# Learning from Data Through Models

Alberto Bisin, Guillaume Frechette and Jimena Galindo

July 13, 2023

**Abstract**

TBW

# 1 Introduction

# 2 Literature Review

# 3 Framework and Predictions

A finite number of observable random variables, $X = (x_1, x_2, ..., x_N)$, determine a binary state, $y$. The observable variables are independently distributed and each follows a Normal distribution with corresponding mean $\mu_i$ and standard deviation $\sigma_i$. The state is *Red* (**R**) or *Blue* (**B**) and is determined in the following way:

$$y = \begin{cases} \mathbf{B} & \text{if } a_1 x_1 + ... + a_N x_N + K \geq 0 \\ \mathbf{R} & \text{otherwise} \end{cases}$$

An observable variable $x_i$ is *relevant* if the associated coefficient, $a_i$, is non-zero. Likewise, variable $i$ is *irrelevant* if $a_i = 0$. We call the set of relevant variables $\mathcal{R}$ and the set of irrelevant variables $\mathcal{I}$.

1

In each period, an agent observes the realization of all observable variables and predicts the state. Their flow payoff in period $t$ is 1 if they predict the state correctly, and it is zero otherwise. Denote the prediction by $\hat{y}$.

In order to make a prediction, the agent uses a map $M_t \rightarrow \{\mathbf{B}, \mathbf{R}\}$ where $M_t \subseteq \{x_1, x_2, ..., x_N\}$ is the set of variables that they choose to consider in period $t$. Our main object of study is the set of variables that the agent decides to consider in a given period, and the procedure by which the maping or the set are updated. We refer to the set of variables that are taken into consideration, $M_t$, as the agent's *mental model*, and assume that the mapping is linear.

A mental model is *simple* if it does not include all the relevant variables, *i.e.* if $M \subset \mathcal{R}$. And it is *complex* if it includes all relevant variables, $M = \mathcal{R}$. In addition, mental models can be ranked in terms of their simplicity by the number of relevant variables that they consider.

Whenever agents use simple models, they ignore some of the relevant information that is available to them. In contrast, whenever agents use complex models, they use all the relevant information available. This means that when agents with complex mental models are restricted in terms of how many variables they can consider, their ability to predict the state will be impaired. On the contrary, restricting the number of variables that an agent with a simple model can use, does not necessarily affect their ability to predict the state. In particular, if they are considering fewer variables than what the restriciton allows, their predictions should be as accurate as without the restriciton. These two effects are captured in predictions 1 and 1B respectively.

**Prediction 1:** If subjects use complex models, then allowing them to use an additional variable will always increase the frequency with which they make correct predictions.

**Prediction 1B:** If subjects use a simple mental model that considers $K$ variables, allowing them to use $K + 1$ variables will not increase the frequency with which they make correct predictions.

## 3.1  Optimal Rules

Conditional of using a particular mental model, the most accurate prediction rule that the agent can use is one that predicts the state to be $\mathbf{B}$ whenever, given the observables under consideration, $\mathbf{B}$ is more likely. And it predicts $\mathbf{R}$ whenever $\mathbf{R}$ is more likely. Because we abstract away from the procedure through which the rule is updated, we will make predictions under the assumption that the agent has learned the parameters perfectly. That is, they have seen enough information that, under an updating procedure that converges to the truth—for example, correctly specified regression model or Bayes rule—they would have converged to the true parameters already. This provides an upper bound for how well any model can perform.

We refer to a prediction rule as an *Optimal Rule* if it is the rule that a fully rational agent would converge to with perfect Bayesian learning. In this context, Bayesian learning corresponds to using Bayes rule to learn the coefficients $a = (a_1, a_2, ..., a_N)$ as well as the constant $K$. This means that in an $N$ dimensional space, the agent must learn $N + 1$ parameters which corresponds to learning the halfspace in which the state is $\mathbf{B}$. The optimal rule for the complex model coincides with the true data-generating process when learning is perfect and the parameters are identified.

An optimal rule conditional on model $M$ is the rule that, when restricted to the variables being considered, does the best at predicting the state. In this case, the agent is learning only the parameters that pertain the variables in $M$ plus a constant term. The following proposition illustrates exactly what the optimal prediction rules look like when the agent is able to learn perfectly.

**Proposition 1** *Fix any model that considers $m \leq N$ variables and relabel the variable so that $x_1, ..., x_m$ are the variables that the model considers and $x_{m+1}, ..., x_N$ are the ones that the model does not consider. Relabel the coefficients $a = (a_1, ..., a_N)$ accordingly as well.*

*Proof.* Let $M := a_1 x_1 + ... + a_m x_m$ and $M^c = a_{m+1} x_{m+1} + ... + a_N x_N + k$ and define the

latent variable $y^L := M + M^c$. The optimal prediction rule for model $M$ predicts the state to be **B** whenever $P[y^L|M] \geq \frac{1}{2}$ and predicts **R** otherwise.

Using the fact that for a random variable $z \sim N(\mu_z, \sigma_z)$ we have that $P[z \geq \mu_z] = \frac{1}{2}$ and noticing that $M^c$ and $y^L|M$ are Normally distributed with means $\mathbb{E}[M^c]$ and $M + \mathbb{E}[M^c]$ respectively, it is easy to see that $P[y^L \geq 0|M] \geq P[y^L \geq M + \mathbb{E}[M^c|M]] = \frac{1}{2}$ whenever $M + \mathbb{E}[M^c] \geq 0$. Similarly, whenever $M + \mathbb{E}[M^c] < 0$ we will have that $P[y^L \geq 0|M] < P[y^L \geq M + \mathbb{E}[M^c|M]] = \frac{1}{2}$. Therefore, the optimal rule is to predict **B** whenever $a_1 x_1 + ... + a_k x_k \geq -\mathbb{E}[M^c]$ and **R** otherwise. $\square$

With limited data, it is not guaranteed that the agent is able to perfectly learn the parameters required for the optimal rule. However, the notion of optimal rules is useful for determining a benchmark for how well each model can perform, as well as the scenarios in which we can expect polarization to arise. The next section uses the notion of optimal rules to formalize the concept of polarization and to establish the main prediction of the theory.

## 3.2 Polarization

We take the definition of polarization from Haghtalab et al. [2021]: two agents are said to be polarized if, when observing the realization of X, they predict the state to be different. That is, they both observe the same draw of $(x_1, ..., x_N)$ but one agent predicts the state to be **B** and the other predicts the state to be **R**.

Within out framework—as in Haghtalab et al. [2021]—polarization may persist even when agents have access to unlimited data. There are two factors that dare necessary for polarization to arise: the use of different models and the realization of the observables must be such that the models will make different predictions. It is important to observe that even when agents are using different models, not all realiations of $X$ will produce polarized predictions. The following 2-dimensional example illustrates why the use of different models

is necessary but not sufficient for polarization to arise.

**Example.** *The optimal rule for the complex model coincides with the truth, it predicts $\boldsymbol{B}$ if $x_2 - 1 \geq x_1$. In addition, there are two simple models: one that considers only $x_1$—call it $M_1$—and one that considers only $x_2$—call it $M_2$. The optimal rules for $M_1$ and $M_2$ are threshold rules that predict the state is $\boldsymbol{B}$ if the value of $x_i$ is below (above) the threshold $t_i := \mu_{-i} + k$ and the predicted state is $\boldsymbol{R}$ otherwise. These prediction rules are illustrated in . As Figure 1 makes clear, both rules agree on the predicted state for some values of $(x_1, x_2)$ and disagree for other values. The regions in which the rules make different predictions are the disagreement zones.*
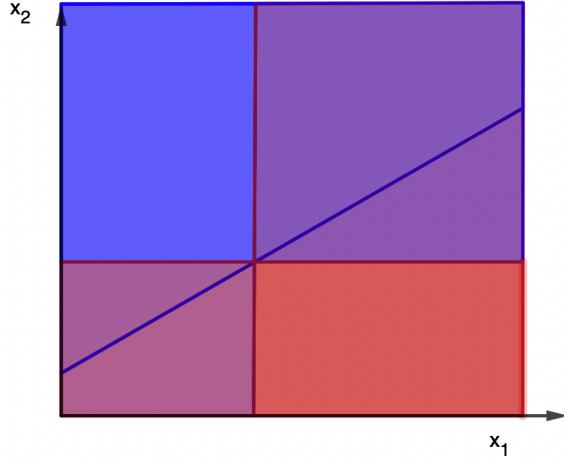


Figure 1: Polarization regions for two single variable models

For any pair of agents who observe the same realization of the observables and use different models, it is possible to determine whether they will be polarized or not by looking at the predictions made by the optimal rule corresponding to their models. Polarizing observations—those for which the predictions are different—are different for each pair of models. Also notice that we determine whether an observation is polarizing or not by using the optimal rules. We do not expect subjects in our experiment to converge to the optimal rule, this is a theoretical construct that allows us to formalize Prediction 2.

**Prediction 2 :** Polarization will arise more often when the theory predicts so. That is,

when two agents who have different models of the world, are facing a realization of observables that polarizing for their models.

# 4    Experimental Design

Subjects for the experimet were recruited from NYU's CESS-lab subject pool. We recruited 30 undergraduate students who participated in this experiment in person during the summer and fall of 2023. The experiment was coded using oTree (Chen et al. [2016]) and it consisted of two parts: Part 1 was meant to expose the subjects to the data generating process and allow them to learn how to predict the state; part 2 was meant to elicit the models of the world that subjects might have developed in part 1.

Part 1 consited of 20 rounds in which subjects observed the realizations of all of the observables and were asked to predict the state. We had 5 observable variables called variable 1, 2, 3, 4 and 5. The variables were independent and Normally distributed with means and variances as described in Table 1.

The state was determined by the following linear rule: $11x_1 + 6x_2 + 4.5x_3 + 2.4x_4 + k \geq 0$. The values of the coefficients were chosen so that all variables had a similar but unique level of informativeness when considered on their own. The constant $k$ was chosen so that the state was **B** in 50% of the cases. The last column of Table 1 shows the probability of the optimal prediction rule making a correct prediction when for each single-variable model. This is the measure that we use to determine the informativeness of each variable.

Table 1: Parameters for the Data Generating Process

| Variable | Mean | Variance | Modified Variable | Modified Variance | Informativeness |
|---|---|---|---|---|---|
| $x_1$ | 0 | 1 | $11x_1$ | 11 | .65 |
| $x_2$ | -10 | 2 | $6x_2$ | 12 | .67 |
| $x_3$ | 5 | 3 | $4.5x_3$ | 13.5 | .69 |

| Variable | Mean | Variance | Modified Variable | Modified Variance | Informativeness |
|---|---|---|---|---|---|
| $x_4$ | -5 | 4 | $2.4 * x_4$ | 12 | .63 |
| $x_5$ | 100 | 5 | $0x_5$ | 0 | 0 |

In each round, subjects were asked to predict the state. They were given feedback on whether their prediction was correct or not and they had access to the entire history of the game. Throughout the experiment, subjects were paired with another participant at random. Both subjects in a pair were shown exactly the same realization of the variables and had to predict the same state. This allows us to determine whether the pair was polarized for that particular observation. Subjects were not told that there was someone else observing the same realizations as them.

Part 2 started immediatly after part 1 and had 70 rounds[1]. In instead of showing subjects the realizations of all 5 variables, they were told that they can disclose up to $m$ variables, where $m$ was drawn at random from the set $\{1, 2, ..., 5\}$. m was redrawn independently across rounds and across subjects.

We use the random assignment of $m$ to investigate whether the subjects are using siumple or complex models. If allowing subjects to disclose an additional variable improves their prediction accuracy, then it must be that they are in fact using the additional information for prediciton. If, on the other hand, they do not perform better when they have access to more information, it must be that, although they are disclosing an additional variable, they do not correctly incorporate the infomation into their prediction rule. Thus, if we observe that the performance of our subjects does not improve when they are allowed to disclose $m + 1$ variables, relative to their performance when they were allowed to disclose only $m$ variables, their model of the world must be of size $m$ or lower. This feature design allows us to test predictions 1 and 1B and determine the number of variables that subjects are using.

---

[1]In the first 2 sessions we had 40 rounds in part 2 and the experiment was done in under 30 minutes. Because the experiment concluded in less time than expected, we increased the number of rounds to 70 for the following sessions. Only 12 of our subjects participated in the sessions with 40 rounds, all others were in sessions with 70 rounds for part 2

In order to explore the polarization prediction, we use the fact that subjects are in fixed pairs throughout the experiment. Both subjects in a pair observe the same realization of all the variables in each round. Therefore, we can determine whether they are polarized or not by looking at whether they make the same prediction or not. We can then determine whether the theory predicts that they should be polarized or not by looking at what the optimal rules for those models dictate. We use the optimal rules that correcpond to the variables that they chose to reveal in that round. It could be that although they are revealing certain variables, they are not using them for prediction. Since we do not have a good method to determine which variables they actually pay attention to, we use the variables that they reveal as a proxy for the variables that they are using for prediction and apply the optimal rules for those models.

# 5   Results

In this section we present the results from the experiment. We start by exploring the general learning patterns of the subjects as well as the model choices. We then look at the effect of allowing subjects to disclose an additional variable and lastly we look at the polarization results.

## 5.1   Learning

In the first part of the experiment, subjects were asked to predict the state given the realization of all the variables. Having access to all the information they could predict the state perfectly if they managed to learn the parameters of the hyperplane that determines the state. Figure 2 shows the share of correct predictions across rounds in part 1. We see that subjects are able to learn how to predict the state at a rate that is higher than random. They initially predict the state corectly about 50% of the time, which is consistent with random guessing. The share of correct guesses increases to 59% by the end of part 1. And this last

number is significantly higher than 50% (p-value = 0.00031). They continue to learn even after the first 20 rounds when we look at the rounds in which subjects had access to all the information. By the end of the experiment the rate of correct predictions is 59% which is still significantly higher than 50% (p-value = 0.00000019).
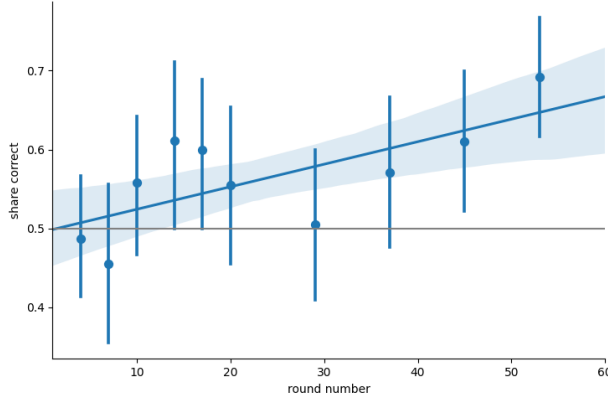


Figure 2: Share of correct predictions across rounds by performance in rounds 20 to 40

We also observe that there is heterogeneity in how our subjects learn. In particular, in the first 20 rounds of part 2 we are able to identify two distinct groups of subjects: those who guess more than 50% of the states correctly and those who guess 50% of fewer of them correctly. As can seen in 3, these two groups have very different learning patterns throughout the experiment. The group that does better than random in rounds 20 to 40 (the first 20 rounds of part 2) consists of 20 subjects, which account for 67% of the sample. These subjects improved their performance from 50% to 64% in the first 20 rounds of part 2, and continued to improve to get to 69% by the end of part 2. Meanwhile, the group that does worse than random in rounds 20 to 40 consists of 10 subjects and they do not seem to learn how to predict the state more accurately even after all 60 rounds of the experiment. Whenever it is relevant, we will show results separately for these two groups.

## 5.2   Simplicity

In part 2, subjects were allowed to disclose up to $m$ variables before making their predictions. and $m$ was assigned at random for every subject in every round. We use this random
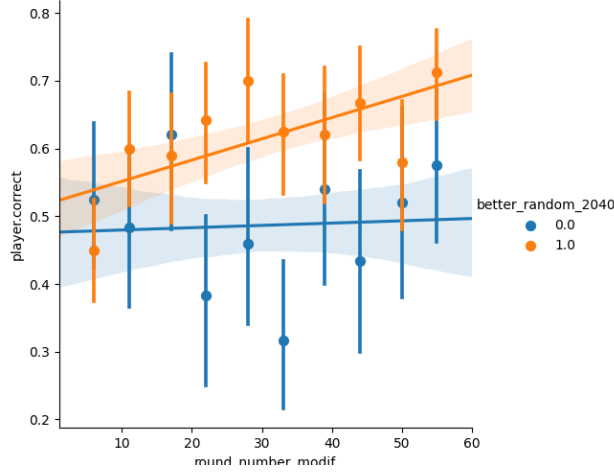
Figure 3: Share of correct predictions across rounds

assignment to understand if allowing them to reveal an additional variable improves their prediction accuracy. If it does, it must be that they are using the additional information for prediction. Figure 4 shows the share of correct predictions by the number of variables that subjects chose to disclose.
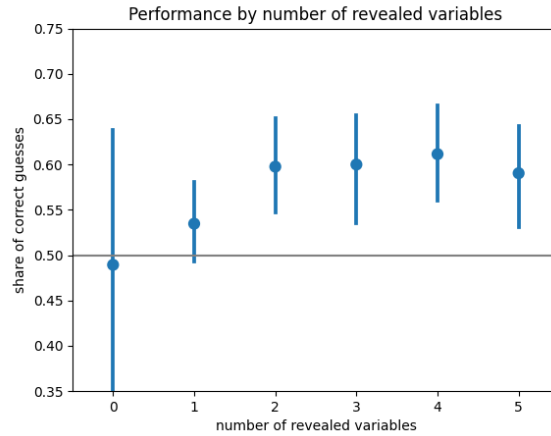


Figure 4: Share of correct predictions by number of disclosed variables

We see that there is a sharp increase when going from 1 to 2 variables, but no significant increase when going from 2 to 3 and beyond. This suggests that on average subjects are using mental models that consider 2 variables. However, since subjects were allowed to reveal fewer variables than what was assigned to them, the estimates in 4 might be biased. To account for this, we look at the share of correct predictions by the number of variables that subjects

10

were allowed to disclose. The results are presented in Table

|  | Disclosed | Allowed |
|---|---|---|
| Zero variables |  | 0.489 *** |
|  |  | (0.072) |
| One variable | 0.534 *** | 0.535 *** |
|  | (0.027) | (0.023) |
| Two variables | 0.602 *** | 0.598 *** |
|  | (0.027) | (0.027) |
| Three variables | 0.584 *** | 0.600 *** |
|  | (0.027) | (0.029) |
| Four variables | 0.599 *** | 0.611 *** |
|  | (0.025) | (0.028) |
| Five variables | 0.576 *** | 0.590 *** |
|  | (0.027) | (0.030) |
| N. obs. | 1740 | 1740 |
| R squared | 0.580 | 0.581 |

*** $p < 0.001$; ** $p < 0.01$; * $p < 0.05$.

## 5.3 Model Choices

## 5.4 Polarization Results

# 6 Conclusion

# Appendix

Daniel L. Chen, Martin Schonger, and Chris Wickens. otree—an open-source platform for laboratory, online, and field experiments. *Journal of Behavioral and Experimental Finance,*

9:88–97, 3 2016. ISSN 2214-6350. doi: 10.1016/J.JBEF.2015.12.001.

Nika Haghtalab, Matthew O. Jackson, and Ariel D. Procaccia. Belief polarization in a complex world: A learning theory perspective. *PNAS*, 118, 2021.