

# PRIOR SYMMETRY, CATEGORIZATION, AND SIMILARITY-BASED REASONING\*

MARCIN PEŃSKI\*\*

ABSTRACT. This paper presents a rational theory of categorization and similarity-based reasoning. I study a model of sequential learning in which the decision maker infers unknown properties of an object from information about other objects. The decision maker may use the following heuristics: Divide objects into categories with similar properties, then predict that a member of a category has a property if some other member of this category has this property. The environment is *symmetric*: The decision maker has no reason to believe that the objects and properties are a priori different. In symmetric environments, categorization is an optimal solution to an inductive inference problem. Any optimal solution looks *as if* the decision maker categorizes. Various experimental observations about similarity-based reasoning coincide with the optimal behavior in the model.

" ... we conclude [that] because *a* resembles *b* in one or more properties, that it does so in a certain other property." (Mill 1874)

## 1. INTRODUCTION

In a class of discrete prediction problems, the decision maker (henceforth called the DM) is supposed to predict whether an object  $o \in O$  has a property  $p \in P$ . Quite often, the prediction is based on some notion of similarity: the DM predicts that  $o$  has property  $p$  if a similar object  $o'$  has this property. Table 1 presents an example of such reasoning. The DM uses information about chickens and tigers to predict that a falcon, which is more similar to chicken than to a tiger, does not mix oxygen with carbon dioxide while breathing. Arguably, the knowledge that a chicken is similar to a falcon is not inborn to the DM, but arises as a result of learning that both correspond to objects with wings, feathers and beaks and share other properties. From this perspective, similarity-based reasoning is an application of the Principle of Similarity: a presumption that similarity of some properties of two objects

---

This paper was previously presented under the title "Categorization." I am grateful to Jeff Ely, Sham Kakade, Philip Reny, John Conlon, Willie Fuchs, Ali Hortacsu, Motty Perry, Bart Lipman, Eddie Dekel, Alberto Bisin, Jorg Stoye, and Nabil al-Najar for helpful conversations and participants at seminars at the University of Toronto, University of Chicago, MIT/Harvard Theory Seminar, Center of Rationality, Tel Aviv University and New York University for discussion and comments. P. Reny suggested a simplification in the proof of Theorem 1. The author bears full responsibility for the remaining errors.

\*\*The University of Chicago, Department of Economics. E-mail: mpeski@uchicago.edu.

Premise:	Chicken does not mix oxygen with carbon dioxide while breathing Tiger mixes oxygen with carbon dioxide while breathing
Conclusion:	Falcon does not mix oxygen with carbon dioxide while breathing

TABLE 1. Example of similarity-based reasoning.

indicated the similarity of other properties. This principle, in various forms, is recognized as one of the most salient features of human reasoning and the foundation of any inductive argument (for example, (Mill 1874), (Hume 1777)).

Categorization is a model of similarity-based reasoning (see also (Tversky 1977), (Gilboa and Schmeidler 1995) or (Osherson, Wilkie, Smith, Lopez, and Shafir 1990)). It has three defining characteristics. First, the DM divides objects and properties into a finite number of groups called *categories* considering two objects to be similar if they are assigned to the same category. Thus, similarity has a particularly simple binary form. Second, categories are few (relative to the number of objects) and large. Third, prediction is category-based: All members of each category are presumed to share properties, and an observation that a member of a category has a certain property is generalized to the other members. In the above example, the DM's reasoning can be interpreted as predicting a falcon's properties given that it falls into the category "birds."

In the paper, I argue that both similarity-based reasoning and categorization have a rational explanation as optimal behavior in a model of sequential learning. There are infinite sets of objects  $O$  and properties  $P$ . Each period  $t$ , the DM is asked whether object  $o_t \in O$  has property  $p_t \in P$ . After making a prediction, the DM learns the correct answer  $\theta(o_t, p_t) \in \{0, 1\}$ , where function  $\theta : O \times P \rightarrow \{0, 1\}$  is called the *state of the world*. While making a prediction, the DM uses the information that she acquired in the previous periods.

The key assumption is the symmetry of the distribution  $\omega$  from which state of the world  $\theta$  is drawn. The assumption says that the prior is invariant with respect to relabelling of objects and properties. This means that a priori all objects and properties are considered by the DM to be perfectly symmetric and exchangeable. In particular, the assumption eliminates a possibility of categorization before the first period because before the first period the DM cannot distinguish which objects are more similar to one another than to the rest.

The symmetry assumption has a sound interpretation from the subjective point of view, where  $\omega$  is treated as the beliefs of the Bayesian DM. The symmetry is a reflection of her ignorance about objects and properties before the experiment. Given the lack of prior information, she has no reason to treat any two objects differently. This interpretation of symmetry is also known as the Laplacian Principle of Insufficient Reason, or Keynesian Principle of Indifference ((Keynes 1921), (Jaynes 1988), (Savage 1972) chapters 3.7 and 4.5,

(Kreps 1988) chapter 11, (Gilboa, Postlewaite, and Schmeidler 2004)). Alternatively, there are many situations in which the symmetry assumption makes sense from the objective point of view. In particular, suppose that names of objects and properties are assigned to the objects randomly, without any attention to the state of the world. No relabelling should affect the observed frequency of outcomes. From the perspective of the DM, this looks *as if* the distribution of the states of the world were invariant with respect to any relabelling. I present examples of such situations.

Section 5 constructs a learning algorithm that has the features of the psychological model of categorization: The DM assigns objects into a few large categories and makes category-based predictions. The assignments into categories minimizes the inner-category entropy and it maximizes the informational content of the categories. I compare two types of DMs. The Bayesian type knows the distribution  $\omega$  and uses Bayes formula to make her predictions. The non-Bayesian type does not know  $\omega$  and uses the categorization algorithm. In Section 5.3, I show that the non-Bayesian type makes asymptotically the same predictions as the Bayesian one, *no matter what the symmetric distribution  $\omega$* . In other words, the categorization algorithm achieves the quality of Bayesian prediction *uniformly* across all symmetric distributions. This leads to the first argument of this paper:

**Argument 1.** Categorization is an optimal tool of prediction in symmetric environments.

There is an important statistical reason why the size of the categories is large. Having fewer and larger categories helps with the problem of *overfitting*. Recall that the statistical literature warns against using high-dimensional models to "explain" limited observations, the reason being that we risk losing all the predictive power as a price for fitting the past data precisely. The use of fewer and larger categories alleviates this problem: Classification of an object into the correct category is easier, given that there are many objects to compare it with; also, predictions in noisy categories are more precise, if based on a larger number of observations.

The categorization algorithm is not the only optimal solution to the prediction problem. In Section 5.4, I discuss a model in which the DM uses predictions to make decisions and to obtain payoffs. This allows the prediction problem to be cast as an exercise in payoff maximization under uncertainty. I demonstrate that all optimal solutions asymptotically lead to the same behavior. Hence,

**Argument 2.** In the long-run, any optimal behavior in symmetric environments is behaviorally indistinguishable from categorization.

This implies that a behavior of the rational DM should look *as if* she categorized, even if she is really following some other optimal learning rule.

The above two arguments are concerned with the long-run optimality of categorization. In particular, it is optimal in the long-run to apply the Principle of Similarity and to predict that objects with similar properties observed in the past are going to have similar properties in the future. Can one make any statements about the short-run behavior? The question is clearly relevant, as the various properties of similarity-based reasoning listed in the psychological literature ((Rips 1975), (Osherson, Wilkie, Smith, Lopez, and Shafir 1990)) come from experiments in which subjects have access to only very limited data samples. To address this question, I consider a Bayesian DM whose beliefs are symmetric. In Section 6, I describe the qualitative implications of prior symmetry and Bayesian rationality for predictions. The result is that

**Argument 3.** The qualitative properties of Bayesian updating in symmetric environments coincide with experimental observations about similarity-based reasoning.

It needs to be emphasized that the symmetry assumption is motivated by the Principle of Insufficient Reason, not by any similarity considerations. This makes the coincidence between theoretical analysis and empirical observations somewhat unexpected. In my opinion, the coincidence between theory and empirics is indirect evidence that the Principle of Insufficient Reason is strongly embedded in human reasoning.

There are three insights from the results of this paper. First, the fact that humans categorize is itself *not* evidence for bounded rationality. It is often a temptation to denounce heuristics as irrational and attribute using them to a mishandling of available information. On the contrary, categorization arises as an optimal statistical procedure. Second, categorization, if done optimally, does *not* lead to a persistent bias. The fully rational DM dynamically and endogenously adapts her categories to observations. There might be only a temporary bias, which is a consequence of insufficient data. Finally, it seems plausible that, through evolution, Nature equipped us with a tool for making predictions. The evolutionary pressure should have optimized this tool for standard environments. The above results imply that if the standard environments are symmetric, then any such a tool should look like categorization.

## 2. RELATED LITERATURE

The literature on heuristics is divided between two strands. The first strand is concerned with the biases resulting from the use of heuristics. The most influential paper in this literature is (Tversky and Kahneman 1974). (Tversky 1977) presents a model of similarity-based reasoning. More specifically, biases of categorization have been studied in recent papers ((Mullainathan 2002), (Lam 2000), (Jackson and Fryer Jr 2005)). (Jehiel 2005) and

(Jehiel and Samet 2007) analyze the behavior of categorizing agents in games. (Azrieli and Lehrer 2004) develops an axiomatic characterization of categorization. The second strand argues that heuristics are efficient ("fast and frugal") tools for processing information. For an example, see (Gigerenzer and Todd 1999). The current paper belongs to the second strand of the literature.

I. Gilboa and D. Schmeidler's theory of case-based reasoning is an axiomatic approach to questions that are similar to the ones I ask here ((Gilboa and Schmeidler 1995), (Gilboa and Schmeidler 2000), (Gilboa and Schmeidler 2001), (Gilboa and Schmeidler 2002), (Billot, Gilboa, Samet, and Schmeidler 2005)). These papers identify axioms under which the DM's behavior looks as if her prediction were guided by similarity between objects. The most important of these is the combination axiom, which says that if two different databases lead to the same prediction, their union should also lead to the same prediction. The combination axiom is controversial because many prediction rules (including the Bayes formula) do not satisfy it (see (Gilboa and Schmeidler 1995) for a discussion). In particular, neither the categorization algorithm nor the Bayesian prediction of my model satisfy the combination axiom.<sup>1</sup> Thus, the combination axiom is not necessary for similarity-based reasoning.

### 3. MODEL

Let  $X = X^1 \times X^2$  be a set of *instances* (inputs, independent variables, decision problems), where each instance  $(x^1, x^2)$  is a pair of two features  $x^1$  and  $x^2$ . Assume that sets of features  $X_i$  are infinite. Let  $\{0, 1\}$  be a set of *outcomes* (outputs, dependent variables, solutions). A *state of the world* is an assignment of an outcome to every instance,  $\theta : X \rightarrow \{0, 1\}$ . Let  $\Theta = \{0, 1\}^X$  be the space of states of the world;  $\Theta$  is a compact space in product topology, and it is a measurable space with Borel  $\sigma$ -field. A state of the world  $\theta$  is chosen from distribution  $\omega \in \Delta\Theta$ .

Consider the following examples.

**Example 1** (Objects and Properties). *Let  $X^1 = O$  be a space of objects and  $X^2 = P$  be a space of properties. Interpret instance  $(o, p) \in O \times P$  as a query "Does object  $o$  have property  $p$ ?" with an answer  $\theta(o, p) \in \{0, 1\}$ . Say that objects  $o$  and  $o'$  share property  $p$  if  $\theta(o, p) = \theta(o', p)$ .*

**Example 2** (Students and Grades). *An undergraduate advisor helps students to predict grades. Each problem (instance, in my terminology) is described as a pair of two features  $x = (x^1, x^2) = (\text{student}, \text{course})$ . An outcome  $\theta(x^1, x^2)$  is equal to 1 if student  $x^1$  receives a good grade in course  $x^2$  and  $\theta(x^1, x^2) = 0$  if student  $x^1$  receives a bad grade.*

<sup>1</sup>(Gilboa, Lieberman, and Schmeidler 2005) argue that, as reasonable as it seems, the combination axiom should not hold when the DM "uses both inductive and deductive reasoning" at the same time.

**Example 3** (Recommendation algorithms). *Netflix.com is a DVD rental company. Among numerous services, Netflix helps consumers choose movies using a recommendation algorithm: a customer rates movies, and Netflix uses the ratings of that customer and other customers to predict a rating for the movies that the customer has not yet watched. Let  $X^1 = C$  be a space of customers and  $X^2 = M$  be a space of movies. Each instance  $(c, m) \in X$  can be interpreted as a query "Is customer  $c$  interested in movie  $m$ ?"*

A sequence of instances  $\bar{x} = x_1, x_2, \dots \in X^\infty$  is called an *instance process*. To avoid trivial cases, I assume that the DM never observes the same instance twice,  $x_s \neq x_t$  for  $s \neq t$ . Each period, the DM observes an instance  $x_t$ , makes a prediction and subsequently observes outcome  $y_t = \theta(x_t)$ . *Learning rule*  $l : \bigcup_t (X \times Y)^{t-1} \times X \rightarrow \Delta\{0, 1\}$  is a complete description of the predictive behavior of the DM. The predicted probability that the outcome of  $x_t$  is equal to  $y$  is denoted as  $l(\{x_s, y_s\}_{s < t}, x_t)(y)$ . In particular, each distribution  $\omega$  induces a Bayesian learning rule

$$l_\omega((x_s, y_s)_{s < t}, x_t)(y) := \omega(\theta(x) = y | \{x_s, y_s\}_{s < t}).$$

Let  $A$  be a compact and normed space of actions and  $u : A \times \{0, 1\} \rightarrow R$  be a continuous utility function, such that the solution to the optimization problem

$$\max_a (1 - p)u(a, 0) + pu(a, 1)$$

exists, and it is unique and continuous in  $p$ .<sup>2</sup> Denote the solution as  $a_{\max}(p)$ . Let  $a : \bigcup_t (X \times Y)^{t-1} \times X \rightarrow R$  designate a *behavioral rule*. Let

$$u(a((x_s, \theta(x_s))_{s < t}, x_t), \theta(x_t))$$

be a payoff in period  $t$  from behavioral rule  $a$  in the state of the world  $\theta$ . For any distribution  $\omega$  and any instance process  $\bar{x}$ , let

$$U(a; \omega, \bar{x}) := \liminf_{t \rightarrow \infty} E \frac{1}{t} \sum_{s < t} u(a((x_s, \theta(x_s))_{s < t}, x_t), \theta(x_t))$$

denote a long-run expected quality of behavioral rule  $a$ , where the expectation is taken with respect to the distribution over databases of past observations induced by  $\omega$  and instance process  $\bar{x}$ . Let

$$a^l := a_{\max} \circ l$$

denote the behavioral rule induced by learning rule  $l$ .

---

<sup>2</sup>For example, suppose that  $A = [0, 1]$  and  $u(a, y) = ay - \frac{1}{2}a^2$ .

A *database*  $d$  is any finite subset of observations  $d \subseteq X \times Y$ . The size of database  $d$  is denoted with  $|d|$ . For any database  $d$ , let

$$d^i = \{x^i \in X^i : (x^i, x^j, y) \text{ for some}\}$$

denote the number of distinct features  $i$  in database  $d$ . In particular, given an instance process  $\bar{x}$  and state of the world  $\theta$ , the period  $t$  database of past observations is defined as  $d_t = \{(x_s, \theta(x_s))\}_{s < t}$  and  $|d_t| = t$ . If the value of the learning rule does not depend on the order of past observations, write  $l(d, x)$  for database  $d$  and an instance  $x$ . For example, Bayesian learning rule  $l_\omega$  does not depend on the order of observations.

**Definition 1.** *Instance process  $\bar{x}$  satisfies sufficient data condition if*

$$\frac{t}{d_t^1 + d_t^2} \rightarrow \infty.$$

The sufficient data condition implies that the number of observations grows quicker than the number of distinct features in the database of past observations. In the Netflix example (Example 3), this means that the number of observations per customer and per movie increases to infinity.

#### 4. PRIOR SYMMETRY

This section discusses the key assumption of this paper.

**4.1. Symmetric distributions.** A *permutation of instances*  $\pi$  is a bijection of  $X$  onto itself. Denote the set of permutation of instances as  $\Pi$ . A *permutation  $\pi^i$  of dimension  $i$*  is a bijection of  $X^i$  onto itself. Denote the set of all permutations of dimension  $i$  as  $\Pi_i^F$ . A *permutation of features* is a mapping  $\pi = \pi^1 \times \pi^2$ , where  $\pi^i \in \Pi_i^F$  for both  $i$  and

$$(\pi^1 \times \pi^2)(x^1, x^2) = (\pi^1(x^1), \pi^2(x^2)).$$

Thus, the permutation of features is a product of permutations of each dimension separately. Denote the set of all permutations of features with  $\Pi^F$ . Note that  $\Pi^F \subseteq \Pi$ , but  $\Pi^F \neq \Pi$ : Not all permutations of instances are also permutations of features. (For example, suppose that  $x \neq x'$ ; permutation  $\pi_{x,x'} \in \Pi$ , which exchanges instance  $x$  with  $x'$  and keeps all other instances constant, is not a permutation of features,  $\pi_{x,x'} \notin \Pi^F$ .)

For any permutation of instances  $\pi$  and any state of the world  $\theta$ , define  $\pi\theta \in \Theta$  as a state of the world, such that  $(\pi\theta)(x) = \theta(\pi(x))$  for each instance  $x$ .

**Definition 2.** *Distribution  $\omega \in \Delta\Theta$  is symmetric (with respect to renaming features) if, for any permutation of features  $\pi \in \Pi^F$ , for any measurable subset  $E \subseteq \Theta$ ,*

$$\omega(\theta \in E) = \omega(\pi\theta \in E).$$

The symmetry condition generalizes exchangeability of (de Finetti 1964) to two dimensions. It was introduced in (Aldous 1981) and (Hoover 1982) (see also (Kallenberg 2005)). The condition says that, a priori, all features are symmetric. In particular, no Bayesian DM considers any two features as a priori more similar to each other than to any other feature.

**4.2. Interpretation and examples.** From the subjective point of view, symmetry is a restatement of the Laplacian Principle of Insufficient Reason. It should be satisfied by any beliefs of the DM in a hypothetical state of perfect ignorance. In such a state, a Bayesian DM hasn't yet observed any instances, outcomes, or anything that might be correlated with the state of the world. She has no reason to treat any two features differently a priori.

From the objective point of view, symmetry should not be interpreted as an assumption about the distribution from which Nature draws the state of the world (which would be quite restrictive) but about the DM's perception of it (which is not so restrictive). Imagine that (a) there is an objective state of the world, and (b) Nature randomly and uniformly mixes features before letting the DM observe outcomes. This leads to two labels of features: "original" and "perceived." If the mixing is truly uniform, then, from the point of view of the DM, the "perceived" feature  $x^i$  looks like the "original"  $x^i$  with the same probability as it looks like the "original"  $x^{i'}$ .

To see this argument more clearly, consider Example 2: Suppose that the registrar office randomly assigns ID numbers to students and courses. The DM knows the IDs but not the individual names of students or courses. If the assignment is completely random, then, from the point of view of the DM, it looks *as if* the state of the world is drawn from a symmetric distribution.

As another example, consider the Netflix problem from Example 3. From the point of view of Netflix, no renaming of its 5 million customers should change the correlations among customers, movies and their preferences. Analogously, no renaming of their 65 000 movie titles should affect the distribution of the states of the world.<sup>3</sup>

Since a priori all features are exchangeable, the assumption precludes any possibility of non-empirical categorization. However, the assumption allows for a wide range of theories

---

<sup>3</sup>Recently, Netflix announced a public competition for a recommendation algorithm that improves on its own (see [www.netflixprize.com](http://www.netflixprize.com)). A database of 100 million customer-movie rankings is available for any contestant. In order to protect the confidentiality of ratings, Netflix replaced the customers' and movies' names by randomly drawn IDs. In other words, from the perspective of the contestant, the Netflix database looks as if it were drawn from an invariant distribution.

about correlations between outcomes of instances. These correlations may open a possibility of ex post categorization. Consider the following examples. In the first example, all correlations are eliminated. Two subsequent examples are more sophisticated.

**Example 4** (Idiosyncratic preferences). *In the Netflix example, suppose that each outcome  $\theta(x)$  is chosen i.i.d. from uniform distribution on  $\{0, 1\}$ , independently across instances.*

**Example 5** (Bad and good movies). *Suppose that each movie  $m \in X^2$  is independently chosen to be good with probability  $p \in [0, 1]$  or bad with probability  $1 - p$ . State of the world  $\theta$  depends deterministically on the quality of movies: for any instance  $(c, m) \in X$ , let  $\theta(c, m) = 1$  if movie  $m$  is good; otherwise, let  $\theta(c, m) = 0$ .*

**Example 6** (Two types of movies and customers). *There are two types of customers, Men and Women and two types of movies, Action and Romance. Each customer is chosen independently to be Man or Woman with probability  $\frac{1}{2}$ ; similarly, each movie is chosen to be Action or Romance with equal probability. State of the world  $\theta$  depends on the types of customers and movies:*

$$\theta(c, m) = \begin{cases} 1, & \text{if } c \text{ is Man and } m \text{ is Action or } c \text{ is Woman and } m \text{ is Romance,} \\ 0, & \text{otherwise.} \end{cases}$$

In the examples, customers and movies are divided into types (categories). In the first example, there is only one category for customers and movies; in the second, there is one category for customers and two categories for movies; in the last, there are two categories for customers and movies. The outcomes depend on the category assignment either probabilistically (as in the first example) or deterministically (as in the two subsequent examples.) Appendix B.1 contains the Representation Theorem for symmetric distribution. The Theorem shows that any symmetric distribution is a mixture of distributions generated as in the examples, but with, possibly, infinitely many categories.

It is also good to give an example of distribution that does not satisfy symmetry. In general, any logical relationship eliminates symmetry.

**Example 7.** *Consider Example 1 and suppose that property  $p$  is known to be a logical negation of property  $p'$ , i.e., if  $\theta(o, p) = 1$  then  $\theta(o, p') = 0$  for any object  $o$ . Then, no distribution that respects this relation is symmetric.*

The results of the paper remain true under assumptions that are somewhat weaker than symmetry (and which apply to Example 7). This is discussed in the last section.

**4.3. Principle of Similarity.** Recall that the Principle of Similarity says that if two objects were observed to have similar properties, their unobserved properties should be expected to

be similar. Without any further analysis, it is unclear what the Principle has to do with symmetric distribution. Nevertheless, the connection can be illustrated with a simple result. Consider Example 1. Suppose that the Bayesian DM with symmetric beliefs  $\omega$  observes properties  $p \in P'$  of two objects  $o_1, o_2 \in O$ , where  $P'$  is finite set. The DM wonders whether these objects share an unobserved property  $p^* \notin P'$ . The next Proposition says that the probability of such an event increases in the number of shared properties  $p \in P'$ .

**Proposition 1.** *For any symmetric  $\omega$ , any two sets  $P_1, P_2 \subseteq P'$ , if  $|P_1| \leq |P_2|$ , then*

$$\begin{aligned} \omega(\theta(o_1, p^*) = \theta(o_2, p^*) | \theta(o_1, p) = \theta(o_2, p) \text{ if and only if } p \in P_1) \\ \leq \omega(\theta(o_1, p^*) = \theta(o_2, p^*) | \theta(o_1, p) = \theta(o_2, p) \text{ if and only if } p \in P_2). \end{aligned}$$

*Proof.* Consider a random variable  $s : P \rightarrow \{0, 1\}$  defined as a function of the state of the world: for any property  $p$ ,  $s(p) = 1$  if and only if objects  $o_1$  and  $o_2$  share property  $p$ , i.e.  $\theta(o_1, p) = \theta(o_2, p)$ . Let  $\varpi \in \Delta \{0, 1\}^P$  be the distribution of variable  $s$  induced by symmetric distribution  $\omega$ . The symmetry of  $\omega$  implies that distribution  $\varpi$  is invariant with respect to permutation of properties:  $\varpi(s \in E) = \varpi(\pi^P s \in E)$  for any bijection  $\pi^P : P \rightarrow P$  and any measurable  $E \subseteq \{0, 1\}^P$ . By de Finetti's Theorem,  $\varpi$  can be interpreted as a distribution of infinitely many Bernoulli draws indexed with  $p \in P$ , with a parameter that is stochastically drawn once from some  $\mu \in \Delta [0, 1]$ . It is a simple consequence of the representation that the conditional probability of  $s(p^*) = 1$  increases with the number of 1s observed so far,

$$\begin{aligned} \omega(s(p^*) = s(p^*) | s(p) = s(p) \text{ iff } p \in P_1) \\ \leq \omega(s(p^*) = s(p^*) | s(p) = s(p) \text{ iff } p \in P_2). \end{aligned}$$

This yields the Proposition. □

The Proposition provides the main intuition for the connection between symmetry and similarity-based reasoning. All the subsequent results should be seen as generalizations of this intuition.

## 5. CATEGORIZATION

**5.1. Categorization algorithm.** Next, I construct two learning rules. In both rules, the DM divides instances into finitely many categories. The number of categories is fixed in the first learning rule, and it increases with the number of observations in the second one. The rules share stylized characteristics with the categorizing behavior discussed in the psychology literature: (a) categorization is endogenous and dynamic, (b) the number of categories is small compared to the number of objects, and (c) the prediction is category-based: (i.e., all objects in the same category are predicted to share similar properties).

Any categorization process must solve two difficulties: how to allocate instances into categories and how to find a prediction of an outcome conditional on the category. Both difficulties are addressed simultaneously. There are  $k$  possible "bins" for features  $i$ , and each feature is assigned into only one bin. If features  $x^1$  and  $x^2$  are assigned to categories  $k^1$  and  $k^2$ , respectively, then the outcome  $\theta(x^1, x^2)$  is predicted to be 1 with probability  $\rho(k^1, k^2) \in [0, 1]$ . Initially, the DM is uncertain which assignment into bins and which prediction function  $\rho$  are the best, i.e., the most helpful in facilitating predictions. She acts as a Bayesian: She starts with a uniform prior over all possible assignments and functions  $\rho$ . When new information comes, she updates her prior through Bayes' formula.

Formally, let  $k$ -(category) assignment of feature  $i$  be a map  $c^i : X^i \rightarrow \{1, \dots, k\}$ . I refer to  $c^i(x^i)$  as a category of feature  $x^i$ . Let  $\mathcal{C}_i^k = \{1, \dots, k\}^{X^i}$  be a set of  $k$ -assignments of feature  $i$  and let  $\mathcal{C}^k = \mathcal{C}_1^k \times \mathcal{C}_2^k$  be a set of  $k$ -assignments. For any  $k$ -assignment  $c \in \mathcal{C}^k$ , any instance  $x = (x^1, x^2) \in X$ , write

$$c(x) = (c^1(x^1), c^2(x^2)) \in \{1, \dots, k\}^2$$

and call  $c(x)$  a category of instance  $x$  with respect to assignment  $c$ . Hence,  $k$ -assignment divides instances into  $k^2$  categories.

A (category-based) prediction is a function  $p : \{1, \dots, k\}^2 \rightarrow \Delta\{0, 1\}$  with the following interpretation: If the DM decides to assign instance  $x$  to category  $\mathbf{k} \in \{1, \dots, k\}^2$ , then she predicts that the outcome of  $x$  is equal to  $y$  with probability  $p(\mathbf{k})(y)$ . Define the space of prediction functions as

$$\mathcal{R}^k := \{\rho : \{1, \dots, k\}^2 \rightarrow \Delta\{0, 1\}\} = [0, 1]^{k^2}.$$

A couple of an assignment and a prediction function is called a *theory*. Let

$$\mathcal{T}^k = \mathcal{C}^k \times \mathcal{R}^k$$

be a space of theories.

The DM starts with prior beliefs over theories. Let  $\Psi_i^k \in \Delta\mathcal{C}_i^{k_i}$  be the "uniform" measure over assignments of  $i$ . It is formally defined as a measure such that for any feature  $x^i \in X^i$ , category  $c^i(x^i)$  is drawn independently and uniformly from set  $\{1, \dots, k\}$ . Let

$$\Psi_C^k = \Psi_1^k \otimes \Psi_2^k \in \Delta\mathcal{C}^k$$

be the independent product of two measures. Distribution  $\Psi_C^k$  is the uniform measure over a space of assignments  $\mathcal{C}^k$ . Let  $\Psi_R$  be the Lebesgue measure on the space of prediction functions  $\mathcal{R}^k$ . Let

$$\Psi^k = \Psi_C^k \otimes \Psi_R^k \in \Delta\mathcal{T}^k$$

be the independent product of measures  $\Psi_C^k$  and  $\Psi_R^k$ . Distribution  $\Psi^k$  is the uniform measure over an infinitely dimensional space of theories, and it is treated as the prior beliefs.

The posterior density of theory  $(c, \rho)$  is defined through Bayes' formula

$$\psi^k(c, \rho|d) = \frac{\prod_{(x,y) \in d} \rho(c(x))(y)}{\int_{\mathcal{T}^k} \prod_{(x,y) \in d} \rho'(c'(x))(y) d\Psi^k(c', \rho')},$$

where

$$\prod_{(x,y) \in d} \rho(c(x))(y)$$

is equal to the probability of realization of database  $d$  given theory  $(c, \rho)$ .

Finally, define learning rule  $l^k$  as the prediction that would be made by a DM with "beliefs"  $\psi(\cdot|d)$ :

$$l^k(d, x)(y) = \int_{\mathcal{T}^k} [\rho(c(x))(y)] \psi^k(c, \rho|d) d\Psi^k(c, \rho).$$

I refer to learning rule  $l^k$  as a *k-categorization algorithm*.

So far, I have kept the number of categories constant. The second learning rule increases the number of categories together with the size of the data sample. Let  $K_t$  be defined as

$$K_t := \left\lfloor \sqrt[3]{\frac{t}{d_t^1 + d_t^2}} \right\rfloor, \quad (5.1)$$

where  $\lfloor x \rfloor$  is the largest natural number not larger than  $x$ .  $K_t$  will determine the number of categories used in period  $t$ . I refer to  $K_t$  as a *switching rule*. If the sufficient data condition (Definition 1) is satisfied, then  $\lim_{t \rightarrow \infty} K_t = \infty$ .

Let  $d_t(k)$  be defined as follows:

$$d_t(k) = \{(x_s, y_s) : s < t, K_s = k\}.$$

This is a database of observations before period  $t$  in which the value of the switching rule is equal to  $k$ . Define an (*adaptive*) *categorization algorithm*  $l^C$  as

$$l^C(d_t, x) := l^{K_t}(d_t(K_t), x).$$

In period  $t$ , the DM applies the  $K_t$ -categorization algorithm to database  $d_t(K_t)$ .

Notice that both learning rules satisfy two defining properties of categorization that have been described in the beginning of this Section. The assignment into categories depends on past data. The size of categories grows to infinity as the number of observations increases. (This is because

$$\lim_{t \rightarrow \infty} \frac{t}{k} = \infty \text{ and } \lim_{t \rightarrow \infty} \frac{t}{K_t} \geq \lim_{t \rightarrow \infty} t^{2/3} = \infty.$$

Thus, the categorizing behavior does not disappear in the limit.) Finally, the predictions of outcomes are based on category assignments.

**5.2. Categorization as entropy minimization.** It is instructive to reinterpret the categorization algorithm  $l^k$  as entropy minimization. For each database  $d$  and category assignment  $c$ , define the number of instances assigned to category  $(k^1, k^2)$  and the frequency of outcome 1 among these instances as

$$n(k^1, k^2|c, d) = \#\{(x^1, x^2, y) \in d : c(x^1, x^2) = (k^1, k^2), y \in \{0, 1\}\},$$

$$\phi(k^1, k^2|c, d) := \frac{\#\{(x^1, x^2, 1) \in d : c(x^1, x^2) = (k^1, k^2)\}}{n(k^1, k^2|c, d)},$$

if  $n(k^1, k^2|c, d) \neq 0$  and  $\frac{1}{2}$  otherwise. Define the *entropy* of assignment  $c$  as

$$E(c|d) = -\frac{1}{|d|} \sum_{(k^1, k^2) \in \{1, \dots, k\}^2} n(k^1, k^2|c, d) h(\phi(k^1, k^2|c, d)),$$

where  $h(\cdot)$  is the entropy function

$$h(\phi) = \phi \log \phi + (1 - \phi) \log (1 - \phi).$$

Entropy  $E(c|d)$  measures the informational content of assignment  $c$  in database  $d$ . In particular, if the entropy is close to 0, then, for most categories,  $\phi(k^1, k^2)$  is very close to 0 or very close to 1, i.e., predictions inside categories are almost deterministic. On the other hand, if the entropy is close to  $\log 2$  (which is the maximal possible value), then  $\phi(k^1, k^2)$  is close to  $\frac{1}{2}$ , and categories are quite useless in facilitating prediction. Define the minimal value of entropy in database  $d$  as

$$E_{\min}(d) = \min_{c \in \mathcal{C}^k} E(c|d).$$

The next Proposition says that, if the sufficient data condition is satisfied, then, asymptotically, the categorization algorithm puts probability close to 1 on the set of assignments with entropy close to the minimal entropy.

**Proposition 2.** *Suppose that instance process  $\bar{x}$  satisfies the sufficient data condition. Then, for any  $\varepsilon > 0$ ,*

$$\lim_{t \rightarrow \infty} \int_{\{c: E(c|d_t) \leq E_{\min}(d_t) + \varepsilon\} \times \mathcal{R}^k} \psi^k(c, \rho|d_t) d\Psi^k(c, \rho) = 1.$$

The (standard) proof is contained in Appendix A. I will use the Proposition to compare the categorization algorithm from this paper with a heuristics proposed in (Jackson and Fryer Jr 2005). There, the decision maker looks for a categorization assignment that minimizes variance inside categories. Since entropy minimization is not the same as minimization of inner-category variance, (Jackson and Fryer Jr 2005)'s algorithm is not going to satisfy the optimality results that are described next.

**5.3. Optimality of categorization.** This section shows that categorization is an optimal solution to the prediction problem. Consider two types of DMs. A Bayesian DM knows which symmetric distribution  $\omega$  generates the state of the world. The best prediction she can make is to use Bayes' learning rule  $l_\omega$ . A non-Bayesian DM does not know the true distribution  $\omega$ , but nonetheless believes that  $\omega$  is symmetric.

The next Proposition says that if  $k$  is sufficiently high, then the non-Bayesian DM who uses  $k$ -categorization algorithm  $l^k$  makes asymptotically similar predictions to those made by the Bayesian DM. For any  $p, q \in \Delta\{0, 1\}$ , let  $\|p - q\| = \sum_{y \in \{0, 1\}} |p(y) - q(y)|$  be the  $L^1$ -distance between  $p$  and  $q$ . Let

$$E \left\| l^k(d_t, x_t) - l_\omega(d_t, x_t) \right\|$$

be the expected difference between predictions made in period  $t$  by the Bayesian DM and by the  $k$ -categorization algorithm, where the expectation is taken with respect to the distribution over past databases induced by  $\omega$  and instance process  $\bar{x}$ .

**Proposition 3.** *Suppose that process  $\bar{x}$  satisfies the sufficient data condition. For any symmetric  $\omega$ ,*

$$\lim_{k \rightarrow \infty} \limsup_{t \rightarrow \infty} \frac{1}{t} E \sum_{s < t} \left\| l^k(d_s, x_s) - l_\omega(d_s, x_s) \right\| = 0.$$

It is a consequence of the Proposition that any symmetric  $\omega$  can be approximated by a model with a finite number of categories. For a  $k$  high enough, additional categories do not substantially increase the predictive power of the categorization algorithm. The key step in the proof uses the Representation Theorem from Appendix B.1 to show that any symmetric distribution can be approximated as a mixture of distributions generated as in Examples 4, 5 and 6 but with  $k$  categories for each feature. Given sufficient data, the categorization algorithm  $l^k$  behaves as the Bayesian updating on the approximating distribution. Details can be found in Appendix C.

The Proposition is not satisfactory because the number of categories  $k$  needed depends on distribution  $\omega$ . That, in particular, is not known by the non-Bayesian DM. This is a problem because the DM should choose a learning rule before making any observations and before learning anything about distribution  $\omega$ .

The DM faces a fundamental trade-off. On one hand, a higher  $k$  allows for a better limit approximation, as described in Proposition 3. On the other hand, too high a  $k$  creates an *overfitting* problem: The more categories she has, the more difficult it is to use finite data to choose the right category assignment and the right prediction function. The solution to this problem must balance the growth of  $k$  with overfitting risks. I use the adaptive categorization

algorithm. In this algorithm, the DM changes the number of categories as time goes by, and the database of past observations increases.<sup>4</sup>

**Theorem 1.** *Suppose that process  $\bar{x}$  satisfies the sufficient data condition. For any symmetric  $\omega$ ,*

$$\limsup_{t \rightarrow \infty} \frac{1}{t} E \sum_{s < t} \|l^C(d_s, x_s) - l_\omega(d_s, x_s)\| = 0. \quad (5.2)$$

The Theorem says that the adaptive categorization  $l^C$  makes, on average, the same predictions as the Bayesian DM, *uniformly* across all symmetric beliefs  $\omega$ .<sup>5</sup> The Theorem fixes the flaw of Proposition 3: Adaptive categorization does not depend on the initial choice of the number of categories and increases the number of categories together with the data sample. The proof can be found in Appendix C.

To interpret the Theorem, consider first the subjective point of view on  $\omega$ . Any Bayesian DM with symmetric beliefs expects to predict *as if*, approximately and asymptotically, she were using the categorization algorithm. Therefore, any Bayesian DM is indifferent between Bayesian updating and categorization.

From the objective viewpoint, Nature draws the state of the world from symmetric distribution  $\omega$ . The DM may understand that the distribution is symmetric, even if she does not know  $\omega$ . The categorization algorithm guarantees payoffs as high as if the DM knew the true  $\omega$ . This is very good news for any ambiguity-averse decision maker.

Note that the Theorem does not guarantee that the DM will predict all outcomes correctly, only that the DN will predict as well as the Bayesian. In Example 4, all outcomes are i.i.d. equal to 1 with probability  $\frac{1}{2}$ . This is the prediction of the Bayesian DM no matter how much past data she observes. By the Theorem, this is also the asymptotic prediction of the categorization algorithm  $l^C$ . On the other hand, one can show that if the sufficient data condition is satisfied, then, asymptotically, the categorization algorithm predicts correctly almost all outcomes in Examples 5 and 6.

**5.4. Uniqueness of optimal solution.** Theorem 1 shows that there is a learning rule that guarantees uniformly good predictions. Moreover, this rule, by construction, can be interpreted as a categorization. The Theorem does not guarantee that an optimal rule is unique. In fact, there are infinitely many optimal learning rules. This is a consequence of my

---

<sup>4</sup>In general, this method of solving the overfitting problem is known in the statistical decision theory as the *capacity control method* (see, among others, (Vapnik 1998) and (Bousquet, Boucheron, and Lugosi 2004)).

<sup>5</sup>The statement of the Theorem is related to the literature on the Bayesian merging of opinions ((Blackwell and Dubins 1962), (Lehrer and Smorodinsky 2000); also (Jackson, Kalai, and Smorodinsky 1999)). Here and there, the predictions of two learning rules converge only on average. Note, however, that the categorization algorithm  $l^C$  is not, strictly speaking, Bayesian, as it is not based on updating one "uniform prior" over all symmetric distributions. (In fact, such a prior does not exist.)

optimality criterion: Any learning rule that behaves differently from  $l^C$  in the first  $t$  periods and then follows the same predictions as  $l^C$  satisfies the formula (5.2).

**Definition 3.** *Behavior  $a$  is uniformly optimal if  $U(a; \omega, \bar{x}) \geq U(a'; \omega, \bar{x})$  for any symmetric distribution  $\omega$  and for any other behavior  $a'$ .*

The definition of uniformly optimal behavior is very strong. It requires the behavior to be (weakly) better than any other behavior for any other distribution over states of the world. Any uniformly optimal behavior is robust to misspecification of prior beliefs. The notion of robustness is stronger if the behavior was simply optimal with respect to the minimax preferences of (Gilboa and Schmeidler 1989). In the latter, the DM cares only about the worst-case payoff. Here, the DM achieves the optimal payoff given any distribution  $\omega$ .

**Corollary 1.** *Categorization behavior  $a^{lc}$  is uniformly optimal. For any process  $\bar{x}$ , for any uniformly optimal behavior  $a$ ,*

$$\limsup_{t \rightarrow \infty} \frac{1}{t} E \sum_{s < t} \|a^{lc}(d_s, x_s) - a(d_s, x_s)\| = 0.$$

The Corollary says that all uniformly optimal behavioral rules are asymptotically equal, and, in particular, all of them are equal to the categorization algorithm. The idea is very simple. The best prediction possible, given any symmetric  $\omega$ , is the Bayesian prediction. Because categorization makes the Bayesian prediction asymptotically, any uniformly optimal behavior must do the same. Therefore, any two uniformly optimal behaviors are asymptotically equal. In particular, any uniformly optimal behavior is asymptotically indistinguishable from categorization.

*Proof.* Let  $a_\omega := a_{\max} \circ l_\omega$  denote Bayesian behavioral rule. By standard arguments, for any instance process  $\bar{x}$ , any symmetric  $\omega$  and any behavioral rule  $a$ ,

$$U(a; \omega, \bar{x}) \leq U(a_\omega; \omega, \bar{x})$$

with strict inequality if

$$\limsup_{t \rightarrow \infty} \frac{1}{t} E \sum_{s < t} \|a_\omega(d_s, x_s) - a(d_s, x_s)\| > 0.$$

Together with Theorem 1, this implies the thesis of the Corollary.  $\square$

## 6. SIMILARITY-BASED REASONING

In this section, I discuss a series of laboratory observations and stylized facts about the similarity-based reasoning. In each case, the empirical observation is compared with a Bayesian prediction of the DM with symmetric beliefs. In all cases but the last one, the

Premise	Robins use serotonin as a neurotransmitter. Bluejays use serotonin as a neurotransmitter.
Conclusion	Sparrows use serotonin as a neurotransmitter.
Premise	Robins use serotonin as a neurotransmitter. Bluejays use serotonin as a neurotransmitter.
Conclusion	Geese use serotonin as a neurotransmitter.

TABLE 2. Premise-conclusion similarity

theoretical predictions coincide with the laboratory observations (the empirical evidence in the last case seems to be quite weak, though).

**6.1. Premise-conclusion similarity.** (Osherson, Wilkie, Smith, Lopez, and Shafir 1990) present subjects with pairs of inductive arguments. Each argument consists of a premise, which is followed by a conclusion. Subjects are asked to choose the more credible argument in a pair.

An example of such a comparison is presented in Table 2. The first argument was chosen as more credible by 59 out of 80 subjects. The interpretation is that a category of sparrows is more similar to categories of robins and bluejays and the inductive argument seems more appropriate.

Note that the careful choice of objects and properties for the experiments ensures that the subjects have "no reason to believe" that the objects and properties are a priori different. (It is rare for non-professionals to be able to differentiate prior probabilities of "using" and "not using serotonin as a neurotransmitter.") Hence, if the Principle of Insufficient Reason holds, then the subjects should have symmetric beliefs. This allows me to treat the coincidence between empirics and theory as an indirect positive verification of the model, or precisely, the fact that the DM is rational and that she has symmetric beliefs. Under these two assumptions, the Principle of Similarity is an implication of Proposition 1.

Next, (Osherson, Wilkie, Smith, Lopez, and Shafir 1990) argue that an inductive argument is more credible if the conclusion is more specific (and more similar to the premise.). Consider an example in Table 3. Here, 75 out of 80 subjects point to the second argument as more credible.

This experiment differs from the previous one because the conclusion concerns a class rather than a specific object. Nevertheless, conclusion specificity is an application of the Principle of Similarity (formalized as, possibly a variation of, Proposition 1).

Premise	Bluejays require Vitamin K for their liver to function. Falcons require Vitamin K for their liver to function.
Conclusion	All animals require Vitamin K for their liver to function.

Premise	Bluejays require Vitamin K for their liver to function. Falcons require Vitamin K for their liver to function.
Conclusion	All birds require Vitamin K for their liver to function.

TABLE 3. Premise diversity

Premise	Hippopotamuses have a higher sodium concentration in their blood than humans. Hamsters have a higher sodium concentration in their blood than humans.
Conclusion	All mammals have a higher sodium concentration in their blood than humans.

Premise	Hippopotamuses have a higher sodium concentration in their blood than humans. Rhinoceroses have a higher sodium concentration in their blood than humans.
Conclusion	All mammals have a higher sodium concentration in their blood than humans.

TABLE 4. Premise diversity

**6.2. Premise diversity.** The next experiment indicates that an inductive argument is more credible if the premise is more diverse. Consider an example in Table 4. Here, 76 out of 80 subjects point to the first argument as more credible.<sup>6</sup>

The premise diversity is somehow surprising, given what has been said so far about similarity-based reasoning. Nevertheless, it has a consistent explanation. Given that rhinoceroses are known a priori to be similar to hippopotamuses, *it is not unexpected* that rhinoceroses and hippopotamuses have similar amounts of sodium. In particular, the fact about rhinoceroses does not add to what is already known from the analogous statement about hippopotamuses. On the other hand, the same statement about hamsters is more informative: it signals that "high sodium concentration" is shared by other mammals.

I formalize this argument. Suppose that the set of objects is divided into infinitely many categories,  $O = C \times I$ , where  $C$  is an infinite set of categories and  $I$  is an infinite set of identifiers of individual objects. The division into categories can be endogenous, as a result of the past learning, or it might come from some other source of information (such

<sup>6</sup>It is instructive to compare premise diversity with an observation reported in (Glazer and Rubinstein 2001) that subjects believe that the counterargument to some thesis is more credible if the object is more similar to the object in the original argument. It is difficult to explain Glazer and Rubinstein's phenomenon in a Bayesian setting and (Glazer and Rubinstein 2001) propose a game-theoretic explanation.

as instructions given in the experiment). Let  $\omega \in \{0, 1\}^O$  be the beliefs of the DM about whether objects  $o \in O$  have certain fixed property.

I will make two symmetry assumptions about  $\omega$ : that objects inside each category are exchangeable and that categories are exchangeable. Formally, for each  $c$ , for all permutations  $\pi : I \rightarrow I$ ,  $\pi_C : C \rightarrow C$ , define permutations  $\pi_c^*, \pi_C^* : O \rightarrow O$ , as

$$\pi_c^*(c', i) = \begin{cases} (c', i), & \text{if } c' \neq c, \\ (c', \pi(i)) & \text{if } c' = c, \end{cases},$$

$$\pi_C^*(c, i) = (\pi_C(c), i).$$

Thus, permutation  $\pi_C^*$  exchanges categories and permutation  $\pi_c^*$  exchanges objects inside a particular category  $c$ . I assume that for all permutations  $\pi : I \rightarrow I$ ,  $\pi_C : C \rightarrow C$ , any category  $c \in C$ , any measurable set  $E \subseteq \{0, 1\}^O$ ,

$$\omega(\theta \in E) = \omega(\pi_c^*\theta \in E) \text{ and } \omega(\theta \in E) = \omega(\pi_C^*\theta \in E).$$

These two assumptions seems appropriate for the experiment in Table 4. Suppose that  $c_0$  is a category of large African mammals, that includes rhinoceroses, hippopotamuses, and, possibly elephants. The first assumption implies that the probability that elephants have the property given that rhinoceroses have it is the same as given hippopotamuses have the property. Unless a DM has biology training, there is no reason why the former should be more informative about the functioning of elephant liver than the latter, and vice versa. Additionally, suppose that  $c_1$  and  $c_2$  are categories of rodents and dogs, respectively. The second assumption says that the probability that dogs have the property is the same, conditionally on either the fact about rodents or large African mammals. Again, there is no reason to expect rodents bring more information about dogs than hippopotamuses.

**Proposition 4.** *Suppose that  $\omega$  satisfies the above assumptions. Then, for any objects  $(c_0, i_0), (c_1, i_1), (c_2, i_2)$ , and  $(c_0, i'_0)$  such that  $c_k \neq c_l$  for  $k \neq l$  and  $i_0 \neq i'_0$ ,*

$$\omega(\theta(c_2, i_2) = \theta(c_0, i_0) \mid \theta(c_0, i_0) = \theta(c_0, i'_0))$$

$$\leq \omega(\theta(c_2, i_2) = \theta(c_0, i_0) \mid \theta(c_0, i_0) = \theta(c_1, i_1)).$$

The conditioning event in the LHS of the above inequality contains information about two objects from the same category. Thus, it can be interpreted as a *similar premise*. The conditioning event in the RHS contains information about two objects from different categories, and it is interpreted as a *diverse premise*. The Proposition says that the probability that new object has given property increases when the premise is diverse. The proof of the Proposition can be found in Appendix D.

Mice have a lower body temperature than humans  


---

Bats have a lower body temperature than humans

Bats have a lower body temperature than humans  


---

Mice have a lower body temperature than humans

TABLE 5. Premise-conclusion asymmetry

**6.3. Premise-conclusion asymmetry.** Finally, I present an example of similarity-based reasoning that cannot be captured in the model of this paper. Consider the inductive arguments in Table 5. (Osherson, Wilkie, Smith, Lopez, and Shafir 1990) report that a majority of students select the first argument as more credible.<sup>7</sup> Following (Rips 1975), they argue that "mice" are more typical animals; hence, it is more informative about properties shared by a general category of similar animals. "Bats" are exotic, and it is not surprising that they have exotic properties.

No Bayesian model exhibits consistent asymmetry of this form. To see it, suppose that the DM has a probability distribution over whether "mice" and "bats" have property  $p$ :  $\mu \in \Delta(\{0, 1\}^{\{M, B\}})$ . Then the conditional probabilities conditional on the premise in the first and second arguments are equal, respectively, to

$$\frac{\mu(y(M) = 1, y(B) = 1)}{\mu(y(M) = 1)} \quad \text{and} \quad \frac{\mu(y(M) = 1, y(B) = 1)}{\mu(y(B) = 1)}.$$

In the benchmark case, the probabilities that "mice" and "bats" have property  $p$  should be equal,  $\mu(y(M) = 1) = \mu(y(B) = 1)$ . (One can think about it as an application of the Principle of Insufficient Reason: Notice that property  $p$  "have lower body temperature than humans" seems to be exchangeable with its logical negation  $\neg p$  "have higher body temperature than humans.") But then both conditional probabilities must be equal.

## 7. CONCLUDING REMARKS

The psychological literature mentions two major functions of categorization (Smith 1995). First, as in this paper, categorization is a tool of inductive inference. Second, categorization serves as a device to code experience without being too demanding on our memory. It is

---

<sup>7</sup>The results seem to be very weak. In the first round of the experiment, 41 out of 80 students pointed to the first argument, and 39 pointed to the second. Only in the second round, when the subjects were explicitly instructed "Although the arguments may seem similar, there is always a difference in how much reason the facts of an argument give to believe its conclusions," did 40 out of 60 students select the first argument.

probably the second role of the categorization that is more connected to bounded rationality and biases. A complete theory should take into account both functions of categorization.

There are various possible extensions of the model of learning through categories. Details of some of the generalizations can be found in the previous version of the current paper (Peski 2006).

*Finite space of outcomes.* So far, a state of the world has been defined as a mapping  $\theta : X \rightarrow Y$ , where the space of outcomes  $Y = \{0, 1\}$  is binary. As a simple extension, consider any finite space  $Y$ . The results do not change, and the proofs change in a predictable way. For example, the prediction function in Section 5 should be redefined as a mapping  $\rho : \{1, \dots, k\}^2 \rightarrow \Delta Y$ .

*Stochastic instance process.* Suppose that instances are drawn from a stochastic distribution  $\mu_X \in \Delta X^\infty$ . Say that the sufficient data condition is satisfied for  $\mu_X$  if it is satisfied almost surely for each of the realizations. If one assumes that the path of instances  $x_1, x_2, \dots$ , is drawn independently from the realization of the state of the world  $\theta$ , then all the results hold  $\mu_X$ -surely.

The independence of instance process and the state of the world achieve the following goal. Consider a scientist who designs experiments, i.e., chooses the instance process, and whose goal is not to predict well, but to find interesting observations. Such a scientist will be interested mostly in outcomes of instances that are difficult to predict. This is because such outcomes are probably most interesting and studying them will increase the knowledge of the scientist. The assumption makes such experiments impossible. I believe that the assumption is sufficient for empirical (or any non-experimental) sciences.

The assumption also eliminates self-selection. For example, one can imagine that students ask about their grades only if they are hopeful of getting a positive grade. However, allowing for a possibility of self-selection should not affect any of the results of this paper. Notice that, if the DM cares only about predictions, self-selection, if biased in a consistent way, only helps: It adds an additional possibility of inference of an outcome from the fact of the query.

*Multiple dimensions of features.* Suppose that the space of instances is equal to  $X = X^1 \times \dots \times X^D$ , where  $X_D$  is infinite. So far I have assumed that  $D = 2$ . When  $D > 2$ , an adequate version of the main result of this paper still holds. In particular, the categorization algorithm categorizes not only each of  $D$  features but also pairs of features, triples, ..., and  $(D - 1)$ -tuples of features.

*Finite symmetry.* One can also to some degree relax the symmetry assumption. Say that distribution  $\omega \in \Delta \Theta$  is finitely symmetric if there are finite partitions  $X^i = \bigcup_{k \leq K} X^{i,(k)}$ , such

that for each  $k, l \leq K$ , the marginal distribution

$$\text{marg}_{\{0,1\}^{X^{1,(k)} \times X^{1,(l)}}} \omega$$

is symmetric. Such a distribution allows for limited a priori categorization of objects and properties where the number of prior categories is bounded by the size of the partition. Up to minor modifications of the proof, Theorem 1 holds for any finitely symmetric  $\omega$ . Thus, adaptive categorization produces predictions asymptotically close to Bayesian even if the distribution over states of the world is finitely symmetric. Corollary 1 follows.

In particular, consider Example 7 and suppose that  $P = P^0 \cup P^1$  is an union of two disjoint sets  $P^0$  and  $P^1$ . There is a bijection  $i : P^0 \rightarrow P^1$  such that  $i(p)$  is interpreted as the property that is a logical negation of property  $p$ . Then, any distribution that is symmetric on  $P^0 \times O$  but otherwise respects logical relationships satisfies the above assumptions.

## REFERENCES

- ALDOUS, D. (1981): “Representations for Partially Exchangeable Arrays of Random Variables,” *Journal of Multivariate Analysis*, 11, 581–598. [8](#), [26](#)
- AZRIELI, Y., AND E. LEHRER (2004): “Categorization Generated by Prototypes - An Axiomatic Approach,” [. 5](#)
- BILLOT, A., I. GILBOA, D. SAMET, AND D. SCHMEIDLER (2005): “Probabilities as similarity-weighted frequencies,” *Econometrica*, 73(4), 1125–1136. [5](#)
- BLACKWELL, D., AND L. DUBINS (1962): “Merging of Opinions with Increasing Information,” *Annals of Mathematical Statistics*, 33, 882–886. [15](#)
- BOUSQUET, O., S. BOUCHERON, AND G. LUGOSI (2004): “Introduction to Statistical Learning Theory,” in *Advanced Lectures in Machine Learning*, ed. by O. Bousquet, U. Luxburg, and G. Rätsch, pp. 169–207. Springer. [15](#)
- DE FINETTI, B. (1964): “La Prevision: Ses Lois Logiques, Ses Sources Subjectives”, in *Studies in Subjective Probability*, ed. by H. E. J. Kyburg, and H. E. Smokler. John Wiley and Sons, New York, translation from French. [8](#)
- GIGERENZER, G., AND P. M. TODD (1999): *Simple Heuristics That Make Us Smart*. Oxford University Press, Oxford. [5](#)
- GILBOA, I., O. LIEBERMAN, AND D. SCHMEIDLER (2005): “Empirical Similarity,” Cowles Foundation Discussion Papers. [5](#)
- GILBOA, I., A. POSTLEWAITE, AND D. SCHMEIDLER (2004): “Rationality of Belief. Or Why Bayesianism is Neither Necessary Nor Sufficient for Rationality,” PIER Working Paper. [3](#)
- GILBOA, I., AND D. SCHMEIDLER (1989): “Maxmin Expected Utility with Non-Unique Priors,” *Journal of Mathematical Economics*, 18, 141–153. [16](#)

- GILBOA, I., AND D. SCHMEIDLER (1995): “Case-Based Decision Theory,” *Quarterly Journal of Economics*, 110(3), 605–639. [2](#), [5](#)
- (2000): “Case-based knowledge and induction,” *Ieee Transactions on Systems Man and Cybernetics Part a-Systems and Humans*, 30(2), 85–95. [5](#)
- GILBOA, I., AND D. SCHMEIDLER (2001): *A Theory of Case-Based Decisions*. University Press, Cambridge, UK. [5](#)
- GILBOA, I., AND D. SCHMEIDLER (2002): “Cognitive foundations of probability,” *Mathematics of Operations Research*, 27(1), 65–81. [5](#)
- GLAZER, J., AND A. RUBINSTEIN (2001): “Debates and Decisions: On a Rationale of Argumentation Rules,” *Games and Economic Behavior*, vol. 36(2), 158–173. [18](#)
- HEWITT, E., AND J. L. SAVAGE (1955): “Symmetric Measures on Cartesian Products,” *Tran. Amer. Math. Soc.*, 80, 470–501. [34](#)
- HOOVER, D. (1982): *Row-Column Exchangeability and a Generalized Model for Probability* North-Holland. [8](#), [26](#)
- HUME, D. (1777): *An Enquiry Concerning Human Understanding*. London. [2](#)
- JACKSON, M. O., AND R. G. FRYER JR (2005): “Categorical Cognition: A Psychological Model of Categories and Identification in Decision Making,” Discussion paper. [4](#), [13](#)
- JACKSON, M. O., E. KALAI, AND R. SMORODINSKY (1999): “Bayesian Representation of Stochastic Processes under Learning: De Finetti Revisited,” *Econometrica*, 67, 875–894. [15](#)
- JAYNES, E. T. (1988): “How Does the Brain Do Plausible Reasoning?,” in *Maximum-Entropy and Bayesian Methods in Science and Engineering*, ed. by G. J. Erickson, and C. R. Smith, p. 1. Kluwer, Dordrecht. [2](#)
- JEHIEL, P. (2005): “Analogy-Based Expectation Equilibrium,” *Journal of Economic Theory*, 123, 81–104. [4](#)
- JEHIEL, P., AND D. SAMET (2007): “Valuation Equilibrium,” *Theoretical Economics*, 2. [5](#)
- KALLENBERG, O. (2005): *Probabilistic Symmetries and Invariance Principles*, Probability and Its Applications. Springer, New York. [8](#), [26](#)
- KEYNES, J. M. (1921): *A Treatise on Probability*. Macmillan, London. [2](#)
- KINGMAN, J. F. C. (1978): “Uses of Exchangeability,” *Ann. Probability*, 6, 183–197. [34](#)
- KREPS, D. M. (1988): *Notes on the Theory of Choice*. Westview Press, Boulder. [3](#)
- LAM, R. (2000): “Learning Through Stories,” Yale University Ph.D. Thesis. [4](#)
- LEHRER, E., AND R. SMORODINSKY (2000): “Relative Entropy in Sequential Decision Problems,” *Journal of Mathematical Economics*, 33, 425–439. [15](#)
- MILL, J. S. (1874): *A System of Logic: Ratiocinative and Inductive*. Harper, New York. [2](#)
- MULLAINATHAN, S. (2002): “Thinking Through Categories,” NBER and MIT Working Paper. [4](#)
- OSHERSON, D. N., O. WILKIE, E. E. SMITH, A. LOPEZ, AND E. SHAFIR (1990): “Category-Based Induction,” *Psychological Review*, 97, 185–200. [2](#), [4](#), [17](#), [20](#)
- PESKI, M. (2006): “Categorization,” University of Chicago, working paper, <http://home.uchicago.edu/mpeski/learning.pdf>. [21](#)
- RIPS, L. J. (1975): “Inductive Judgements About Natural Categories,” *Journal of Verbal learning and Verbal Behavior*, 14, 665–681. [4](#), [20](#)
- SAVAGE, L. J. (1972): *The Foundations of Statistics*. Dover Publications, Toronto. [2](#)
- SMITH, E. E. (1995): “Concepts and Categorization,” in *An Invitation to Cognitive Science, Vol.*, ed. by D. N. Osherson, chap. 1, pp. 3–33. MIT Press, Cambridge. [20](#)

TVERSKY, A. (1977): “Features of Similarity,” *Psychological Review*, 84, 327–352. [2](#), [4](#)

TVERSKY, A., AND D. KAHNEMAN (1974): “Judgment under Uncertainty: Heuristics and Biases,” *Science*, 185, 1124–1131. [4](#)

VAPNIK, V. N. (1998): *Statistical Learning Theory*. Wiley-Interscience. [15](#)

## APPENDIX A. PROOF OF PROPOSITION [2](#)

The proof relies on a well-known connection between entropy and Bayesian updating. Some preliminary remarks are needed. Observe that

$$\begin{aligned}
& \prod_{(x,y) \in d} \rho(c(x))(y) \\
&= \prod_{k^1, k^2} \left( (\rho(k^1, k^2)(1))^{\phi(k^1, k^2|c, d)} (\rho(k^1, k^2)(0))^{(1-\phi(k^1, k^2|c, d))} \right)^{n(k^1, k^2|c, d)} \\
&= \exp \left( \sum_{k^1, k^2} n(k^1, k^2|c, d) \left( \begin{array}{c} \phi(k^1, k^2|c, d) \log(\rho(k^1, k^2)(1)) \\ + (1 - \phi(k^1, k^2|c, d)) \log(\rho(k^1, k^2)(0)) \end{array} \right) \right) \\
&\leq \exp(-|d| E(c|d)),
\end{aligned}$$

because the expression in the third line is maximized when  $\rho(k^1, k^2)(1) = \phi(k^1, k^2|c, d)$ . For any database  $d$ , and any assignment  $c$ , let  $C(c|d)$  be the set of all assignments that coincide with  $c$  on all observations in  $d$ :

$$C(c|d) := \{c' \in \mathcal{C}^k : c(x) = c'(x) \ \forall (x, y) \in d\}.$$

Sets  $C(c|d)$  induce partition space  $\mathcal{C}^k$  into exactly  $k^{d^1+d^2}$  disjoint sets. By the definition of the uniform distribution  $\Psi_C^k$ , for any assignments  $c$  and  $c'$ ,

$$\Psi_C^k(C(c|d)) = k^{-d^1-d^2}.$$

Also, note that  $\psi^k(c, \rho|d)$  depends on assignment  $c$  only up to the observations in database  $d$ : For any  $c' \in C(c|d)$ , any  $\rho \in \mathcal{R}^k$ ,

$$\psi^k(c, \rho|d) = \psi^k(c', \rho|d).$$

By the above,

$$\int_{\{c: E(c|d_t) > E_{\min}(d_t) + \varepsilon\} \times \mathcal{R}^k} \prod_{(x,y) \in d_t} \rho(c(x))(y) d\Psi^k(c, \rho) \leq \exp(-tE_{\min}(d_t) - \varepsilon t).$$

Find an assignment  $c_{\max}$  that minimizes the entropy in database  $d_t$ , i.e.,  $E(c_{\max}|d_t) = E_{\min}(d)$ . Define prediction function  $\rho_{\max}$  such that for each category  $k^1, k^2$ ,

$$\rho_{\max}(k^1, k^2)(1) = \phi(k^1, k^2|c_{\max}, d_t).$$

Denote the set of prediction functions

$$\mathcal{R}_t : \left\{ \rho \in \mathcal{R}^k : \forall_{k^1, k^2} \forall_y \rho(k^1, k^2)(y) \geq e^{-\frac{\varepsilon}{2}} \rho_{\max}(k^1, k^2)(y) \right\}.$$

By the definition of the uniform distribution  $\Psi_R^k$ ,

$$\Psi_R^k(\mathcal{R}_t) \geq \left( \frac{1}{2} (1 - e^{-\frac{\varepsilon}{2}}) \right)^{k^2}.$$

Using the above calculations, for any prediction function  $\rho \in \mathcal{R}_t$ ,

$$\prod_{(x,y) \in d_t} \rho(c_{\max}(x))(y) \geq \exp\left(-t \left(E_{\min}(d_t) + \frac{\varepsilon}{2}\right)\right).$$

Hence,

$$\begin{aligned} & \int_{\{c: E(c|d_t) \leq E_{\min}(d_t) + \varepsilon\} \times \mathcal{R}^k} \prod_{(x,y) \in d_t} \rho(c(x))(y) d\Psi^k(c, \rho) \\ & \geq \int_{\{c \in C(c_{\max}|d)\} \times \mathcal{R}_t} \prod_{(x,y) \in d_t} \rho(c(x))(y) d\Psi^k(c, \rho) \\ & \geq \exp\left(-t \left(E_{\min}(d_t) + \frac{\varepsilon}{2}\right)\right) k^{-d_t^1 - d_t^2} \left(\frac{1}{2} (1 - e^{-\frac{\varepsilon}{2}})\right)^{k^2}. \end{aligned}$$

Observe that

$$\begin{aligned} & \lim_{t \rightarrow \infty} \frac{\int_{\{c: E(c|d_t) \leq E_{\min}(d_t) + \varepsilon\} \times \mathcal{R}^k} \psi^k(c, \rho|d_t) d\Psi^k(c, \rho)}{\int_{\{c: E(c|d_t) > E_{\min}(d_t) + \varepsilon\} \times \mathcal{R}^k} \psi^k(c, \rho|d_t) d\Psi^k(c, \rho)} \\ & = \lim_{t \rightarrow \infty} \frac{\int_{\{c: E(c|d_t) \leq E_{\min}(d_t) + \varepsilon\} \times \mathcal{R}^k} \prod_{(x,y) \in d_t} \rho(c(x))(y) d\Psi^k(c, \rho)}{\int_{\{c: E(c|d_t) > E_{\min}(d_t) + \varepsilon\} \times \mathcal{R}^k} \prod_{(x,y) \in d_t} \rho(c(x))(y) d\Psi^k(c, \rho)} \\ & \geq \left(\frac{1}{2} (1 - e^{-\frac{\varepsilon}{2}})\right)^{k^2} \lim_{t \rightarrow \infty} \frac{\exp\left(-t \left(E_{\min}(d_t) + \frac{\varepsilon}{2}\right)\right) k^{-d_t^1 - d_t^2}}{\exp\left(-t E_{\min}(d_t) - \varepsilon t\right)} \\ & = \left(\frac{1}{2} (1 - e^{-\frac{\varepsilon}{2}})\right)^{k^2} \lim_{t \rightarrow \infty} \exp\left(t \left(\frac{\varepsilon}{2} - \frac{(d_t^1 + d_t^2)}{t} \log k\right)\right) = \infty, \end{aligned}$$

where the limit is a consequence of the sufficient data condition. This finishes the proof of the Proposition.

## APPENDIX B. REPRESENTATION OF SYMMETRIC DISTRIBUTIONS

**B.1. Representation theorem.** In this part of the Appendix, I present a useful representation of a symmetric distribution. This can be described as follows. For any symmetric  $\omega$ , there exists a measurable function  $q : [0, 1]^3 \rightarrow \Delta \{0, 1\}$ , which can be used to generate the state of the world in the following procedure:

- draw variable  $\xi^\emptyset$  from the uniform measure on interval  $[0, 1]$ ,
- for each feature  $i$ , for any  $x^i \in X^i$ , draw independently  $\xi^i(x^i)$  from measure  $U[0, 1]$ ,
- for each instance  $(x^1, x^2) \in X$ , draw independently  $\theta(x^1, x^2)$  from distribution  $q(\xi^\emptyset, \xi^1(x^1), \xi^2(x^2))$

The Representation Theorem says that the distribution of  $\theta$  generated in such a procedure is equal to  $\omega$ . Variables  $\xi^i(x^i)$  are interpreted as categories of features  $x^i$ , or, alternatively, shocks to the outcomes of instances that are feature specific. Variable  $\xi^\emptyset$  can be interpreted as an aggregate shock to the outcomes of all instances.

Formally, let  $\Xi := [0, 1] \times [0, 1]^{X^1} \times [0, 1]^{X^2}$ . I refer to  $\Xi$  as the space of auxiliary variables with a typical element  $\xi \in \Xi$ . For any  $x \in X$ , denote

$$\xi(x) = (\xi^\emptyset, \xi^1(x^1), \xi^2(x^2)) \in [0, 1]^3.$$

$\xi(x)$  is equal to a triple of an auxiliary variable  $\xi^\emptyset$ , an auxiliary variable assigned to feature  $x^1$  and an auxiliary variable assigned to feature  $x^2$ .

**Theorem 2** (Representation of Invariant Distributions). *For any symmetric distribution  $\omega \in \Delta\Theta$ , there a distribution  $\omega^* \in \Delta(\Theta \times \Xi)$ , such that*

- (1)  $\omega = \text{marg}_\Theta \omega^*$ ;
- (2)  $\text{marg}_\Xi \omega^*$  is a product of independent uniform measures on the interval  $[0, 1]$ ;
- (3) there exists a measurable function  $q : [0, 1]^3 \rightarrow \Delta \{0, 1\}$  such that for any  $x$

$$\begin{aligned} q(\xi(x)) &= \omega^*(\theta(x) | \xi(x)) \\ &= \omega^*\left(\theta(x) | \xi, \{\theta(x')\}_{x' \neq x}\right), \end{aligned}$$

*i.e., conditional distribution of  $\theta(x)$ , conditional on the realization of all auxiliary variables  $\xi$ , and all other outcomes  $\theta(x')$ ,  $x' \neq x$ , depends only on the realization of  $\xi(x)$ .*

*Proof.* The Theorem is a restatement of Corollary 7.23 of (Kallenberg 2005). This result was originally proven in (Aldous 1981) and (Hoover 1982).  $\square$

Say that measure  $\omega^*$  represents distribution  $\omega$ . The representing measure is not unique, but the choice of representation is not important for the proof as long as it is fixed. From

now on, instead of writing  $\omega^*$ , I always write  $\omega$ . The second property says that variables

$$\xi^\emptyset, \xi^1(x^1)_{x^1 \in X^1}, \xi^2(x^2)_{x^2 \in X^2}$$

are independent and uniformly distributed on the interval  $[0, 1]$ . The third property says that, conditional on the realization of  $\xi(x)$ , no additional information (apart from observing outcome  $\theta(x)$  itself) affects the prediction of outcome  $\theta(x)$ . In other words, variable  $\xi(x)$  is a *sufficient statistic* for outcome  $\theta(x)$ .

In order to shorten the notation, write  $E_\theta, E_\xi, E_{\theta|\xi}$  to denote expectations with respect to  $\theta \in \Theta, \xi \in \Xi$ , and  $\theta$ , conditional on the realization of  $\xi$ . In particular,

$$E_{\theta, \xi} = E_\xi E_{\theta|\xi}.$$

**B.2. Approximation.** By the above Theorem, each symmetric distribution can be represented by infinitely many categories from interval  $[0, 1]$ . One shows that each such distribution can be approximated by distributions generated only with finitely many categories. This part of the Appendix develops the appropriate notation. Divide  $[0, 1]$  into  $k$  intervals of equal length and for any  $z \in [0, 1]$ , let  $A^k(z) \in \{1, \dots, k\}$  be the index of the covering interval:  $z \in \left[ \frac{A^k(z)-1}{k}, \frac{A^k(z)}{k} \right]$ . For any  $(z_0, z_1, z_2) \in [0, 1]^3$ , define

$$q^k(z_0, z_1, z_2) := E_{z'_1, z'_2} (q(z_0, z'_1, z'_2) | A^k(z_i) = A^k(z'_i) \text{ for } i = 1, 2),$$

where the expectation is taken with respect to the uniform measure on  $[0, 1]^2$ . Hence,  $q^k(z_0, z_1, z_2)$  is equal to the expectation of  $q(z_0, z'_1, z'_2)$  with respect to i.i.d. uniformly distributed  $z'_1$  and  $z'_2$ , conditional on the fact that  $A^k(z_i) = A^k(z'_i)$  for  $i = 1, 2$ .

Let the expected difference between  $q(z)$  and  $q^k(z)$  be denoted as

$$\Delta^k := E_z \|q^k(z) - q(z)\|, \tag{B.1}$$

where the expectation is taken with respect to the uniform measure on  $[0, 1]^3$ . By standard arguments based on the Martingale Convergence Theorem,

$$\lim_{k \rightarrow \infty} \Delta^k = 0. \tag{B.2}$$

Given a realization of auxiliary variables  $\xi$ , define *approximate k-assignment*, defined as  $c^k \in \mathcal{C}^k$ , where

$$c^*(x) := (A^k(\xi^1(x^1)), A^k(\xi^2(x^2))).$$

Define *approximate prediction function*  $q^{k*} \in \mathcal{R}^k$ : for any  $k_i \in \{1, \dots, k\}$

$$\rho^*(k_1, k_2) := q^k \left( \xi^\emptyset, \frac{k_1}{k}, \frac{k_2}{k} \right).$$

Both approximate assignment  $c^*$  and prediction function  $\rho^*$  are defined as functions of  $\xi$  and, thus, they are random variables.

## APPENDIX C. PROOFS OF OPTIMALITY RESULTS

In this part of the Appendix, I present proofs of the optimality results in Section 5.3. The goal is to show that the predictions of categorization algorithms  $l^k$  and  $l_C$  are close to the prediction of  $l_\omega$ . This is shown in three steps. First, I show that the period  $t$  predictions of the categorization algorithm are close to the prediction of  $q(\xi(x_t))$ . Next, I show that  $l_\omega$  is close to  $q(\xi(x_t))$ . The proofs of Proposition 3 and Theorem 1 follow.

C.1. **Convergence**  $l^k \rightarrow q$ . For any  $p, p' \in \Delta\{0, 1\}$ , let

$$D(p, p') = \sum_y p(y) \log \frac{p(y)}{p'(y)}.$$

$D(p, q)$  measures the distance between  $p$  and  $p'$ . In particular,  $D(p, p') = 0$  if and only if  $p = p'$ .

**Proposition 5.** *For any instance process  $\bar{x}$ , for any symmetric  $\omega$ ,*

$$\begin{aligned} E_{\theta, \xi} \sum_{s < t} D(q(\xi(x_s)), l^k(d_s, x_s)) \\ \leq 4 \left( \sqrt{\Delta^k} \log \frac{1}{\Delta^k} + 6\sqrt{\Delta^k} \right) t + k^2 \log k + k(d_t^1 + d_t^2). \end{aligned}$$

(Recall that constant  $\Delta^k$  is defined in equation (B.2).) The Proposition derives a bound on the distance between prediction  $q(\xi(x_s))$  from the Representation Theorem and the prediction of the  $k$ -categorization algorithm. The rest of this part of the Appendix proves the Proposition.

C.1.1. *Notation.* For any database  $d$  and any assignment  $c$ , let  $C(c|d)$  be the set of assignments that coincide with  $c \in \mathcal{C}^k$  on all observations in  $d$  (recall the definition from the proof of Proposition 2). For any  $\delta < \frac{1}{2}$ , any  $\rho \in \mathcal{R}^k$ , let

$$B(\rho, \delta) = \left\{ \rho' \in \mathcal{R}^k : \begin{array}{l} \sup_{\mathbf{k} \in \{1, \dots, k\}^2} \|\rho(\mathbf{k}) - \rho'(\mathbf{k})\| \leq 2\delta, \\ \inf_{\mathbf{k} \in \{1, \dots, k\}^2} \inf_y \rho'(\mathbf{k})(y) \geq \delta. \end{array} \right\}.$$

$B(\rho, \delta)$  is the set of prediction functions  $\rho'$  such that for each category  $\mathbf{k} \in \{1, \dots, k\}^2$ ,

- the predictions of  $\rho$  and  $\rho'$  differ by at most  $2\delta$ , and
- $\rho'$  assigns probability at least  $\delta$  to each outcome.

The "prior" probability of set  $B(\rho, \delta)$  is bounded from below by

$$\Psi_R^k(B(\rho, 2\delta)) \geq \delta^{k^2}.$$

Fix a realization of auxiliary variables  $\xi$ . Recall that  $\psi^k(\cdot|d)$  denotes the posterior density of the DM's beliefs from the  $k$ -categorization algorithm. Let  $\Psi(\cdot|d)$  denote the corresponding

distribution. Also, recall that  $d_s$  is a database of observations made before period  $s$ . Further denote

$$G(\xi, \delta) := \{(c, \rho) : c \in C(c^*|d_t), \rho \in B(\rho^*, \delta)\},$$

$$l_s^\delta(y) = \frac{1}{\Psi(G(\xi, \delta)|d_s)} \int_{G(\xi, \delta)} \rho(c(x_s))(y) d\Psi(c, \rho|d_s),$$

Set  $G(\xi, \delta)$  consists of theories  $(c, \rho)$  such that

- assignment  $c$  agrees with approximate assignment  $c^*$  on database  $d_t$ , and
- prediction function  $\rho$  is close to prediction  $\rho^*$ .

Prediction  $l_s^\delta$  is equal to the prediction of the  $k$ -categorization algorithm, conditional on the fact that the theory belongs to set  $G(\xi, \delta)$ . In particular,

$$\|q^k(\xi(x_s)) - l_s^\delta\| = \|\rho^*(c^*(x_s)) - l_s^\delta\| \leq 3\delta.$$

### C.1.2. Intermediary lemmas.

#### Lemma 1.

$$E_{\theta, \xi} \sum_{s < t} D(q(\xi(x_s)), l_s^\delta) \leq 2 \left( \frac{\Delta^k}{\delta} \log \frac{1}{\delta} + 6\delta \right) t$$

*Proof.* Notice that

$$\begin{aligned} & E_{\theta, \xi} \sum_{s < t} D(q(\xi(x_s)), l_s^\delta) \\ & \leq E_{\theta, \xi} \sum_{s < t} \mathbf{1} \{ \|q(\xi(x_s)) - q^k(\xi(x_s))\| \geq \delta \} \max_{p \in \Delta\{0,1\}} D(p, l_s^\delta) \\ & \quad + E_{\theta, \xi} \sum_{s < t} \mathbf{1} \{ \|q(\xi(x_s)) - q^k(\xi(x_s))\| \leq \delta \} D(q(\xi(x_s)), l_s^\delta). \end{aligned}$$

To bound the first term, recall that by construction  $l_s^\delta(y) \geq \delta$ ; hence,  $D(q, l_s^\delta) \leq 2 \log \frac{1}{\delta}$  for any  $q \in \Delta\{0,1\}$ . By Chebyshev's inequality,

$$E_{\theta, \xi} \sum_{s < t} \mathbf{1} \{ \|q(\xi(x_s)) - q^k(\xi(x_s))\| \geq \delta \} \leq t \frac{\Delta^k}{\delta},$$

where  $\Delta^k$  is defined in equation (B.1).

For the second term, notice that if  $\|q(\xi(x_s)) - q^k(\xi(x_s))\| \leq \delta$ , then  $\|q(\xi(x_s)) - l_s^\delta\| \leq 3\delta$ , and

$$\begin{aligned} D(q(\xi(x_s)), l_s^\delta) &\leq \sum_y (l_s^\delta(y) + 3\delta) \log \left(1 + \frac{3\delta}{l_s^\delta(y)}\right) \\ &\leq 3\delta \sum_y \frac{l_s^\delta(y) + 3\delta}{l_s^\delta(y)} \leq 12\delta. \end{aligned}$$

The last inequality uses the fact that  $l_s^\delta(y) \geq \delta$  by construction.  $\square$

**Lemma 2.**

$$\begin{aligned} E_{\theta, \xi} \sum_{s < t} \sum_y q(\xi(x_s))(y) \log \frac{l_s^\delta(y)}{l^k(d_s, x_s)(y)} \\ \leq k^2 \log \frac{1}{\delta} + k(d_t^1 + d_t^2). \end{aligned}$$

*Proof.* Notice that for each  $y$ ,

$$\begin{aligned} &\frac{l_s^\delta(y)}{l^k(d_s, x_s)(y)} \\ &= \frac{1}{\Psi(G(\xi, \delta) | d_s)} \frac{\int_{G(\xi, \delta)} \rho(c(x))(y) d\Psi(c, \rho | d_s)}{\int_{T^k} \rho(c(x))(y) d\Psi(c, \rho | d_s)} \\ &= \frac{1}{\Psi(G(\xi, \delta) | d_s)} \frac{\Psi(\theta(x_s) = y, G(\xi, \delta) | d_s)}{\Psi(\theta(x_s) = y | d_s)} \\ &= \frac{\Psi(G(\xi, \delta) | d_s, \theta(x_s) = y)}{\Psi(G(\xi, \delta) | d_s)}. \end{aligned}$$

Here,  $\{d_s, \theta(x_s) = y\} = d_{s+1}$  is a database of past observations in period  $t + 1$ . Hence,

$$\begin{aligned} &\sum_y q(\xi(x_s))(y) \log \frac{l_s^\delta(y)}{l^k(d_s, x_s)(y)} \\ &= \sum_y q(\xi(x_s))(y) \log \frac{\Psi(G(\xi, \delta) | d_{s+1})}{\Psi(G(\xi, \delta) | d_s)} \\ &= E_{d_{s+1} | d_s, \xi} \log \frac{\Psi(G(\xi, \delta) | d_{s+1})}{\Psi(G(\xi, \delta) | d_s)}, \end{aligned}$$

where  $E_{d_{s+1} | d_s, \xi}$  is the expectation of the realization of outcome  $\theta(x_s)$ , conditional on the realization of outcomes  $\theta(x_{s'})$ ,  $s' < s$ , and auxiliary variables  $\xi$ . (The last equality follows from the conditional independence of point 3 of the Representation Theorem.) By the Law

of Iterated Expectations,

$$\begin{aligned}
 & E_{\theta, \xi} \sum_{s < t} \sum_y q(\xi(x_s))(y) \log \frac{l_s^\delta(y)}{l^k(d_s, x_s)(y)} \\
 &= E_\xi \left[ E_{d_1 | \xi} \left[ \log \frac{\Psi(G(\xi, \delta) | d_1)}{\Psi(G(\xi, \delta) | d_0)} + \dots \left[ \dots + E_{d_t | d_{t-1}, \xi} \left[ \log \frac{\Psi(G(\xi, \delta) | d_t)}{\Psi(G(\xi, \delta) | d_{t-1})} \right] \right] \right] \right] \\
 &= E_{d_t, \xi} \log \frac{\Psi(G(\xi, \delta) | d_t)}{\Psi(G(\xi, \delta))} \leq -E \log \Psi(G(\xi, \delta)) \\
 &= -E \log \Psi_C(C(c^* | d_t)) - E \log \Psi_R(B(\rho^*, 2\delta)).
 \end{aligned}$$

The thesis of the Lemma is a consequence of the fact, that for any assignment  $c \in \mathcal{C}^k$  and for any prediction function  $\rho \in \mathcal{R}^k$ ,

$$\begin{aligned}
 \log \Psi_C(C(c|d)) &\geq -k(d_t^1 + d_t^2) \quad \text{and} \\
 \log \Psi_R(B(\rho, 2\delta)) &\geq k^2 \log \delta.
 \end{aligned}$$

□

C.1.3. *Proof of the Proposition.* Notice that for any  $\rho' \in \Delta\{0, 1\}$ ,

$$D(q, p) = D(q, p') + \sum_y q(y) \log \frac{p'(y)}{p(y)}. \tag{C.1}$$

By (C.1) and the Lemmas,

$$\begin{aligned}
 & E_{\theta, \xi} \sum_{s < t} D(q(\xi(x_s)), l^k(d_s, x_s)) \\
 &= E_{\theta, \xi} \sum_{s < t} D(q(\xi(x_s)), l_s^\delta) \\
 &+ E_{\theta, \xi} \sum_{s < t} \sum_y q(\xi(x_s))(y) \log \frac{l_s^\delta(y)}{l^k(d_s, x_s)(y)} \\
 &\leq 2 \left( \frac{\Delta^k}{\delta} \log \frac{1}{\delta} + 6\delta \right) t + k^2 \log \frac{1}{\delta} + k(d_t^1 + d_t^2).
 \end{aligned}$$

Take  $\delta = \max\left(\sqrt{\Delta^k}, \frac{1}{k}\right)$  to finish the proof of the Proposition.

C.2. **Convergence**  $l_\omega \rightarrow q$ .

**Proposition 6.** *Suppose that instance process  $x_1, x_2, \dots$  satisfies the sufficient data condition. Then,*

$$\lim_{t \rightarrow \infty} \frac{1}{t} E_{\theta, \xi} \sum_{s < t} \|q(\xi(x_s)) - l_\omega(d_s, x_s)\| = 0. \tag{C.2}$$

Denote behavioral rules (see Section 5.4)

$$a^k = a_{\max} \circ l^k \text{ and } a_\omega = a_{\max} \circ l_\omega.$$

By Proposition 5 and convergence (B.2),

$$\lim_{k \rightarrow \infty} U(a^k; \omega, \bar{x}) = \lim_{t \rightarrow \infty} \inf_{E_{\xi, \theta}} E_{\xi, \theta} \frac{1}{t} \sum_{s < t} u(q(\xi(x_s)), \theta(x_t)).$$

Let  $E_{\xi|d}$  denote the expectation with respect to realization of the auxiliary variable  $\xi$ , conditional on the observed database of past cases  $d_t$ . Notice that

$$E_{\xi|d} q(\xi(x_s))(y) = \omega(\theta(x_s) | d_s) = l_\omega(d_s, x_s).$$

By standard convexity arguments, this implies that

$$U(a_\omega; \omega, \bar{x}) \leq \lim_{t \rightarrow \infty} \inf_{E_{\xi, \theta}} E_{\xi, \theta} \frac{1}{t} \sum_{s < t} u(q(\xi(x_s)), \theta(x_t))$$

with equality if and only if (C.2) holds. Because for each  $k$

$$U(a^k; \omega, \bar{x}) \leq U(a_\omega; \omega, \bar{x}),$$

Proposition 5 together with (B.2) imply that (C.2) must be true.

**C.3. Proof of Proposition 3.** The Proposition is a consequence of Propositions 5, 6, and convergence (B.2).

**C.4. Proof of Theorem 1.** Due to Proposition 6, it is enough to show that

$$\limsup_{t \rightarrow \infty} \frac{1}{t} E_{\theta, \xi} \sum_{s < t} \|l^{K_s}(d_s(K_s), x_s) - q(\xi(x_s))\| = 0.$$

Because  $D(p, p')$  is continuous in  $p, p' \in \Delta\{0, 1\}$  and  $D(p, p') = 0$  if and only if  $p = p'$ , it is enough to show that

$$\begin{aligned} 0 &= \limsup_{t \rightarrow \infty} \frac{1}{t} E_{\theta, \xi} \sum_{s < t} D(q(\xi(x_s)), l^{K_s}(d_s(K_s), x_s)) \\ &= \limsup_{t \rightarrow \infty} \sum_k \frac{m_t^k}{t} E_{\theta, \xi} \frac{1}{m_t^k} \sum_{s < t: K_s = k} D(q(\xi(x_s)), l^k(d_s(K_s), x_s)) \end{aligned} \quad (\text{C.3})$$

where

$$m_t^k = \#\{s < t : K_s = k\}.$$

Denote the last period in which  $k$ -assignment is used before period  $t$

$$s_t^k = \max\{s < t : K_s = k\}$$

and let  $s_t^k = \infty$  if  $k$  assignment is never used before period  $t$ . Then, by Proposition 5, (C.3) is not higher than

$$\begin{aligned} & \limsup_{t \rightarrow \infty} \frac{1}{t} E_{\theta, \xi} \sum_{s < t} D(q(x_s), l^{K_s}(d_s(K_s), x_s)) \\ & \leq \limsup_{t \rightarrow \infty} 4 \sum_k \sqrt{\Delta^k} \left( \log \frac{1}{\Delta^k} + 6 \right) \frac{m_t^k}{t} \end{aligned} \quad (\text{C.4})$$

$$+ \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{k: s_t^k < \infty} \left[ k^2 \log k + k \left( d_{s_t^k}^1 + d_{s_t^k}^2 \right) \right] \quad (\text{C.5})$$

By (B.2),

$$\lim_{k \rightarrow \infty} \sqrt{\Delta^k} = 0 \text{ and } \lim_{k \rightarrow \infty} \sqrt{\Delta^k} \log \frac{1}{\Delta^k} = 0.$$

Because the process satisfies the sufficient data condition,  $\lim_{t \rightarrow \infty} K_t = \infty$  and, for each  $k$ ,  $\lim_{t \rightarrow \infty} \frac{m_t^k}{t} = 0$ . Because  $\sum_k \frac{m_t^k}{t} = 1$ , term (C.4) of the inequality above is equal to 0.

By the definition of  $d^i$ ,

$$d_t^1 + d_t^2 \geq \sqrt{t}.$$

By the definition of switching rule  $K_t$  in equation (5.1), if  $s_t^k < \infty$ , then,

$$\begin{aligned} d_{s_t^k}^1 + d_{s_t^k}^2 & \leq \frac{s_t^k}{k^3} \text{ and} \\ \sqrt{s_t^k} & \leq d_{s_t^k}^1 + d_{s_t^k}^2 \leq \frac{s_t^k}{k^3}. \end{aligned}$$

Therefore, if  $s_t^k < \infty$ , then

$$k \leq \sqrt[6]{s_t^k} \leq \sqrt[6]{t} := k_t^{\max}.$$

This permits to bound term (C.5). Because  $\log k \leq k$ ,

$$\begin{aligned} & \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{k: s_t^k < \infty} \left[ k^2 \log k + k \left( d_{s_t^k}^1 + d_{s_t^k}^2 \right) \right] \\ & \leq \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{k=1}^{k_t^{\max}} k^3 + \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{k=1}^{K_t} k \left( d_{s_t^k}^1 + d_{s_t^k}^2 \right) + \limsup_{t \rightarrow \infty} \frac{1}{t} \sum_{k=K_t+1}^{k_t^{\max}} k \frac{s_t^k}{k^3} \\ & \leq \lim_{t \rightarrow \infty} \frac{1}{t} (k_t^{\max})^4 + \limsup_{t \rightarrow \infty} \frac{K_t^2 (d_t^1 + d_t^2)}{t} + \limsup_{t \rightarrow \infty} \sum_{k=K_t+1}^{\infty} \frac{1}{k^2} \\ & \leq \lim_{t \rightarrow \infty} t^{-\frac{1}{3}} + \limsup_{t \rightarrow \infty} \frac{1}{K_t^3} + \limsup_{t \rightarrow \infty} \frac{1}{K_t} = 0. \end{aligned}$$

## APPENDIX D. PROOF OF PROPOSITION 4

Let  $\Delta_I = \Delta \{0, 1\}^I$  be the distribution of outcomes inside each category. Distribution  $\omega$  can be stated as a distribution over product of  $C$  elements of  $\Delta_I : \omega \in \Delta (\Delta_I)^C$ . Because of exchangeability of categories,  $\omega$  is De Finetti exchangeable. By an extension of the De Finetti's Theorem to Polish spaces ((Hewitt and Savage 1955), (Kingman 1978)), there is a distribution  $\mu_\omega \in \Delta (\Delta_I)$  measurable function  $q : [0, 1] \rightarrow \Delta_I$  such that, for any finite set of objects  $O' \subseteq O$ , any assignment  $y : O' \rightarrow \{0, 1\}$ ,

$$\begin{aligned} & \omega (\theta (o) = y (o) \text{ for } o \in O') \\ &= \int_{\Delta_I} d\mu_\omega (\omega_I) \prod_{c \in C} \omega_I (\theta (c, i) = y (c, i) \text{ for all } i \text{ st. } (c, i) \in O'). \end{aligned}$$

Because of exchangeability inside categories, measure  $\mu_\omega$  assigns probability 1 to exchangeable measures  $\omega_I \in \Delta_I = \Delta \{0, 1\}^I$ . By the De Finetti's theorem, any exchangeable measure  $\omega_I$  is represented by a distribution  $\mu_I \in \Delta [0, 1]$  so that

$$\begin{aligned} & \omega (\theta (o) = y (o) \text{ for } o \in O') \\ &= \int_{\Delta_I} d\mu_\omega (\mu_I) \prod_{c \in C} \int_{[0,1]} d\mu_I (p) \prod_{i:(c,i) \in O'} p^{y(c,i)} (1-p)^{1-y(c,i)}. \end{aligned}$$

Then,

$$\begin{aligned} A &:= \omega (\theta (c_2, i_2) = \theta (c_0, i_0) = \theta (c_0, i'_0)) \\ &= \int_{\Delta_I} \left( \begin{aligned} & \left[ \int_{[0,1]} p^2 d\mu_I (p) \right] \left[ \int_{[0,1]} p d\mu_I (p) \right] \\ & + \left[ \int_{[0,1]} (1-p)^2 d\mu_I (p) \right] \left[ \int_{[0,1]} (1-p) d\mu_I (p) \right] \end{aligned} \right) d\mu_\omega (\mu_I) \\ &= \int_{\Delta_I} \left( \left[ \int_{[0,1]} d\mu_I (p) p \right]^3 + \left[ \int_{[0,1]} d\mu_I (p) (1-p) \right]^3 \right) d\mu_\omega (\mu_I) \\ &+ \int_{\Delta_I} \left( \text{Var}_{\mu_I} p \left[ \int_{[0,1]} p d\mu_I (p) \right] + \text{Var}_{\mu_I} p \left[ \int_{[0,1]} (1-p) d\mu_I (p) \right] \right) d\mu_\omega (\mu_I) \\ &= \int_{\Delta_I} \left( \left[ \int_{[0,1]} d\mu_I (p) p \right]^3 + \left[ \int_{[0,1]} d\mu_I (p) (1-p) \right]^3 \right) d\mu_\omega (\mu_I) \\ &+ \int_{\Delta_I} (\text{Var}_{\mu_I} p) d\mu_\omega (\mu_I), \end{aligned}$$

where  $\text{Var}_{\mu_I}$  is the variance of random variable  $p$  with respect to distribution  $\mu_I$ . Similarly,

$$\begin{aligned} B &:= \omega(\theta(c_0, i_0) = \theta(c_0, i'_0)) \\ &= \int_{\Delta_I} \left( \left[ \int_{[0,1]} d\mu_I(p) p \right]^2 + \left( \int_{[0,1]} d\mu_I(p) (1-p) \right)^2 \right) d\mu_\omega(\mu_I) \\ &\quad + 2 \int_{\Delta_I} (\text{Var}_{\mu_I} p) d\mu_\omega(\mu_I), \end{aligned}$$

$$\begin{aligned} C &:= \omega(\theta(c_2, i_2) = \theta(c_0, i_0) = \theta(c_1, i_1)) \\ &= \int_{\Delta_I} \left( \left[ \int_{[0,1]} d\mu_I(p) p \right]^3 + \left[ \int_{[0,1]} d\mu_I(p) (1-p) \right]^3 \right) d\mu_\omega(\mu_I), \end{aligned}$$

$$\begin{aligned} D &:= \omega(\theta(c_0, i_0) = \theta(c_1, i_1)) \\ &= \int_{\Delta_I} \left( \left[ \int_{[0,1]} d\mu_I(p) p \right]^2 + \left( \int_{[0,1]} d\mu_I(p) (1-p) \right)^2 \right) d\mu_\omega(\mu_I). \end{aligned}$$

Denote  $x = \int_{\Delta_I} (\text{Var}_{\mu_I} p) d\mu_\omega(\mu_I) \geq 0$ . Then,

$$\begin{aligned} \omega(\theta(c_2, i_2) = \theta(c_0, i_0) \mid \theta(c_0, i_0) = \theta(c_0, i'_0)) &= \frac{A}{B} = \frac{C+x}{D+2x}, \\ \omega(\theta(c_2, i_2) = \theta(c_0, i_0) \mid \theta(c_0, i_0) = \theta(c_1, i_1)) &= \frac{C}{D}. \end{aligned}$$

Finally, observe that

$$\begin{aligned} &2C - D \\ &= 2 \int_{\Delta_I} \left( 1 - 3 \left[ \int_{[0,1]} d\mu_I(p) p \right] + 3 \left[ \int_{[0,1]} d\mu_I(p) p \right]^2 \right) d\mu_\omega(\mu_I) \\ &\quad - \int_{\Delta_I} \left( 1 - 2 \left[ \int_{[0,1]} d\mu_I(p) p \right] + 2 \left[ \int_{[0,1]} d\mu_I(p) p \right]^2 \right) d\mu_\omega(\mu_I) \\ &= \int_{\Delta_I} \left( 1 - 4 \left[ \int_{[0,1]} d\mu_I(p) p \right] + 4 \left[ \int_{[0,1]} d\mu_I(p) p \right]^2 \right) d\mu_\omega(\mu_I) \\ &= \int_{\Delta_I} \left( 1 - 2 \left[ \int_{[0,1]} d\mu_I(p) p \right] \right)^2 d\mu_\omega(\mu_I) \geq 0. \end{aligned}$$

Hence,  $\frac{C}{D} \geq \frac{1}{2}$  and

$$\frac{A}{B} = \frac{C+x}{D+2x} \leq \frac{C}{D}.$$

UNIVERSITY OF CHICAGO, DEPARTMENT OF ECONOMICS

*E-mail address:* [mpeski@uchicago.edu](mailto:mpeski@uchicago.edu)