal of Business Research* 65.5, pp. 648–657.
- Helsen, K. and D.C. Schmittlein (1992). “Some characterizations of stockpiling behavior under uncertainty”. In: *Marketing Letters* 3, pp. 5–16.
- Hendel, I. and A. Nevo (2003). “The post-promotion dip puzzle: What do the data have to say?” In: *Quantitative Marketing and Economics* 1, pp. 409–424.
- (2006). “Measuring the implications of sales and consumer inventory behavior”. In: *Econometrica* 74, pp. 1637–1673.
- Hosken, D. and D. Reiffen (2004). “How retailers determine which products should go on sale: Evidence from store-level data”. In: *Journal of Consumer Policy* 27, pp. 141–177.
- In, Y. and J. Wright (2014). “Loss-leader pricing and upgrades”. In: *Economics Letters* 122.1, pp. 19–22.

- K., Gurumurthy and Russell S. W. (1995). "Empirical Generalizations from reference price research". In: *Marketing Science* 14.3 Supplement, G161–G169.
- Kahneman, D., J. Knetsch, and R. Thaler (1990). "Experimental tests of the endowment effect and the Coase Theorem". In: *Journal of Political Economy* 98.6, pp. 1325–1348.
- Keller, K. L. and D. Lehmann (2006). "Brands and branding: Research findings and future priorities". In: *Marketing Science* 25.6, pp. 551–765.
- Kőszegi, B. and M. Rabin (2007). "Reference-dependent risk attitudes". In: *American Economic Review* 97.4, pp. 1047–1073.
- Lattin, J. M. and R. E. Bucklin (1989). "Reference effects of price and promotion on brand choice behavior". In: *Journal of Marketing Research* 26.3, pp. 299–310.
- Lee, S., S. Shimojo, and J. O'Doherty (2014). "Neural computations underlying arbitration between model-based and model-free learning". In: *Neuron* 81.3, pp. 687–699.
- Marom, O. and A. Seidmann (2011). "Using "last-minute" sales for vertical differentiation on the Internet". In: *Decision Support Systems* 51.4, pp. 894–903.
- Marr, D. (1982). *Vision: A computational investigation into the human representation and processing of visual information*. Cambridge, Mass: MIT Press.
- Mazumdar, T., S. P. Raj, and I. Sinha (2005). "Reference price research: Review and propositions". In: *Journal of Marketing* 69.4, pp. 84–102.
- Mela, C., S. Gupta, and D. Lehmann (1997). "The long-term impact of promotion and advertising on consumer brand choice". In: *Journal of Marketing Research* 34.2, pp. 248–261.
- Mela, C., K. Jedidi, and D. Bowman (1998). "The long-term impact of promotions on consumer stockpiling behavior". In: *Journal of Marketing Research* 35.2, pp. 250–262.
- Monroe, K. B. and R. Krishnan (1985). "The effect of price on subjective product evaluation". In: *The perception of merchandise and store quality*. Ed. by J. Jacoby and J. Olson. Lexington: Lexington Book, pp. 209–232.
- O'Brien, E. (2021). "A mind stretched: The psychology of repeat consumption". In: *Consumer Psychology Review* 4, pp. 42–58.
- Oliver, R. (1999). "Whence consumer loyalty?" In: *Journal of Marketing* 63, pp. 33–44.
- Osborne, M. (2018). "Frequency versus depth: How changing the temporal process of promotions impacts demand for a storable good". In: *Japanese Economic Review* 69, pp. 258–283.
- Pauwels, K., D. M. Hanssens, and S. Siddarth (2002). "The long-term effects of price promotions on category incidence, brand choice, and purchase quantity". In: *Journal of Marketing Research* 39.4, pp. 421–439.

- Pesendorfer, M. (2002). "Retail sales: A study of pricing behavior in supermarkets". In: *Journal of Business* 75.1, pp. 33–66.
- Putler, D. S. (1992). "Incorporating reference price effects into a theory of consumer choice". In: *Marketing Science* 11.3, pp. 287–309.
- Ray, D., M. Shum, and C. F. Camerer (2015). "Loss aversion in post-sale purchases of consumer products and their substitutes". In: *American Economic Review* 105.5, pp. 376–80.
- Romaniuk, J. and M. Nenycz-Thiel (2013). "Behavioral brand loyalty and consumer brand associations". In: *Journal of Business Research* 66.1, pp. 67–72.
- Rust, J. (1994). "Structural estimation of Markov decision processes". In: *Handbook of Econometrics*. Ed. by Engle R. and McFadden D. Elsevier.
- Seiler, S. (2013). "The impact of search costs on consumer behavior: A dynamic approach". In: *Quantitative Marketing and Economics* 11.2, pp. 155–203.
- Slonim, R. and E. Garbarino (1999). "The effect of price history on demand as mediated by perceived price expensiveness". In: *Journal of Business Research* 45, pp. 1–14.
- Sun, B. (2005). "Promotion effect on endogenous consumption". In: *Marketing Science* 24.3, pp. 430–443.
- Tam, L., W. Wood, and M. F. Ji (2009). "Brand loyalty is not habitual". In: *Handbook of brand relationships*. Ed. by D. J. MacInnis, C. W. Park, and J. R. Priester. M E Sharpe, pp. 43–62.
- Tversky, A. and D. Kahneman (1979). "Prospect theory: An analysis of decision under risk". In: *Econometrica* 47.2, pp. 263–291.
- Valette-Florence, P., H. Guizani, and D. Merunka (2011). "The impact of brand personality and sales promotions on brand equity". In: *Journal of Business Research* 64.1, pp. 24–28.
- Wisker, Z. L., D. Kadirov, and J. Nizar (2020). "Marketing a destination brand image to Muslim tourists: Does accessibility to cultural needs matter in developing brand loyalty?" In: *Journal of Hospitality & Tourism Research*.
- Yeh, C-H., Y-S. Wang, and K. Yieh (2016). "Predicting smartphone brand loyalty: Consumer value and consumer-brand identification perspectives". In: *International Journal of Information Management* 36.3, pp. 245–257.
- Zizzo, D.J. (2010). "Experimenter demand effects in economic experiments". In: *Experimental Economics* 13, pp. 75–98.

BIBLIOGRAPHY

- Abeler, J. et al. (2011). "Reference points and effort provision". In: *American Economic Review* 101.2, pp. 470–92.
- Acland, D. and M.R. Levy (2015). "Naiveté, projection bias, and habit formation in gym attendance". In: *Management Science* 61.1, pp. 146–160.
- Adams, C. D. and A. Dickinson (1981). "Instrumental responding following reinforcer devaluation". In: *Quarterly Journal of Experimental Psychology* 33B, pp. 109–121.
- Adda, J. and F. Cornaglia (2006). "Taxes, cigarette consumption, and smoking intensity". In: *American Economic Review* 96.4, pp. 1013–1028.
- Adriaanse, M. A. et al. (2014). "Effortless inhibition: Habit mediates the relation between self-control and unhealthy snack consumption". In: *Frontiers in Psychology* 5, p. 444.
- Agranov, M. and A. Buyalskaya (2021). "Deterrence effects of enforcement schemes: An experimental study". In: *Management Science*, forthcoming.
- Ailawadi, K. L., K. Gedenk, et al. (2007). "Decomposition of the sales impact of promotion-induced stockpiling". In: *Journal of Marketing Research* 44.3, pp. 450–467.
- Ailawadi, K. L., D. Lehmann, and Scott A. Neslin (2003). "Revenue premium as an outcome measure of brand equity". In: *Journal of Marketing* 67.4, pp. 1–17.
- Aldrich, J. H., J. M. Montgomery, and W. Wood (2011). "Turnout as a Habit". In: *Political Behavior* 33.4, pp. 535–563.
- Allen, E. and P. Dechow (2020). "Gender differences in the strategies used for task completion: An analysis of marathon runners". In: *Available at SSRN: <https://ssrn.com/abstract=3536645>*.
- Anderson, E. T., G. J. Fitzsimons, and D. Simester (2006). "Measuring and mitigating the costs of stockouts". In: *Management Science* 52.11, pp. 1751–1763.
- Anderson, E. T. and D. Simester (2004). "Long-run effects of promotion depth on new versus established customers: Three field studies". In: *Marketing Science* 23.1, pp. 4–20.
- Auld, M. C. and P. Grootendorst (2004). "An empirical analysis of milk addiction". In: *Journal of Health Economics* 23.6, pp. 1117–1133.
- Balleine, B. and A. Dezfouli (2019). "Hierarchical Action Control: Adaptive Collaboration Between Actions and Habits". In: *Frontiers in Psychology* 10.1, p. 2735.

- Balleine, B. and J. O'Doherty (2010). "Human and rodent homologies in action control: Corticostriatal determinants of goal-directed and habitual action." In: *Neuropsychopharmacology : Official Publication of the American College of Neuropsychopharmacology* 35.1, pp. 48–69.
- Bar-Isaac, H., G. Caruana, and V. Cuñat (2012). "Search, design, and market structure". In: *American Economic Review* 102.2, pp. 1140–60.
- Bargh, J. (1994). "The four horsemen of automaticity: Awareness, intention, efficiency, and control in social cognition". In: *Handbook of Social Cognition: Basic Processes; Applications*. Ed. by R. S. Wyer and T.K. Srull. Lawrence Erlbaum Associates, Inc., pp. 1–40.
- Batra, R., A. Ahuvia, and R. P. Bagozzi (2012). "Brand love". In: *Journal of Marketing* 76.2, pp. 1–16.
- Baucells, M., M. Weber, and F. Welfens (2011). "Reference-point formation and updating". In: *Management Science* 57.3, pp. 506–519.
- Baumeister, R., T. Heatherton, and D. Tice (1994). *Losing control*. San Diego: Academic Press.
- Bayley, P., J. Franscino, and L. R. Squire (2005). "Robust habit learning in the absence of awareness and independent of the medial temporal lobe". In: *Nature* 436.7050, pp. 550–553.
- Becker, G. S. and K. M. Murphy (1988). "A Theory of Rational Addiction". In: *Journal of Political Economy* 96.4, pp. 675–700.
- Bedolla, L. G. and M. R. Michelson (2012). *Mobilizing Inclusion: Transforming the Electorate Through Get-Out-the-Vote Campaigns*. Yale University Press.
- Belk, R. W. (1975). "Situational variables and consumer behavior". In: *Journal of Consumer Research* 2.3, pp. 157–164.
- Benjamin, J. and L. et al Li (1996). "Population and familial association between the D4 dopamine receptor gene and measures of Novelty Seeking". In: *Nature Genetics* 12, pp. 81–84.
- Bernheim, B.D. and A. Rangel (2004). "Addiction and cue-triggered decision processes". In: *American Economic Review* 94, pp. 1558–1590.
- Black, J. (2020). "How Google Got Its Employees to Eat Their Vegetables". In: *OneZero, a Medium publication*.
- Blain, B. et al. (2019). "Neuro-computational impact of physical training overload on economic decision-making". In: *Current Biology* 29.19, pp. 3289–3297.
- Blaisdell, B. (2012). *The Wit and Wisdom of Abraham Lincoln*. Dover Publications.
- Brady, H. E. and J. E. McNulty (2011). "Turning out to vote: The costs of finding and getting to the polling place". In: *The American Political Science Review* 105.1, pp. 115–134.

- Breakwell, G. M. and D. Rose (2006). “Preface to the Handbook”. In: *Research Methods in Psychology*. Ed. by G. M. Breakwell et al. Sage Publications, pp. 2–23.
- Brody, R. A. and P. M. Sniderman (1977). “From life space to polling place: The relevance of personal concerns for voting behavior”. In: *British Journal of Political Science* 7.3, pp. 337–360.
- Brown, A., C. F. Camerer, and D. Lovallo (2013). “Estimating structural models of limited strategic thinking in the field: The case of missing movie critic reviews”. In: *Management Science* 59.3, pp. 733–747.
- Buyalskaya, A., M. Gallo, and C. F. Camerer (2021). “The golden age of social science”. In: *Proceedings of the National Academy of Sciences* 118.5.
- Buyalskaya, A., H. Ho, et al. (2021). “Predicting context-sensitivity of behavior in field data using machine learning”. In: *Under Review at Science*.
- Camerer, C. F. (2010). “The promise and success of lab-field generalizability in experimental economics: A critical reply to Levitt and List”. In: *The Methods of Modern Experimental Economics*. Ed. by G. Frechette and A. Schotter. Oxford University Press.
- (2011). “Prospect theory in the wild: Evidence from the field”. In: *Advances in Behavioral Economics*. Ed. by C. F. Camerer, G. Loewenstein, and M. Rabin. Princeton University Press, pp. 148–161.
- Camerer, C. F., A. Dreber, et al. (2016). “Evaluating replicability of laboratory experiments in economics”. In: *Science* 351.6280, pp. 1433–1436.
- Camerer, C. F., P. Landry, and R. Webb (2021). “The neuroeconomics of habit”. In: *The State of Mind in Economics*. Ed. by A. Kirman and M. Teschi.
- Camerer, C. F. and D. Mobbs (2017). “Differences in behavior and brain activity during hypothetical and real choices”. In: *Trends in Cognitive Sciences* 21.1, pp. 46–56.
- Campbell, J. Y. and J. H. Cochrane (1999). “By force of habit: A consumption-based explanation of aggregate stock market behavior”. In: *Journal of Political Economy* 107.2, pp. 205–251.
- Chan, T., C. Narasimhan, and Q. Zhang (2008). “Decomposing purchase elasticity with a dynamic structural model of flexible consumption”. In: *Journal of Marketing Research* 45.4, pp. 487–498.
- Chandon, P., B. Wansink, and G. Laurent (2000). “A benefit congruency framework of sales promotion effectiveness”. In: *Journal of Marketing* 64.4, pp. 65–81.
- Charness, G. and U. Gneezy (2009). “Incentives to exercise”. In: *Econometrica* 77.3, pp. 909–931.
- Chen, Z. and P. Rey (2012). “Loss Leading as an Exploitative Practice”. In: *American Economic Review* 102.7, pp. 3462–82.

- Chevalier, J. A., K. K. Anil, and P. E. Rossi (2003). “Why don’t prices rise during periods of peak demand? Evidence from scanner data”. In: *American Economic Review* 93.1, pp. 15–37.
- Ching, A. T. and M. Osborne (2020). “Identification and estimation of forward-looking behavior: The case of consumer stockpiling”. In: *Marketing Science* 39.4, pp. 669–848.
- Coelho, P. S., P. Rita, and Z. Raposo Santos (2018). “On the relationship between consumer-brand identification, brand community, and brand loyalty”. In: *Journal of Retailing and Consumer Services* 43, pp. 101–110.
- Constantinides, G. M. (1990). “Habit formation: A resolution of the equity premium puzzle”. In: *Journal of Political Economy* 98.3, pp. 519–543.
- Coppock, A. and D. Green (2016). “Is voting habit forming? New evidence from experiments and regression discontinuities”. In: *American Journal of Political Science* 60.4, pp. 1044–1062.
- Cravens, M. (2020). “Measuring the strength of voter turnout habits”. In: *Electoral Studies* 64, pp. 102–117.
- Crawford, I. (2010). “Habits revealed”. In: *The Review of Economic Studies* 77.4, pp. 1382–1402.
- Dai, H., K. L. Milkman, et al. (2015). “The impact of time at work and time off from work on rule compliance: The case of hand hygiene in health care”. In: *Journal of Applied Psychology* 100.3, pp. 846–862.
- Dai, H., K.L. Milkman, and J. Riis (2014). “The fresh start effect: Temporal landmarks motivate aspirational behavior”. In: *Management Science* 60.10, pp. 2563–2582.
- Danner, U., N. de Vries, and H. Aarts (2008). “Habit vs. intention in the prediction of future behaviour: The role of frequency, context stability and mental accessibility of past behaviour”. In: *British Journal of Social Psychology* 47.2, pp. 245–265.
- Daw, N. et al. (2011). “Model-Based Influences on Humans’ Choices and Striatal Prediction Errors”. In: *Neuron* 69.6, pp. 1204–1215.
- Dayan, P. and K. Berridge (2014). “Model-based and Model-free pavlovian reward learning: Revaluation, revision and revelation”. In: *Cognitive Affective Behavioral Neuroscience* 14.2, pp. 473–492.
- De Witt Huberts, J., C. Evers, and D. De Ridder (2012). “License to sin: Self-licensing as a mechanism underlying hedonic consumption”. In: *European Journal of Social Psychology* 42.4, pp. 490–496.
- (2013). ““Because I am worth it” A theoretical framework and empirical review of a justification-based account of self-regulation failure”. In: *Personality and Social Psychology Review* 18.2, pp. 119–138.

- De Witt Huberts, J., C. Evers, and D. De Ridder (2014). “Thinking before sinning: Reasoning processes in hedonic consumption”. In: *Frontiers in Psychology* 5, p. 1268.
- Deaton, A. (1992). *Understanding Consumption*. Oxford University Press.
- DellaVigna, S. and U. Malmendier (2006). “Paying not to go to the gym”. In: *American Economic Review* 96.3, pp. 694–719.
- DelVecchio, D., D. H. Henard, and T. H. Freling (2006). “The effect of sales promotion on post-promotion brand preference: A meta-analysis”. In: *Journal of Retailing* 82.3, pp. 203–213.
- DelVecchio, D., H. S. Krishnan, and D. C. Smith (2007). “Cents or percent? The effects of promotion framing on price expectations and choice”. In: *Journal of Marketing* 71.3, pp. 158–170.
- Denny, K. and O. Doyle (2009). “Does voting history matter? Analysing persistence in turnout”. In: *American Journal of Political Science* 53.1, pp. 17–35.
- Dick, A. and K. Basu (1994). “Customer loyalty: Toward an integrated conceptual framework”. In: *JAMS* 22, pp. 99–113.
- Dickinson, A., D. J. Nicholas, and C. D. Adams (1983). “The effect of the instrumental training contingency on susceptibility to reinforcer devaluation”. In: *The Quarterly Journal of Experimental Psychology* 35.1, pp. 35–51.
- Dolan, P. and M. Galizzi (2014). “Because I’m worth it. A lab-field experiment on the spillover effects of incentives in health”. In: *LSE CEP Discussion Paper CEPDP 1286, London*.
- (2015). “Like ripples on a pond: Behavioral spillovers and their implications for research and policy”. In: *J. Econ. Psychol* 47, pp. 1–16.
- Downs, A. (1957). *An Economic Theory of Democracy*. New York: Harper & Row.
- Dubé, J.-P., G. J. Hitsch, and P. E. Rossi (2009). “Do switching costs make markets less competitive?” In: *Journal of Marketing Research* 46.4, pp. 435–445.
- (2010). “State dependence and alternative explanations for consumer inertia”. In: *The RAND Journal of Economics* 41.3, pp. 417–445.
- Dubois, P., R. Griffith, and M. O’Connell (2018). “The effects of banning advertising in junk food markets”. In: *Review of Economic Studies* 85.1, pp. 396–436.
- Duckworth, A., K. Milkman, and D. Laibson (2018). “Beyond willpower: Strategies for reducing failures of self-control”. In: *Psychological Science in the Public Interest* 19.3, pp. 102–129.
- Duesenberry, J. S. (1949). *Income, Saving and the Theory of Consumption Behavior*. Cambridge, Mass.: Harvard University Press.
- Dutcher, G., T. Salmon, and K. J. Saral (2015). “Is ‘Real’ Effort More Real?” In: Available at SSRN: <http://dx.doi.org/10.2139/ssrn.2701793>.

- Eales, T. (2016). *Price and Promotion in Western Economies*. Tech. rep. IRI.
- Erdem, T., S. Imai, and M. P. Keane (2003). “Brand and quantity choice dynamics under price uncertainty”. In: *Quantitative Marketing and Economics* 1, pp. 5–64.
- Ersche, K. et al. (2017). “Creature of Habit: A self-report measure of habitual routines and automatic tendencies in everyday life”. In: *Personality and Individual Differences* 116, pp. 73–85.
- Faruk, G. and W. Pesendorfer (2008). “The case for mindless economics”. In: *The Foundations of Positive and Normative Economics*. Ed. by A. Caplin and A. Schotter. Oxford University Press, pp. 3–40.
- Ferguson, S. G. and S. Shiffman (2009). “The relevance and treatment of cue-induced cravings in tobacco dependence”. In: *Journal of Substance Abuse Treatment* 36.3, pp. 235–243.
- Fishbach, A., R. K. Ratner, and Y. Zhang (2011). “Inherently loyal or easily bored? Nonconscious activation of consistency versus variety-seeking behavior.” In: *Journal of Consumer Psychology* 21, pp. 38–48.
- Fisher, I. (1930). *The theory of interest : As determined by impatience to spend income and opportunity to invest it*. New York: Macmillan Co.
- Fournier, M. et al. (2017). “Effects of circadian cortisol on the development of a health habit”. In: *Health Psychology* 36.11, pp. 1059–1064.
- Fox, H. C. et al. (2007). “Enhanced sensitivity to stress and drug/alcohol craving in abstinent cocaine-dependent individuals compared to social drinkers”. In: *Neuropsychopharmacology* 33.4, pp. 796–805.
- Franklin, M. N. and S. B. Hobolt (2011). “The legacy of lethargy: How elections to the European Parliament depress turnout”. In: *Electoral Studies* 30.2, pp. 67–76.
- Fransen, A. (2019). “Experimental considerations on habit formation in humans”. MA thesis. Maastricht University.
- Franzen, A. and S. Pointner (2013). “The external validity of giving in the dictator game: A field experiment using the misdirected letter technique”. In: *Experimental Economics* 16, pp. 155–169.
- Galizzi, M. and D. Navarro-Martinez (2019). “On the external validity of social preference games: A systematic lab-field study”. In: *Management Science* 65.3, pp. 976–1002.
- Gardner, B. (2015). “A review and analysis of the use of ‘habit’ in understanding, predicting and influencing health-related behavior”. In: *Healthy Psychology Review* 9.3, pp. 277–295.
- Gardner, B., C. Abraham, et al. (2012). “Towards parsimony in habit measurement: Testing the convergent and predictive validity of an automaticity subscale of the Self-Report Habit Index”. In: *International Journal of Behavioral Nutrition and Physical Activity* 9.102.

- Gardner, B. and P. Lally (2018). “Modelling Habit Formation and Its Determinants”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 207–229.
- Garr, E. and A. Delamater (2019). “Exploring the relationship between actions, habits, and automaticity in an action sequence task”. In: *Learning and Memory* 26, pp. 128–132.
- Gillebaart, M. and M. A. Adriaanse (2017). “Self-control predicts exercise behavior by force of habit, a conceptual replication of Adriaanse et al. (2014)”. In: *Frontiers in Psychology* 8, p. 190.
- Gläscher, J. et al. (2010). “States versus rewards: Dissociable neural prediction error signals underlying model-based and model-free reinforcement learning”. In: *Neuron* 66.4, pp. 585–595.
- Gonsalves, B. D. et al. (2005). “Memory strength and repetition suppression: Multimodal imaging of medial temporal cortical contributions to recognition”. In: *Neuron* 47.5, pp. 751–761.
- Graybiel, A. M. (1998). “The basal ganglia and chunking of action repertoires”. In: *Neurobiology of Learning and Memory* 70.1-2, pp. 119–136.
- Green, D. and A. S. Gerber (2002). “The downstream benefits of experimentation”. In: *Political Analysis* 10.4, pp. 394–402.
- Gullo, K. et al. (2019). “Does Time of Day Affect Variety-Seeking?” In: *Journal of Consumer Research* 46.1, pp. 20–35.
- Gupta, S. (1988). “Impact of sales promotion on when, what, and how much to buy”. In: *Journal of Marketing Research* 25, pp. 342–355.
- Hanseens, D.M. (2018). *Long-Term Impact of Marketing*. World Scientific.
- Harrington, N. (2017). “Commentary: Why it doesn’t pay to ask consumers about habitual behaviors”. In: *Journal of the Association for Consumer Research: The Habit-Driven Consumer* 2.3.
- He, H., Y. Li, and L. Harris (2012). “Social identity perspective on brand loyalty”. In: *Journal of Business Research* 65.5, pp. 648–657.
- Heckman, J. (1981). “Heterogeneity and State Dependence”. In: *Studies in Labor Markets*. Ed. by Sherwin Rosen. National Bureau of Economic Research, Inc., pp. 91–140.
- Helsen, K. and D.C. Schmittlein (1992). “Some characterizations of stockpiling behavior under uncertainty”. In: *Marketing Letters* 3, pp. 5–16.
- Hendel, I. and A. Nevo (2003). “The post-promotion dip puzzle: What do the data have to say?” In: *Quantitative Marketing and Economics* 1, pp. 409–424.
- (2006). “Measuring the implications of sales and consumer inventory behavior”. In: *Econometrica* 74, pp. 1637–1673.

- Hendel, I. and A. Nevo (2013). “Intertemporal price discrimination in storable goods markets”. In: *American Economic Review* 103.7, pp. 2722–2751.
- Henrich, J. et al. (2005). ““Economic man” in cross-cultural perspective: Behavioral experiments in 15 small-scale societies”. In: *Behavioral and Brain Sciences* 28.6, pp. 795–815.
- Hogenelst, K., R. A. Schoevers, and M. Rot (2015). “Studying the neurobiology of human social interaction: Making the case for ecological validity”. In: *Social Neuroscience* 10.3, pp. 219–229.
- Hopkins, M. et al. (2011). “The relationship between substrate metabolism, exercise and appetite control: Does glycogen availability influence the motivation to eat, energy intake or food choice?” In: *Sports Medicine* 41.6, pp. 507–21.
- Hosken, D. and D. Reiffen (2004). “How retailers determine which products should go on sale: Evidence from store-level data”. In: *Journal of Consumer Policy* 27, pp. 141–177.
- In, Y. and J. Wright (2014). “Loss-leader pricing and upgrades”. In: *Economics Letters* 122.1, pp. 19–22.
- Jacoby, J. and R. Chestnut (1978). *Brand Loyalty: Measurement and Management*. John Wiley and Sons, New York.
- Jasny, B. R. et al. (2011). “Again, and Again, and Again . . .” In: *Science* 334.6060, pp. 1225–1225.
- Ji, M. F. and W. Wood (2007). “Purchase and consumption Habits: Not necessarily what you intend”. In: *Journal of Consumer Psychology* 17.4, pp. 261–276.
- Job, V., V. Sieber, et al. (2018). “Age differences in implicit theories about willpower: Why older people endorse a nonlimited theory”. In: *Psychology and Aging* 33.6, pp. 940–952.
- Job, V., G. M. Walton, et al. (2013). “Beliefs about willpower determine the impact of glucose on self-control”. In: *PNAS* 110.37, pp. 14837–14842.
- K., Gurumurthy and Russell S. W. (1995). “Empirical Generalizations from reference price research”. In: *Marketing Science* 14.3 Supplement, G161–G169.
- Kahn, B. E. (1995). “Variety-Seeking Among Goods and Services: An Integrative Review”. In: *Journal of Retailing and Consumer Services* 2.3, pp. 139–148.
- Kahneman, D., J. Knetsch, and R. Thaler (1990). “Experimental tests of the endowment effect and the Coase Theorem”. In: *Journal of Political Economy* 98.6, pp. 1325–1348.
- Karlan, D. and J. Appel (2016). “Part 1: Leading causes of research failures”. In: *Failing in the Field*. Princeton University Press, pp. 17–70.
- Karmarkar, U. and B. Bollinger (2015). “BYOB: How bringing your own shopping bags leads to treating yourself and the environment”. In: *Journal of Marketing* 79.4, pp. 1–15.

- Karni, E. (2008). "State-dependent preferences". In: *The New Palgrave Dictionary of Economics*. Ed. by Palgrave Macmillan. Palgrave Macmillan, London.
- Kaushal, N. and R. Rhodes (2015). "Exercise habit formation in new gym members: a longitudinal study". In: *Journal of Behavioral Medicine* 38, pp. 652–663.
- Keane, M. P. (1997). "Modeling heterogeneity and state dependence in consumer choice behavior". In: *Review of Economics and Statistics* 15.3, pp. 310–327.
- Keller, K. L. and D. Lehmann (2006). "Brands and branding: Research findings and future priorities". In: *Marketing Science* 25.6, pp. 551–765.
- Kessler, J. and L. Vesterlund (2015). "The external validity of laboratory experiments: The misleading emphasis on quantitative effects". In: *Handbook of Experimental Economic Methodology*. Ed. by G. Frechette and A. Schotter. Oxford University Press, pp. 391–406.
- Khan, U. and R. Dhar (2006). "Licensing effect in consumer choice". In: *Journal of Marketing Research* 43.2, pp. 259–266.
- Kirchner, T. et al. (2013). "Geospatial Exposure to Point-of-Sale Tobacco: Real-Time Craving and Smoking Cessation Outcomes". In: *American Journal of Preventative Medicine* 45.4.
- Kivetz, R. and T. Simonson (2002). "Earning the right to indulge: Effort as a determinant of customer preferences toward frequency program rewards". In: *Journal of Marketing Research* 39.2, pp. 155–170.
- Knowlton, B. and T. K. Patterson (2016). "Habit formation and the striatum". In: *Behavioral Neuroscience of Learning and Memory. Current Topics in Behavioral Neurosciences, vol 37*. Ed. by R.E. Clark and S. Martin. Springer.
- Koopmans, T. C. (1960). "Stationary ordinal utility and impatience". In: *Econometrica* 28.2, pp. 287–309.
- Kőszegi, B. and M. Rabin (2007). "Reference-dependent risk attitudes". In: *American Economic Review* 97.4, pp. 1047–1073.
- Kuehn, A. (1962). "Consumer brand choice as a learning process". In: *Journal of Advertising Research* 2, pp. 10–17.
- Kwon, H. E. et al. (2016). "Excessive dependence on mobile social apps: A rational addiction perspective". In: *Information Systems Research* 27.
- Laibson, D. (2001). "A Cue-Theory of Consumption". In: *Quarterly Journal of Economics* 116.1, pp. 81–119.
- Lally, P. et al. (2010). "How are habits formed: Modelling habit formation in the real world". In: *European Journal of Social Psychology* 40.6, pp. 998–1009.
- Lattin, J. M. and R. E. Bucklin (1989). "Reference effects of price and promotion on brand choice behavior". In: *Journal of Marketing Research* 26.3, pp. 299–310.

- Lee, S., S. Shimojo, and J. O'Doherty (2014). "Neural computations underlying arbitration between model-based and model-free learning". In: *Neuron* 81.3, pp. 687–699.
- Levitt, S. and J. List (2007). "What do laboratory experiments measuring social preferences reveal about the real world?" In: *Journal of Economic Perspectives* 21.2, pp. 153–174.
- (2008). "Homo economicus Evolves". In: *Science* 319.5865, pp. 909–910.
- Liu, Y. and S. Balachander (2014). "How long has it been since the last deal? consumer promotion timing expectations and promotional response". In: *Quantitative Marketing and Economics* 12.1, pp. 85–128.
- Marom, O. and A. Seidmann (2011). "Using "last-minute" sales for vertical differentiation on the Internet". In: *Decision Support Systems* 51.4, pp. 894–903.
- Marr, D. (1982). *Vision: A computational investigation into the human representation and processing of visual information*. Cambridge, Mass: MIT Press.
- Mazar, A. and W. Wood (2018). "Defining habit in psychology". In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 13–29.
- Mazumdar, T., S. P. Raj, and I. Sinha (2005). "Reference price research: Review and propositions". In: *Journal of Marketing* 69.4, pp. 84–102.
- McAlister, L. and E. Pessemier (1982). "Variety Seeking Behavior: An Interdisciplinary Review". In: *Journal of Consumer Research* 9.3, pp. 311–322.
- McNutt, M. (2014). "Reproducibility". In: *Science* 343.229.
- Mela, C., S. Gupta, and D. Lehmann (1997). "The long-term impact of promotion and advertising on consumer brand choice". In: *Journal of Marketing Research* 34.2, pp. 248–261.
- Mela, C., K. Jedidi, and D. Bowman (1998). "The long-term impact of promotions on consumer stockpiling behavior". In: *Journal of Marketing Research* 35.2, pp. 250–262.
- Meredith, M. (2009). "Persistence in political participation". In: *Quarterly Journal of Political Science* 4.3, pp. 187–209.
- Middleton, J. C. et al. (2010). "Effectiveness of policies maintaining or restricting days of alcohol sales on excessive alcohol consumption and related harms". In: *American journal of Preventive Medicine* 39.6, pp. 575–589.
- Monroe, K. B. and R. Krishnan (1985). "The effect of price on subjective product evaluation". In: *The perception of merchandise and store quality*. Ed. by J. Jacoby and J. Olson. Lexington: Lexington Book, pp. 209–232.
- Moors, A. and J. De Houwer (2006). "Automaticity: A theoretical and conceptual analysis". In: *Psychological Bulletin* 132.2, pp. 297–326.

- Mullan, B. and E. Novoradovskaya (2018). “Habit Mechanisms and Behavioural Complexity”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 71–90.
- Neal, D., W. Wood, and J. Quinn (2006). “Habits—A repeat performance”. In: *Current Directions in Psychological Science* 15.4, pp. 198–202.
- Neal, D., W. Wood, M. Wu, et al. (2011). “The Pull of the past: When do habits persist despite conflict with motives?” In: *Personality and Social Psychology Bulletin* 37.11, pp. 1428–1437.
- O’Brien, E. (2021). “A mind stretched: The psychology of repeat consumption”. In: *Consumer Psychology Review* 4, pp. 42–58.
- Oliver, R. (1999). “Whence consumer loyalty?” In: *Journal of Marketing* 63, pp. 33–44.
- Orbell, S. and B. Verplanken (2010). “The automatic component of habit in health behavior: Habit as cue-contingent automaticity”. In: *Health Psychology* 29.4, pp. 374–383.
- Osborne, M. (2018). “Frequency versus depth: How changing the temporal process of promotions impacts demand for a storable good”. In: *Japanese Economic Review* 69, pp. 258–283.
- OSC (2012). “An open, large-scale, collaborative effort to estimate the reproducibility of psychological science”. In: *Perspectives on Psychological Science* 7.6.
- Östling, R. et al. (2011). “Testing game theory in the field: Swedish LUPI lottery games”. In: *American Economic Journal: Microeconomics* 3.3, pp. 1–33.
- Ouellette, J. A. and W. Wood (1998). “Habit and intention in everyday life: The multiple processes by which past behavior predicts future behavior”. In: *Psychological Bulletin* 124.1, pp. 54–74.
- Parry, D. et al. (2021). “A systematic review and meta-analysis of discrepancies between logged and self-reported digital media use”. In: *Nature Human Behaviour* (forthcoming).
- Pauwels, K., D. M. Hanssens, and S. Siddarth (2002). “The long-term effects of price promotions on category incidence, brand choice, and purchase quantity”. In: *Journal of Marketing Research* 39.4, pp. 421–439.
- Peltzman, S. (1975). “The Effects of Automobile Safety Regulation”. In: *Journal of Political Economy* 83.4, pp. 677–726.
- Perez, O. D. and A. Dickinson (2020). “A theory of actions and habits: The interaction of rate correlation and contiguity systems in free-operant behavior”. In: *Psychological Review* 127.6, pp. 945–971.
- Pesendorfer, M. (2002). “Retail sales: A study of pricing behavior in supermarkets”. In: *Journal of Business* 75.1, pp. 33–66.

- Pires, T. (2016). “Costly search and consideration set in storable goods markets”. In: *Quantitative Marketing and Economics* 14.3, pp. 157–193.
- Plott, C. R. and V. L. Smith (2008). “Preface to the Handbook”. In: *Handbook of Experimental Economics Results*. Ed. by Charles R. Plott and Vernon L. Smith. Handbook of Experimental Economics Results. Elsevier, pp. 1–2.
- Pool, E. et al. (2021). “Determining the effects of training duration on the behavioral expression of habitual control in humans: A multi-laboratory investigation”. In: *PsyArXiv*: <https://psyarxiv.com/z756h>.
- Potthoff, S. et al. (2018). “Creating and breaking habit in healthcare professional behaviours to improve healthcare and health”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 247–265.
- Putler, D. S. (1992). “Incorporating reference price effects into a theory of consumer choice”. In: *Marketing Science* 11.3, pp. 287–309.
- Quinn, J. and W. Wood (2021). “Habits Across the Lifespan”. In: *Working Paper*.
- Ray, D., M. Shum, and C. F. Camerer (2015). “Loss aversion in post-sale purchases of consumer products and their substitutes”. In: *American Economic Review* 105.5, pp. 376–80.
- Rebar, A. et al. (2018). “The measurement of habit”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 31–49.
- Rhodes, R. and A. Rebar (2018). “Physical activity habit: Complexities and controversies”. In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 91–109.
- Roberts, W. (Published on January 6, 2014). “Why do runners gain weight?” In: *Runner’s World Magazine*.
- Romaniuk, J. and M. Nenycz-Thiel (2013). “Behavioral brand loyalty and consumer brand associations”. In: *Journal of Business Research* 66.1, pp. 67–72.
- Roth, A. (1986). “Laboratory experimentation in economics”. In: *Economics and Philosophy* 2, pp. 245–273.
- Royer, H., M. Steher, and J. Sydnor (2015). “Incentives, commitments, and habit formation in exercise: Evidence from a field experiment with workers at a Fortune-500 company”. In: *American Economic Journal: Applied Economics* 7.3, pp. 51–84.
- Rozen, K. (2010). “Foundations of intrinsic habit formation”. In: *Econometrica* 78.4, pp. 1341–1373.
- Rust, J. (1994). “Structural estimation of Markov decision processes”. In: *Handbook of Econometrics*. Ed. by Engle R. and McFadden D. Elsevier.
- Ryder, H. E. and G. M. Heal (1973). “Optimal growth with intertemporally dependent preferences”. In: *The Review of Economic Studies* 40.1, pp. 1–31.

- Samuelson, P. (1948). "Consumption theory in terms of revealed preference". In: *Economica* 15, pp. 243–253.
- Schwabe, L. and O. Wolf (2009). "Stress Prompts Habit Behavior in Humans". In: *The Journal of Neuroscience* 29.22, pp. 7191–7198.
- Seger, C. A. and B. J. Spiering (2011). "A critical review of habit learning and the basal ganglia". In: *Frontiers in Systems Neuroscience* 5.66.
- Seiler, S. (2013). "The impact of search costs on consumer behavior: A dynamic approach". In: *Quantitative Marketing and Economics* 11.2, pp. 155–203.
- Shen, K. and D. E. Giles (2006). "Rational exuberance at the mall: Addiction to carrying a credit card balance". In: *Applied Economics* 38.5, pp. 587–592.
- Sherman, W. M. (1995). "Metabolism of sugars and physical performance". In: *The American Journal of Clinical Nutrition* 62, 228S–241S.
- Shi, S. and L. Epstein (1993). "Habits and time preference". In: *International Economic Review* 34.1, pp. 61–84.
- Simonson, I. (1990). "The Effect of Purchase Quantity and Timing on Variety-Seeking Behavior". In: *Journal of Marketing Research* 27.2, pp. 150–162.
- Sinha, R. (2009). "Modeling stress and drug craving in the laboratory: Implications for addiction treatment development". In: *Addiction Biology* 14.1, pp. 84–98.
- Slonim, R. and E. Garbarino (1999). "The effect of price history on demand as mediated by perceived price expensiveness". In: *Journal of Business Research* 45, pp. 1–14.
- Smith, K. S. et al. (2012). "Reversible online control of habitual behavior by optogenetic perturbation of medial prefrontal cortex". In: *Proceedings of the National Academy of Sciences* 109.46, pp. 18932–18937.
- Staats, B. R. et al. (2017). "Motivating process compliance through individual electronic monitoring: An empirical examination of hand hygiene in healthcare". In: *Management Science* 63.5, pp. 1563–1585.
- Steinfeld, M. R. and M. E. Bouton (2020). "Context and renewal of habits and goal-directed actions after extinction". In: *Journal of Experimental Psychology: Animal Learning and Cognition* 46.4, pp. 408–421.
- Sun, B. (2005). "Promotion effect on endogenous consumption". In: *Marketing Science* 24.3, pp. 430–443.
- Sundaresan, S. M. (1989). "Intertemporally dependent preferences and the volatility of consumption and wealth". In: *The Review of Financial Studies* 2.1, pp. 73–89.
- Tam, L., W. Wood, and M. F. Ji (2009). "Brand loyalty is not habitual". In: *Handbook of brand relationships*. Ed. by D. J. MacInnis, C. W. Park, and J. R. Priester. M E Sharpe, pp. 43–62.

- Tanaka, T., C. F. Camerer, and Q. Nguyen (2010). "Risk and Time Preferences: Linking Experimental and Household Survey Data from Vietnam". In: *American Economic Review* 100.1, pp. 557–571.
- Thaler, R. and S. Benartzi (2004). "Save more tomorrow (TM): Using behavioral economics to increase employee saving". In: *Journal of Political Economy* 112.S1, S164–S187.
- Tricomi, E., B. Balleine, and J. O'Doherty (2009). "A specific role for posterior dorsolateral striatum in human habit learning". In: *The European journal of neuroscience* 29.11, pp. 2225–2232.
- Tversky, A. and D. Kahneman (1979). "Prospect theory: An analysis of decision under risk". In: *Econometrica* 47.2, pp. 263–291.
- (1992). "Advances in prospect theory: Cumulative representation of uncertainty". In: *Journal of Risk and Uncertainty* 5, pp. 297–323.
- Valette-Florence, P., H. Guizani, and D. Merunka (2011). "The impact of brand personality and sales promotions on brand equity". In: *Journal of Business Research* 64.1, pp. 24–28.
- Verplanken, B. (2018). "Introduction". In: *The Psychology of Habit*. Ed. by B. Verplanken. Springer, pp. 1–10.
- Verplanken, B., O. Friberg, et al. (2007). "Mental habits: Metacognitive reflection on negative self-thinking". In: *Journal of Personality and Social Psychology* 92.3, pp. 526–541.
- Verplanken, B. and S. Orbell (2003). "Reflections on past behavior: A self-report index of habit strength". In: *Journal of Applied Social Psychology* 33.6, pp. 1313–1330.
- Volpp, K. G. and G. Loewenstein (2020). "What is a habit? Diverse mechanisms that can produce sustained behavior change". In: *Organizational Behavior and Human Decision Processes* 161, pp. 36–38.
- Volpp, K. G. et al. (2009). "A randomized, controlled trial of financial incentives for smoking cessation". In: *The New England Journal of Medicine* 360.7, pp. 699–709.
- Walker, M. and J. Wooders (2001). "Minimax play at Wimbledon". In: *The American Economic Review* 91.5, pp. 1521–1538.
- Wertenbroch, K. (1998). "Consumption self-control by rationing purchase quantities of virtue and vice". In: *Marketing Science* 17.4, pp. 317–337.
- Wilcockson, Thomas D., D. A. Ellis, and H. Shaw (2018). "Determining typical smartphone usage: What data do we need?" In: *Cyberpsychology, Behavior, and Social Networking* 21.6, pp. 395–398.

- Wisker, Z. L., D. Kadirov, and J. Nizar (2020). "Marketing a destination brand image to Muslim tourists: Does accessibility to cultural needs matter in developing brand loyalty?" In: *Journal of Hospitality & Tourism Research*.
- Wit, S. de et al. (2018). "Shifting the balance between goals and habits: Five failures in experimental habit induction". In: *Journal of Experimental Psychology* 147.7, pp. 1043–1065.
- Wood, W. and D. Neal (2007). "A new look at habits and the habit-goal interface". In: *Psychological Review* 114.4, pp. 843–863.
- (2009). "The habitual consumer". In: *Journal of Consumer Psychology* 19, pp. 579–592.
- Wood, W., J. Quinn, and D. A. Kashy (2002). "Habits in everyday life: Thought, emotion, and action". In: *Journal of Personality and Social Psychology* 83.6, pp. 1281–1297.
- Wood, W. and D. Runger (2016). "Psychology of Habit". In: *Annual review of psychology* 67, pp. 289–314.
- Wood, W., M. G. Witt, and L. Tam (2005). "Changing Circumstances, Disrupting Habits". In: *Journal of Personality and Social Psychology* 88.6, pp. 918–933.
- Yeh, C-H., Y-S. Wang, and K. Yieh (2016). "Predicting smartphone brand loyalty: Consumer value and consumer-brand identification perspectives". In: *International Journal of Information Management* 36.3, pp. 245–257.
- Yin, H. H. and B. Knowlton (2006). "The role of the basal ganglia in habit formation". In: *Nature Reviews Neuroscience* 7, pp. 464–476.
- Zizzo, D.J. (2010). "Experimenter demand effects in economic experiments". In: *Experimental Economics* 13, pp. 75–98.