

From predictions to insights: an empirical study of risky choice behavior

Ana B. Barcenas & Yolanda Diao

Duke University

Master of Interdisciplinary Data Science

Table of Contents

- I. Introduction
- II. Literature Review on CPC18 & Data Description
- III. Aggregate-Level Prediction
 - A. Data Preprocessing and Exploratory Data Analysis
 - B. Traditional Machine Learning Models
 - C. Neural Network Model
- IV. Individual-Level Prediction
 - A. Data Preprocessing and Exploratory Data Analysis
 - B. Naïve model and Factorization Machines
 - C. Proposed approach
- V. Applications in Financial Industry
 - A. Next steps
- VI. Supplemental Material
 - A. Theoretical insights
 - B. Sampling tools
 - C. Predictions
 - D. Individual variation in predictions

I. Introduction

“Risky choices, such as ... whether or not to go to war, are made without advance knowledge of their consequences. Because the consequences of such actions depend on uncertain events ... the choice of an act may be interpreted as the acceptance of a gamble that can yield various outcomes with different probabilities” (Kahneman & Tversky, 1983, pp. 341). A risky choice occurs whenever we choose between multiple uncertain prospects. Within an uncertain prospect, there are different possible outcomes with associated probabilities. It is well known in behavioral decision research that individuals tend to deviate from maximizing expected values when facing risky choices. Thus, predicting decision-making has become a field in itself due to the complexity of incorporating the effect of behavioral anomalies.

In the present work, we aim to predict how groups and individuals make choices in high-risk situations. To elucidate these, it is first necessary to distinguish between two types of risky choice prediction. At the aggregate-level, risky choice prediction entails predicting the rate at which an uncertain prospect will be chosen among the population. The aggregated-level predictions will be made based on the characteristics of the uncertain prospect (associated probabilities, payoffs, etc.). At the individual-level, risky choice prediction aims to predict a decision-maker’s choice based on some information about that individual (previous choices, psychological and demographic information, etc.).

Our first goal is to replicate and outperform the performance in aggregate-level predictions of the Choice Prediction Competition 2015 and 2018 (CPC15 & CPC18) winner models (<https://cpc-18.com/>). Our second goal is to replicate the individual-level predictions of the CPC18 winners and build new models more applicable in the financial industry such as

commercial banks or insurance companies. To doing so, we integrate psychological theories proposed by Plonsky et.al. (2017) as additional variables to raise performance.

I. Literature Review on CPC18 & Data Description

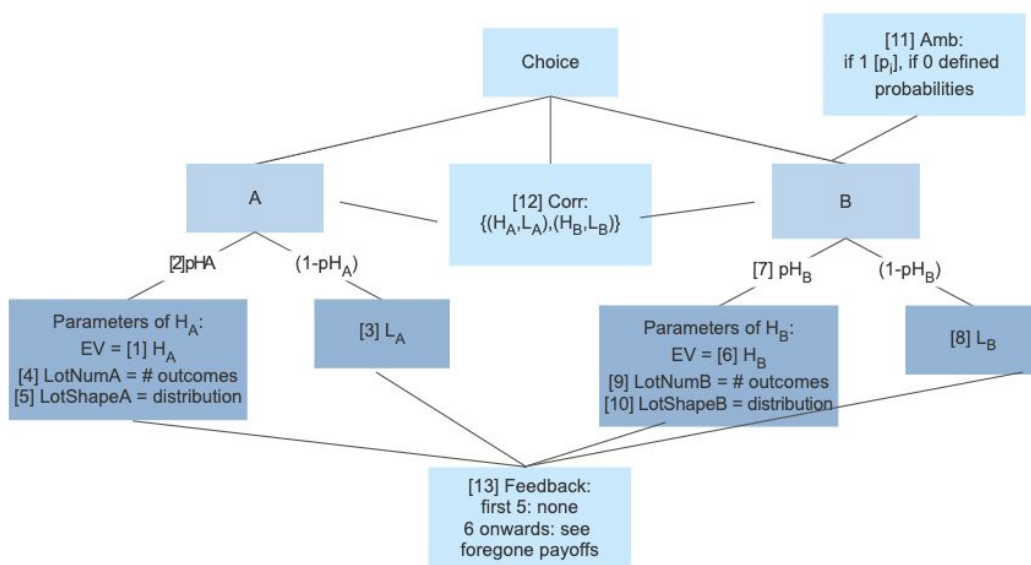
The present work is inspired by and started from the findings of two tournaments recently hosted by researchers in Israel. These Choice Prediction Competitions CPC15 and CPC18 (Erev et al., 2017; Plonsky et al., 2017; Plonsky et al., 2018; Plonsky et al., 2019) challenged participants to provide predictive models that could accurately predict both aggregate and individual behavior in risky choice tasks. One of the shortcomings of theories regarding risky choices is that they tend to explain only a fraction of the observed paradoxes. The organizers of CPC15 (Erev et al., 2017) set out hoping to unify what we have learned in over 50 years of risky choice research into a single framework. To do this they initially outlined 14 risky choice paradoxes. These included the four Kahneman & Tversky (1979) addressed in prospect theory; the certainty effect, the reflection effect, overweighting rare events and loss aversion, but also phenomena ranging from Thaler & Johnson's (1990) breakeven effect to the St Petersburg paradox (Bernoulli, 1954).

To construct the dataset for CPC15, the organizers first hand-selected 30 problems that would manifest each of these 14 behavioral anomalies. They then collected data from participants at two universities in Israel to capture human behavior in these tasks (Erev et al., 2017). They added the results of a calibration study (60 additional problems) to the original training set, withholding the results of a separate test set study (another 60 additional problems). The results of both of these studies were subsumed into the training set of CPC18. Furthermore, the final training data of CPC18 included the results of a new experiment (a

final, further 60 additional problems). This brought the total training data for the aggregate-level prediction track to 530,000 choices. The test set for this track was 180,000 unseen choices in 60 new problems.

It is necessary to be concrete about the dimensions of the problem space of CPC18. At a basic level, each problem in the space is a choice between Option A and Option B. Between CPC15 (Erev et al., 2017) and CPC18 (Plonsky et al., 2018) the organizers added one parameter - the number of outcomes in lottery A. This used to be fixed at 2. This brings the total of dimensions of the problem space in CPC18 to 12 (see Figure 1). Ten out of twelve of the dimensions can be described as parameters of the payoff distributions of the risky prospects A & B. The index [13] makes salient the temporal sequence of seeing the choices with and without feedback. Feedback plays an important role in the CPC18 paradigm. Some risky choice paradoxes emerge without feedback - these are decisions from description. However, other phenomena arise when we rely on feedback - these are decisions from experience. Participants in the studies faced each problem first without feedback and then with full feedback (seeing obtained and foregone outcomes following each choice).

Figure 1



Both options A & B offer a low but certain payoff with a fixed probability and a gamble-within-gamble with a fixed probability. In a gamble-within-gamble there are parameters for the expected value, the number of different outcomes and the skew of probabilities over this set. The distribution of outcomes in Lottery A [5] defines whether the distribution is symmetric around its mean, right-skewed, left-skewed or undefined (if LotNum = 1). The 11th parameter is ambiguity. This is an important parameter to consider. In decisions under risk, we do not know what outcome will be realized, because there is uncertainty over states of the world. This is simplified into two extreme cases Amb = 1 (no initial info concerning these probabilities) and Amb = 0 (complete info and no ambiguity). Finally, the Correlation Parameter [12] captures whether there is a correlation (positive, negative, or zero) between the payoffs of the two options (A & B).

CPC18 consisted of two tracks: Track I, the prediction of aggregate-level behavior and Track II, the prediction of individual-level behavior. For Track II, 30/240 people from Experiment 1 (Plonsky et al., 2018), participants were asked to predict their average rate of choosing B in 5 blocks of 5 problems. The data available to make these predictions were; the choices those 30 people made in the 5 blocks (of 5 trials) for the other 25 problems, the aggregate choices rates of all other decision-makers (at least 60 people, as at least 90 decision-makers answered each problem) in those problems in each block and the choices decision-makers made in alternative problems.

The organizers (Erev et al., 2017) constructed a process model 'BEAST' (Best Estimate and Sampling Tools) based on psychological theories and sampling tools, which is set as the baseline model. Theoretical insights and summaries of the BEAST model is put in Supplemental Materials, Section A.

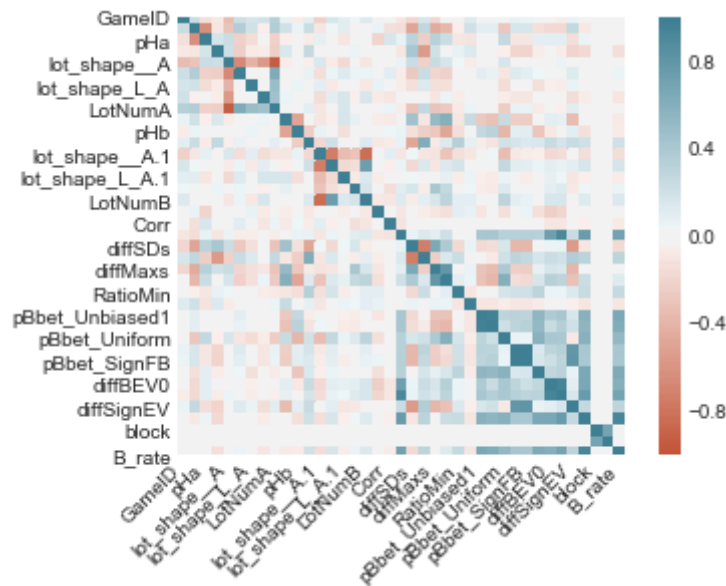
II. Aggregate-level Prediction

A. Data Preprocessing & Exploratory Data Analysis

The dataset for aggregate-level prediction, as described in the last section, is transformed from the raw data by aggregating the population's response over each gamble problem and block. The data set is of size 1050 (210 problems * 5 block) * 39 variables with no missing value or outlier.

Categorical variables are transformed into numerical type through one-hot encoding. Each column of data is scaled into a standard normal distribution. Variables with high correlation (0.8) are removed to eliminate potential multicollinearity, as shown in Figure 2.

Figure 2 Variable Correlation Heat Map



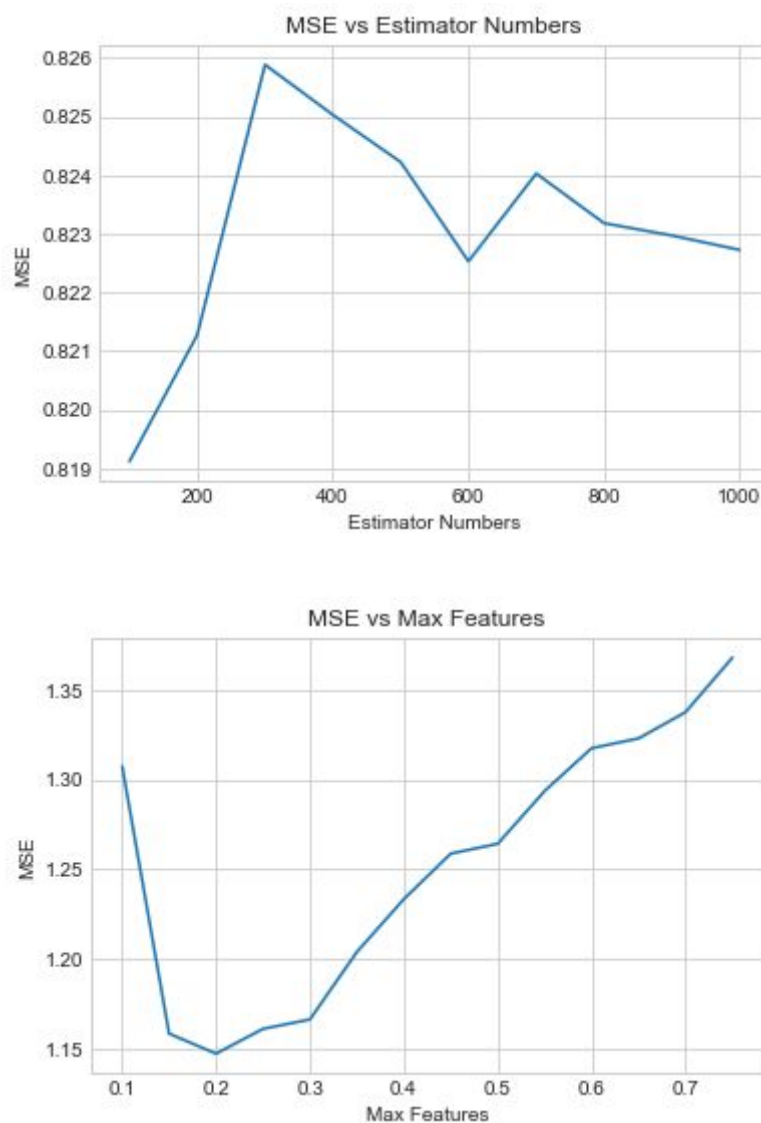
B. Traditional Machine learning Models

The competition participants tested many kinds of algorithms, including; random forest, neural nets (with various architectures), support vector machines with radial and polynomial kernels and k-nearest neighbors (with various specifications). They found a

simple random forest algorithm using psychological inputs to be the only method that improved performance compared to the BEAST model.

We further fine-tuned the parameters of the random forest and extreme gradient descent (Xg-Boosting) algorithms using grid search. By assuming the loss function is continuous in the parameter space, we successively iterated through the possible range of each parameter. The parameter with the smallest mean squared error (MSE) was chosen, as demonstrated in Figure 3.

Figure 3 Fine-Tuning Xg-Boosting Model Hyperparameters with Grid Search



By fine-tuning the hyper-parameters, we were able to achieve a better MSE compared to the competition winner, as shown in Table 1.

Table 1 Performance Comparison btw CPC18 Winner and Our Model

| Algorithm | Features | | |
|---------------|--|---|-------------|
| | Obj. + Naive + Psych. (Ori, et al.) | Obj. + Naive + Psych. (our implementation) | |
| | MSE * 100 | MSE * 100 | R - Squared |
| Random Forest | 0.87 | 0.815 | 0.89 |
| XGBoost | - | 0.84 | 0.89 |

C. Neural Network Model

One surprising result of the two choice prediction competitions (Erev et al., 2017; Plonsky et al., 2018; Plonsky et al., 2019) was the failure of deep learning approaches. Previous work has shown that deep learning could outperform behavioral models in the prediction of boundedly rational behavior in economic games (Hartford, Wright & Leyton-Brown, 2016). The failure in performance is mainly due to data scarcity and a team of researchers (Bourgin et al., 2019) came up with an innovative methodology to overcome the problem by creating synthetic data, which we followed after.

To overcome data scarcity, 3 steps are taken: (1) create synthetic choice predictions from the BEAST model that (2) will serve as cognitive priors to pre-train machine learning models. Once the models are pre-trained, (3) the small-sized real human choice dataset will be employed to fine-tune the parameters of the model to better capture real human behavior.

To create the synthetic choice dataset, we used the problem selection algorithm that Plonsky et al. (2018) developed to draw random problems from the sample space detailed above. The algorithm consists of (1) constraining each of the 11 dimensions, (2) ensuring that no problem was duplicated, and (3) removing those problems where both gambles have the same distribution and payoffs or at least one gamble had no variance, but the payoffs were correlated. After 85 thousand problems were generated by this algorithm, we input them to the BEAST model and got as a result 85,000 synthetic choice predictions that followed the most common behavioral phenomena.

Then, we used these 85,000 synthetic data as a cognitive model prior to pre-train a neural network model as Bourgin et al. (2019) proposed. The neural network model consists of 5 layers with 3 hidden layers. Our input layer had 13 neurons because our problems had 13 dimensions. The output layer had 1 neuron because we were predicting a 1-dimension probability of choosing gamble B over gamble A. The initial hyperparameters were set as in table 2 based on literature reference and our knowledge.

Table 2 Initial Hyperparameters of the Neural Network Model

| Initial Hyperparameters | |
|-------------------------|---------|
| Activation Functions | ReLU |
| Layer-wise Dropout Rate | 0.15 |
| Learning Rate | 0.001 |
| Optimizer | RMSProp |
| Batch Size | 100 |
| Epochs | 200 |

After the model was pre-trained, the small-sized real human predictions were used to tune the parameters. Although this procedure let us leverage state-of-the-art machine learning

algorithms to predict choice behavior without sacrificing accuracy due to data scarcity, the neural network model failed to yield satisfying results, with its $MSE \times 100$ as large as 2.087.

III. Individual-Level Prediction

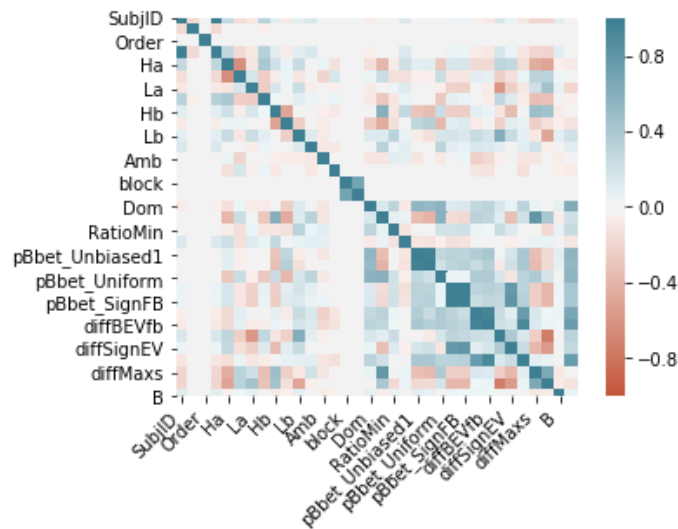
Following the structure proposed by the organizers of CPC18, the individual-level track aims to predict the progression of behavior over time in an individual level.

A. Data Preprocessing and Exploratory Data Analysis

The dataset for individual-level prediction, as described in section I, is transformed from the raw data by averaging the response of the same participant in the same block. The size of the data is 102,150 (20,430 problems * 5 blocks) * 38 variables and there are no missing values or outliers.

We first transformed the categorical variables into numerical types using one-hot encoding. Then, we standardized each column of data into standard normal distribution. Variables with high correlation (0.8) are removed to eliminate potential multicollinearity (Figure 4).

Figure 4



Surprisingly, we did not find a significant relationship between the demographic features (gender, location, and age) and their preference rate over gamble choices. This indicates that the individual differences in risky choices lie in the different choice history of each individual but not the demographic features themselves. We believe this why almost all individual-level models fell behind the naïve model, which will be introduced in the next section.

B. Naïve Model and Factorization Machine

As baseline models, the organizers of the competition proposed two models and its associated MSEs (Table 3): a naïve approach and the use of factorization machines (FM). The naïve baseline predicts that each individual decision-maker in each block of its individual target problems behaves the same as the average population behaves in the same block of that problem. This baseline was hard to defeat by more sophisticated models. The second baseline is factorization machines, which improves on the performance of the naïve baseline by a small margin.

Table 3.

| Baseline model | Features | MSE |
|-----------------------|-----------------------|------------|
| Naïve | Game ID Subject ID | 0.1038 |
| Factorization machine | Game ID Subject ID | 0.0976 |

Factorization machines combine Support Vector Machines (SVM) with factorization models (Rendle, 2012). Similar to SVM models, FMs are a general predictor working with any real value feature vector. In contrast to SVMs, FMs model all interactions between

variables using factorized parameters. Thus they are able to estimate interactions even in problems with huge sparsity, which made it better for our case compared to SVMs.

In this case, each observation supplied to the baseline implementation of the FM is composed of a long binary feature vector with only two non-zero elements that correspond to the active decision-maker and the active block within an active problem. The response is the observed choice rate of the active decision-maker in the active block of the active problem (first transformed to imply the maximization rate of the problem, and then after making the prediction transformed back to implying the choice rate of Option B). This means the FM model did not directly use the knowledge that behavior across different blocks of the same problem is likely correlated.

C. Proposed approach

As mentioned before, the baseline methodologies do not use other variables apart from the subjectID (active decision-maker) and gameID (active choice problem). We proposed a novel approach that leverages the predictive power of psychological features and we also use the features with greater correlation to the target variable such as gender and characteristics of the decision problem (probabilities, payoffs, etc.). The models we proposed, its most important features, and its associated MSEs are summarized in Table 4.

Table 4.

| Algorithm | Most important features | MSE |
|---------------|---|--------|
| Fixed effects | (1) Sign heuristic and (2) tendency to minimize immediate regret. | 0.1224 |
| Random forest | (1) Difference on EVs. | 0.1044 |
| XGBoost | (1) Driven by stochastic dominance, and (2) difference on EVs. | 0.0988 |

The psychological features we use are proposed and implemented in this dataset by Plonsky et.al. (2007) and consist of 13 variables that aim to capture direct research made by social scientists on decision making and the psychology of choice. Some of the behavioral anomalies captured by these features are the sensitivity to the payoff sign and sensitivity to the difference between the gambles' expected values. As shown in Figure 6, the variables with higher predictive value are psychological features. In particular, these features are related to (1) high sensitivity to the payoff sign, (2) the tendency to minimize immediate regret, that implies a preference for the option that produces a higher outcome most of the time, (3) the difference between expected values, and (4) the preference for the option that dominates in terms of payoffs and probabilities when the choice problem is trivial.

Lastly, it is worth noting that the XGBoost model outperforms the naïve baseline and is slightly greater than the factorization machine.

IV. Applications in Financial Industry

Many financial institutions rely on predicting its users' decision to provide the right product to the right person at the right time. Leveraging the predictive power of psychological features to better understand their customers could be a breakthrough for banks and other financial institutions.

A. Next steps

As part of the efforts we made to apply the insights we get from the abstract decision-making predictions, we reached out to banks and insurance companies offering to work with their customer's data. Right now, we are in collaboration with The Standard Bank of South Africa, which eventually will provide us with real data of their customer decisions

on the products they offer. We are currently working on collecting a small set of applied decisions from Amazon's Mechanical Turk that will simulate the data we could get from The Standard Bank of South Africa and the psychological features to predict customers' decisions.

Data collected from Mechanical Turk, compared to the data used in our research, would have a better representation of the real-world financial decisions people face. Instead of using gamble games to capture behaviour tendencies, participants would be asked to choose between real financial products and portfolios. Models built upon this data set are expected to better capture people's choices in the real world.

Our work laid a solid foundation in predicting these choices in 2 ways: 1) we recognized the key factors in risky decisions, as listed in *Section III.c*. These factors, combined with heuristic theories, can guide the design and layout of the questions participants would face. For example, our model indicates people's over emphasis on reward signs. Therefore, portfolios with even one single possible negative return rate should be linked with much higher expected outcomes. 2) The data processing techniques, models we built and the fine-tuning methods could be applied to new data sets with rigorous flexibility and interpretability.

Besides, another interesting topic would be detecting any difference between people's choice behaviours in gamble games and in real finance decisions.

B. Application Scenarios

Companies can take advantage of the ability to predict people's financial choices based on their demographics, historical choices and characteristics of each financial alternative, and make their products more attractive. For both insurance companies and

finance companies, the results can predict clients' preferences over different terms and products, and therefore help with design and recommendations.

IV. Supplemental material

A. Theoretical insights

In their proposition of CPC15, the organizers (Erev et al., 2017) constructed a process model 'BEAST' (Best Estimate and Sampling Tools) to explain the broad panel of risky choice paradoxes they included in their dataset. This sought to incorporate four behavioral tendencies; the equal weighting of possible outcomes, sensitivity to payoff signs, pessimism and the minimization of regret. These tendencies were considered in addition to a (normative) sensitivity to expected value. A key modeling choice the authors made was to assume that the heuristic-based decision tools that people use provide subjective estimates of the value of prospects, rather than deterministic rules with which to make choices.

B. Sampling tools

The four sampling tools mentioned above can be grouped into unbiased and biased behavioral tendencies. The unbiased sampling tool captures the preference for the option or the gamble that produces the higher outcome most of the time. The other three sampling tools constitute a mental draw from distributions that differ from the expected distributions. However, biased behavior decreases when the agent receives full feedback on the payoffs of their decisions. The uniform sampling tool captures the tendency to weight all of the outcomes as equally likely. The contingent pessimism sampling tool represents the belief that

the worst outcomes are more likely than they are. Lastly, the sampling tool sign implies high sensitivity to the payoff sign.

C. Predictions

The BEAST model produces predictions regarding behavior by returning an attractiveness score. This attractiveness is calculated by adding the best estimate of the expected value of the prospect and the average value generated by the use of the four sampling tools. For example, the value of option j based on the use of sampling tools (ST_j) equals the average of k_i (where $i=1,2,3,4$) outcomes that are each generated by using one sampling tool. People are presumed to choose option A if it scores more highly on attractiveness than option B (see Figure 2). If one of the options dominates the other, the error term in this equation is 0. In all other cases, the error term is drawn from a normal distribution with a mean of 0 and a standard deviation of σ_i (a property of agent i). Note that when payoff distributions are known, the best estimation of expected value is the actual expected value of a prospect.

$$[BEV_A(r) - BEV_B(r)] + [ST_A(r) - ST_B(r)] + e(r) > 0$$

(Erev et al., 2017, p.17)

A further improvement was made on the original BEAST model in the BEAST.sd (BEAST subjective dominance) model based on the assumption that the estimation noise is reduced when the problems are perceived as “trivial” (Plonsky et al., 2018). Specifically, according to BEAST.sd, a problem is likely to be perceived as trivial if both the EV rule and the equal weighting rule of sampling tools favor the same prospect, and the choice of that prospect does not lead to immediate regret. While the estimation noise is reduced in “trivial” problems, it is increased in “complex” problems (one option has at least 2 possible outcomes

and the other has at least 3 possible outcomes). What's more, the BEAST.sd model assumes faster learning from feedback in ambiguous problems.

D. Individual variation in predictions

The process by which the BEAST model generates predictions is complicated. One could be forgiven for believing that the above equation should produce a prediction of either A or B for each risky choice, when in fact the output is the probability of choosing option B. This is because the BEAST model relies on six free parameters that delineate the upper limits of individual level draws from uniform distributions (see Figure 3). These allow for heterogeneity within the sample. The process by which individual differences are incorporated into the model is an area for potential improvement. Whilst the parameters are estimated across the sample as a whole using a grid search approach to maximize predictive performance, for any particular individual, their parameter is a random draw.

uniform distributions between 0 and the model's parameters: $\sigma_i \sim U(0, \sigma)$, $\kappa_i \sim (1, 2, 3, \dots, \kappa)$, $\beta_i \sim U(0, \beta)$, $\theta_i \sim U(0, \theta)$, $\gamma_i \sim U(0, \gamma)$, and $\varphi_i \sim U(0, \varphi)$. That is, the model has six free parameters: σ , κ , β , γ , φ , and θ . Note that only four of these parameters (σ , κ , β , and γ) are needed to capture decisions under risk without feedback (the class of problems addressed by prospect theory). The parameter φ captures attitude toward ambiguity, and θ abstracts the reaction to feedback.

(Erev et al., 2017, p.46)

Briefly, it is necessary to consider the mechanics of the BEAST model under ambiguity. When the probabilities of the gamble-within-a-gambles are ambiguous, participants are assumed to estimate them with a pessimistic bias. The initial expected value of the prospect is estimated as a weighted average of three terms; the expected value of prospect A, the minimum payoff of prospect B and the estimated expected value of prospect

B under the assumption of uniform likelihood. The two expected value estimates are equally weighted, and the minimum payoff term is weighted by the parameter ϕ_i , which captures the ambiguity aversion trait of agent i . When the probabilities of a prospect are unknown, feedback enables the adjustment of expected value estimates. Put simply, each trial with feedback moves the best estimate of a prospect's expected value towards its true expected value.

References

- Bernoulli, D. (1954). Exposition of a New Theory on the Measurement of Risk (original 1738). *Econometrica*, 22(1), 22–36.
- Bourgin, D. D., Peterson, J. C., Reichman, D., Griffiths, T. L., & Russell, S. J. (2019). Cognitive Model Priors for Predicting Human Decisions. *ArXiv:1905.09397 [Cs, Stat]*. Retrieved from <http://arxiv.org/abs/1905.09397>
- Erev, I., Ert, E., Plonsky, O., Cohen, D., & Cohen, O. (2017). From anomalies to forecasts: Toward a descriptive model of decisions under risk, under ambiguity, and from experience. *Psychological Review*, 124(4), 369–409. <https://doi.org/10.1037/rev0000062>
- Hartford, J. S., Wright, J. R., & Leyton-Brown, K. (2016). Deep Learning for Predicting Human Strategic Behavior. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, & R. Garnett (Eds.), *Advances in Neural Information Processing Systems* 29 (pp. 2424–2432). Retrieved from <http://papers.nips.cc/paper/6509-deep-learning-for-predicting-human-strategic-behavior.pdf>
- Kahneman, D., & Tversky, A. (1979). Prospect Theory: An Analysis of Decision under Risk. *Econometrica*, 47(2), 263–292.
- Kahneman, D., & Tversky, A. (1983). *Choices, Values and Frames*. Presented at the APA Award Addresses. Retrieved from <http://web.missouri.edu/~segerti/capstone/choicesvalues.pdf>
- Pan, W., & Chen, Y.-S. (2018). Network approach for decision making under risk—How do we choose among probabilistic options with the same expected value? *PLOS ONE*, 13(4), e0196060. <https://doi.org/10.1371/journal.pone.0196060>

- Plonsky, O., Apel, R., Erev, I., Ert, E., & Tennenholtz, M. (2018). *When and how can social scientists add value to data scientists? A choice prediction competition for human decision making*. Presented at the Society of Judgment and Decision-Making, New Orleans. Retrieved from <https://cpc-18.com/wp-content/uploads/2018/03/cpc18-white-paper-march-update.pdf>
- Plonsky, O., Apel, R., Ert, E., Tennenholtz, M., Bourgin, D., Peterson, J. C., ... Erev, I. (2019). Predicting human decisions with behavioral theories and machine learning. *ArXiv:1904.06866 [Cs]*. Retrieved from <http://arxiv.org/abs/1904.06866>
- Plonsky, O., Hazan, T., & Tennenholtz, M. (2017). Psychological Forest: Predicting Human Behavior. *SSRN Electronic Journal*. <https://doi.org/10.2139/ssrn.2816450>
- Rendle, S. (2012). Factorization Machines with libFM. *ACM Transactions on Intelligent Systems and Technology*, 3(3), 1–22. <https://doi.org/10.1145/2168752.2168771>
- Thaler, R. H., & Johnson, E. J. (1990). Thaler_and_johnson.pdf. *Management Science*, 36(6), 643–660.