# Pre-analysis plan: Learning about one's self

Yves Le Yaouanq[1]        Peter Schwardmann[2]

## 1   Introduction

Individuals with present bias overweigh more imminent consumption events. When choosing whether to work today or tomorrow, they allocate more work to tomorrow. However, come tomorrow, they will work less than they would have wanted to from yesterday's perspective. Experimental studies document both present bias (e.g. Read and Van Leeuwen, 1998; Augenblick et al., 2015) and people's naiveté about being subject to present bias (e.g. Augenblick and Rabin, 2015; Fedyk, 2016). Unfortunately, it is in the presence of naiveté that one's present bias becomes particularly costly (e.g. DellaVigna and Malmendier, 2006).

However, it is puzzling that naiveté that pertains to everyday behavior and a stable trait of the individual can persist. Why do people not learn about their bias? Our experiment seeks to provide a first step to resolving this puzzle empirically. We investigate whether people are able to learn about their present bias if they have the opportunity to engage in an unpleasant task repeatedly. To this end, we develop a framework that allows us to quantify what it means to learn optimally from one's own behavior and to compare individuals' observed learning to this benchmark.

In the experiment, subjects can complete two unpleasant tasks at two separate dates in the future. We elicit subjects' beliefs before and after the first task. During the first elicitation, we elicit a subject's joint prior distribution over her behavior at each of the two dates. After she either completes or fails to complete the first task, we elicit her beliefs about the likelihood of completing the second task. In a treatment, we vary whether the second task is the same task subjects completed at the first date, or a different one.

The resulting dataset allows us to study the extent to which people are able to learn from

---

[1]Department of Economics, University of Munich (LMU), Ludwigstr. 28, D-80539 Munich, Germany; email: yves.leyaouanq@econ.lmu.de.

[2]Department of Economics, University of Munich (LMU), Ludwigstr. 28, D-80539 Munich, Germany; email: peter.schwardmann@econ.lmu.de.

their past behavior. Moreover, we can investigate three potential drivers of people's failure to learn. First, people may underestimate how informative their behavior tomorrow is about their behavior the day after tomorrow. That is, they may not understand that behavior is driven by deep time preferences and, hence, underappreciate the correlation between their behavior at future dates.

Second, people may interpret what task completion today implies for their task completion tomorrow in an upwardly biased way. To investigate this hypothesis, we ask whether a subject's belief after observing her task completion at the first date are compatible with her prior and its implied Bayesian posterior.

Third, people may struggle to transport what they learn about their self-control problem in one environment into a (slightly) different environment. By varying the nature of the second task, we ask whether the previous two impediments to learning are more severe if the decision-making environment changes.

# 2 Experimental design

## 2.1 Main experiment

The experiment will feature two unpleasant tasks. Task A is a slider task. Subjects can complete up to 40 screens of 40 sliders each. Task B is a counting zeroes task. Subjects will be afforded the opportunity to complete up to 40 screens that each feature 12 matrices that each contain 40 ones and zeroes in random order.

In the control condition, subjects will face task A twice. In the treatment condition, they will first face task A and then task B. The payment for each task is a function of the number of screens a subject completes, with decreasing marginal benefit. Subjects decide how many screens they complete. The timing of the experiment is as follows.

- t=1 (in the lab)

    - Detailed instructions about the experiment

    - Familiarization with tasks A and B, i.e. subjects have to complete 5 screens of each

- Instructions on the Becker-DeGroot-Marschak (BDM) mechanism used to elicit beliefs.

- t=2 (online):

  - Commitment choice: subjects decide on action at t=3 and at t=5. Their choice is implemented with 5 percent probability

  - Belief elicitation about joint probability distribution over task completion at t=3 and t=5

- t=3 (online):

  - Task A

- t=4 (online):

  - Commitment choice: previously non-committed subjects decide on action at t=5 and their choice is implemented with 5 percent probability

  - Belief elicitation about probability of task completion at t=5

- t=5 (online):

  - Task A (or task B in treatment)

  - Payoff announcement

Belief elicitations at $t = 2$ and $t = 4$ take place two days after the familiarization with and completion of the tasks respectively. This is done in order to avoid that tired subjects state beliefs that are affected by projection bias. Beliefs are about the probability of completing more than 20 screens in the task(s). At $t = 2$, we will ask about the following four joint probabilities (with one randomly chosen scenario being payoff relevant):

- The probability of completing at least 20 screens at $t = 3$ and at least 20 screens at $t = 5$.

- The probability of completing at least 20 screens at $t = 3$, but less than 20 screens at $t = 5$.

- The probability of completing less than 20 screens at $t = 3$, but at least 20 screens at $t = 5$.

- The probability of completing less than 20 screens at $t = 3$, but less than 20 screens at $t = 5$.

At $t = 4$, we will ask subjects about their probability of completing at least 20 screens at $t = 5$.

**Incentives in belief elicitation.** Suppose a subject states a subjective probability about an event of $X$ percent. Then the BDM mechanism will randomly select an integer $Y$ between 0 and 100. If $X < Y$, then the subject receives a prize (of 3 euros) with probability Y. If $X \geq Y$, then the subject receives a prize if the event occurs.

**Show-up fee and incentivization of task** A subject receives 25 euros for completing the entire experiment in addition to the payoffs from the belief elicitations and the tasks. If a subject fails to log in at any one day, she will lose all payments. If she logs in, she is allowed to complete zero screens (at $t = 3$ and $t = 5$). However, failure to complete either training, elicitations or choices at dates $t = 1$, $t = 2$, and $t = 4$ will result in exclusion from the experiment with no payment. This design feature is intended to keep attrition at a minimum.

Subjects will be paid for every 5 screens they complete. The cumulative payment schedule is as follows: 5 screens = 5 euros ; 10 screens = 9 euros ; 15 screens = 12 euros; 20 screens = 14 euros; 25 screens = 15.5 euros; 30 screens = 16.5 euros; 35 screens = 17 euros; 40 screens = 17.10 euros;.

**Commitment choice.** At $t = 2$ and $t = 4$, we let subjects commit to future effort choices. At $t = 2$, subjects can commit to effort choices at $t = 3$ and $t = 5$. At $t = 4$ subjects can commit to effort choice at $t = 5$. Each commitment choice is implemented with 5 percent probability. If it is implemented, then the pre-committed number of screens cannot be exceeded and in case of a failure to complete the pre-committed number of screens, a subject will not be paid for a particular task, regardless of how many screens she did complete. The

commitment choice is not crucial for studying naiveté and its evolution, but helps provide a richer picture of the setting: commitment choices that exceed effort choice on the day of the task imply present bias.

**Belief elicitation as a soft commitment.** Subjects may use the belief elicitations at $t = 2$ and $t = 4$ as a soft commitment device to induce themselves to exert effort at $t = 3$ and $t = 5$ respectively. Subjects may thus state a higher belief than their truly held belief and thereby exhibit apparent naiveté. While this is theoretically possible, data in Augenblick and Rabin (2015) suggests that experimental subjects are unlikely to be this sophisticated.

Nonetheless, we will randomize whether the belief elicitations are payoff-relevant or not, with the uncertainty realizing after beliefs have been elicited but before the task. If there is no difference in task completion between states of the world in which the beliefs are payoff-relevant and states of the world in which they are not, then the belief elicitation is not used as a commitment device. As an additional step, we will include a question in the post-experimental survey that asks whether the belief elicitation was used as a soft commitment.

**Randomization of dates.** Relative to the starting date at $t = 1$ the dates of the task will be randomized on the individual level. This is done to avoid correlated shocks to subject's private information about task completion. $t = 1$ will take place on a Monday or Tuesday, $t = 2$ takes place 2 days after $t = 1$. $t = 3$ takes place the next week on Monday, Tuesday, or Wednesday, with equal probabilities. $t = 4$ will happen two days after the realized $t = 3$. $t = 5$ will happen the week after $t = 3$ on Monday, Tuesday and Thursday, with equal probabilities and independently from the realization of $t = 3$. The exact dates are realized before the experiment begins and communicated to the subjects during the initial session in the lab ($t = 1$).

**Short post-experimental survey** We will ask students about their high-school math grade, gender, age, and their parent's income.

**Setting and sample size.** The respective first session a subject participates in will take place at the Munich Experimental Laboratory for the Social Sciences (MELESSA), whereas

later sessions take place online. There will be 10 initial sessions, for which we will invite 22 subjects each. Every subject that shows up will be allowed to participate in the experiment. Given the average show-up rate, we thus expect a sample size of 200 subjects.

## 2.2 Pilot experiment

In order to study learning, we need enough identifying variation in task completion and beliefs about task completion. This makes it essential that our experimental parameters are well calibrated. We ran a pilot with 40 subjects that featured only periods $t = 1$, $t = 2$ and $t = 3$. The pilot could **not** be used to study learning or pre-test any of our main hypotheses. It merely allowed us to select a threshold that hopefully yields enough variation in task completion and beliefs about task completion to allow for an investigation of learning.

# 3 Hypotheses

Denote the binary effort level on dates $t = 3$ and and $t = 5$ as $a_3$ and $a_5$ respectively, where $a_i = 1$ if the subject performs 20 screens or more, and $a_i = 0$ otherwise. Moreover, let $\mathbb{Q}(x)$ denote the frequency of behavior $x$ in the data and $\mathbb{P}_t(x)$ the participants' expectation of that behavior at time $t$. Our dataset consists of:

- each participant's prior beliefs $\mathbb{P}_2(a_3, a_5)$ about the 4 possible scenarios ($(a_3, a_5) \in \{0, 1\}^2$);

- each participant's posterior beliefs $\mathbb{P}_4(a_5 = 1)$;

- each participant's commitment choice $e_3^{\text{commitment}}$ at date 2, and $e_5^{\text{commitment}}$ at date 4; from these measures we construct the binary variables $a_3^{\text{commitment}}$ and $a_5^{\text{commitment}}$, equal to 1 if the number of tasks chosen exceeds the threshold, and 0 otherwise;

- each participant's effort level $e_3$ at date 3, and $e_5$ at date 5; from these measures we construct the binary variables $a_3$ and $a_5$.

The following measures will also be constructed and used in the analysis:

- at the individual level, the perceived marginal probabilities $\mathbb{P}_2(a_3)$ and $\mathbb{P}_2(a_5)$, and the perceived conditional probabilities $\mathbb{P}_2(a_3 \mid a_5)$ for all values of $a_3, a_5$. These objects are not elicited directly but constructed from the primitive dataset. For instance, the conditional probability $\mathbb{P}_2(a_3 = 1 \mid a_5 = 1)$ is defined as

$$\frac{\mathbb{P}_2(a_3 = 1, a_5 = 1)}{\mathbb{P}_2(a_3 = 1, a_5 = 1) + \mathbb{P}_2(a_3 = 0, a_5 = 1)}.$$

- the empirical distribution $\mathbb{Q}(a_3, a_5)$ of the 4 possible scenarios, that is, the fraction of participants whose behavior is given by $(a_3, a_5)$; similarly, marginal and conditional probabilities are constructed from these 4 probabilities;

- the aggregate posterior beliefs $\mathbb{P}_4(a_5 = 1 \mid a_3 = 1)$ defined as the mean value of $\mathbb{P}_4(a_5 = 1)$ elicited among individuals who passed the threshold at date 3 ($a_3 = 1$); similarly, $\mathbb{P}_4(a_5 = 1 \mid a_3 = 0)$ is elicited among individuals who did not pass the threshold at date 3.

Our hypotheses are divided into preliminary and main hypotheses, where the former are based on the literature and, to some extent, constitute prerequisites for the study of learning.

**Preliminary hypothesis 1** *Individuals are present biased. That is, $\mathbb{Q}(a_3^{commitment} = 1) > \mathbb{Q}(a_3 = 1)$ and $\mathbb{Q}(a_5^{commitment} = 1) > \mathbb{Q}(a_5 = 1)$, meaning that participants are more likely to choose an effort level larger than the threshold ex-ante than on the spot. We will also test this hypothesis with the non-binary effort level (see below).*

**Preliminary hypothesis 2** *Individuals are initially naive or optimistic about their future effort. That is, $\mathbb{P}_2(a_3 = 1) > \mathbb{Q}(a_3 = 1)$.*

**Preliminary hypothesis 3** *Behavior at date 3 is informative about behavior at date 5 and the relationship is positive, i.e.*

$$\frac{\mathbb{Q}(a_3 = 1 \mid a_5 = 1)}{\mathbb{Q}(a_3 = 1 \mid a_5 = 0)} > 1.$$

*This inequality means that passing the threshold at date 3 predicts participants' ability to pass the threshold at date 5. It is equivalent to*

$$\frac{\mathbb{Q}(a_3 = 0 \mid a_5 = 1)}{\mathbb{Q}(a_3 = 0 \mid a_5 = 0)} < 1,$$

*which means that failing to pass the threshold at date 3 positively predicts failing to pass the threshold at date 5.*

As a first step in the analysis of learning, we will then test whether the absolute degree of naiveté, as defined in Ahn et al. (2017), decreases over time.

**Main hypothesis 1**   *The absolute degree of naiveté is decreasing over time, i.e.*

$$\mathbb{P}_2(a_3 = 1) - \mathbb{Q}(a_3 = 1) > \mathbb{P}_4(a_5 = 1) - \mathbb{Q}(a_5 = 1).$$

Our second main hypothesis is that individuals' prior beliefs misperceive the correlation in their future behavior.

**Main hypothesis 2**   *Individuals underestimate the informativeness of their date-3 behavior, i.e.,*

$$\frac{\mathbb{P}_2(a_3 = 1 \mid a_5 = 1)}{\mathbb{P}_2(a_3 = 1 \mid a_5 = 0)} < \frac{\mathbb{Q}(a_3 = 1 \mid a_5 = 1)}{\mathbb{Q}(a_3 = 1 \mid a_5 = 0)}$$

*and*

$$\frac{\mathbb{P}_2(a_3 = 0 \mid a_5 = 1)}{\mathbb{P}_2(a_3 = 0 \mid a_5 = 0)} > \frac{\mathbb{Q}(a_3 = 0 \mid a_5 = 1)}{\mathbb{Q}(a_3 = 0 \mid a_5 = 0)}.$$

Next, we gauge the internal consistency of the learning process. One possible inconsistency in updating behavior is that there is a systematic movement between prior and (average) posterior beliefs, which would contradict the law of iterated expectations.

**Main hypothesis 3**  *Beliefs after observing the task completion are upwardly biased relative to the prior beliefs, i.e.,*

$$\mathbb{P}_4(a_5 = 1) > \mathbb{P}_2(a_5 = 1).$$

Crucially, to understand the mechanism behind hypothesis 3, we will also test for an upward bias in $t = 4$ beliefs *conditioning* on $t = 3$ behavior. This also allows us to ask whether subjects' updating process conditional on the information received is (on average) Bayesian or not. We will analyze whether individuals respond differently to the good signal (task completion) and the bad (failure to complete the task). In particular, subjects may exhibit biases akin to asymmetric updating and conservatism, documented in the literature on updating over ego-relevant characteristics (Eil and Rao, 2011; Mobius et al., 2014; Buser et al., 2017).

Next, we turn to the question of whether subjects are able to transport what they learn about their preferences in one setting into another setting. To this end, we compare your two treatments.

**Main hypothesis 4**  *The underestimation of the informativeness is more severe in the treatment condition, where date 5 involves a different task. Formally, the biases in the perception of the likelihood ratio, equal to*

$$\left| \frac{\mathbb{P}_2(a_3 \mid a_5 = 1)}{\mathbb{P}_2(a_3 \mid a_5 = 0)} - \frac{\mathbb{Q}(a_3 \mid a_5 = 1)}{\mathbb{Q}(a_3 \mid a_5 = 0)} \right|$$

*are larger in treatment $T = 2$ (with two different tasks) than in treatment $T = 1$ (with the same task twice), for $a_3 = 0$ and $a_3 = 1$.*

If the learning process appears to be upwardly biased or conservative (Main Hypothesis 2 or 3), we will test whether this bias is more severe in the treatment condition.

**Main hypothesis 5**  *The bias in the learning process is more severe in the treatment condition, i.e. the biases*

$$|\mathbb{P}_4(a_5 = 1) - \mathbb{P}_2(a_5 = 1)|$$

*are larger in treatment $T = 2$ than in treatment $T = 1$ (for $a_3 = 0$ and $a_3 = 1$). Hypothesis 5 will again be tested separately for the cases of $a_3 = 1$ and $a_3 = 0$.*

# 4  Analysis

For the analysis we will use only the data from the main experiment. We will necessarily exclude all participants for whom the commitment choice was binding at date 3 or at date 5 (approx. 20 out of 200).[3] This leaves us with 180 subjects for whom we observe commitment choices, effort decisions on the spot, and beliefs.

**Preliminary hypotheses.**   To test for present bias, we test whether $\mathbb{Q}(a_3^{\text{commitment}} = 1) > \mathbb{Q}(a_3 = 1)$ using Fischer's exact test. We will also consider the continuous variables and test whether $e_3^{\text{commitment}} > e_3$ by means of a one-sided t-test. To test for naiveté at $t = 2$ and $t = 4$, we will test whether $\mathbb{P}_2(a_3 = 1) - \mathbb{Q}(a_3 = 1) > 0$ and $\mathbb{P}_4(a_5 = 1) - \mathbb{Q}(a_5 = 1) > 0$ respectively using a one-sided t-test. To test for the empirical informativeness of $t = 3$ behavior for $t = 5$ behavior, we test whether $\mathbb{Q}(a_5 = 1 \mid a_3 = 1) - \mathbb{Q}(a_5 = 1 \mid a_3 = 0) > 0$ with Fischer's exact test.

**Main hypotheses.**   As a first step in the analysis of learning, we will test whether $\mathbb{P}_2(a_3 = 1) - \mathbb{Q}(a_3 = 1)$ exceeds $\mathbb{P}_4(a_5 = 1) - \mathbb{Q}(a_5 = 1)$ by means of a one-sided t-test.

To test hypothesis 2 (misperception of the informativeness of behavior) we will calculate the bias in the perceived likelihood ratio, given by

$$Bias_1 = \frac{\mathbb{P}_2(a_3 \mid a_5 = 1)}{\mathbb{P}_2(a_3 \mid a_5 = 0)} - \frac{\mathbb{Q}(a_3 \mid a_5 = 1)}{\mathbb{Q}(a_3 \mid a_5 = 0)}$$

and test whether $Bias_1 > 0$ by means of a one-sided t-test over the pooled data. Note that $Bias_1$ will be calculated for $a_3 = 0$ and $a_3 = 1$.

To test hypothesis 3 (consistency of the updating process), we will construct

$$Bias_2 = \mathbb{P}_4(a_5 = 1) - \mathbb{P}_2(a_5 = 1)$$

---

[3]Except for the case of hypothesis 2, which relies only on $t = 2$ data, so that only approx. 10 observations will need to be excluded.

and test whether $Bias_2 > 0$ with a one-sided t-test over the whole sample. In a Bayesian model, the law of iterated expectations implies that $Bias_2 = 0$. It will also be crucial to construct two measures of the bias conditional on date 3-behavior, i.e. $\mathbb{P}_4(a_5 = 1 \mid a_3 = 1) - \mathbb{P}_2(a_5 = 1 \mid a_3 = 1)$ and $\mathbb{P}_4(a_5 = 1 \mid a_3 = 0) - \mathbb{P}_2(a_5 = 1 \mid a_3 = 0)$ and test whether they are significantly different from zero by means of one-sided t-tests. This constitutes the cleanest test of an upward deviation from Bayesian updating conditional on the information received.[4]Finally, we test for differences between these two measures by means of two-sided t-tests.

Finally, to test hypotheses 4 and 5 we will use one-sided t-tests to test whether the biases defined above are larger in $T = 2$ than in $T = 1$.

**Statistical power.** We test main hypothesis 1 with a one-sided t-test and a sample of 180 observations. We have power of 0.8 to detect biases of a size of 0.185 standard deviations at $\alpha = 0.05$.

We test main hypothesis 2 using a one-sided t-test and a an expected sample of 170 observations.[5] We have power of 0.8 to detect biases of a size of 0.191 standard deviations at $\alpha = 0.05$.

We test main hypothesis 3 with a one-sided t-test and sample of 180 observations, which gives us power of 0.8 to detect biases of a size of 0.185 standard deviations at $\alpha = 0.05$. The updating conditional on realized behavior at $t = 3$ can be tested with an expected sample size of 85.[6] We have power of 0.8 to detect biases of a size of 0.269 standard deviations at $\alpha = 0.05$.

Hypotheses 4 and 5 rely on the treatment comparison and half the samples used for the tests of hypotheses 2 and 3 respectively. For main hypotheses 4 we therefore have power

---

[4]The unconditional $Bias_2$ maybe non-zero simply because a subject has wrong beliefs about how often different $t = 3$ behaviors are realized. Our measures will allow us to see whether this is the case or not.

[5]Hypothesis 2 is based on belief data from $t = 2$, where the first commitment choice will have eliminated 5% of the observations (10 out of 200). Moreover, since likelihood ratios cannot be calculated when there is a zero in the denominator, we expect to lose up to 20 more observations. This leaves us with a sample of 170 observations.

[6]We half the initial 180 observations when we condition on the realized behavior. We expect to discard another 5 observations in each condition for whom the Bayesian posterior is not defined (e.g. because a subject expected to complete the task with probability 1 and then failed to complete the task).

of 0.8 to detect treatment effects of a size of 0.269 standard deviations at $\alpha = 0.05$. For main hypothesis 5, at $\alpha = 0.05$, we have power of 0.8 to detect treatment effects of a size of 0.262 standard deviations in the unconditional data and 0.381 standard deviations when we condition on $t = 3$ behavior.

**Baseline balance**   We will test for the baseline balance between our two treatment groups according to the following characteristics: gender, age, math score and $t = 3$ task completion. If the sample is unbalanced on a variable that correlates with the bias under investigation, we will run OLS regression of the respective bias on a treatment dummy, while controlling for the unbalanced correlate.

# References

Ahn, D. S., R. Iijima, Y. Le Yaouanq, and T. Sarver (2017). Behavioral characterizations of naiveté for time-inconsistent preferences. Working paper.

Augenblick, N., M. Niederle, and C. Sprenger (2015). Working over time: dynamic inconsistency in real effort tasks. *Quarterly Journal of Economics 130*(3), 1067–1115.

Augenblick, N. and M. Rabin (2015). An experiment on time preference and misprediction in unpleasant tasks. Working paper.

Buser, T., L. Gerhards, and J. van der Weele (2017). Responsiveness to feedback as a personal trait. *Journal of Risk and Uncertainty (forthcoming)*.

DellaVigna, S. and U. Malmendier (2006, June). Paying not to go to the gym. *American Economic Review 96*(3), 694–719.

Eil, D. and J. M. Rao (2011). The good news-bad news effect: asymmetric processing of objective information about yourself. *American Economic Journal: Microeconomics 3*(2), 114–38.

Fedyk, A. (2016). Asymmetric naivete: Beliefs about self-control. Working paper.

Mobius, M. M., M. Niederle, P. Niehaus, and T. S. Rosenblat (2014). Managing self-confidence. Working paper.

Read, D. and B. Van Leeuwen (1998). Predicting hunger: The effects of appetite and delay on choice. *Organizational behavior and human decision processes 76* (2), 189–205.