# Habit Formation, Naiveté, and Projection Bias in Gym Attendance\*

Dan Acland<sup> $\dagger$ </sup> and Matthew Levy<sup> $\ddagger$ </sup>

April 13, 2011

<sup>\*</sup>The authors would like to thank Stefano DellaVigna, Gary Charness, Uri Gneezy, Teck Hua Ho, Shachar Kariv, Botond Koszegi, Ulrike Malmendier, Matthew Rabin, and seminar participants at UC Berkeley and Harvard for their helpful comments. Financial support was provided by the National Institute on Aging through the Center on the Economics and Demography of Aging at UC Berkeley [grant number P30 AG12839].

<sup>&</sup>lt;sup>†</sup>University of California, Berkeley

 $<sup>^{\</sup>ddagger} \rm Corresponding author.$  Harvard University. Email mattlevy@rwj.harvard.edu. Telephone 617.495.5365. Fax 617.496.1636

#### Abstract

We develop a model that parsimoniously captures habit formation, projection bias, and present bias in an intertemporal-choice setting, and conduct a field experiment to calibrate the parameters of the model. Building on the gym-attendance study of Charness and Gneezy (2009), we incentivize subjects to attend the gym for a month, observe pre- and post-treatment attendance relative to a control group, and elicit pre- and post-treatment predictions of posttreatment attendance. Present-biased subjects will want their future selves to exercise more than they actually will, and naivete about this bias will cause them to over-estimate their future gym attendance. Subjects with projection bias expect their future preferences will too closely resemble their current ones, and therefore under-estimate any habit-formation effect of our treatment *ex-ante*. We find that subjects form a short-run habit, although we also find evidence of substantial decay caused by the shock of the semester break. Subjects appear not to embed habit formation into their ex-ante predictions. We additionally find that subjects greatly over-predict future attendance, which we interpret as evidence of partial naivete with respect to present bias. We estimate a structural model to identify the key parameters and interpret the economic significance of the habit. We find that approximately one-third of subjects formed a habit equivalent, on average, to a \$2.60 per-visit subsidy, while their predictions correspond to 90% projection bias over the habit formation. Subjects are also substantially naive about their self-control problems: the structural results indicate they expect their future selves to be two-thirds less "present biased" than they currently are. This experiment provides novel estimates of important economic parameters, and identifies a pattern of procrastination and underinvestment in health.

JEL Codes: B49, C93, D03, I19

# 1 Introduction

Individuals routinely make decisions involving predictions of how their preferences, costs, and beliefs will unfold in the future. It is commonly assumed that individuals have rational expectations, and while exact preferences, costs, and beliefs may not be known, people know the range of possibilities and make accurate predictions based on expected values. If people's predictions are incorrect, however, their decisions may not achieve long-run optimality. In this paper, we investigate the possibility that people make systematically biased predictions in the domain of physical exercise. We choose this domain both because of its intrinsic importance, and because the choices one makes about exercise are potentially susceptible to two well-established cognitive biases: "present bias" and "projection bias". The former may cause people to procrastinate on going to the gym even when the benefits are large, while the latter may lead people to under-estimate the positive habit formation associated with gym exercise. While these biases have previously been documented in other contexts, we leverage unique features of our experiment to calibrate the extent of both biases simultaneously, as well as the size of the associated welfare costs.

We explore projection bias by building on the consensus in the behavioral health literature that habit formation plays an important role in physical exercise.<sup>1</sup> Optimal choices about exercise therefore require correct beliefs about this habit formation process. Becker and Murphy (1988), in their "Theory of Rational Addiction", incorporate rationality for both positive and negative habits—and hence preserve welfare optimality—by modeling addicts as perfectly forward-looking with respect to the habit-forming effects of current and future consumption.<sup>2</sup> Loewenstein, O'Donoghue and Rabin (2003) demonstrate the importance of prediction of preferences for Becker and Murphy's results, and show how misprediction of habit formation can lead to long-term welfare losses. Much attention is given to negative habits such as smoking, but equally, if people do not foresee the way that healthy behaviors such as exercise can become more enjoyable after a period of habit formation, they may miss out on a lifetime of health benefits. Loewenstein, O'Donoghue and Rabin (2003) model a form of systematic misprediction in which individuals correctly foresee the direction in which their preferences will change, but underappreciate the magnitude of the change.

<sup>&</sup>lt;sup>1</sup>see Valois, Dersharnais and Godin (1988), Dzewaltowski, Noble and Shaw (1990), Reynolds, Killen, Bryson, Maron, Taylor, Maccoby and Farquhar (1990), Godin, Valois and Lepage (1993), Godin (1994)

 $<sup>^{2}</sup>$ This approach has been relaxed to allow imperfect information and ex-post sub-optimal choices, as in Orphanides and Zervos (1995), but still generates a prediction of ex-ante optimality.

They refer to this kind of misprediction as "projection bias" because people are thought of as projecting their current preferences (in the case of exercise, their current level of habituation) onto their future selves. Current beliefs about future preferences are represented by an alpha-mixture of current preferences and true future preferences. Thus, a single parameter,  $\alpha$ , captures the degree of misprediction, with  $\alpha = 0$  corresponding to correct beliefs and  $\alpha = 1$  corresponding to complete projection bias. In this paper, we relax the rational-expectations hypothesis and explore whether projection bias plays a role in people's predictions of their own future exercise behavior.

A second potential source of concern for exercise is "present bias", wherein people have a timeinconsistent taste for immediate gratification. We follow the model of Laibson (1997), and consider subjects with a long-run exponential discount factor  $\delta$  and a short-run discount factor  $\beta$ . Because gym exercise tends to have up-front costs and delayed rewards, a person with quasi-hyperbolic time preferences may always want their future selves to go to the gym, but may never want to go in the present. A "sophisticate" understands her future self-control problems and will therefore seek out commitment devices to help her overcome them. Equally important is the possibility of "naivete", whereby people do not fully appreciate that they will be just as impatient tomorrow as they are today. We follow O'Donoghue and Rabin (1999a), and model agents as believing that their future selves will discount according to some  $\hat{\beta} \in [\beta, 1]$ , with  $\hat{\beta} = 1$  corresponding to the case of complete naivete. This can lead to procrastination on exercise, with people never going to the gym but always believing they will go tomorrow. Such procrastination in this setting can lead to a much greater welfare loss than that which sophisticated agents would experience. Moreover, naivete will cause a person to undervalue the habit induced by regular gym exercise, insofar as it could help act as a self-control device in the future.

The experimental design of this paper draws on Charness and Gneezy (2009), who paid subjects to attend the gym for several weeks and found significantly higher gym attendance in the period after the payment ended than in the pre-intervention period. In particular, we build upon the interpretation of this change as evidence that being paid to attend for a month had led to habit formation. Their subjects were university undergraduates who were randomized into three groups. A "low-incentive" group were offered \$25 to attend the gym once during the initial week of the study. A "high-incentive" group received the same \$25 offer, and were additionally offered \$100 to attend the gym another eight times in the subsequent four weeks for a total of nine visits over five weeks. A third group, who received no offers for gym attendance, served as controls. Gym-attendance data was collected for all students for a period beginning eight weeks before the treatment and ending seven weeks after.<sup>3</sup> They show that subjects in the high-incentive group continue to attend the gym significantly more after the incentive period ends: 0.67 visits per week more than the control group, and 0.58 visits per week more than the low-incentive group. Our goals in this paper are to build on this study not only to better understand the economic significance of this habit formation, but also to use it as a setting in which to estimate systematic biases in subjects' beliefs about their own future preferences.

To test for misprediction of future gym preferences, we first built on Charness and Gneezy's high-incentive and low-incentive treatments. We recruited 120 subjects who were self-reported non-regular gym attenders. This allows us to plausibly assume that none are ex-ante habituated, and therefore to attribute the entire habit to our intervention. In addition to the \$25 and \$100 attendance incentives, we elicited subjects' predictions of their post-treatment gym attendance. conducting an elicitation both immediately before and immediately after the treatment period. The timeline of our experiment was as follows. Attendance incentives were offered at a first meeting, one week prior to the beginning of the treatment period. At the end of that week, pre-treatment predictions of future attendance were elicited. Subjects first indicated their willingness to pay for a certificate that paid off according to actual gym attendance in one of five post-treatment "target weeks". They then stated an unincentivized prediction of their attendance during that target week, should they actually receive the certificate.<sup>4</sup> Thus we have both incentivized and unincentivized measures of predicted future attendance. Next came the five weeks of the treatment period, at the end of which, post-treatment predictions were elicited for the same five target weeks, using the same procedures. Attendance was then monitored during the remaining post-treatment period, comprising a one-week buffer, the five target weeks for which attendance predictions had been elicited, and weeks with no further subsidies which carried over into the following semester.

Conditional on treated subjects developing a gym-attendance habit during the treatment period, we can test whether they predicted this habit by comparing the change in their predictions—from

 $<sup>^{3}</sup>$ In the same paper they conducted a second study with a slightly different design and 13 weeks of post-intervention data, which yielded largely similar results.

 $<sup>^{4}</sup>$ These elicitation procedures, and the endogeneity issues attendant upon our incentivized elicitations, are described in detail in Section 3.

pre- to post-treatment elicitations—with the change for control subjects. If treated subjects fail to foresee all or part of the habit, their attendance predictions should go up more than their controlgroup counterparts. Moreover, by comparing subjects' post-treatment predictions of attendance with their actual post-treatment attendance in general, we can estimate the extent of naivete with regard to self-control problems. If either control or treated subjects (who by this point have already developed any habit) systematically over-predict their gym attendance, we will interpret this as direct evidence that they are naive about their self-control problems. Finally, by offering small attendance incentives in some of the post-treatment weeks, we are able to manipulate subjects' weekly attendance directly. We are therefore able to estimate the costs and benefits associated with attendance, as well as to calibrate both the habit-formation effect and any systematic misprediction, in dollar terms.

We find a significant short-run habit formation effect among our subjects of 0.256 visits per week, which is smaller than, but statistically indistinguishable from, Charness and Gneezy's result. The effect appears to largely decay during the semester break, however, suggesting that this type of habit formation can be short-lived.<sup>5</sup> Moreover, the treatment effect is highly concentrated in the upper tail of the post-treatment attendance distribution. We find no evidence that subjects predicted this habit-formation effect overall, which we interpret as evidence of projection bias. We also find that all subjects substantially over-predict their future gym attendance in general: even in our simplest elicitation task, subjects over-predicted attendance by roughly a factor of three. We interpret this as evidence of naive present bias. Predictions are closer to actual attendance after the treatment period than before. By fixing the delay between the week in which predictions are made and the week about which they are made, we rule out intertemporal discounting as an explanation for this shift. The secular trend is instead interpreted as a general tendency among our subjects, at the start of the semester, to over-predict their level of free time later in the semester.

Finally, we estimate a structural model that extends the analysis beyond the reduced-form evaluation of the intervention effect. We are therefore able to estimate two generalizable behavioral parameters that have previously proven difficult to directly identify: the degree of projection bias

<sup>&</sup>lt;sup>5</sup>Indeed, Kane, Johnson, Town and Butler (2004) in a review find that monetary incentives are generally effective at generating short-run behavioral changes, but the literature typically does not find evidence of long-run effects that extend even as far as those we identify in this study. Our structural estimates will suggest that our intervention did not induce the steady-state level of habituation in subjects.

and the degree of present-bias naivete. We estimate that subjects over-estimate their willingness to attend the gym tomorrow by approximately \$3, and hence greatly over-estimate their likelihood of gym attendance. This is slightly larger than the value of the habit, which we estimate at \$2.60 and find in roughly 1/3 of treated subjects. In contrast, subjects only predict a habit value of less than \$0.25, which corresponds to a degree of projection bias  $\alpha = 0.9$ . This is considerably greater than the  $\alpha \in [0.31, 0.50]$  range found by Conlin, O'Donoghue and Vogelsang (2007) for cold-weather clothing catalog sales (except vests) and the  $\alpha \in [0.41, 0.49]$  range found by Levy (2009) for young smokers, although both sets of estimates lie within our 95% confidence interval. We cannot say without further research whether this difference is driven by differences in setting or differences in methodology, and we are open to the interpretation that the ability to understand changes in preferences is context-dependent.

Our second structural parameter of interest re-parameterizes naivete over present bias so that it follows a parallel construction to projection bias. While a present-biased agent currently would trade off utility in two future periods  $(\tau, \tau + s)$  at a marginal rate of substitution of  $\delta^s$ , in period  $\tau$ she will trade them off at a rate of  $\beta \cdot \delta^s$ . A partially naive agent believes their future marginal rate of substitution will be  $\hat{\beta} \cdot \delta^s$ , and we simply index her beliefs by  $\hat{\beta} = \omega \cdot 1 + (1 - \omega) \cdot \beta$ .<sup>6</sup> Thus  $\omega = 0$ corresponds to full sophistication, and  $\omega = 1$  to fully naive beliefs. By using the "commitment value" embedded in subjects' valuations for a contract that rewards future gym attendance, we are able to estimate a value of  $\omega = 0.666$ : subjects are two-thirds naive about their future self-control problems. If one uses the value  $\beta = 0.7$  typically found in other studies (DellaVigna 2009), this corresponds to  $\hat{\beta} = 0.9$ . Given the importance of naivete in the theoretical literature, it is surprising to note the lack of published estimates of  $\hat{\beta}$ . Skiba and Tobacman (2008), the only other estimate we could find, use a sample of payday loan borrowers to estimate an almost-identical  $\hat{\beta} = 0.9$ ; more work must be done, however, to confirm the regularity of this result.<sup>7</sup>

The remainder of this paper is organized as follows. Section two presents a simple model of habit formation which nests the rational-addiction model within the projection-bias framework.

<sup>&</sup>lt;sup>6</sup>It is tempting, but not correct, to infer that naivete in present bias is merely a case of projection bias where a subject's state is given by the current period. Our goal in introducing  $\omega$  is not to unify these two biases, but rather simply to provide a means of characterizing naivete in present bias that is independent of the underlying level of time-inconsistency.

<sup>&</sup>lt;sup>7</sup>Their estimate unfortunately comes alongside an atypically low  $\beta = 0.53$  in addition to an annual long-run discount factor  $\delta = 0.45$ , suggesting their sample put an non-representatively low weight on future consumption.

In section three we describe the experimental design, and in section four we present our results. Section five discusses our findings, and concludes.

## 2 Model

In this section we develop a simple model of gym attendance that incorporates habit formation, projection bias, and present-biased preferences. The experimental design described in Section 3 will then build on this to identify the key parameters of this model.<sup>8</sup> Consider a model with three periods. Initially all subjects are non-habituated, and are randomly assigned to a control or treated group.

At the beginning of the first period treated subjects are told that, at the end of the period, they will undergo a treatment which endows them with a gym-attendance habit. All subjects then give an incentive-compatible valuation of a "p-coupon": a contingent payment contract that rewards gym attendance during period three. They also give a direct, but unincentivized, prediction of how many times they would go to the gym in period three if they received the coupon. At the end of this period, subjects in the treated group are endowed with the habit.<sup>9</sup>

In the second period two things happen. First, subjects once again value the third-period p-coupon and predict their third-period attendance. Then, after the elicitation, all subjects are actually given a p-coupon about which they had no foreknowledge.<sup>10</sup> We assume there is sufficient time between periods two and three to act as a buffer, ensuring that subjects consider period three to be "in the future" when all predictions are elicited. In period three, subjects receive p-coupon rewards according to their contemporaneous gym attendance. There are seven opportunities to attend the gym during this period, corresponding to the seven days of a "target week" in our experiment. Finally, at the end of period three, subjects receive the delayed health benefit of

<sup>&</sup>lt;sup>8</sup>By presenting the model before the experimental design, we hope to underscore the extent to which our empirical strategy is designed to provide identification of key underlying economic parameters. By tightly linking the model to the field experiment, we are following in the "structural behavioral economics" tradition of, among others, DellaVigna, List and Malmendier (2009).

 $<sup>^{9}</sup>$ We are therefore modeling the habit value associated with receiving the treatment offer in our experiment. In practice, the offer will be a \$100 reward for attending the gym twice-weekly for one month. Approximately 80% of treated subjects met the 8-visit threshold for earning the reward, while none of the control subjects visited the gym 8 times during this same month.

<sup>&</sup>lt;sup>10</sup>In the model we are ignoring the fact that the elicitation process requires one or two subjects to wind up with two coupons. In practice, because there were multiple target weeks, most of the auction winners did not end up holding multiple p-coupons for the same week. The two subjects who did wind up with two p-coupons for the same target week simply received double the reward and are counted in the analysis as such.

whatever gym attendance they have engaged in.

We model utility as quasi-linear in money. Without loss of generality, utility from all non-gym sources will be normalized to zero. Let the immediate utility of gym attendance on day d be given by  $(-c + \varepsilon_d)$  with c > 0 and i.i.d.  $\varepsilon_d$ , and let the present value of the delayed benefits of gym attendance be b > 0. Thus we model gym attendance as an "investment good" in the language of DellaVigna and Malmendier (2004), meaning that costs are immediate while rewards are delayed. We abstract from the model of Becker and Murphy (1988) and O'Donoghue and Rabin (1999b) by modeling habituation as a binary state variable rather than a stock variable with geometric decay. While we do not explicitly model the habit formation process, we assume that a single gym visit or absence is not sufficient to change a subject's state.

When subjects are habituated they receive additional, immediate utility for gym attendance of  $\eta_i \geq 0$ , so that the immediate utility of gym attendance for a habituated subject is  $\eta_i - c + \varepsilon_d$ . To capture habit formation heterogeneity parsimoniously, the habit value will take one of two values. With probability  $\pi$ , a subject has  $\eta_i = \overline{\eta}$  strictly greater than zero, and with probability  $1 - \pi$ , they have  $\eta_i = 0$ . For simplicity, we present the case in which subjects know their own type ex ante, but neither our modeling nor empirical strategy distinguishes this from the case where subjects merely know  $\overline{\eta}$  and  $\pi$  ex ante, but not their own type.

On top of these basic preferences, we allow subjects to exhibit any degree (including none) of two psychological biases.<sup>11</sup> First, individuals may have a time-inconsistent taste for immediate gratification: they discount all future periods relative to the present. Such "present bias" has been observed in a wide range of contexts, from long-term savings behaviors (Angeletos, Laibson, Repetto, Tobacman and Weinberg 2001) to daily caloric intake (Shapiro 2005).<sup>12</sup> Following the formulation of Laibson (1997), the degree of present bias is captured by an extra discount factor  $\beta < 1$  applied uniformly to all future periods, in addition to a standard exponential discount factor  $\delta$ . Because the experimental period is brief we normalize the long-run discount factor  $\delta$  to one. Moreover, subjects may be naive about their present bias. Rather than assuming subjects use the correct  $\beta$ , we follow O'Donoghue and Rabin (1999a) and endow subjects with a belief that their future short-run discount factor will be given by  $\hat{\beta} \in [\beta, 1]$ . The lower and upper bounds on  $\hat{\beta}$  refer

<sup>&</sup>lt;sup>11</sup>Thus we embed the standard model of no bias, allowing us to simultaneously test the standard model and estimate behavioral parameters.

<sup>&</sup>lt;sup>12</sup>An overview of the literature, with many additional examples, is available in DellaVigna (2009).

to full sophistication and full naivete, respectively, while intermediate values correspond to partial naivete. We note that while the present bias itself may generate an under-investment in exercise relative to one's long-run preferences, naivete is necessary for subjects to hold systematically biased beliefs (including those that can lead to procrastination).

The second source of bias we consider is "projection bias", whereby subjects do not appreciate the extent to which their future preferences may differ from their current ones as a result of changes to their exercise-habit "state". In our setting, this implies that individuals will correctly foresee the direction of the habit-formation process, but may partially or fully "project" their current level of habit onto their future selves. Such projection appears in a wide variety of settings. For example, Read and van Leeuwen (1998) show that people who are currently hungry act as though their future selves will also be relatively hungry, and people who are currently sated act as through their future selves will also be relatively sated. Similar effects have been shown for the "endowment effect" (Loewenstein and Adler 1995), sexual arousal (Ariely and Loewenstein 2005), cold weather clothing catalog orders (Conlin et al. 2007), and drug addiction (Badger, Bickel, Giordano, Jacobs, Loewenstein and Marsch 2007). We model subjects as having "simple projection bias" as defined by Loewenstein, O'Donoghue and Rabin (2003), using  $\alpha \in [0,1]$  to index the strength of the bias. That is, when considering future consumption decisions, subjects believe that their future utility function will be an alpha-mixture of their current and future utility functions, with a weight of  $\alpha$  on the current utility function and  $1 - \alpha$  on the future utility function. Thus  $\alpha = 0$  refers to the case of no projection bias, in which subjects correctly foresee the actual future instantaneous utility function, and  $\alpha = 1$  refers to the case of full projection bias, in which subjects believe that their instantaneous utility function will not change with their state of habituation.<sup>13</sup>

We define a p-coupon as a contract that pays \$p concurrently with each day the holder attends the gym. Let  $V_{t,g}(p)$  refer to the valuation of a p-coupon during the elicitation session  $t \in \{pre, post\}$ (i.e. pre-treatment and post-treatment) of a subject in group  $g \in \{C, T\}$  (i.e. control or treated). Because we allow some subjects' habit to have a value of zero, we can treat group assignment as equivalent to habit assignment in the post-treatment period. Let  $Z_{d,g}^t(p)$  be an indicator for

<sup>&</sup>lt;sup>13</sup>Because of the linearity embedded in our model, the simple projection bias of Loewenstein, O'Donoghue and Rabin (2003) is indistinguishable from subjects systematically holding the incorrect belief that the value of the gym habit is  $(1 - \alpha)$  times the actual habit value. This ambiguity could in principle be resolved by a modification of our experiment that shocked people out of the habit rather than into it.

whether a subject in group g actually attends the gym on day in period t, so that, for example,  $Z_g^{post}(p) = \sum_{d=1}^{7} Z_{d,g}^t(p)$  is the number of gym visits during a given week post-treatment week for a subject in group q.

If a subject holding a \$p coupon attends the gym on a given day during the target week, her utility for that day will be  $p + \beta b + \mathbb{G}\eta_i - c + \varepsilon_d$ , where  $\mathbb{G} = \mathbb{1} \cdot \{g = T\}$  is simply an indicator for assignment to the treated group. She will attend the gym if this is greater than zero. Thus  $Z_{d,g}^{post}(p) = \mathbb{1} \cdot \{\varepsilon_d > p + \beta b + \mathbb{G}\eta_i - c\}$ , and  $Z_g^{post}(p) = \sum_{d=1}^7 \mathbb{1} \cdot \{\varepsilon_d > p + \beta b + \mathbb{G}\eta_i - c\}$ . Denote the unobserved distribution of the daily shock by  $F(\varepsilon)$ . In expectation, total target-week gymattendance will be,

(1) 
$$\sum_{d=1}^{7} \Pr\left(Z_{d,g}^{post}(p) = 1\right) = 7 \times \int_{c-\beta b - \mathbb{G}\eta_i - p}^{\infty} dF(\varepsilon)$$

and the habit-formation effect, the increase in attendance caused by habituation, will be,

(2) 
$$\sum_{d=1}^{7} \Pr\left(Z_{d,g}^{post}(p) = 1\right) - \sum_{d=1}^{7} \Pr\left(Z_{d,g}^{pre}(p) = 1\right) = 7 \times \int_{c-\beta b - \mathbb{G}\eta_i - p}^{c-\beta b - p} dF(\varepsilon)$$

Equation (2) highlights both a simple reduced-form test of habit formation, and an intuitive means to calibrate the strength of the effect. Any relative change in attendance between control and treated subjects' attendance is directly captured by  $\eta$ . That is, any test for the significance or sign of  $\eta$  — noting that it is not constrained to be positive — is equivalent to testing the difference-indifferences in attendance for control and treated subjects. Moreover, the quasi-linearity of money means that the external subsidy that equates control subjects' attendance with that of unincentivized treated subjects is directly interpretable as the utility equivalent of habit formation.

We now turn from subjects' actual attendance to their predictions. The perceived probability of target-week gym-attendance, from the perspective of any previous period, depends upon the subject's belief about future self-control,  $\hat{\beta}$  and on her projection-bias parameter,  $\alpha$ . She believes she will attend on any given day of the target week if  $\varepsilon_d > p + \hat{\beta}b + \mathbb{G}(1-\alpha)\eta_i - c$ , and therefore her expected attendance for a future week is given by:

(3) 
$$7 \times \int_{c-\widehat{\beta}b-\mathbb{G}(1-\alpha)\eta_i-p}^{\infty} dF(\varepsilon)$$

Equation (3) is the key to identifying present bias. For simplicity, consider a subject's prediction in period 2. The only difference between their prediction in (3) and their actual behavior in (1) is the use of a possibly-biased  $\hat{\beta}$  in the former and their actual discount factor  $\beta$  in the latter.<sup>14</sup> This suggests a simple test of naive present bias: if subjects' predictions are systematically greater than their actual attendance, we can simultaneously reject both time-consistency ( $\beta = \hat{\beta} = 1$ ) and complete sophistication ( $\hat{\beta} = \beta$ ).

A subject's ex-ante prediction of her total utility for the target-week, given that she holds a p-coupon, is,

(4) 
$$7 \times \beta \times \int_{c-\widehat{\beta}b-\mathbb{G}(1-\alpha)\eta_i-p}^{\infty} (p+b+\mathbb{G}(1-\alpha)\eta_i-c+\varepsilon) dF(\varepsilon)$$

The lower limit of integration reflects the perceived probability of a future self going to the gym, and therefore uses the perceived future discount factor  $\hat{\beta}$ . Because the utility flows during that target week and the delayed benefits of attendance are both in the future from the point of view of the current period, however, we use her actual discount factor  $\beta$  to discount the entire expression. Setting p = 0 gives us the perceived utility without a coupon. The perceived value of the p-coupon is simply the difference between expected utility with a p-coupon and expected utility without. To compare this with a payment delayed until the target week, we eliminate the  $\beta$  outside the integral. In the pre-treatment elicitation session, the perceived value of a p-coupon therefore is:

(5)

$$V_{0,g}(p) = \left[7 \times \int_{c-\widehat{\beta}b-\mathbb{G}(1-\alpha)\eta_i-p}^{\infty} p \, dF(\varepsilon)\right] + \left[7 \times \int_{c-\widehat{\beta}b-\mathbb{G}(1-\alpha)\eta_i-p}^{c-\widehat{\beta}b-\mathbb{G}(1-\alpha)\eta_i} (b + \mathbb{G}(1-\alpha)\eta_i - c + \varepsilon) \, dF(\varepsilon)\right].$$

<sup>&</sup>lt;sup>14</sup>We note that a partially naive present-biased agent uses this  $\hat{\beta}$  to predict how her future self will *behave*, but would use her actual discount factor  $\beta$  to evaluate the utility *resulting* from any choices her future self makes as in (4) below.

And in the post-treatment elicitation session, when the full habit-formation effect is known to the subject, it is:

(6) 
$$V_{1,g}(p) = \left[7 \times \int_{c-\widehat{\beta}b-\mathbb{G}\eta_i-p}^{\infty} p \, dF(\varepsilon)\right] + \left[7 \times \int_{c-\widehat{\beta}b-\mathbb{G}\eta_i-p}^{c-\widehat{\beta}b-\mathbb{G}\eta_i} (b+\mathbb{G}\eta_i-c+\varepsilon) \, dF(\varepsilon)\right].$$

The first term in both (5) and (6) is the expected redemption value of the coupon, which is always weakly positive. We note that the face value incorporates an estimate of the behavioral response to the p-coupon's implicit subsidy. The second term is the subject's valuation of the behavioral change that results from holding the coupon, which we will call the "commitment value". This is the change in *utility* caused by those gym-visits that the subject would not have made in the absence of the p-coupon. The sign depends on the subject's ex-ante belief about future self-control problems. If the subject believes that she will not have self-control problems in the target week then the commitment value is negative because the subject believes that the p-coupon will make her attend the gym at times when she would ex-ante prefer not to. If the subject believes that she will have self-control problems in the target week then the commitment value may be positive because she foresees that the p-coupon will make her more likely to attend the gym and gain a long-term benefit that she would otherwise have foregone.

The total ex-ante value of the p-coupon is always non-negative. Intuitively, the endogeneity of the attendance decision gives positive option value to even a p-coupon one does not expect to use. We prove in appendix A.1 that this holds for subjects with both present-biased and projectionbiased preferences.

We next turn to our reduced-form test for projection bias. Suppose we compare the differencein-differences in valuations of a p-coupon for the treated and control subjects, before and after the intervention. For subjects with no projection bias, regardless of their level of present bias, we would expect this difference-in-differences to be zero. Dividing by 7 for simplicity, this double difference is given by:

$$(7) \frac{[\bar{V}_{1,T} - \bar{V}_{0,T}] - [\bar{V}_{1,C} - \bar{V}_{0,C}]}{7} = \pi \cdot \int_{c-\hat{\beta}b-\bar{\eta}-p}^{c-\hat{\beta}b-(1-\alpha)\bar{\eta}-p} p \, dF(\varepsilon) + \pi \cdot \left[\int_{c-\hat{\beta}b-\bar{\eta}-p}^{c-\hat{\beta}b-\bar{\eta}-p} (b+\bar{\eta}-c+\varepsilon) \, dF(\varepsilon) - \int_{c-\hat{\beta}b-(1-\alpha)\bar{\eta}-p}^{c-\hat{\beta}b-(1-\alpha)\bar{\eta}} (b+(1-\alpha)\bar{\eta}-c+\varepsilon) \, dF(\varepsilon)\right]$$

The first term in (7) is the more intuitive one, namely the misprediction of gym attendance caused by projection bias. It is weakly positive for projection-biased agents and zero for agents without projection bias, regardless of self-control problems. The latter two terms, however, reflect the difference in perceived incentive value from before the treatment to after. For any given value of  $\hat{\beta}$ , the difference in incentive values depends exclusively on the distribution of  $\varepsilon_d$ , and could take on either positive or negative values.<sup>15</sup> The overall effect is therefore ambiguous, but still provides a test inasmuch as any observed value of the double-difference that is significantly different from zero indicates projection bias. A more powerful test uses only unincentivized predictions, and therefore reflects only the first term of (7) without incorporating the change in commitment value. We present both in Section 4.

### 3 Design

We recruited one hundred and twenty subjects from the students and staff of UC Berkeley and randomly assigned them to treated and control groups.<sup>16</sup> Since we want to move subjects from an unhabituated state to a habituated state, we screened for subjects who self-reported that they had not ever regularly attended any fitness facility.<sup>17</sup> Treated and control subjects met in separate

<sup>&</sup>lt;sup>15</sup>This result follows from the unknown density  $F(\varepsilon)$ , which may differ in the regions around the threshold values of  $\varepsilon$  for the actual habit and the projection-biased predicted habit. If, for example, a subject believed that the coupon would make her much more likely to attend the gym if the habit were  $(1 - \alpha)\overline{\eta}$  but would be entirely infra-marginal with a habit of  $\overline{\eta}$ , this latter difference would be negative. The opposite case would yield a positive difference. A noteworthy special case occurs when  $\varepsilon$  has constant density over the full region traced by (7). In this case, there is no relative change in the perceived commitment value and the term in brackets reduces to zero.

<sup>&</sup>lt;sup>16</sup>Due to attrition and missing covariates, our final sample includes 54 treated subjects and 57 control subjects. Details of the sample appear in appendix A.2.

<sup>&</sup>lt;sup>17</sup>Our screening mechanism is described in appendix A.3. We expected that using a sample of students who did not regularly attend the gym would increase the power of our habit-formation test, as Charness and Gneezy found

sessions on the same day, at the beginning of the second week of the fall semester. Both treatment and control subjects were asked to complete a questionnaire, and were then given an offer of \$25 to attend the gym once during the following week.<sup>18</sup> We call this the "learning week" offer, and it is identical to Charness and Gneezy's low-incentive condition. Our control group is therefore comparable to Charness and Gneezy's low-incentive group. We chose this as our control in order to separate the effect of overcoming the one-time fixed cost of learning about the gym from the actual habit formation that occurs after multiple visits.<sup>19</sup>

At the same initial meeting, the treatment group received an additional offer of \$100 to attend the gym twice a week in each of the four weeks following the learning week. We call this the treatment-month offer, and it is the same as Charness and Gneezy's high-incentive offer, except that they did not require the eight visits to be evenly spaced across the four weeks. This difference was intended to limit the potential for procrastination so that naive present-biased subjects in the treated group would be more likely to meet the eight-visit threshold, although our compliance rate was not distinguishable from the less-restrictive design. The other difference between this offer and Charness and Gneezy's high-incentive offer is that we made our offer at the first meeting, at the same time as the \$25 learning-week offer, whereas Charness and Gneezy made their high-incentive offer at their second meeting, a week later. We made our treatment-month offer earlier because we wanted treated subjects to have a week to contemplate the idea of going to the gym twice weekly for a month before making predictions. Moreover, if subjects have reference-dependent preferences for money then suddenly announcing a gain of \$100 to one group but not the other could introduce systematic bias into the incentive compatible procedure we used to elicit predictions. Waiting a week after treatment subjects learn they will earn \$100 will help us overcome a potential "house money effect" on predictions, and should not affect actual gym attendance decisions.

At the end of the learning week both groups of subjects again met separately and completed pencil-and-paper tasks (described in detail below) designed to elicit their predictions of gym attendance during each of five post-treatment "target weeks". Both groups were reminded of the offers

the treatment effect to be concentrated among subjects whose average pre-treatment attendance was less than 1 visit per week.

<sup>&</sup>lt;sup>18</sup>For this and all subsequent offers, subjects were told that a visit needed to involve at least 30 minutes of some kind of physical activity at the gym. We were not able to observe actual behavior at the gym and did not claim that we would be monitoring activity.

<sup>&</sup>lt;sup>19</sup>We also paid the \$10 gym-membership fee for all students, and filed the necessary membership forms for those who were not already members.

they had received. Four weeks later, at the end of the treatment month, both groups again met separately, completed an additional questionnaire, and completed the same elicitation tasks as in the second session. The target weeks were separated from this second elicitation session by one week, to ensure that present-biased subjects would see the target weeks as being "in the future" from the perspective of both elicitation sessions. The timeline of the experiment is illustrated in Figure 1.

Gym attendance data were collected for a 17-month period stretching from 37 weeks before the learning week to 33 weeks after it. This period includes summer and winter breaks as well as three full semesters. This attendance data was based on recorded ID card swipes required for gym entry.<sup>20</sup>

#### 3.1 Elicitation procedures

To elicit predictions of target-week gym attendance we created what we call a "p-coupon", which is a contingent payment contract that rewards the holder with p for each day that he or she attends the gym during a specified "target week". The value of p, which ranged from 1 to 7, was printed on the coupon, along with the beginning and end dates of the target-week. A subject with a 7 coupon, for example, could potentially earn up to 35 during the target week. We used an incentive-compatible multiple price listing mechanism to elicit subjects' valuations for p-coupons of various values with various target weeks.<sup>21</sup> A sample p-coupon is included in Appendix A.4, along with the pencil-and-paper task we used to elicit valuations for p-coupons, the instructions we gave them for completing the task, and further description of how the elicitation mechanism worked. Each subject completed this incentive-compatible elicitation task for four of the five target weeks in our design, and for a different value of p-coupon in each of those four weeks. The values of the p-coupons for the different weeks was randomized among subjects, as was the order in which those

<sup>&</sup>lt;sup>20</sup>Because swipes were necessary to enter the gym but not exit, we cannot always determine the length of a visit. In some cases where a person swiped multiple times, e.g. to enter the locker room and the pool, we can form a lower bound on the length of the visit. We acknowledge that many of the recorded swiped during the treatment month may be subjects swiping to receive the reward but not exercising. However, there is no reason to continue to engage in such false visit behavior in the post-treatment period, when the incentives to attend are removed. To the extent that some people did swipe without exercising, our estimates in Section 4 reflect a lower bound on the treatment effect of the program. One may also interpret some proportion of the subjects who do not form a habit (i.e. for whom  $\eta_i = 0$ ) as being in the category of false swipers.

<sup>&</sup>lt;sup>21</sup>Subjects made a series of choices between a p-coupon and an incrementally increasing fixed amount of money. We infer their valuation from the indifference point between the coupon and the fixed sum. The elicitation mechanism is described in detail in Appendix A.4.

weeks were presented.<sup>22</sup> Because at most one p-coupon would be awarded as part of the elicitation, subjects' valuations are not confounded by uncertainty about how many p-coupons they would receive.

Subjects' willingness to pay for a coupon that pays out as a function of their future behavior will of course not be based entirely on their underlying predictions. Risk-aversion alone implies we would only observe subjects' certainty equivalents, even for an exogenous event.<sup>23</sup> But for an endogenous event like gym attendance, there is the additional confound that the p-coupon itself incentivizes the subject to go to the gym, thus influencing the very behavior we are asking them to predict. There is an important distinction to be made between the effect of the incentive on the underlying behavior – the "incentive effect" – and the distortion the incentive induces for a subject bidding on their own p-coupon. The incentive effect alone would still allow for a direct comparison between p-coupon valuations and attendance, as both incorporate the behavioral response equally. The latter effect, which we call the "commitment value", drives a wedge between the face value of a coupon and a subject's personal value from holding it and complicates such a comparison. The direction of the distortion depends on the type of present bias a subject experiences. In general, time-consistent and naive subjects will view the incentive as creating a costly distortion and lower their valuations relative to the expected face value. Sophisticated subjects, however, will value the p-coupon as a commitment device for their future selves, and will raise their valuation over the expected face value. We use this later case in naming the effect, but note that this "commitment value" may increase or decrease subjects' demand for a p-coupon, and care must be taken not to interpret subjects' valuations as directly proportional to their beliefs.

As a check on this mechanism, we also directly asked subjects to state how many times they thought they would go to the gym during the specified target weeks if they had been given the p-coupon they just bid on in the incentive-compatible task. Thus they were making unincentivized *predictions* of hypothetical future attendance under the same set of *attendance* incentives as in the incentivized task.<sup>24</sup> This unincentivized mechanism also allowed us to ask subjects how often they

 $<sup>^{22}</sup>$ Among each subject-group/target-week intersection, subgroups of fifteen subjects received \$1, \$2, and \$3 coupons, ten received \$5 coupons, and five received \$7 coupons.

 $<sup>^{23}</sup>$ An alternative design which would have allowed us to sidestep assumptions about the linearity of money utility, would have been to have the coupons pay off not with a dollar sum per visit, but with a per-visit increment in the cumulative probability of winning some fixed-sum prize. We believe our design is more intuitive for subjects, and easier for them to understand.

<sup>&</sup>lt;sup>24</sup>It is important to note that the p-coupons incentivize both target-week attendance and accurate predictions of

thought they would go to the gym during the one target week for which they were not presented with a p-coupon, the so-called "zero week" (because it is equivalent to a p of zero). The zero week gives us an additional unincentivized prediction of behavior in the absence of any effect of attendance incentives. More importantly, by comparing the unincentivized prediction with the valuation of the corresponding p-coupon, we have an estimate of the commitment value provided by that coupon. We use this difference in one of the structural model specifications to calibrate the extent of naivete, and thereby leverage what would otherwise be a confound to provide additional identification.

Subjects went through the same set of elicitation tasks in both the pre-treatment and posttreatment elicitation sessions. Then, at the end of the second elicitation session, after all of the elicitation tasks had been completed, each subject was given one of the four coupons they had been presented with during the elicitation process.<sup>25</sup> We therefore have two target weeks for each subject in which we can compare their predictions with their actual gym attendance under the same conditions, the first being the zero-week, and the second being the week for which they received a p-coupon in the giveaway. The giveaway was a surprise to the subjects—having been conducted unannounced only after the second elicitation session—and thus did not affect their bids or unincentivized responses during the elicitation tasks. We discuss compliance with the treatment incentive, attrition, and our randomization procedure in Appendix A.5.

# 4 Results

Of the 54 subjects in our final treatment sample, 43 completed the eight necessary bi-weekly visits in order to earn the \$100 incentive: a compliance rate of 80%. In Charness and Gneezy's (2009) high-incentive group the compliance rate was approximately 83%, suggesting that our more restrictive twice-weekly attendance requirement did not have a significant effect on subjects' ability to make the required number of visits.

target-week attendance.

 $<sup>^{25}</sup>$ We used a block-random design to assign coupons to 12 control and 12 treated subjects in each of the five target weeks. Within each treatment group, 15 subjects received a \$1 coupon, 15 received a \$2 coupon, 15 received a \$3 coupon, 10 received a \$5 coupon, and 5 received a \$7 coupon.

#### 4.1 Habit formation

Figure 2 shows average weekly attendance for the treated and control groups over the duration of the study period.<sup>26</sup> In the pre-treatment period, attendance in the two groups moves together tightly. In the treatment period, treated subjects attend much more than control subjects. In the two months immediately following the treatment period, leading up to, but not including the semester break, the treatment group consistently attends the gym more than the control group. In the four months after the semester break the graph suggests persistence of the increased treatment-group attendance, but the difference is not as striking.

We estimate a linear difference-in-differences panel regression model to determine if these patterns are statistically significant. Each observation in the panel is a specific individual on a specific week of the study, and we therefore cluster all standard errors throughout by subject.<sup>27</sup> We regress weekly gym attendance on a treated-group dummy, week-of-study dummies, and the interactions of the treated-group dummy with dummies for the treatment period and each of the two posttreatment periods. We also control for individual characteristics, including demographics and demand shifters such as travel time to the gym.<sup>28</sup> The results of this regression appear in the first column of Table 1.

The coefficient on the treated-group dummy indicates no statistically significant difference in pre-treatment gym attendance between treated and control subjects. The coefficient on the interaction of the treated-group and treatment-period dummies, roughly the product of the twice-weekly incentive target and the 80% compliance rate, reiterates that the treatment-incentive was effective. The remaining two interaction terms tell us the effect of the treatment on treated-group attendance in the two post-treatment periods. The point estimate is 0.256 additional visits per week for the immediate post-treatment period, representing approximately a doubling of average attendance in our sample. While significantly different from zero, it is also misleadingly small as it combines a large mass of unaffected subjects with a smaller mass of subjects for whom the habit formation

 $<sup>^{26}</sup>$ We have removed observations for target weeks when subjects received p-coupons to make the graph easier to interpret.

<sup>&</sup>lt;sup>27</sup>We again exclude observations for the one target week for each subject for which they received an actual p-coupon. <sup>28</sup>When we omit the individual characteristics, the main effect is no longer significant at standard levels. A Hausman test between the two specifications obtains a p-value of 0.051, suggesting that we may be correcting for some lumpiness in our randomization. Following Gelbach (2009), we find that demographic covariates, naivete proxies, and attitudes about gym attendance explain three-quarters of the change in the coefficient, but do not enter significantly.

appears substantial. We will address this heterogeneity below, and in our structural model.

Because not all subjects in the treatment group made the requisite eight visits to the gym, the results in the first column represent the "intention to treat" effect, or ITT. To see the effect on those who complied with the treatment we instrument for compliance with the treated-group dummy, including our vector of individual covariates in the first stage. This gives us the average "treatment effect on the treated", or ATT, controlling for observable differences between compliers and non-compliers. This analysis assumes there is no effect on subjects who did not meet the 8-visit threshold, which is not implausible given the average of only two visits during the treatment period for such subjects. These results are reported in the second column of Table 1. Not surprisingly, the ATT is larger than the ITT. We now see an increase in immediate post-treatment gym attendance for the treated-group of a third of a visit per week.

In the later post-treatment period we see no statistically significant increase, despite the apparent visual difference between treated and control attendance in Figure 2. We estimate a later post-treatment ITT effect of 0.045 visits per week, and an ATT effect of 0.061. Neither effect is significant, suggesting that the habit induced by four weeks of exogenous gym attendance was dissipated by a similar period of quasi-exogenous non-attendance. It is noting that this long-run decay supports our interpretation of the short-run effect as habit formation over alternatives such as learning, for which one would not expect to find decay.

To compare our results with the results from Charness and Gneezy's first study we ran the same regression on their data, the results of which constitute the final column of Table 1.<sup>29</sup> The double difference in average weekly attendance between their high-incentive and low-incentive subjects in the immediate post-treatment period was 0.585 visits per week. Stacking their data with ours allows us to conduct a Chow test of the equality of their habit-formation coefficient with the one in our column-one specification. The p-value, reported in square brackets, is 0.186. Thus we cannot reject that the habit-formation effect in our sample was the same as the habit-formation effect in their sample.

To get a better picture of the heterogeneity of treatment effect in the immediate post-treatment period, we can also compare the empirical CDFs of average post-treatment attendance in the treated

<sup>&</sup>lt;sup>29</sup>This specification differs from the one they report, which uses pre- and post-treatment averages rather than the full panel of weeks.

and control groups. Following Frandsen (2010), we can estimate exact finite-sample confidence intervals for the distributions of post-treatment attendance for compliers. Figure 5 shows the 90% confidence intervals for immediate post-treatment attendance, by treatment assignment.<sup>30</sup> We find a large mass of subjects—roughly 60%—did not go to the gym prior to the experiment and did not go to the gym after. The treatment effect is instead concentrated in the upper portion of the distribution, with a significant difference above the 85th percentile. A quantile regression (not shown) confirms a significant treatment effect over this range.

Another straightforward test is to compare individual post-treatment attendance against their predicted attendance were they to be in the control group. We imputed this counterfactual based on a regression of attendance on week dummies and covariates using control group data for all weeks and treated group data for the pre-treatment period, and while it clearly includes noise (for both controls and treated subjects) it is the best prediction of post-treatment attendance in the absence of any intervention. Similar to Charness and Gneezy, we identify as "habit formers" those subjects in each group for whom average attendance in the immediate post-treatment period was at least one visit per week greater than their predicted attendance. Shown in Figure 4, this applies to 8 of 54 treated subjects and 3 of 57 control subjects, the latter serving as an estimate of false positives due to noise. A one-sided test of equal proportions rejects the null that there are more habit formers in the control group at a p-value of 0.046.<sup>31</sup>

It is not surprising that we find heterogeneity in our treatment effect. One possibility, which we cannot fully address, is that some subjects in the treated group swiped their ID cards at the gym but did not actually exercise. We would not expect such subjects to form any habit, and our estimates of the treatment effect in Table 1 would be biased towards zero by their presence. An alternative explanation would be that some subjects would have formed a habit, but our month-long treatment was too short for them to do so. This interpretation is consistent with recent findings such as Lally, van Jaarsveld, Potts and Wardle (2010), who estimate a range of 18 to 254 days in their subjects' time for habit formation for various tasks. Finally, it is possible that some subjects

<sup>&</sup>lt;sup>30</sup>That is, we compare the distribution of post-treatment attendance among compliers in the treatment group with that of compliers assigned to the control group. Given that no control subject met the 8-visit threshold during the treatment month, we restrict the proportion of always-takers to be zero and therefore estimate a population of "compliers" and "never-takers".

 $<sup>^{31}</sup>$ Relaxing the threshold for habit formation to 0.5 visits/week, which introduces more noise, also yields a significant result at a p-value of 0.066.

simply do not find exercising at the gym to be sufficiently habit-forming to alter their behavior.

#### 4.2 Predictions

#### 4.2.1 Naivete over self-control

We next turn our attention to subjects' predictions. Figure 5 shows predicted versus actual gym attendance, first for the weeks that subjects actually received a p-coupon in the giveaway at the end of the experiment, and then for weeks when no p-coupon was offered—so-called "zero-weeks". The two panels break the subjects into control and treated groups. Within each group we separate observations into p-coupon weeks and zero-weeks.<sup>32</sup> Finally, we separate subjects predictions by when they were elicited. We show only subjects' unincentivized predictions for clarity, but Tables 2 and 3 confirm that incentivized and unincentivized predictions follow similar patterns. We find in general that the valuations of the p-coupons are less than would be implied by the unincentivized predictions, indicating that the commitment value of the p-coupons is typically negative for subjects.

In both the pre- and post-treatment elicitation sessions, both the treated and control groups predicted future gym attendance that substantially exceeds their actual gym attendance. This pattern holds for both p-coupon weeks and zero-weeks. Furthermore, introducing a p-coupon seems to increase both actual and predicted attendance, as we would expect. Finally, there is a consistent pattern of less over-prediction in the later elicitation session.

Table 2 tests differences between predicted and actual attendance for the different groups and elicitation sessions, pooled over values of the p-coupon. The first column of each panel looks at predictions as captured by subjects' p-coupon valuations.<sup>33</sup> The second and third refer to their unincentivized predictions, for p-coupon weeks and zero-weeks. In all cases subjects significantly over-predict future gym attendance, by as much as two visits per week. It is particularly striking that subjects substantially over-predict gym attendance in weeks with no p-coupon, suggesting that the overprediction is not driven by the p-coupon incentives. On the basis of these results we can

 $<sup>^{32}</sup>$ We group all non-zero values of p-coupon together here for simplicity — the effect of each separate p-coupon value is shown in Table 3.

 $<sup>^{33}</sup>$ We include subjects' valuations for a p-coupon divided by the subsidy, which clearly does not account for the commitment value of the p-coupon. While we will exploit this in Section 4.3, we present the simple inferences here and caution readers not to interpret the valuations literally as predictions. For this section, we prefer the un-incentivized predictions which are directly measuring predicted attendance.

rule out, in our model, both time consistency ( $\beta = 1$ ) and full sophistication ( $\hat{\beta} = \beta$ ) if, after the treatment, subjects have rational expectations over their future costs.<sup>34</sup> Furthermore, while it does not rule them out, this result indicates that our data cannot be explained solely by other models of self control problems such as the "temptation utility" of Gul and Pesendorfer (2001, 2004) which also embed rational expectations.

In Table 3 we explore the effect of p-coupon value, and the change in predictions over time. The first column regresses actual attendance on dummies for the various values of p-coupon.<sup>35</sup> The point estimates on the p-value dummies indicate a nearly monotonic effect of monetary incentives, and pairwise comparisons of the coefficients do not reject monotonicity. This is reassuring, as it indicates the upward-sloping supply curve for exercise that we would expect. The second and third columns regress normalized coupon valuations and un-incentivized predictions on the same p-coupon dummies, plus a dummy for the post-treatment elicitation session. Subjects appear to predict the slope of their labor-supply curve relatively accurately, despite consistently over-predicting its intercept.

The session dummy implies that, between the first and second elicitation sessions, subjects reduce their predictions by roughly two-thirds of a visit per week. These sessions differ in two ways: they are a month apart in time, and the second session is closer to the target weeks than the first. One possibility is that subjects' discount factors decrease smoothly over time rather than abruptly as in the beta-delta model. If so, we would see a change in mispredictions merely because the temporal proximity of the target weeks is greater in the post-treatment elicitation session.<sup>36</sup> We can examine this by comparing first-session predictions for the first target week with second-session predictions for the fifth target week. This comparison holds temporal proximity constant. Columns (4) and (5) report the results of this regression. The coefficients on the session dummy, for both coupon valuations and unincentivized predictions, still show a substantial decrease in over-prediction over time. Such a secular drift in misprediction suggests that subjects may begin the

 $<sup>^{34}</sup>$ These two cases require rational expectations, which appears to be strongly violated here. It is possible that other un-modeled biases could cause the mis-prediction in Table 2, but we feel that many such alternatives would, unlike naive present bias, be ad hoc for this result.

<sup>&</sup>lt;sup>35</sup>The omitted category is p = \$7 throughout this table. This is so that we can compare coefficients across 'Actual' and 'Un-incentivized' (for each of which the lowest value is p = \$0), and 'Coupon Value' (where the lowest value is p = \$1). In addition, all specifications in this table include individual covariates.

<sup>&</sup>lt;sup>36</sup>This would be the case, for example, if the correct model of present bias was a true hyperbolic model rather than the familiar, quasi-hyperbolic  $\beta$ ,  $\delta$  model.

semester with overly optimistic beliefs about their amount of free time, and grow more realistic.<sup>37</sup> It is worth noting that typical models of learning would suggest that treated subjects, who have more experience with the gym, should by at least as much as control subjects do. Differences in this learning would therefore bias against our test of projection bias, which uses the relative increase in predicted attendance among treated subjects.

#### 4.2.2 Projection bias

Our test for projection bias is implemented with a difference-in-differences regression of predictions, shown in Table 4. While the GMM estimator in the next section will accommodate the complex effect of projection bias at different values of p-coupons, we restrict our attention here to the simplest cases. The first column utilizes unincentivized predictions about unsubsidized weeks, so as not to conflate mistakes about habit formation with mistakes about responsiveness to gymattendance subsidies. We find that treated subjects revise their predictions upwards by 0.458 visits per week relative to control subjects. This revision is in fact larger than the treatment effect itself, although not significantly. We therefore interpret this as consistent with full projection bias. Indeed, this double-difference exceeding the actual habit effect could be used as a test of projection bias against certain informational explanations if post-treatment subjects are temporarily above their steady-state level of habituation.

We replicate the test using subjects' valuations of p-coupons in the second column of Table 4. Because no incentives can be obtained for unsubsidized weeks, we instead look at predictions for the smallest subsidy (p=\$1), for which any distortion is minimized. We find a similar upward revision among treated subjects of 0.325 visits per week, but the coefficient is no longer significant. Because the change in perceived commitment value may offset learning about one's habit, however, we prefer the first column as a more powerful test of projection bias. We take this commitment value into account in the following section, however, and obtain a precise and highly significant estimate of projection bias.

<sup>&</sup>lt;sup>37</sup>See, e.g. Bénabou and Tirole (2002) for why subjects may begin the semester with overly optimistic beliefs.

#### 4.3 Structural Estimation

We next develop a vector of moment restrictions derived from our structural model and estimate the parameters of the model using the generalized method of moments. Let  $\overline{Z}_{t,g}(p)$  denote the average attendance of group g during period  $t \in \{pre, post\}$ , when holding coupon p. Let  $g \in \{C, T\}$  denote control and treated subjects, respectively, and  $\mathcal{T} = \{i|g_i = T\}$  denote the set of treated subjects. Let  $\overline{Y}_{t,g}(p)$  correspond to the analogous average unincentivized predictions for group g during the pre- and post-treatment elicitation sessions. Let the corresponding observations for an individual i be denoted by  $Z_t^i(p)$  and  $Y_t^i(p)$ . For any individual i, let  $p_{i,w}$  denote the wth p-coupon that she was offered in the elicitation session ( $w \in \{1, 2, 3, 4\}$ ) and let  $V_{w,i}^t(p_{i,w})$  denote her valuations of these coupons. Then we use the following moment restrictions:

(8) 
$$\bar{Z}_{post,T}(0) = 7 \cdot \left[ \pi \left( 1 - F \left( c - \beta b - \eta; \sigma_{\varepsilon} \right) \right) + (1 - \pi) \left( 1 - F \left( c - \beta b; \sigma_{\varepsilon} \right) \right) \right]$$

(9) 
$$\bar{Z}_{post,C}(0) = 7 \cdot \left[1 - F(c - \beta b; \sigma_{\varepsilon})\right]$$

(10) 
$$\bar{Z}_{post,T}(p>0) = 7 \cdot \left[\pi \sum_{i \in \mathcal{T}} \sum_{w=1}^{4} (1 - F(c - \beta b - \eta - p_{iw}; \sigma_{\varepsilon}))\right]$$

$$+ (1-\pi) \sum_{i \in \mathcal{T}} \sum_{w=1}^{4} (1 - F(c - \beta b - p_{iw}; \sigma_{\varepsilon})) \Big]$$

(11) 
$$\bar{Z}_{pre,C\cup T}(0) = 7 \cdot \left[1 - F(c - \beta b; \sigma_{\varepsilon})\right]$$

(12) 
$$\bar{Y}_{pre,T}(0) = 7 \cdot \left[ \pi \cdot (1 - F(c - \widehat{\beta}b - (1 - \alpha)\eta; \sigma_{\varepsilon})) + (1 - \pi) \cdot (1 - F(c - \widehat{\beta}b; \sigma_{\varepsilon})) \right]$$

(13) 
$$\bar{Y}_{pre,C}(0) = 7 \cdot \left[1 - F(c - \hat{\beta}b; \sigma_{\varepsilon})\right]$$

(14) 
$$\bar{Y}_{post,C}(0) = 7 \cdot \left[1 - F(c - \hat{\beta}b; \sigma_{\varepsilon})\right]$$

(15) 
$$\sum_{i \in \mathcal{T}} \mathbb{1}\{\bar{Z}_{pre,i}(0) < \bar{Z}_{post,i}(0)\} = \sum_{i \in \mathcal{T}} \left(\pi + \frac{1}{2}(1-\pi)\right)$$

The habit-formation effect itself is identified by equations (8) and (9), relying on the assumption that the habit is the only systematic driver of post-treatment differences in attendance among control and treated subjects in unincentivized weeks. In these and all the subsequent restrictions, we rely on our experimental randomization and do not condition on observable characteristics of the subjects. The rest of the attendance parameters are identified by (10) and (11), which establish the responsiveness to coupons and the pre-treatment baseline. Equation (12) mirrors (8), but uses pre-treatment predictions of attendance rather than actual attendance in order to identify the ex-ante expectations of the habit value. Finally, in (13) and (14), we use control subjects' pre- and post-treatment predictions to identify the general over-confidence driven by naivete about self-control. For these three sets of expectations, we use unincentivized weeks to avoid embedding an assumption that subjects correctly predict their response to small monetary incentives.

We introduce (15) to estimate the fraction of subjects developing a strictly positive habit. As the number of pre- and post-treatment periods grows large, the probability that a subject with a positive habit will have higher average attendance in the post-treatment period converges to 1. The corresponding probability for a subject who did not form a habit converges to 0.5. We use this limit as our moment for estimation, but note that (15) gives a conservative estimate of  $\pi$  due to the finite pre- and post-period samples. Calibrations at the estimated coefficients suggest the approximation is good.

In addition, in Model 2 we use the difference between coupon valuations and the valuation implicit in the unincentivized predictions to gain an additional moment. Letting  $\mathbb{G}_i = \mathbb{1} \cdot \{g_i = T\}$ and  $\mathbb{T} = \mathbb{1} \cdot \{t = post\}$  for convenience, we write:

$$\sum_{t \in \{pre, post\}} \sum_{i} \frac{1}{4} \sum_{w=1}^{4} \left( V_{w,i}^{t}(p_{i,w}) - Y_{t}^{i}(p_{i,w}) \cdot p_{i,w} \right) =$$

$$(16) \qquad \sum_{t} \sum_{i} \frac{1}{4} \sum_{w=1}^{4} 7 \cdot \left[ \pi \int_{c-\widehat{\beta}b-\mathbb{G}_{i}(1-\alpha(1-\mathbb{T}))\eta}^{c-\widehat{\beta}b-\mathbb{G}_{i}(1-\alpha(1-\mathbb{T}))\eta} (b + \mathbb{G}_{i}(1-\alpha(1-\mathbb{T}))\eta - c + \varepsilon) dF(\varepsilon) + (1-\pi) \int_{c-\widehat{\beta}b-p_{i,w}}^{c-\widehat{\beta}b} (b - c + \varepsilon) dF(\varepsilon) \right]$$

While Equation (16) may appear complex, it has a straightforward interpretation. The lefthand side is the average difference between incentivized and unincentivized predictions, both preand post-treatment, across both groups. The right-hand side is the commitment value described in (4) and (5), taking into consideration the different coupons offered to subjects and their different beliefs at the time of each prediction. The linearity that simplified the previous moments does not extend into the distribution of  $\varepsilon$ , and so we must write out the averages using summations. Because this moment relies strongly on the assumption that the only difference between the unincentivized and incentivized predictions comes through the commitment value of a p-coupon, we present the results from estimating a model both without (Model 1) and with (Model 2) making use of Equation (16).

The results of estimating the structural models using GMM are presented in table 5. We assume a Type 1 extreme value distribution for  $\varepsilon$  with a zero mean and scale parameter  $\sigma_{\varepsilon}$ .<sup>38</sup> Panel A presents those parameters estimated directly, while additional parameters derived from these are presented in Panel B. The structural parameters confirm the reduced-form results in Sections 4.1 and 4.2. On average, the immediate utility cost of gym attendance exceeds the discounted future benefits by \$4.71. Naivete about their future self-control problems causes people to under-estimate their future impatience about gym attendance by \$3.10, however. This corresponds directly to the significant over-prediction of future attendance relative to actual attendance found in the previous results. This naivete also explains why unincentivized predictions lie above the normalized valuations of p-coupons, as subjects underestimate their need for commitment and view the p-coupons as including a costly distortion.

Turning to the habit formation, we find that 32% of treated subjects formed a habit equivalent to a \$2.60 daily gym-attendance subsidy.<sup>39</sup> This is a significant habit — our \$100 treatment would be recouped after only 50 visits. It is still substantially smaller than the net daily cost, however, so that even a habituated subject would not on average enjoy going to the gym. This helps explain why we estimate one-third of subjects forming the habit but far fewer actually regularly going to the gym – many of the habituated subjects still did not receive sufficiently good shocks to push them over the attendance threshold.

In contrast to the not-inconsequential actual habit, subjects' predicted habit value was trivial. We estimate that subjects expected a habit worth only \$0.16, and the coefficient is not significantly different from zero. Because (12) assumes that subjects are fully aware of  $\pi$ , it is worth noting that any overoptimism about the probability of habit formation would imply that even this \$0.16 is biased upwards. By combining the predicted habit with the actual habit, we can estimate a large and highly significant degree of projection bias  $\alpha = 0.94$ . That is, while we can strongly reject the

<sup>&</sup>lt;sup>38</sup>The extreme value distribution is the most appropriate choice because we have aggregated the daily attendance decision to a weekly attendance observation. As a robustness check, we estimate the same models in Table A.4 using a normally distributed error term.

<sup>&</sup>lt;sup>39</sup>It is also possible to estimate a model with a homogeneous treatment effect assumption by restricting  $\pi = 1$ . In this case the overall habit value is \$1.33 (s.e. 0.45), although the overidentifying restrictions are now rejected at p = 0.031.

null of "no projection bias" ( $\alpha = 0$ ), we cannot reject "complete projection bias" ( $\alpha = 1$ ).

The second column of Table 5 presents the results from incorporating the additional moment restriction on the commitment value of our p-coupons. While the effect on most of the preceding parameters is small, and qualitatively unimportant, this additional restriction allows us to estimate an additional parameter,  $(1 - \hat{\beta})b$ . Because this reduces to zero in the case of complete beta-naivete, our significant estimate of 1.500 allows us to reject this null. We can combine this estimate with the our estimate of  $(\hat{\beta} - \beta)b$  in order to estimate the extent of beta-naivete. We show this in Panel B, finding a tightly estimated  $\omega = 0.666$ . This means that subjects believe their future discount factor  $\hat{\beta}$  will be a weighted average between 1 and their actual discount factor, unfortunately. If we use the typical value of  $\beta = 0.7$  found in other studies, this would imply that subjects hold the partially-naive  $\hat{\beta} = 0.9$ .

Finally, Figure 6 uses the estimated structural parameters to plot the probability of nonzero weekly gym attendance as a function of an individual's habit value. This figure makes clear the link between utility shifters and behavioral changes: in particular, our estimated habit value of \$2.60 increases the probability that a subject will attend the gym at least once in a given week by 47 percentage points. On the one hand, this can be viewed as a large effect. On the other hand, such a habit value is still only expected to lead to gym attendance in 63% of weeks. This is another way to understand why we estimate a larger proportion of "habit formers" than we find evidence of using behavior changes in the data. It also prompts us to expect that the habit effect may have decayed on its own in the absence of the quasi-exogenous break imposed by the semester break, as weeks with insufficient utility shocks gradually cause subjects to de-habituate.<sup>40</sup> More optimistically, however, we note that we estimate that a habit value of \$4.30 would generate a 95% probability of weekly gym attendance, which we view as the steady-state level of habituation.<sup>41</sup>This leaves as an open question for future work whether such a habit can be achieved by varying the treatment, or by complementing the habit with long-term small subsidies.

<sup>&</sup>lt;sup>40</sup>Although we have examined this point in our data, we have insufficient power to detect a downward trend in post-treatment attendance in addition to the the semester break effect.

 $<sup>^{41}</sup>$ Given our assumption of an unbounded error term, there is of course no finite habit level that would guarantee weekly attendance. We choose a 95% cutoff both as a "rule of thumb" and to avoid over-inferring from the tail of our assumed error distribution.

# 5 Conclusion

We find that incentivizing gym-attendance creates a significant short-run habit, albeit one which decays substantially as the result of a quasi-exogenous break in attendance. Among a sample of students who did not initially attend the gym regularly, a \$100 incentive induced 80% of our treatment group to visit twice-weekly for a month and raised their post-treatment attendance relative to a control group. In addition, smaller financial incentives significantly increased post-treatment attendance among both groups.<sup>42</sup> The overall habit we estimate is statistically indistinguishable in size from Charness and Gneezy's (2009) effect, and we find significant underlying heterogeneity: two-thirds of treated subjects appear unaffected, while one-third of subjects developed a gym attendance habit equivalent in utility to a \$2.60 per-visit subsidy. Despite this significant habit. subjects did not appear to predict any habit formation ex ante—an effect we interpret through the framework of projection bias. Furthermore, we find that subjects are greatly overoptimistic about their subsequent gym attendance in general, which we interpret as over-optimism about their own self-control. Even in weeks with no p-coupon to complicate the prediction task, subjects over-predict attendance by a factor of about three. This is a sufficient degree of mis-prediction to explain the result in DellaVigna and Malmendier (2006) that people purchase monthly health club memberships when their actual attendance only justifies the purchase of single-visit passes.<sup>43</sup> Thus, somewhat counterintuitively, we find that subjects under-value exercise and at the same time over-predict their willingness to engage in it, both of which may lead to under-attendance.

In our calibrated structural model, we find that subjects exhibit a significant amount of both projection bias and naivete regarding present bias. Subjects appear to only predict \$0.23 worth of the \$2.62 habit, corresponding to a degree of projection bias  $\alpha = 0.91$ . Leveraging the differences between their predicted and actual attendance, and unincentivized and incentivized predictions, we estimate that subjects are two-thirds naive with respect to their future self-control problems. Thus, while we model the exercise habit in the spirit of "rational addiction", we find that a substantial

 $<sup>^{42}</sup>$ This addresses the concern of, among others, Gneezy and Rustichini (2000) that small financial incentives may actually reduce a desirable behavior by crowding out subjects' intrinsic motivation. Although the scope of such motivation is limited in our sample – subjects did not initially attend the gym on their own – we find in no case did the provision of incentives lower attendance or was the attendance for a higher p-coupon significantly lower than that for a smaller p-coupon.

 $<sup>^{43}</sup>$ DellaVigna and Malmendier (2006) consider a different population, of course, so we do not claim that this is driving their result.

relaxation of the rational expectations assumption is necessary to organize the data.

We can also use the structural parameter estimates to address welfare effects and cost-effectiveness of the main experimental intervention. If we accept the  $\beta = 0.7$  typically found in the literature, then the undiscounted future benefits of gym attendance are  $b = \frac{1.5}{(1 - \hat{\beta})} = 15$ . Given the average increase of 0.256 visits/week and the 80% compliance rate, the increased attendance would have to last approximately 20 weeks for the program to be cost-effective. In our sample of students, however, we see significant decay after winter break, suggesting that exogenous interruptions in attendance may undermine the intervention before this threshold is met. It is possible, however, that a smaller incentive would have been sufficient to obtain compliance for a significant subset of our subjects, which would lower the duration necessary to justify the intervention. It is also worth noting that we have identified mistakes which people would themselves recognize as such at least ex-post—and thus lower the normative bar for intervention.<sup>44</sup> One potential difficulty for policy implementations, however, is that the distorted beliefs inherent to naivete and projection bias are such that the relevant population will not recognize the value of an intervention ex ante, nor properly make use of most forms of self-control devices. The policies developed in the face of naivete-driven mistakes may therefore be necessarily paternalistic, but still a Kaldor-Hicks improvement over status quo. One must exercise caution in extrapolating our results to other populations, of course, where compliance, habit formation, and the extent of biased beliefs might all be quite different.<sup>45</sup>

To that end, future research should explore the habit-formation and habit-decay effects in a more policy-relevant population. We also believe it is crucial to test the implications of this study for successfully translating short-run, fragile habits into longer-run, durable ones. Our results suggest that the same biases that prevent the initiation of an exercise habit may also subvert the remainder of the behavior change process.

 $<sup>^{44}</sup>$ See Bernheim and Rangel (2007) for a discussion of the difficulties in applying standard welfare criteria to nonstandard preferences.

<sup>&</sup>lt;sup>45</sup>For example, at the risk of over-inferring from our results we note that the mistakes we find subjects making, coupled with a shift in physical activity from work to leisure, could serve as an explanation of the rising obesity trend. While this hypothesis shares many similarities to "rational obesity" explanations such as Philipson and Posner (2003), it differs greatly in the welfare analysis.

## References

- Angeletos, George-Marios, David Laibson, Andrea Repetto, Jeremy Tobacman, and Stephen Weinberg, "The Hyperbolic Consumption Model: Calibration, Simulation, and Empirical Evaluation," *The Journal of Economic Perspectives*, Summer 2001, 15 (3), 47–68.
- Ariely, Dan and George Loewenstein, "The Heat of the Moment: the Effect of Sexual Arousal on Sexual Decision Making," *Journal of Behavioral Decision Making*, July 2005, 19 (2), 87–98.
- Badger, Gary, Warren Bickel, Louis Giordano, Eric Jacobs, George Loewenstein, and Lisa Marsch, "Altered states: The impact of immediate craving on the valuation of current and future opioids," *Journal of Health Economics*, September 2007, 26 (5), 865–876.
- Becker, Gary and Kevin Murphy, "A Theory of Rational Addiction," Journal of Political Economy, August 1988, 96 (4), 675–700.
- Bénabou, Roland and Jean Tirole, "Self-Confidence and Personal Motivation," The Quarterly Journal of Economics, August 2002, 117 (3), 871–915.
- Bernheim, B. Douglas and Antonio Rangel, "Behavioral Public Economics," in Peter Diamond and Hannu Vartiainen, eds., *Behavioral Economics and Its Applications*, Princeton University Press, 2007.
- Charness, Gary and Uri Gneezy, "Incentives to Exercise," *Econometrica*, May 2009, 77 (3), 909–931.
- Conlin, Michael, Ted O'Donoghue, and Timothy J. Vogelsang, "Projection Bias in Catalog Orders," The American Economic Review, September 2007, 97 (4), 1217–1249.
- **DellaVigna, Stefano**, "Psychology and Economics: Evidence from the Field," Journal of Economic Literature, 2009, 47 (2), 315–372.
- **and Ulrike Malmendier**, "Contract Design and Self-Control: Theory and Evidence," *The Quarterly Journal of Economics*, May 2004, *119* (2), 353–402.
- \_\_\_\_ and \_\_\_\_, "Paying Not To Go To The Gym," The American Economic Review, June 2006, 96 (3), 694–719.
- \_\_\_\_\_, John List, and Ulrike Malmendier, "Testing for Altruism and Social Pressure in Charitable Giving," *NBER Working Paper No. 15629*, December 2009.
- Dzewaltowski, David, John Noble, and Jeff Shaw, "Physical activity participation: social cognitive theory versus the theories of reasoned action and planned behavior," *Sport Psychology*, December 1990, *12* (4), 388–405.
- Frandsen, Brigham, "Randomization Inference on Quantiles Under Imperfect Compliance," working paper, 2010.
- Gelbach, Jonah, "When Do Covariates Matter? And Which Ones, and How Much?," Working Paper, June 2009.
- Gneezy, Uri and Aldo Rustichini, "Pay Enough, or Don't Pay at All," The Quarterly Journal of Economics, August 2000, 115 (3), 791–810.

- Godin, Gaston, "Theories of reasoned action and planned behavior: usefulness for exercise promotion," *Medicine and Science in Sports and Exercise*, November 1994, 26 (11), 1391–1394.
- \_\_\_\_\_, Pierre Valois, and Linda Lepage, "The pattern of influence of perceived behavioral control upon exercising behavior: An application of Ajzen's theory of planned behavior," *Journal of Behavioural Medicine*, 1993, 16 (1), 81–102.
- Gul, Faruk and Wolfgang Pesendorfer, "Temptation and Self-Control," Econometrica, November 2001, 69 (6), 1403–1435.
- \_\_\_\_ and \_\_\_\_, "Self-Control and the Theory of Consumption," *Econometrica*, January 2004, 72 (1), 119–158.
- Kane, Robert, Paul Johnson, Robert Town, and Mary Butler, "A Structured Review of the Effect of Economic Incentives on Consumers' Preventive Behavior," American Journal of Preventive Medicine, 2004, 27 (4).
- Laibson, David, "Golden Eggs and Hyperbolic Discounting," The Quarterly Journal of Economics, May 1997, 112 (2), 443–477.
- Lally, Phillppa, Cornelia H. M. van Jaarsveld, Henry W. W. Potts, and Jane Wardle, "How are habits formed: Modelling habit formation in the real world," *European Journal of Social Psychology*, 2010, 40 (6), 998–1109.
- Levy, Matthew, "An Empirical Analysis of Biases in Cigarette Addiction," mimeo, 2009.
- Loewenstein, George and Daniel Adler, "A Bias in the Prediction of Tastes," *Economic Journal*, 1995, 105, 929–937.
- \_\_\_\_\_, Ted O'Donoghue, and Matthew Rabin, "Projection Bias in Predicting Future Utility," *The Quarterly Journal of Economics*, March 2003, 118 (4), 1209–1248.
- O'Donoghue, Ted and Matthew Rabin, "Doing It Now or Later," The American Economic Review, March 1999, 89 (1), 103–124.
- \_\_\_\_ and \_\_\_\_, "Addiction and Self Control," in Jon Elster, ed., Addiction: Entries and Exits, Russel Sage Foundation, 1999.
- **Orphanides, Athanasios and David Zervos**, "Rational Addiction with Learning and Regret," *The Journal of Political Economy*, August 1995, 103 (4), 739–758.
- Philipson, Tomas and Richard Posner, "The Long-Run Growth in Obesity as a Function of Technical Change," *Perspectives in Biology and Medicine*, Summer 2003, 46 (3), S87–S107.
- Read, Daniel and Barbara van Leeuwen, "Predicting Hunger: The Effects of Appetite and Delay on Choice," Organizational Behavior and Human Decision Processes, November 1998, 76 (2), 189–205.
- Reynolds, Kim, Joel Killen, Susan Bryson, David Maron, C. Barr Taylor, Nathan Maccoby, and John Farquhar, "Psychosocial predictors of physical activity in adolescents," *Preventive Medicine*, September 1990, 19 (5), 541–551.
- Shapiro, Jesse, "Is There a Daily Discount Rate? Evidence From the Food Stamp Nutrition Cycle," *Journal of Public Economics*, 2005, 89 (2), 303–325.

- Skiba, Paige and Jeremy Tobacman, "Payday Loans, Uncertainty and Discounting: Explaining Patterns of Borrowing, Repayment, and Default," Vanderbilt Law and Economics Research Paper No. 08-33, August 2008.
- Valois, Pierre, Raymond Dersharnais, and Gaston Godin, "A comparison of the Fishbein and Ajzen and the Triandis attitudinal models for the prediction of exercise intention and behavior," *Journal of Behavioural Medicine*, 1988, 11 (5), 459–472.



## Figure 1: Our Experimental Design



Notes: Average weekly gym attendance, by treatment group status. Weeks in which a subject received a p-coupon for attendance are omitted from this figure.

	(1)	(2)	(Charness
	(1)	(2)	& Gneezy)
Treated	0.045		-0.100
	(0.057)		(0.196)
			$[0.477]^{a}$
Treatment Period X Treated	$1.209^{***}$		$1.275^{***}$
	(0.150)		(0.181)
			$[0.780]^{\rm a}$
Imm. Post-Treatment X Treated <sup>b</sup>	$0.256^{**}$		$0.585^{***}$
	(0.122)		(0.217)
			$[0.186]^{a}$
Later Post-Treatment x Treated <sup>b</sup>	0.045		_
	(0.098)		
Complied w/ treatment		0.057	
		(0.071)	
Treatment Period X Complied		$1.582^{***}$	
		(0.180)	
Imm. Post-Treatment X Compliance <sup>b</sup>		$0.338^{**}$	
		(0.154)	
Later Post-Treatment x Compliance <sup>b</sup>		0.061	
		(0.126)	
Week Efffects	Yes	Yes	Yes
Controls	Yes	Yes	_
IV	_	Yes	_
Observations	7433	7433	1520
Num Clusters	111	111	80
R-squared	0.21	0.22	0.13

 Table 1: Habit Formation: Regression of average weekly attendance.

Notes: Observations of weekly attendance at the subject-week level. Robust standard errors in parentheses, clustered by individual. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%.

<sup>a</sup>Terms in square brackets are p-values from a Chow test of equal coefficients between our sample (column 1) and Charness and Gneezy (2009)'s sample.

<sup>b</sup> "Immediate" refers to the 8 weeks following the intervention (excluding the "dead week" for columns (1) and (2). "Later" refers to the 19 weeks of observations in the following semester (excluding the semester break).



Figure 3: Cumulative Distribution of Immediate Post-Treatment Attendance

Notes: Cumulative distribution functions of post-treatment attendance among compliers in control and treated groups, based on finite-sample randomization inference with imperfect compliance. 90% confidence intervals shown; corresponding interval for fraction of compliers is [0.7083, 0.8083]; fraction of always-takers is constrained to zero. Weeks in which subjects received a p-coupon are omitted.

Figure 4: Actual vs. Predicted Post-Treatment Attendance



Notes: Average weekly gym attendance in the 8 weeks following the treatment month plotted against predicted attendance conditional on receiving no treatment. Weeks in which a subject received a p-coupon for attendance are omitted from this calculation. Predicted attendance is based on an OLS regression of attendance on week dummies and covariates using control group data for all weeks and treated group data for the pre-treatment period.



Figure 5: Predicted versus Actual Attendance

		Control grou	р		r	Freatment grou	ıp	
	Coupon	Un-	Un-	-	Coupon	Un-	Un-	
	Value $p > 0$	Incentivized $p > 0$	Incentivized $p = 0$		Value $p > 0$	Incentivized $p > 0$	Incentivized $p = 0$	
Pre-Treatment Predic	etions							
Predicted attendance	3.868	4.053	1.418		3.63	3.963	1.231	
Actual attendance	1.561	1.561	0.255		1.463	1.463	0.365	
Difference	2.307	2.491	1.164		2.167	2.500	0.865	
St. Error	(0.297)	(0.235)	(0.149)		(0.350)	(0.318)	(0.178)	
No. of observations	57	57	55		54	54	52	
Post-Treatment Predi	ictions							
Predicted attendance	3.395	3.614	1.058		3.185	3.056	1.313	
Actual attendance	1.561	1.561	0.269		1.463	1.463	0.396	
Difference	1.833	2.053	0.788		1.722	1.593	0.917	
St. Error	(0.321)	(0.299)	(0.144)		(0.315)	(0.299)	(0.171)	
No. of observations	57	57	52		54	54	48	

Table 2: Misprediction of attendance

Notes: Coupon value refers to the average valuation of a p-coupon normalized by its subsidy, and includes only observations for the week a subject actually received a p-coupon. Un-incentivized refers to subjects' direct predictions, and is separated into this subsidized week and the unincentivized week for which subjects were only asked to make predictions without a p-coupon.

	Table 3: Predictions: Delay versus Session Effects								
	(1)	(2)	(3)	(4)	(5)				
	Actual	Coupon	Un-	Coupon	Un-				
	1100000	Value	incentivized	Value	incentivized				
Post-Treatment		-0.630***	-0.707***	-0.476**	-0.810***				
		(0.132)	(0.112)	(0.226)	(0.187)				
p=\$0	$-2.275^{***}$		-3.360***		-3.925***				
	(0.611)		(0.498)		(0.598)				
p=\$1	$-1.669^{**}$	-0.924	$-1.650^{***}$	-0.512	$-1.618^{**}$				
	(0.689)	(0.581)	(0.482)	(1.235)	(0.640)				
p=\$2	$-1.304^{*}$	-0.760	-1.288***	-1.522	-2.213***				
	(0.708)	(0.579)	(0.478)	(1.232)	(0.617)				
p=\$3	$-1.440^{**}$	-0.530	$-0.924^{*}$	-0.489	-1.276**				
	(0.714)	(0.580)	(0.472)	(1.233)	(0.634)				
p=\$5	-0.050	-0.081	-0.272	0.027	-0.698				
	(0.808)	(0.623)	(0.523)	(1.241)	(0.648)				
Constant	$2.600^{***}$	$3.865^{***}$	$4.953^{***}$	$3.988^{***}$	$5.405^{***}$				
	(0.609)	(0.613)	(0.497)	(1.233)	(0.590)				
Observations	551	875	1088	176	217				
R-squared	0.20	0.06	0.27	0.11	0.33				
Num Clusters:	111	111	111	110	111				
Sample	Full	Full	Full	5-wk delay	5-wk delay				

Notes: Observations are at the subject-week level. Coupon value refers to the average valuation of a p-coupon normalized by its subsidy, and includes only target weeks associated with a non-zero subsidy. Un-incentivized refers to subjects' direct predictions, and includes all target week predictions. Robust standard errors in parentheses, clustered by individual. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. p = \$7 is the omitted category.

41

	Unincentivized (\$0)	Incentivized (\$1)	
Post-Trmt X Treated	$0.458^{**}$	0.325	
	(0.227)	(0.315)	
Post-Trmt	-0.396**	-0.825***	
	(0.170)	(0.193)	
Treated	-0.283	-0.061	
	(0.233)	(0.346)	
R-squared	0.255	0.293	
Num clusters	107	111	

 Table 4: Difference-in-differences in Predictions

 Unincentivized (\$0)
 Incentivized (\$1)

Notes: Robust standard errors in parentheses, clustered by individual. \* significant at 10%; \*\* significant at 5%, \*\*\* significant at 1%. The dependent variable in column 1 is subjects' subjective predictions for weeks in which they received no p-coupon, and in column 2 is subjects' valuations of p-coupons for p=\$1.

Table 5: GMM Parameters							
Name	Parameter	Model 1	Model 2				
Panel A: Directly Estimate	ed Parameters						
Net daily cost	$C - \beta b$	$4.713^{***}$	$4.582^{***}$				
		(0.417)	(0.421)				
Cost of naivete	$(\widehat{eta}-eta)b$	$3.099^{***}$	$2.993^{***}$				
		(0.289)	(0.290)				
Habit value	$\eta$	$2.602^{***}$	$2.618^{***}$				
		(0.733)	(0.738)				
Predicted habit value	$(1-lpha)\eta$	0.160	0.226				
		(0.747)	(0.757)				
Probability of habituation	$\pi$	$0.320^{**}$	$0.306^{**}$				
		(0.133)	(0.136)				
Demand for commitment	$(1-\widehat{eta})b$	—	$1.500^{***}$				
			(0.250)				
Scale parameter, daily shock	k $\sigma$	$1.528^{***}$	$1.490^{***}$				
		(0.148)	(0.150)				
Panel B: Extended Parame	eters						
Degree of projection bias	lpha	$0.939^{***}$	0.914***				
		(0.285)	(0.286)				
Degree of beta-naivete	ω		$0.666^{***}$				
			(0.037)				

Notes: The daily shock,  $\epsilon$ , is drawn from a mean-zero type-1 extreme value distribution. Standard errors in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. Parameters in Panel B are calculated by transformations of parameters in Panel A, with standard errors implied by the delta rule. Model 2 includes an additional moment restriction on the difference between unincentivized predictions and p-coupon valuations.



Figure 6: Simulated Weekly Attendance

Notes: Simulation of the probability of observing a non-zero weekly attendance as a function of the habit value, based on Model 1 in Table 5. Dashed lines indicate the 95% confidence interval.

# A Appendices

## A.1 Value of a p-coupon

The ex-ante value of a p-coupon is

$$X_2^g = X_6^g = 7 \times \int_{c-\widehat{\beta}b+g\cdot\eta-P}^{\infty} P \, dF(\varepsilon) + 7 \times \int_{c-\widehat{\beta}b+g\cdot\eta-P}^{c-\widehat{\beta}b-g\cdot\eta} (b+g\cdot\eta-c+\varepsilon) \, dF(\varepsilon).$$

To see that this is weakly positive, note that the first integral is always non-negative, and the second integral is bounded below by

$$\int_{c-\widehat{\beta}b-g\cdot\eta-P}^{c-\widehat{\beta}b-g\cdot\eta} dF(\varepsilon)\cdot\left[(1-\widehat{\beta})b-g\cdot\eta-P\right].$$

Thus, dividing by 7 for notational convenience, and noting that by assumption  $b \ge 0$ :

$$\begin{aligned} \frac{X_2^C}{7} &= \frac{X_6^C}{7} > \int\limits_{c-\widehat{\beta}b-g\cdot\eta-P}^{\infty} P \, dF(\varepsilon) + \int\limits_{c-\widehat{\beta}b-g\cdot\eta-P}^{c-\widehat{\beta}b-g\cdot\eta} dF(\varepsilon) \cdot \left[ (1-\widehat{\beta})b - g\cdot\eta - P \right] \\ &= \int\limits_{c-\widehat{\beta}b-g\cdot\eta-P}^{c-\widehat{\beta}b-g\cdot\eta} P \, dF(\varepsilon) + \int\limits_{c-\widehat{\beta}b-g\cdot\eta}^{\infty} P \, dF(\varepsilon) - \int\limits_{c-\widehat{\beta}b-g\cdot\eta-P}^{c-\widehat{\beta}b-g\cdot\eta} P \, dF(\varepsilon) + \int\limits_{c-\widehat{\beta}b-g\cdot\eta-P}^{c-\widehat{\beta}b-g\cdot\eta} (1-\widehat{\beta})b \, dF(\varepsilon) \\ &= \int\limits_{c-\widehat{\beta}b-g\cdot\eta}^{\infty} P \, dF(\varepsilon) + \int\limits_{c-\widehat{\beta}b-g\cdot\eta-P}^{c-\widehat{\beta}b-g\cdot\eta} (1-\widehat{\beta})b \, dF(\varepsilon) \ge 0 \end{aligned}$$

#### A.2 Sample

Our initial sample consisted of 120 subjects, randomly assigned to treated and control groups of 60 subjects each. Table A.1 provides a comparison of the treated and control groups. Due to attrition and missing covariates the final number of treated subjects in our analysis is 54 and of control subjects 57. Comparing the two groups on the covariates that we used in all of our analysis we find no significant differences in means, and the F-test of joint significance of the covariates in a linear regression of the treatment-group dummy on covariates is 0.387. In addition to basic demographic variables we included discretionary budget and the time and money cost of getting to campus in order to control for differences in the cost of gym attendance and the relative value of monetary incentives. The pre-treatment Godin Activity Scale is a self-reported measure of physical activity in a typical week prior to the treatment. The self-reported importance of physical fitness and physical appearance were included as a proxy for subjects' taste for the outcomes typically associated with gym-attendance. The naivete proxy covariates are subjects answers to a series of questions that we asked in order to get at their level of sophistication about self-control problems. Answers were given on a four-point scale from "Disagree Strongly" to "Agree Strongly". The exact wording of these questions is as follows:

Variable	Question
Forget	I often forget appointments or plans that I've made, so that
	I either miss them, or else have to rearrange my plans at the
	last minute.
Spontaneous	I often do things spontaneously without planning.
Things come up	I often have things come up in my life that cause me to change my plans.
Think ahead	I typically think ahead carefully, so I have a pretty good idea what I'll be doing in a week or a month.
Procrastinate	I usually want to do things I like right away, but put off things that I don't like.

### A.3 Screening mechanism

The webpage we used to screen for non-attenders is shown below. We included three "dummy" questions to make it harder for subjects to return to the site and change their answers in order to be able to join the experiment. Despite this precaution, a handful of subjects did return to the screening site and modify their answers until they hit upon the correct answer to join the experiment. (Which was a "no" on question four.) Out of a total of 497 unique IP addresses in our screening log, we found 5 instances of subjects possibly gaming the system to gain access to the study. We have no way to determine if these subjects wound up in our subject pool.

#### A.4 Elicitation mechanisms

Figure A.2 depicts the sample p-coupon and instructions that subjects saw to prepare them for the incentive-compatible elicitation task. Verbal instructions given at this time further clarified exactly what we were asking subjects to do. Note that the sure-thing values in column A are increments

Table A.1: Co	omparison of	Treated and C	ontrol groups.	
	(1)	(2)	(3)	(4)
	Full sample	Treated group	Control group	T-test p-value
Original sample	120	60	60	
No. of attriters	6	4	2	
No. w/ incomplete controls	3	2	1	
Final sample size	111	54	57	
\$25 learning-week incentive		Yes	Yes	
\$100 treatment-month incentive		Yes	-	
Demographic covariates				
Age	21.919	22.204	21.649	0.639
	(0.586)	(0.990)	(0.658)	
Gender (1=female)	0.685	0.648	0.719	0.425
	(0.044)	(0.066)	(0.060)	
Proportion white	0.36	0.333	0.386	0.568
	(0.046)	(0.065)	(0.065)	
Proportion Asian	0.559	0.63	0.491	0.145
	(0.047)	(0.066)	(0.067)	
Proportion other race	0.081	0.037	0.123	0.01
	(0.026)	(0.026)	(0.044)	
Economic covariates				
Discretionary budget	192.342	208.333	177.193	0.404
	(18.560)	(28.830)	(23.749)	
Travel cost to campus	0.901	0.648	1.14	0.37
	(0.273)	(0.334)	(0.428)	
Travel time to campus (min)	14.662	14.398	14.912	0.811
	(1.071)	(1.703)	(1.335)	
Naivete proxy covariates			i	
$\operatorname{Forget}^{a,b}$	1.595	1.556	1.632	0.573
	(0.067)	(0.090)	(0.099)	
$Spontaneous^{a,b}$	2.486	2.574	2.404	0.281
•	(0.079)	(0.104)	(0.117)	
Things come $up^{a,b}$	2.586	2.611	2.561	0.731
	(0.072)	(0.107)	(0.097)	
Think $ahead^{a,b}$	2.874	2.944	2.807	0.338
	(0.071)	(0.081)	(0.116)	
$Procrastinate^{a,b}$	3.036	3.056	3.018	0.8
	(0.075)	(0.104)	(0.108)	
Exercise experience and attitud	e covariates	. ,	. ,	
Pre-trt Godin Activity Scale	36.05	36.5	35.623	0.855
v	(2.376)	(2.983)	(3.689)	
Fitness is $important^{a,b}$	3.081	2.981	3.175	0.092
Ī	(0.057)	(0.086)	(0.076)	
Appearance is $important^{a,b}$	3.252	3.259	3.246	0.917
	(0.065)	(0.096)	(0.088)	
F-test of joint significance	<pre></pre>	( · · · · /	( /	0.387

Table A 1. ric f Treated and Control  $\mathbf{C}$ 

# Figure A.1: Screening Site

To de	etermine your eligibility for this experiment, please complete this questionnaire and click "submit".
1.	Please enter the verification key supplied in the email.
2.	How many semesters, prior to this one, have you been enrolled at UC Berkeley or another four-year, post-secondary institution? (Include summer session.)
	*
3.	Have you declared a major in the Social Sciences?
	◯ Yes ◯ No ◯ No sure
4.	Do you regularly attend the UC Berkeley Recreational Sports Facility (RSF) or any similar recreational or fitness facility or gym?
	○ Yes ○ No
5.	How frequently do you use the Internet?
	🔘 Several times per day 🔘 Once a day 🔘 A few times each week 🔘 Never

Submit

of P. The line number where subjects cross over from choosing column B to choosing column A bounds their valuation for the p-coupon. We used a linear interpolation between these bounds to create our "BDM" variable. Thus, for example, if a subject chose B at and below line four, and then chose A at and above line five we assigned them a p-coupon valuation of  $P \times 4.5$  In general subjects appear to have understood this task clearly. There were only three subjects who failed to display a single crossing on every task, and all of them appear to have realized what they were doing before the end of the first elicitation session. The observations for which these three subjects did not display a single crossing have been dropped from our analysis.

By randomly choosing only one target week for only one subject we maintain incentive compatibility while leaving all but one subject per session actually holding a p-coupon, and for only one target week. This is important because what we care about is the change in their valuation of a p-coupon from pre- to post-treatment elicitation sessions. Subjects who are already holding a coupon from the first session would be valuing a second coupon in the second session, making their valuations potentially incomparable, rather like comparing willingness-to-pay for a first candy bar to willingness-to-pay for a second candy bar.

The instructions and example for the unincentivized prediction task and the task for prediction of other people's attendance appear as figure A.3.

#### A.5 Compliance, attrition, and randomization.

About 80% of Charness and Gneezy's high-incentive subjects complied with the \$100 treatment incentive by attending the gym eight times during the treatment month. A similar percentage, 75%, of our treatment subjects complied with our treatment incentive by attending the gym twice a week during the treatment month. In our data, a direct comparison of means between treatment and control will only allow us to estimate an "intention to treat" effect (ITT). If compliance were random we could simply inflate this by the inverse of the compliance rate to estimate the average treatment effect. Since compliance is almost certainly not random, we will do our best to estimate an "average treatment effect on the treated" (ATT) by using our rich set of individual covariates to help us control for differences between compliers and non-compliers.

To mitigate attrition over our three sessions we gave subjects two participation payments of \$25 each, in addition to the various gym-attendance offers. The first payment was for attendance at the first session. The second payment required attendance at both the second and third sessions.<sup>46</sup> Despite this titration of rewards, six of the 120 subjects did not complete the study. Two control subjects and two treatment subjects left the study between the first and second sessions, and two more treatment subjects left between the second and third. In order to include an additional handful of subjects who were not able to make the third session, and otherwise would have left the study, we held make-up sessions the following day. Four control subjects and two treatment subjects attended these sessions and we have treated them as having completed the study.

Randomizing subjects into treatment and control presented some challenges. Our design required that treatment and control subjects meet separately. For each of the three sessions we scheduled four timeslots, back-to-back, and staggered them between Control and Treatment. When subjects responded to the online solicitation, and after they had completed the screening questionnaire, they were randomly assigned to either treatment or control and were then asked to choose between the two timeslots allocated to their assigned group. Subjects who could not find a timeslot that fit their schedule voluntarily left the study at this point.<sup>47</sup> As it turned out, subjects assigned

<sup>&</sup>lt;sup>46</sup>Gym-attendance offers were not tied to attendance because this would have created a differential between the treatment and control groups in the incentive to complete the study.

<sup>&</sup>lt;sup>47</sup>Technically they were considered to have never joined the study, and received no payment.

Figure A.2: Sample p-coupon and incentive-compatible elicitation task

# [PRACTICE]

This exercise involves nine questions, relating to the Daily RSF-Reward Certificate shown at the top of the page. Each question gives you two options, A or B. For each question check the option you prefer.

You will be asked to complete this exercise four times, once each for four of the five target weeks. The daily value of the certificate will be different for each of these four target weeks. For one of the five weeks you will not be asked to complete this exercise.

At the end of the session I'll choose one of the five target weeks at random. Then I'll choose one of the nine questions at random. Then I'll choose one subject at random. The randomly chosen subject will receive whichever option they checked on the randomly chosen question for the randomly chosen target week. Thus, for each question it is in your interest to check the option you prefer.

\$ <b>1</b> Daily RSF-Reward Certificate <b>\$1</b>	
This certificate entitles the holder to <b>\$1</b> for every day that he or she attends the RSF during the week	
of Monday, Oct 13 through Sunday, Oct 19.	
\$ 1 \$1	

	S	Μ	Т	W	Т	F	S
SEPT		1	2	3	4	5	6
	7	8	9	10	11	12	13
	14	15	16	17	18	19	20
	21	22	23	24	25	26	27
OCT	28	29	30	1	2	3	4
	5	6	7	8	9	10	11
	12	13	14	15	16	17	18
	19	20	21	22	23	24	25
NOV	26	27	28	29	30	31	1
	2	3	4	5	6	7	8
	9	10	11	12	13	14	15
	16	17	18	19	20	21	22
	23	24	25	26	27	28	29

For each question, check which option you prefer, A or B.

	Option A		Option B
1. Would you prefer	\$1 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
2. Would you prefer	\$2 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
3. Would you prefer	\$3 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
4. Would you prefer	\$4 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
5. Would you prefer	\$5 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
6. Would you prefer	\$6 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
7. Would you prefer	\$7 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
8. Would you prefer	\$8 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.
9. Would you prefer	\$9 for certain, paid Monday, Oct 20.	or	The Daily RSF-Reward Certificate shown above.

# [PRACTICE]

For each target week you will also be asked to complete the following two exercises. Both of these exercises relate to the Daily RSF-Reward Certificate shown at the top of the page, which is the same as the one shown at the top of the preceding page. In addition, there will be one target week for which you will be shown no certificate, and you will be asked to complete only these last two exercises.



	S	Μ	Т	W	Т	F	S
SEPT		1	2	3	4	5	6
	7	8	9	10	11	12	13
	14	15	16	17	18	19	20
	21	22	23	24	25	26	27
OCT	28	29	30	1	2	3	4
	5	6	7	8	9	10	11
	12	13	14	15	16	17	18
	19	20	21	22	23	24	25
NOV	26	27	28	29	30	31	1
	2	3	4	5	6	7	8
	9	10	11	12	13	14	15
	16	17	18	19	20	21	22
	23	24	25	26	27	28	29

Imagine that you have just been given the Daily RSF-Reward Certificate shown above, and that this is the only certificate you are going to receive from this experiment.

How many days would you attend the RSF that week if you had been given that certificate? \_\_\_\_\_

Now imagine that everyone in the room *except you* has just been given the Daily RSF-Reward Certificate shown above, and that this is the only certificate they are going to receive from this experiment.

What do you think would be the average number of days the other people in the room *(not including you)* would go to the RSF that week? \_\_\_\_\_

(Your answer does not have to be a round number. It can be a fraction or decimal.)

Notes: As part of this experiment some subjects will receive real certificates.

I will give a \$10 prize to the subject whose answer to this exercise is closest to the correct, average RSF-attendance for subjects (*other than themselves*) who receive the certificate shown above. The prize money will be paid by check, mailed on Monday,.Oct 20.

to the treatment group were substantially less likely to find a timeslot that worked for them, and as a result the desired number of subjects were successfully enrolled in the control group well before the treatment group was filled. Wanting to preserve the balanced number of Treatment and Control subjects, maintain power to identify heterogeneity within the Treatment group, and stay within the budget for the study, we capped the control group and continued to solicit participants in order to fill the treatment group. Subjects who responded to the solicitation after the Control group was filled were randomly assigned to treatment or control, and those assigned to control were then thanked and told that the study was full. Our treatment group therefore includes subjects who were either solicited later, or responded to the solicitation later than any of the subjects in the control group.<sup>48</sup>

To the extent that these temporal differences are correlated with any of the behaviors we are studying, simple comparisons of group averages may be biased. It appears, however, that the two groups are not substantially different along any of the dimensions we observed in our dataset, as a joint F-test does reject that the two groups were randomly selected from the same population based on observables. A comparison of the two groups appears in a separate appendix. To address the possibility that they may have differed significantly on unobservables we use observable controls in our hypothesis tests.

<sup>&</sup>lt;sup>48</sup>Additionally, the two groups of subjects were available at different times of day. To the extent that what made it hard for Treatment subjects to find a timeslot that fit the schedule may have been correlated with gym-attendance behavior (if, for example, the Treatment timeslots happen to have coincided with the most prefered times for non-gym exercise), then the group averages for some outcome variables may be biased.

10010 11.2. 00	(1)	(2)	(3)	(4)
	Treated Group	Compliers	Non-Compliers	T-test p-value
Demographic covariates	P			p
Age	22.204	22.605	20.636	0.429
0	(0.990)	(1.234)	(0.472)	
Gender (1=female)	0.648	0.651	0.636	0.929
	(0.066)	(0.074)	(0.152)	
Proportion white	0.333	0.349	0.273	0.640
-	(0.065)	(0.074)	(0.141)	
Proportion Asian	0.630	0.651	0.545	0.526
-	(0.066)	(0.074)	(0.157)	
Proportion other race	0.037	0.000	0.182	0.004
	(0.026)	(0.000)	(0.122)	
Economic covariates	· · ·	· ·	· · ·	
Discretionary budget	208.333	222.093	154.545	0.350
	(28.830)	(34.475)	(41.808)	
Travel cost to campus	0.648	0.616	0.773	0.853
	(0.334)	(0.386)	(0.679)	
Travel time to campus (min)	14.398	13.372	18.409	0.237
	(1.703)	(1.790)	(4.564)	
Naivete proxy covariates				
"Forget <sup><math>a,b</math></sup> "	1.556	1.465	1.909	0.047
	(0.090)	(0.096)	(0.211)	
"Spontaneous <sup><math>a,b</math></sup> "	2.574	2.442	3.091	0.011
-	(0.104)	(0.101)	(0.285)	
"Things come $up^{a,b}$ "	2.611	2.558	2.818	0.333
	(0.107)	(0.101)	(0.352)	
"Think ahead <sup><math>a,b</math></sup> "	2.944	2.977	2.818	0.436
	(0.081)	(0.091)	(0.182)	
"Procrastinate $^{a,b}$ "	3.056	2.977	3.364	0.135
	(0.104)	(0.118)	(0.203)	
Exercise experience and attitu	de covariates		· · ·	
Pre-trt Godin Activity Scale	36.500	38.360	29.227	0.221
-	(2.983)	(3.137)	(7.961)	
"Fitness is $important^{a,b}$ "	2.981	2.977	3.000	0.914
-	(0.086)	(0.097)	(0.191)	
"Appearance is important <sup><math>a,b</math></sup> "	3.259	3.256	3.273	0.944
-	(0.096)	(0.095)	(0.304)	
N obs.	54	43	11	
F-test of joint significance				0.635

 Table A.2:
 Comparison of Compliers and Non-Compliers

Notes: <sup>a</sup> 1= Disagree Strongly, 2=Disagree Somewhat; 3=Agree Somewhat; 4=Agree Strongly. <sup>b</sup> Wording of questions in appendix. Standard errors in parentheses.

#### Habit Formers **A.6**

Table A.S.	(1)	(2)	(3)	· (4)
	Treated Group	"Habit-Formers"	Non "Habit-Formers"	T-test p-value
Demographic covariates	incured enoup	110010 1 0111010		1 cost p raide
Age	22.204	19.750	22.630	0.306
	(0.990)	(0.453)	(1.150)	
Gender (1=female)	0.648	0.625	0.652	0.885
	(0.066)	(0.183)	(0.071)	
Proportion white	0.333	0.250	0.348	0.596
1	(0.065)	(0.164)	(0.071)	
Proportion Asian	0.630	0.750	0.609	0.454
1	(0.066)	(0.164)	(0.073)	
Proportion other race	0.037	0.000	0.043	0.557
· F - · · · · · · · · · · · · · · ·	(0.026)	(0.000)	(0.030)	
Economic covariates	. ,	. /	. ,	
Discretionary budget	208.333	181.250	213.043	0.699
	(28.830)	(92.068)	(30.274)	
Travel cost to campus	0.648	0.000	0.761	0.424
I	(0.334)	(0.000)	(0.391)	
Travel time to campus (min)	14.398	9.688	15.217	0.252
r ( )	(1.703)	(1.666)	(1.958)	
Naivete proxy covariates	. ,			
"Forget <sup>a,b</sup> "	1.556	1.500	1.565	0.800
5	(0.090)	(0.327)	(0.091)	
"Spontaneous <sup><math>a,b</math></sup> "	2.574	2.250	2.630	0.198
-	(0.104)	(0.164)	(0.118)	
"Things come $up^{a,b}$ "	2.611	2.375	2.652	0.363
0	(0.107)	(0.263)	(0.117)	
"Think ahead <sup><math>a,b</math></sup> "	2.944	3.000	2.935	0.778
	(0.081)	(0.189)	(0.090)	
"Procrastinate <sup><math>a,b</math></sup> "	3.056	2.875	3.087	0.473
	(0.104)	(0.295)	(0.111)	
Exercise experience and attitu	ude covariates			
Pre-trt Godin Activity Scale	36.500	41.688	35.598	0.474
, i i i i i i i i i i i i i i i i i i i	(2.983)	(3.823)	(3.434)	
"Fitness is important <sup><math>a,b</math></sup> "	2.981	3.500	2.891	0.010
_	(0.086)	(0.189)	(0.089)	
"Appearance is important <sup><math>a,b</math></sup> "	3.259	3.375	3.239	0.620
_	(0.096)	(0.183)	(0.109)	
N obs.	54	8	46	
F-test of joint significance				0.663

Table A 3: Comparison of Habit-Formers and Non Habit-Former

Notes: <sup>a</sup> 1= Disagree Strongly, 2=Disagree Somewhat; 3=Agree Somewhat; 4=Agree Strongly. <sup>b</sup> Wording of questions in appendix. Standard errors in parentheses.

# A.7 GMM Robustness Checks

Name	Parameter	Model 1	Model 2			
Panel A: Directly Estimated Parameters						
Net daily cost	$C - \beta b$	$5.219^{***}$	$5.057^{***}$			
		(0.504)	(0.491)			
Cost of naivete	$(\widehat{eta}-eta)b$	$2.746^{***}$	$2.641^{***}$			
		(0.282)	(0.274)			
Habit value	$\eta$	$2.229^{***}$	$2.250^{***}$			
		(0.731)	(0.740)			
Predicted habit value	$(1-\alpha)\eta$	0.175	0.246			
		(0.825)	(0.839)			
Probability of habituation	$\pi$	$0.320^{**}$	$0.306^{**}$			
		(0.134)	(0.136)			
Demand for commitment	$(1 - \widehat{\beta})b$	_	$1.506^{***}$			
			(0.290)			
Standard deviation, daily shock	$\sigma$	$2.670^{***}$	2.594***			
		(0.264)	(0.258)			
Panel B: Extended Parameters						
Degree of projection bias	$\alpha$	$0.922^{***}$	$0.891^{***}$			
		(0.367)	(0.369)			
Degree of beta-naivete	$\omega$	—	$0.637^{***}$			
			(0.045)			

 Table A.4: GMM Parameters - Alternative Specification

Notes: The daily shock,  $\epsilon$ , is drawn from a mean-zero normal distribution. Standard errors in parentheses. \* significant at 10%; \*\* significant at 5%; \*\*\* significant at 1%. Parameters in Panel B are calculated by transformations of parameters in Panel A, with standard errors implied by the delta rule. Model 2 includes an additional moment restriction on the difference between unincentivized predictions and p-coupon valuations.