# Structured, uncertainty-driven exploration in real-world consumer choice

Eric Schulz[a,1,2], Rahul Bhui[a,1], Bradley C. Love[b,c], Bastien Brier[d], Michael T. Todd[d], and Samuel J. Gershman[a]

[a]Department of Psychology, Harvard University, Cambridge, MA 02138; [b]Department of Experimental Psychology, University College London, London WC1H 0AP, United Kingdom; [c]The Alan Turing Institute, London NW1 2DB, United Kingdom; and [d]Data Science Team, Deliveroo, London EC4R 3TE, United Kingdom

**Making good decisions requires people to appropriately explore their available options and generalize what they have learned. While computational models can explain exploratory behavior in constrained laboratory tasks, it is unclear to what extent these models generalize to real-world choice problems. We investigate the factors guiding exploratory behavior in a dataset consisting of 195,333 customers placing 1,613,967 orders from a large online food delivery service. We find important hallmarks of adaptive exploration and generalization, which we analyze using computational models. In particular, customers seem to engage in uncertainty-directed exploration and use feature-based generalization to guide their exploration. Our results provide evidence that people use sophisticated strategies to explore complex, real-world environments.**

exploration | generalization | reinforcement learning | decision making

**W**hen facing a vast array of new opportunities, a decision maker has two key tasks: to acquire information (often through direct experience) about available options and to apply that information to assess options not yet experienced. These twin problems of exploration and generalization must be tackled by any organism trying to make good decisions, but they are challenging to solve because optimal solutions are computationally intractable (1). Consequently, the means by which humans succeed in doing so—especially in the complicated world at large—have proved puzzling to psychologists and neuroscientists. Many heuristic solutions have been proposed to reflect exploratory behavior (2–4), inspired by research in machine learning (5, 6). However, most studies have used a small number of options and simple attributes (7). To truly ascertain the limits of exploration and generalization requires empirical analysis of behavior outside the laboratory.

We study learning and behavior in a complex environment using a large dataset of human foraging in the "wild"—online food delivery. Each customer has to decide which restaurant to pick out of hundreds of possibilities. How do they make a selection from this universe of options? Guided by algorithmic perspectives on learning, we look for signatures of adaptive exploration and generalization that have been previously identified in the laboratory. This allows us not only to characterize these phenomena in a naturally incentivized setting with abundant and multifaceted stimuli, but also to weigh in on existing debates by testing competing theories of exploratory choice.

We address two broad questions. First, How do people strategically explore new options of uncertain value? Different algorithms have been proposed to describe exactly how uncertainty can guide exploration in qualitatively different ways, such as by injecting randomness into choice or by making choices directed toward uncertainty (8). However, results have been mixed, and these phenomena remain to be studied under real-world conditions. Second, how do people generalize their experiences to other options? Modern computational theories make quantitative predictions about how feature-based similarity should govern generalization, which can in turn guide choice. But again

it is unclear whether these theories can successfully predict real-world choices.

Our results suggest that customers explore (i.e., order from unexperienced restaurants) adaptively based on signals of restaurant quality and make better choices over time. Exploration is indeed risky and leads to worse outcomes on average, but people are more likely to explore in cities where this downside is lower due to higher mean restaurant quality. Moreover, we show that customers' exploratory behavior might take into account not only the prospective reward from choosing a restaurant, but also the degree of uncertainty in their reward estimates. Consistent with an optimistic uncertainty-directed exploration policy, they preferentially sample lesser-known options and are more likely to reorder from restaurants with higher uncertainties.

Importantly, we apply cognitive and statistical modeling to customers' choice behavior and find that their choices are best fitted by a model that includes both an "uncertainty bonus" for unfamiliar restaurants and a mechanism for generalization by function learning (based on restaurant features). People appear to benefit from such generalization, as exploration yields better realized outcomes in cities where features have more predictive power. We also show that people generalize their experiences across different restaurants within the same broad cuisine type, defined both empirically within the dataset and by independent similarity ratings. As predicted by a combination of similarity-based generalization and uncertainty-directed exploration, good experiences encourage selection of other restaurants within the same category, while bad experiences discourage this to an even greater extent.

To set the stage for our analyses of purchasing decisions, we first review the algorithmic ideas that have been developed to explain exploration in the laboratory.

---

## Significance

**We study how people make choices among a large number of options when they have limited experience. In a large dataset of online food delivery purchases, we find evidence for sophisticated exploration strategies predicted by contemporary theories. People actively seek to reduce their uncertainty about restaurants and use similarity-based generalization to guide their selections. Our findings suggest that theories of exploratory choice have real-world validity.**

## Prior Work on the Exploration–Exploitation Dilemma

**Uncertainty-Guided Algorithms.** Most of what we know about human exploration comes from multiarmed bandit tasks, in which an agent repeatedly chooses between several options and receives reward feedback (9, 10). Since the distribution of rewards for each option is unknown at the beginning of the task, an agent is faced with an exploration–exploitation dilemma between two types of actions: Should she exploit the options she currently knows will produce high rewards while possibly ignoring even better options? Or should she explore lesser-known options to gain more knowledge but possibly forego high immediate rewards? Optimal solutions exist only for simple versions of this problem (1). These solutions are in practice difficult to compute even for moderately large problems. Various heuristic solutions have been proposed. Generally, these heuristics coalesce around two algorithmic ideas (8). The first one is that exploration happens randomly, for example by occasionally sampling one of the options not considered to be the best (11), or by so-called soft maximization of the expected utilities for each option—i.e., randomly sampling each option proportionally to its value. The other idea is that exploration happens in a directed fashion, whereby an agent is explicitly biased to sample more uncertain options. This uncertainty guidance is frequently formalized as an uncertainty bonus (5) which inflates an option's expected reward by its uncertainty.

There has been considerable debate about whether or not directed exploration is required to explain human behavior (12). For example, Daw *et al.* (12) have shown that a softmax strategy explains participants' choices best in a simple multiarmed bandit task. However, several studies have produced evidence for a direct exploration bonus (4, 13). Recent studies have proposed that people engage in both random and directed exploration (2, 14). It has also been argued that directed exploration might play a prominent role in more structured decision problems (15). However, evidence for such algorithms is still missing in real-world purchasing decisions, where other mechanisms such as coherency maximization have been observed (7, 16).

**Generalization.** Multiple studies have emphasized the importance of generalization in exploratory choice. People are known to leverage latent structures such as hierarchical rules (17) or similarities between a bandit's arms (18).

Inspired by insights from the animal literature (19), Gershman and Niv (20) investigated how generalization affects the exploration of novel options using a task in which the rewards for multiple options were drawn from a common distribution. Sometimes this common distribution was "poor" (options tended to be nonrewarding), whereas sometimes the common distribution was "rich" (options tended to be rewarding). Participants sampled novel options more frequently in rich environments than in poor environments, consistent with a form of adaptive generalization across options.

Schulz *et al.* (21) investigated how contextual information (an option's features) can aid generalization and exploration in tasks where the context is linked to an option's quality by an underlying function. Participants used a combination of functional generalization and directed exploration to learn the underlying mapping from context to reward (22).

## Results

We looked for signatures of uncertainty-guided exploration and generalization in a dataset of purchasing decisions from the online food delivery service Deliveroo (see *Materials and Methods* for more details), using both statistical and cognitive modeling. Further analyses and details can be found in *SI Appendix*. In the first two sections of *Results*, we provide some descriptive characterizations of the dataset. In particular, we show that customers learn from past experience and adapt their exploratory behavior over time. Moreover, exploration is systematically influenced by restaurant features and hence amenable to quantification. We then turn to tests of our model-based hypotheses. We find that customers' exploratory behavior can be clustered meaningfully, exhibits several signatures of intelligent exploration which have previously been studied in the laboratory, and can be captured by a model that generalizes over restaurant features while simultaneously engaging in directed exploration.

**Learning and Exploration over Time.** We first assessed whether customers learned from past experiences, as reflected in their order ratings over time (Fig. 1*A*). The order rating is defined as customers' evaluation on a scale between 1 (poor) and 5 (great). Customers picked restaurants they liked better over time: There was a positive correlation between the number of a customer's past orders and her ratings ($r = 0.073$; 99.9% CI 0.070, 0.076; see *SI Appendix* for further analyses).

Next, we assessed exploratory behavior by creating a variable indicating whether a given order was the first time a customer had ordered from that particular restaurant—i.e., a signature of pure exploration (20). Fig. 1*B* shows the averaged probability of sampling a new restaurant over time (how many orders a customer had placed previously).

Customers sampled fewer new restaurants over time, leading to a negative overall correlation between the number of past orders and the probability of sampling a new restaurant ($r = -0.139$; 99.9% CI –0.142, –0.136). Exploration also comes at a cost (Fig. 1*C*), such that explored restaurants showed a lower average rating (mean rating = 4.257; 99.9% CI 4.250, 4.265) than known restaurants (mean rating = 4.518; 99.9% CI 4.514, 4.522).

Customers learned from the outcomes of past orders. Fig. 1*D* shows their probability of reordering from a restaurant as a
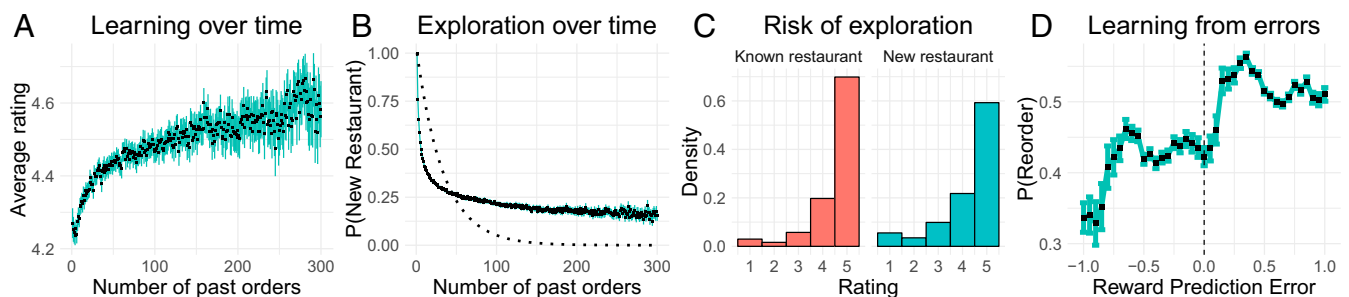


**Fig. 1.** Learning and exploration over time. (*A*) Average order rating by number of past orders. (*B*) Probability of sampling a new restaurant in dependency of the number of past orders. Dashed black line indicates simulated exploratory behavior of agents randomly exploring available restaurants. (*C*) Distribution of order ratings for newly sampled and known restaurants. (*D*) Average probability of reordering from a restaurant as a function of reward prediction error. Means are displayed as black squares and error bars show the 95% confidence interval of the mean.
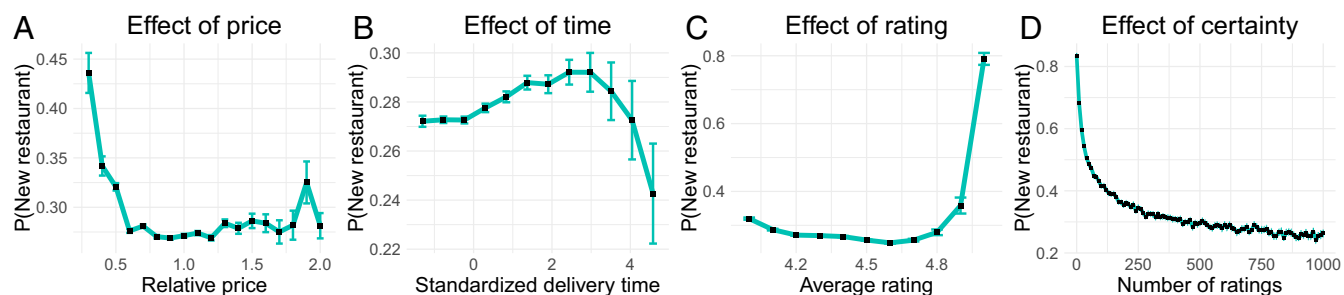
**Fig. 2.** Factors influencing exploration. (*A*) Effect of relative price. The relative price indicates how much cheaper or more expensive a restaurant was compared with an average restaurant in the same city. (*B*) Effect of standardized (z-transformed) estimated delivery time. (*C*) Effect of average rating. (*D*) Effect of a restaurant's number of past ratings (certainty). Means are displayed as black squares and error bars show the 95% confidence interval of the mean.

function of their reward prediction error (RPE, the difference between the expected quality of a restaurant, as measured by the restaurant's average rating at the time of the order, and the actual pleasure customers perceived after they consumed the order, as indicated by their own rating of the order). RPEs are a key component of theories of reinforcement learning (23), and we therefore expected that customers would update their sampling behavior after receiving either a positive or a negative RPE. Confirming this hypothesis, customers were more likely to reorder from a restaurant after an experience that was better than expected (positive RPE: p(reorder) = 0.518, 99.9% CI 0.515, 0.520) than after an experience that was worse than expected (negative RPE: p(reorder) = 0.394, 99.9% CI 0.391, 0.398). The average correlation between RPEs and the probability of reordering was $r = 0.110$ (99.9% CI 0.107, 0.114).

**Determinants of Exploration.** In the next part of our analysis, we focused on what factors were associated with the decision to explore a new restaurant. In particular, we assessed whether exploratory behavior was systematic and therefore looked at the following four restaurant features that were always visible to customers at the time of their order: the relative price (i.e., how much cheaper or more expensive a restaurant is compared with the average within the same country) of a restaurant, its standardized estimated delivery time, the mean rating of a restaurant at the time of the order, and the number of people who had rated the restaurant before.

Customers preferred restaurants that were comparatively cheaper (Fig. 2*A*): The correlation between relative price and the probability of exploration was negative ($r = -0.059$; 99.9% CI $-0.0641$, $-0.0548$). There was a nonlinear relationship between a restaurant's estimated delivery time and its probability of being explored (Fig. 2*B*): Exploration was most likely for standardized delivery times between 1 and 2.5 (0.288; 99.9% CI 0.285, 0.292), and less likely for delivery times below 1 (0.288; 99.9% CI 0.285, 0.292) or above 2.5 (0.252; 99.9% CI 0.229, 0.274). This indicates that customers might have taken into account how long it would take to plausibly prepare and deliver a good meal when deciding which restaurants to explore. The average rating of a restaurant also affected customers' exploratory behavior (Fig. 2*C*): Higher ratings were associated with a higher chance of exploration ($r = 0.038$; 99.9% CI 0.0337, 0.0430). The number of ratings per restaurant also influenced exploration (Fig. 2*D*), with a negative correlation of $r = -0.188$ (99.9% CI $-0.192$, $-0.183$). This may have a mechanical component because restaurants that have been tried more frequently are intrinsically less likely to be explored for the first time. We therefore repeated this analysis for all restaurants that had been rated more than 500 times, yielding a correlation of $r = -0.034$ (99.9% CI $-0.042$, $-0.026$).

We standardized and entered all of the variables into a mixed-effects logistic regression modeling the exploration variable as the dependent variable and adding a random intercept for each customer (see *SI Appendix* for full model comparison). We again found that a smaller number of total ratings ($\beta = -0.475$), a higher average rating ($\beta = 0.086$), and a lower price ($\beta = -0.014$) as well as a quadratic effect of time ($\beta_{Linear} = -0.025$, $\beta_{Quadratic} = 0.015$) were all predictive of customers' exploratory behavior (Table 1). In summary, exploration in the domain of online ordering is systematic, interpretable, and amenable to quantification. We next turned to an examination of our model-based hypotheses concerning directed exploration and generalization.

### Signatures of Uncertainty-Directed Exploration

We probed the data for signatures of uncertainty-directed exploration algorithms that attach an uncertainty bonus to each option. One such signature is that directed and random explorations make diverging predictions about behavioral changes after either a positive or a negative outcome. Whereas random (softmax) exploration predicts no difference between the extent of sampling behavior change following a better-than–expected outcome and that following a worse-than–expected outcome, directed exploration predicts a stronger increase in sampling behavior after a worse-than–expected outcome (*SI Appendix*). This is due to the properties of algorithms that assess an option's utility by a weighted sum of its expected reward and its SD. After a bad experience, the mean and SD both go down, whereas after a good experience the mean goes up but the SD goes down. Thus, there should be greater change in customers' sampling behavior after a bad than after a good outcome.

We verified this prediction by calculating the Shannon entropy of customers' next four purchases after having experienced either a better-than– or a worse-than–expected order. The calculated entropy was higher for negative RPEs (Fig. 3*A*; 1.112; 99.9% CI 1.109, 1.115) than for positive RPEs (1.082; 99.9% CI 1.081, 1.084), in line with theoretical predictions of a directed exploration algorithm.

We calculated each restaurant's relative variance, i.e., how much more variance in its ratings a restaurant possessed

**Table 1. Results of the mixed-effects logistic regression**

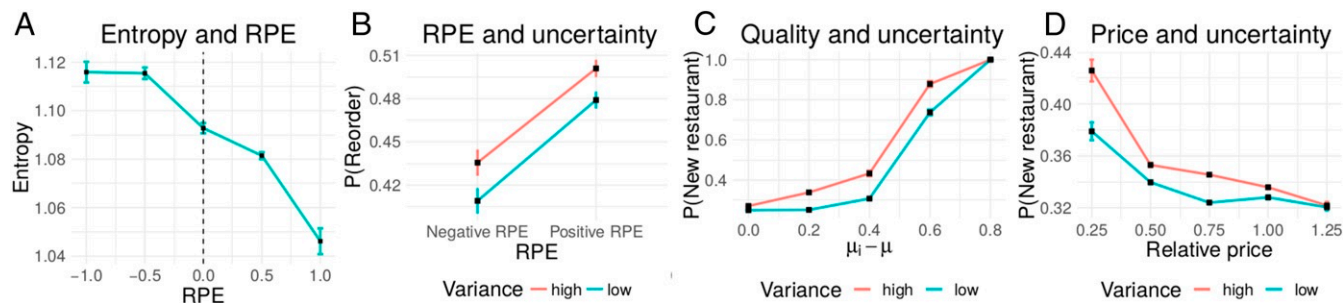|                | Estimate | SE    | z value | Pr($>$ |z|) |
|----------------|----------|-------|---------|-----------|
| Intercept      | $-0.663$ | 0.008 | $-82.01$ | $<0.001$ |
| Relative price | $-0.014$ | 0.006 | $-2.27$  | 0.02     |
| Time linear    | $-0.0246$| 0.008 | $-3.22$  | 0.001    |
| Time quadratic | 0.015    | 0.004 | 3.89     | $<0.001$ |
| Average rating | 0.086    | 0.006 | 13.85    | $<0.001$ |
| No. of ratings | $-0.475$ | 0.007 | $-70.27$ | $<0.001$ |

**Fig. 3.** Signatures of uncertainty-directed exploration. (*A*) Entropy of the next four choices in dependency of RPE. (*B*) Probability of reordering from a restaurant in dependency of RPE, shown for restaurants with high and low relative variance. (*C*) Probability of choosing a novel restaurant in dependency of its difference from an average restaurant within the same cuisine type for restaurants with high and low relative variance. (*D*) Probability of choosing a novel restaurant in dependency of its relative price for restaurants with high and low relative variance.

compared with the average variance per restaurant within the same cuisine type (although customers cannot see the actual estimate of a restaurant's variance in ratings, they can access all past rating as well as a summary that shows the distribution over ratings). We then compared the reorder probability for restaurants with a high vs. low relative rating variance, based on a median split (Fig. 3*B*). This probability was higher for restaurants with high relative variance than for restaurants with low relative variance for both negative and positive RPEs. Thus, customers were more likely to return to restaurants with higher relative uncertainty.

We also assessed customers' exploratory behavior in dependency of the differences in ratings for a given restaurant compared with the average of all restaurants within the same cuisine type (value difference). The probability of exploring a new restaurant increased as a function of the restaurant's value difference (Fig. 3*C*; $r = 0.05$; 99.9% CI 0.045, 0.056). Additionally, a restaurant's relative variance also correlated with its probability of being explored (Fig. 3*C*; $r = 0.05$; 99.9% CI 0.045, 0.056). Comparing restaurants with a high vs. a low relative variance in their ratings revealed a shift of the choice function toward the left. In other words, restaurants with higher relative uncertainty (0.344; 99.9% CI 0.341, 0.349) are preferred to restaurants with lower relative uncertainty (0.319; 99.9% CI 0.317, 0.321), as predicted by uncertainty-directed exploration strategies (2). This difference can also be observed when repeating the same analysis using a restaurant's price (Fig. 3*D*): As restaurants get more expensive, they are less likely to be explored ($r = -0.017$; 99.9% CI $-0.023$, $-0.013$). This function is again shifted for restaurants with higher relative uncertainty: Given a similar price range, relatively more uncertain restaurants are more likely to be explored than less uncertain restaurants.

To further validate these findings, we fitted a mixed-effects logistic regression, using the exploration variable as the dependent variable. For the independent variables, we used the mean difference in ratings between the restaurant and the average restaurant within the same cuisine type, a restaurant's relative price, and its relative uncertainty (Table 2). The average value difference ($\beta = 0.114$), the relative price ($\beta = -0.0876$), and the relative uncertainty ($\beta = 0.084$) all affected a restaurant's probability to be explored. Thus, even when taking into account a restaurant's price and its ratings, customers still preferred more uncertain options. This provides further evidence for a directed exploration strategy.

**Signatures of Generalization.** Having observed how exploratory behavior changes with experience, we investigated how generalization might affect exploration in several ways. First, we looked for evidence of information spillovers by analyzing changes in exploration within cuisine clusters. These seven clusters were

defined in a data-driven manner based on patterns of consecutive explorations, that is, how one exploratory choice predicted the next one (Fig. 4*A* and *Materials and Methods*). This was also related to a subjective understanding of similarity; the frequency of switching between cuisine types was strongly correlated with similarity ratings provided by 200 workers on Amazon Mechanical Turk ($r = 0.78$; Fig. 5*A*). Hinting at strategies of directed exploration as before, we found that bad outcomes had a larger effect than good outcomes compared with a baseline of average switches (Fig. 4*B*)—customers were especially averse to exploring other restaurants in the same cluster after a worse-than–expected outcome ($-5.19\%$), more than they favored such exploration after a better-than–expected outcome ($+2.27\%$). This suggests that uncertainty-modulated exploration takes into account experiences with different restaurants of similar types. Intriguingly, we also observed that customers tended to switch to exploring "unhealthy" cuisines after bad experiences with any other type ($+2.72\%$). This may reflect people balancing differing goals across successive choices (24).

Second, we analyzed how exploration is modulated by the distribution of restaurant quality in a city. Gershman and Niv (20) showed that participants explore novel options more frequently in environments where all options are generally good. We found evidence for this phenomenon in our data (Fig. 5*B*): There was a positive correlation between a city's average restaurant rating and the proportion of exploratory choices in that city ($r = 0.32$; 99.9% CI 0.21, 0.49; see *SI Appendix* for partial correlations). Moreover, there was also a positive correlation between a city's variance of ratings and the proportion of exploratory choices ($r = 0.48$; 99.9% CI 0.37, 0.59), indicating that higher uncertainties in ratings were linked to more exploration.

Third, we examined how the success of exploration depended on the predictability of individual ratings from restaurant features (price, delivery time, mean rating, and number of ratings). Customers gave higher ratings to explored restaurants in cities where ratings were generally more predictable ($r = 0.73$; Fig. 5*C*, 99.9% CI 0.53, 0.84). Thus, exploration seemed to be enhanced by the degree to which features permitted a reduction in uncertainty, similar to findings in contextual bandit tasks (21).

**Table 2. Results of mixed-effects logistic regression**

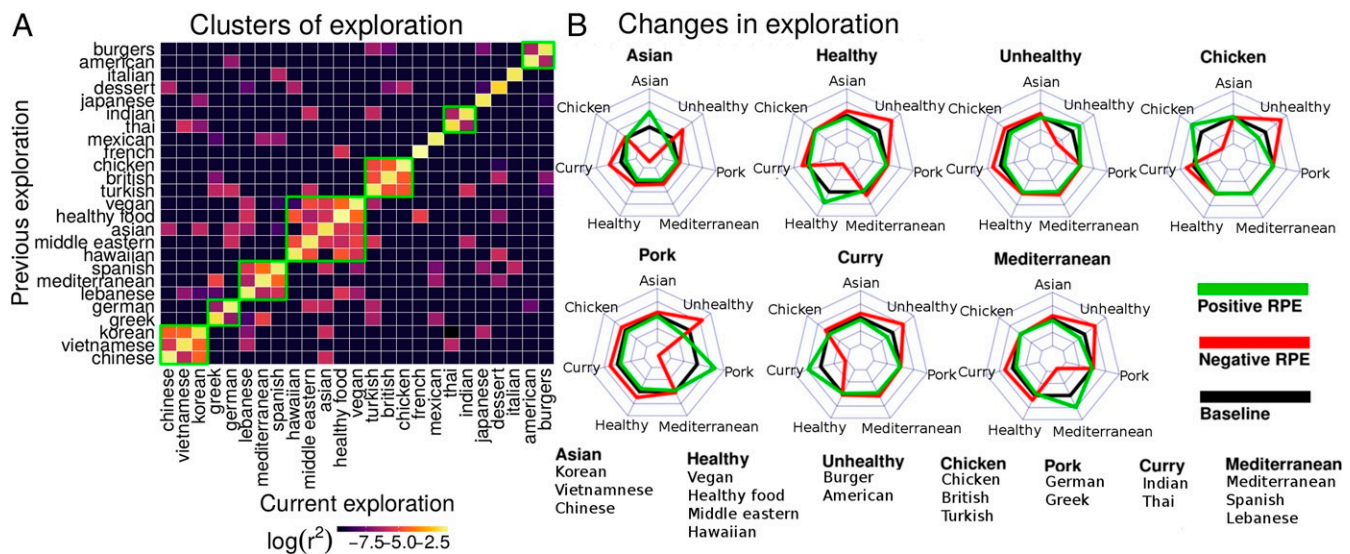|  | Estimate | SE | z value | Pr(> \|z\|) |
|---|---|---|---|---|
| Intercept | −0.342 | 0.007 | 45.81 | <0.001 |
| Value difference | 0.114 | 0.0135 | 8.47 | <0.001 |
| Relative price | −0.087 | 0.007 | −11.67 | <0.001 |
| Variance difference | 0.084 | 0.003 | 24.13 | <0.001 |

**Fig. 4.** Clusters and changes of exploration. (*A*) Clusters of exploration between different cuisine types within customers' consecutive explorations. Green rectangles mark clusters of exploration. (*B*) Moves between clusters after better-than–expected (positive RPE) and worse-than–expected (negative RPE) outcomes compared with a restaurant-specific mean baseline. Centers of radar plots indicate a change of −5%, and outermost lines indicate a change of +5%. A change of 1% roughly translates to 500 orders.

In an attempt to test algorithms of both directed exploration and generalization simultaneously, we compared three models of learning and decision making based on how well they captured the sequential choices of 3,772 new customers who had just started ordering food and who had rated all of their orders. The first model was a Bayesian mean tracker (BMT) that estimates the mean quality for each restaurant independently. The second model was an extension of the BMT model (Gaussian process regression) that estimates mean quality as a function of observable features (price, mean rating, delivery time, and number of past ratings). The shared feature space allows this model to generalize across restaurants. Gaussian process regression is a powerful model of generalization that has been applied to model how participants learn latent functions to guide their exploration (15, 21, 22). It can be seen as a Bayesian variant of similarity-based decision making, akin to economic theories of case-based decision making (25) and psychological formulations of similarity judgments (26). This model was paired with two different policies: stochastic sampling of actions in proportion to their estimated mean quality (Gaussian process with a mean-greedy sampling strategy [GP-M]) or with a directed exploration strategy that sampled based on both the mean and

an uncertainty bonus (formally, an option's upper confidence bound; Gaussian process with an upper confidence bound sampling strategy [GP-UCB]). We treated customers' choices as the arms of a bandit and their order ratings as their utility and then evaluated each model's performance based on its one-step-ahead prediction error, standardizing performance by comparing to a random baseline. Since it was not possible to observe all restaurants a customer might have considered at the time of an order, we compared the different models based on how much higher in utility they predicted a customer's final choice compared with an option with average features out of all of the restaurants available in that customer's city. As Fig. 5*D* shows, the BMT model barely performed above chance ($r^2 = 0.013$; 99.9% CI 0.005, 0.022). Although the GP-M model performed better than the BMT model ($r^2 = 0.231$; 99.9% CI 0.220, 0.241), the GP-UCB model achieved by far the best performance ($r^2 = 0.477$; 99.9% CI 0.465, 0.477). Thus, a sufficiently predictive model of customers' choices required both a mechanism of generalization (learning how features map onto rewards) and a directed exploration strategy (combining a restaurant's mean and uncertainty to estimate its decision value).
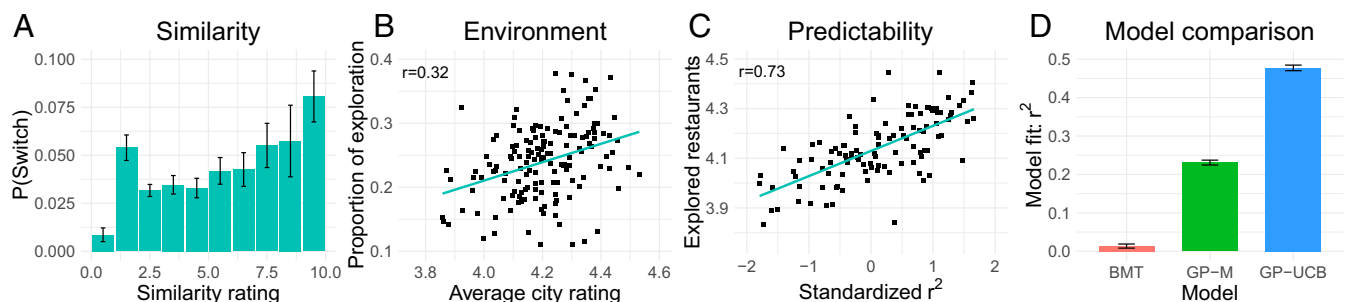


**Fig. 5.** Signatures of generalization. (*A*) Probability of switches between cuisine types and rated similarities between the same types. (*B*) Average rating per city and proportion of exploratory choices. Turquoise line marks least-squares regression line. (*C*) Predictability of a restaurant's quality and average rating of explored restaurants. Turquoise line marks least-squares regression line. (*D*) Results of model comparison for new customers' behavior. Considered models were the Bayesian mean tracker (BMT), a Gaussian process with a mean-greedy sampling strategy (GP-M), and a Gaussian process with an upper confidence bound sampling strategy (GP-UCB).

## Discussion

We investigated customers' exploratory behavior in a large dataset of online food delivery purchases. Customers learned from past experiences, and their exploration was affected by a restaurant's price, average rating, number of ratings, and estimated delivery time. Our results further provide evidence for several theoretical predictions: People engaged in uncertainty-directed exploration, and their exploration was guided by similarity-based generalization. Computational modeling showed that these patterns could be captured quantitatively.

Of course, drawing causal inferences from large datasets is difficult (27, 28). Thus, although we believe that our results provide evidence that people use sophisticated strategies in complex, naturalistic environments, these effects nonetheless deserve further investigation, for example by conducting online experiments.

Furthermore, our model does currently not explain all possible intentions customers might have when ordering food such as maintaining a healthy diet or balancing different goals over successive choices like saving money and trying out expensive food (24). These could hypothetically be incorporated into the kernel function.

Taken together, our results advance our understanding of human choice behavior in complex real-world environments. The results may also have broader implications for understanding consumer behavior. For example, we found that customers frequently switch to unhealthy food options after bad experiences. A potential strategy to increase the exploration of healthy food might thus be to increase healthy restaurants' relative uncertainty by grouping them with other frequently explored options such as Asian restaurants, which showed a comparatively lower relative uncertainty per restaurant.

While we have focused on using cognitive models to predict human choice behavior, the same issues come up for the design of recommendation engines in machine learning. These engines use sophisticated statistical techniques to make predictions about behavior, but do not typically try to pry open the human mind (29). This is a missed opportunity, since one could generate better recommendations of which restaurants to try next, based on a particular customer's estimated values and uncertainties; as models of human and machine learning have become increasingly intertwined, insights from cognitive science may help build more intelligent machines for predicting and aiding consumer choice.

## Materials and Methods

**The Deliveroo Dataset.** The data consisted of a representative random subset of customers ordering food from the online food delivery service "Deliveroo." The dataset contained 195,333 fully anonymized customers. These customers placed 1,613,968 orders over 2 months (February and March 2018) in 197 cities. There were 30,552 restaurants in total, leading to an average of 155 restaurants per city. We arrived at this dataset by filtering out customers with fewer than 5 orders (too few data points to analyze learning) and more than 100 orders (likely multiple people sharing an account). The study was granted ethical approval by the Committee on the Use of Human Subjects at Harvard University. All participants consented to participation through an online consent form at the beginning of the survey.

**Clustering Analysis.** Cuisine tags were manually defined by Deliveroo. We analyzed for each cuisine type how much exploring this type on a time point $t$ was predictive of exploring another cuisine type on a time point $t + 1$, using a linear regression model. Repeating this analysis for every combination of cuisine types led to the graph shown in Fig. 4A. We then analyzed the resulting matrix of $r^2$ values using hierarchical clustering. This clustering excluded the cuisine type "European" as it was found to contain little information about customer choice behavior.

**Similarity Judgments.** To elicit similarity ratings between different cuisine types, we asked 200 participants on Amazon's Mechanical Turk to rate the similarities between two randomly sampled types out of the 20 types used for the clustering analysis reported above. Participants were paid $1 and had to rate 50 pairs of cuisine types on a scale from 0 (not at all similar) to 10 (totally similar).

1. P. Whittle, Multi-armed bandits and the Gittins index. *J. R. Stat. Soc. Ser. B (Methodol.)* **42**, 143–149 (1980).
2. S. J. Gershman, Deconstructing the human algorithms for exploration. *Cognition* **173**, 34–42 (2018).
3. M. Speekenbrink, E. Konstantinidis, Uncertainty and exploration in a restless bandit problem. *Top. Cognit. Sci.* **7**, 351–367 (2015).
4. M. J. Frank, B. B. Doll, J. Oas-Terpstra, F. Moreno, Prefrontal and striatal dopaminergic genes predict individual differences in exploration and exploitation. *Nat. Neurosci.* **12**, 1062–1068 (2009).
5. P. Auer, Using confidence bounds for exploitation-exploration trade-offs. *J. Mach. Learn. Res.* **3**, 397–422 (2002).
6. N. Srinivas, A. Krause, S. M. Kakade, M. W. Seeger, Information-theoretic regret bounds for Gaussian process optimization in the bandit setting. *IEEE Trans. Inf. Theory* **58**, 3250–3265 (2012).
7. P. S. Riefer, R. Prior, N. Blair, G. Pavey, B. C. Love, Coherency-maximizing exploration in the supermarket. *Nat. Hum. Behav.* **1**, 0017 (2017).
8. E. Schulz, S. J. Gershman, The algorithmic architecture of exploration in the human brain. *Curr. Opin. Neurobiol.* **55**, 7–14 (2019).
9. J. D. Cohen, S. M. McClure, J. Y. Angela, Should I stay or should I go? How the human brain manages the trade-off between exploitation and exploration. *Philos. Trans. R. Soc. Lond. B* **362**, 933–942 (2007).
10. K. Mehlhorn et al., Unpacking the exploration–exploitation tradeoff: A synthesis of human and animal literatures. *Decision* **2**, 191–215 (2015).
11. R. S. Sutton et al., *Reinforcement Learning: An Introduction* (MIT Press, 1998).
12. N. D. Daw, J. P. O'Doherty, P. Dayan, B. Seymour, R. J. Dolan, Cortical substrates for exploratory decisions in humans. *Nature* **441**, 876–879 (2006).
13. W. B. Knox, A. R. Otto, P. Stone, B. Love, The nature of belief-directed exploratory choice in human decision-making. *Front. Psychol.* **2**, 398 (2012).
14. R. C. Wilson, A. Geana, J. M. White, E. A. Ludvig, J. D. Cohen, Humans use directed and random exploration to solve the explore–exploit dilemma. *J. Exp. Psychol. Gen.* **143**, 2074–2081 (2014).
15. C. M. Wu, E. Schulz, M. Speekenbrink, J. D. Nelson, B. Meder, Generalization guides human exploration in vast decision spaces. *Nat. Hum. Behav.* **2**, 915–924 (2018).
16. P. M. Todd, Human behaviour: Shoppers like what they know. *Nature* **541**, 294–295 (2017).
17. D. Badre, A. S. Kayser, M. D'Esposito, Frontal cortex and the discovery of abstract action rules. *Neuron* **66**, 315–326 (2010).
18. G. Elliott Wimmer, N. D. Daw, D. Shohamy, Generalization of value in reinforcement learning by humans. *Eur. J. Neurosci.* **35**, 1092–1104 (2012).
19. J. Noble, P. M. Todd, E. Tucif, Explaining social learning of food preferences without aversions: An evolutionary simulation model of Norway rats. *Proc. R. Soc. Lond. Ser. B Biol. Sci.* **268**, 141–149 (2001).
20. S. J. Gershman, Y. Niv, Novelty and inductive generalization in human reinforcement learning. *Top. Cognit. Sci.* **7**, 391–415 (2015).
21. E. Schulz, E. Konstantinidis, M. Speekenbrink, Putting bandits into context: How function learning supports decision making. *J. Exp. Psychol. Learn. Mem. Cognit.* **44**, 927–943 (2018).
22. E. Schulz, N. T. Franklin, S. J. Gershman, Finding structure in multi-armed bandits. bioRxiv:432534 (29 December 2018).
23. W. Schultz, P. Dayan, P. R. Montague, A neural substrate of prediction and reward. *Science* **275**, 1593–1599 (1997).
24. R. Dhar, I. Simonson, Making complementary choices in consumption episodes: Highlighting versus balancing. *J. Market. Res.* **36**, 29–44 (1999).
25. R. Bhui, "Case-based decision neuroscience: Economic judgment by similarity" in *Goal-Directed Decision Making*, Eds. R. W. Morris, A. M. Bornstein, A. Shenhav (Elsevier, 2018), pp. 67–103.
26. R. L. Goldstone, Similarity, interactive activation, and mapping. *J. Exp. Psychol. Learn. Mem. Cognit.* **20**, 3–28 (1994).
27. R. M. Shiffrin, Drawing causal inference from big data. *Proc. Natl. Acad. Sci. U.S.A.* **113**, 7308–7309 (2016).
28. J. Peters, J. M. Mooij, D. Janzing, B. Schölkopf, Causal discovery with continuous additive noise models. *J. Mach. Learn. Res.* **15**, 2009–2053 (2014).
29. T. L. Griffiths, Manifesto for a new (computational) cognitive revolution. *Cognition* **135**, 21–23 (2015).