

# Prediction and Stochastic Choice

## Abstract

In this paper, we study a non-parametric approach to prediction in stochastic choice models in economics. We show that VC complexity characterises the predictability of stochastic choice models. We establish prediction methods and provide corresponding rates of convergence.

Keywords: Statistical Learning; VC dimension; Stochastic Choice

2010 Mathematics Subject Classification: 62F10; 62G05

## 1 Introduction

The study of stochastic choice models in economics seeks to understand aggregate demand behaviour and preference heterogeneity (distribution over preferences) in markets (McFadden [1978], Berry et al. [1995]). A typical model of stochastic choice involves a set of alternatives from which an agent makes a choice and based on a set of underlying characteristics ( $x \in X$ ), the model prescribes the probabilities with which various items ( $a \in A$ ) would be chosen. This is defined in terms of a stochastic choice function  $\mu : X \rightarrow \Delta(A)$ . The probabilities  $\mu(x)$  can be interpreted as the proportion of times a particular alternative is chosen compared to other feasible alternatives. Yet another interpretation involves the choice probabilities resulting from randomisation on part of the agent. The nature of this randomisation depends on the context and details of the decision making process (Manzini and Mariotti [2007], Fudenberg et al. [2015]). Online retail companies such as Amazon, are often interested in assessing consumer data  $(x_i; a_i)_{i=1}^n$ , to predict market shares and learn consumer

preferences, by fitting the data through an economic model, in an effort to improve sales. In this paper,

we study a novel machine learning approach to estimating choice probabilities, which correspond to a model of stochastic choice, based on ideas from statistical learning theory.

We study the problem of learning the map  $\mu$  from finite data on choices. Hence, each data point consists of a pair  $(x_i; a_i)$  where  $a_i$  is the alternative chosen when the characteristics were given by the vector  $x_i$ . We are interested in prediction methods defined on choice data which lead to accurate estimation of choice probabilities and the criterion we require is uniform consistency (Valiant [1984], Vapnik [1998], Kearns and Schapire [1994], Yamanishi [1992]). This requires that for any given precision parameters  $\epsilon; \delta > 0$ , there exists a fixed data size  $N(\epsilon; \delta)$ , such that if the analyst were to apply the prediction method to a data set of size above  $N(\epsilon; \delta)$ , the probability would be at least  $1 - \delta$  that the stochastic choice function conjectured by the prediction method would be close to the true stochastic choice function by at most  $\epsilon$ . Moreover, this holds true irrespective of the process that generates choice problems and the true function  $\mu_0$ . Hence, working with amount of data at least  $N(\epsilon; \delta)$  allows for robust estimates of the choice probabilities, when prediction methods are uniformly consistent. This is a desirable feature of the notion of consistency considered here. In many contexts, it is natural to assume that the analyst knows neither the distribution over characteristics nor the true stochastic choice function, and may wish to work with data sizes that guarantee, in a robust manner, a certain degree of precision in estimation. The minimum number of samples  $N(\epsilon; \delta)$ , which provides such a guarantee is called the sample complexity of the model of stochastic choice and is indeed a central object of study in the present paper.

The estimation strategy developed in this paper relies on the principle of empirical risk minimization (Vapnik [1998], Niyogi [2012], Niyogi [2006], Niyogi et al. [2011]). This involves setting up a risk function to evaluate the goodness of fit of the model. Ideally, these would be constructed in a way that the true choice probabilities would minimise expected risk. Then, since the data generating process is unknown, we approximate the expected risk with the empirical risk on the sample. For large enough samples, this would lead to good estimates for choice probabilities. We define risk in terms of incentive compatible scoring rules which are, in principle, reward schemes to incentivise truthful report of subjective probability judgements (Savage [1971], Gneiting and Raftery [2007]). Hence, when applied as a risk function, it ensures that true choice probabilities minimise expected risk. Furthermore, each scoring rule gives rise to a divergence function which allows us to quantify expected regret in terms of divergence in predictions between the true and estimated choice probabilities.

The criterion for consistent learning, in turn, relies on this divergence converging uniformly to zero.

Hence, scoring rules serve as a natural candidate for defining risk and a key result in the paper establishes that they admit uniform learning with respect to their associated divergence functions (Proposition 3.1). Finally, we obtain sample complexity bounds by applying complexity measures such as Rademacher complexity and Pollard dimension. The details of the approach can be found in Section 2.

There are several reasons as to why the above problem would be of relevance to economists.

We discuss here a few points, highlighting how the present work contributes to the empirical and theoretical literature. The estimation of choice probabilities in discrete choice models has for a long

time been of interest to empiricists and econometricians interested in studying choice behaviour in markets (Manski [1975], McFadden [1978], Berry et al. [1995]). The novelty of the present approach lies in the introduction and assessment of sample complexity. This has two main advantages. Firstly, the sample size  $N(\epsilon; \mathcal{C})$  guarantees robust estimates which are independent of the random process that generates the data. Hence, the analyst/econometrician can safely rely on such a sample size to achieve a predetermined level of accuracy in estimates  $\epsilon; \epsilon > 0$ . This is the key consequence of uniform consistency. Indeed, it also implies the usual consistency notion adopted in econometrics, where only convergence to the true parameter is required without a robustness guarantee in terms of sample sizes. Hence, in the usual definition of consistency, depending on the data generating process, more or less number of samples may be needed for accurate estimation and a priori, the analyst may not know how much data would be needed, for which the present approach would provide guarantees. Sample complexity yields uniform rates of convergence. Secondly, sample complexity also acts as a measure of complexity for stochastic choice models. Different models would have

See, for example, Shalev-Shwartz and Ben-David [2014].  
 Regret is defined simply as the difference between the expected risk of the estimated choice probabilities and the expected risk of the true choice probabilities.

Prediction and Stochastic Choice, Pathikrit Basu

different sample complexity and the models with higher sample complexity, would need more samples for good estimation. Hence, in a sense, we can compare different models based on such measures (say Logit v/s Probit) and this allows the analyst or econometrician to make a formal judgement as to which stochastic choice models are simple and which ones are complex.

We now discuss how the present work pertains to the theoretical literature on stochastic choice.

Typically, in theoretical work, the stochastic choice function is treated as the primitive and is itself interpreted as the data (Luce [1959], Luce [1977]). However, real data sets only involve finitely many data points and typically, each data point represents a choice made by an individual. This suggests a gap between the assumption and the nature of real data sets and the logic that seems to justify the gap is that there is perhaps a "law of large numbers" argument that one could rely on. However, from a formal standpoint, it is not clear whether this logic can always be implemented. For example, it may be that for some instances (such as observing a choice from a menu) there may be enough data points to evoke a law of large numbers argument but not for other instances. Hence, estimation of the entire map  $\mu$  from finite data on choices seems far from obvious especially when keeping in mind the wide variety of decision procedures that the theoretical literature considers and which lead to very sophisticated stochastic choice models. Do all such models admit accurate estimation or predictability of choice probabilities?

This last question also brings us to a point of contrast between the theoretical and empirical literature.

On the one hand, the theoretical models investigate varied contexts and decision procedures. If we assume these models as plausible in understanding choice behaviour, then perhaps a framework such as the present one could provide us a way to reason, in a simple manner, about the estimation and identification of all these models, allowing us to make good predictions about choice. However, on the other hand, empirical work still largely deals with more classical utility models and often for more complicated models, makes several distributional assumptions (on the data generating process) which are needed to obtain consistency results. In the present approach, we show that mild assumptions are needed for estimation. The main result of the paper is that in the context of stochastic preferences, if the preference class has finite VC dimension, then the corresponding stochastic choice model is predictable and allows us to recover the underlying distribution over preferences in the population, uniformly over the data generating process. Under some conditions, with a modified definition of VC dimension with strict preferences, we also obtain a converse result which says that a stochastic preference model is non-predictable if the dimension of the underlying preference class is infinite. This allows us to conclude that for any prediction method and any sample size, the worst case expected distance between true and estimated choice probabilities (in  $\mu(A)$ ) is at least half times the maximum distance i.e  $1 =$

$\frac{1}{2}$

2. In situations involving choice between exactly

two alternatives, this lower bound expected error is above the error level guaranteed by a rule which always makes the prediction that each alternative will be chosen with equal probability.

We provide several applications of our results for specific preference models commonly encountered in Industrial organisation, Decisions under Risk/Uncertainty and Spatial voting. We derive general bounds on rates of convergence based on Rademacher complexity. In particular, when the VC dimension is  $d$  (finite) and there are  $|A|$  many alternatives to choose from, then the Rademacher complexity associated with the stochastic choice model is at most  $O(\frac{1}{\sqrt{n}} \log(|A|)d \log(d))$ ,

$\frac{1}{\sqrt{n}}$

$\log(|A|)d \log(d)$ ,

which gives us an upper bound on the uniform rate of convergence of the estimator, which is defined via empirical risk minimisation with the quadratic scoring rule. We also show that for homothetic preferences specifically, convergence takes place in a stronger sense, the consistent estimates of the stochastic choice function converge in the sup norm, with a proportional rate of convergence. Lastly, we consider non-uniform learning and also address recoverability of cognitive heterogeneity in the setting of the bounded rationality model of choice with consideration sets.

3

Prediction and Stochastic Choice, Pathikrit Basu

**Related Literature** We discuss here, in some more detail, the relationship with prior work in the related literature. Typically, estimation of discrete choice models involves parametric or semiparametric assumptions (McFadden [1984], Manski [1975], Manski [1985]). The most widely used method for estimation is maximum likelihood. In our setting, maximum likelihood arises as a special case when the scoring rule applied is the log rule. We also show that a variant of Manski's maximum score estimation method (Manski [1975]) corresponds to use of an incentive compatible scoring rule. Non-parametric approaches have also been studied in the economics literature (Matzkin [1992], Briesch et al. [2010]). The present paper has a non-parametric approach as well. However, the estimation/learning problem considered here is different in that we consider estimation of choice probabilities and our approach involves techniques and ideas from statistical learning theory at the core. In addition, as noted earlier, the consideration of sample complexity is one aspect that is novel in the present analysis. It is also interesting to compare the results of this paper with McFadden and Train [2000]. The paper establishes that choice probabilities of a random utility model can be approximated by choice probabilities corresponding to a mixed logit specification. In contrast, in the present setting, mixed logit being a subset of linear preferences, would be uniformly predictable (with linear complexity), while all continuous preferences would be non-predictable (with lower bounded predictive error). Hence, the dimensions of preferences ( $V_C$  and  $V_{C^+}$  in this paper) play a critical role and its connection to predictability exhibits a perhaps interesting difference between the mixed logit model and stochastic choice from more general continuous preferences. Relatedly, De Blasi et al. [2010] consider a Bayesian non-parametric approach to the problem of recovering preference heterogeneity in the mixed logit framework via estimation of choice probabilities.

Scoring rules have been applied in prior work to study certain binary classification problems (Buja et al. [2005], Parry et al. [2016]) in other settings with different objectives. In contrast, the problem considered here is of consistent uniform learning for stochastic functions in the sense of PAC (probably approximately correct) learning (Valiant [1984]). This paper is hence, most closely related to prior work in probabilistic concept learning (Kearns and Schapire [1994], Abe et al. [1991]) and learning stochastic rules (Yamanishi [1992], Abe et al. [2001]). While the nature of the learning problem is similar, the present paper is different in that we apply scoring rules to study the learning problem. For instance, one interesting feature is that the related papers mention predictability with respect to specific divergences, which has a very natural connection to scoring rules. We show in this paper how the present approach exploits this connection and this is indeed a key aspect of our estimation strategy. Our approach is also more general and subsumes some of the learning problems studied in previous work. We should also emphasise that we are primarily interested here in stochastic choice models in economics and are motivated by problems in demand estimation.

The application of PAC learning in economics has been studied in different choice contexts in prior work (Kalai [2003], Beigman and Vohra [2006], Zadimoghaddam and Roth [2012], Basu and Echenique [2020], Chambers et al. [2021], Chase and Prasad [2018]). However, the literature has only considered deterministic choice models. The distinguishing feature of the present work in relation to this strand of the literature is the consideration of stochastic choice models. We should also say that the present paper contributes to the growing literature on machine learning methods in economics (Athey and Imbens [2019]).

The outline of the paper is as follows. In section 2, we present the model and the learning problem. In section 3, we discuss consistent prediction methods, which we construct on the basis of the principle of empirical risk minimisation. Finally, in section 4, we apply our approach to a variety of stochastic choice models in economics.

4

Prediction and Stochastic Choice, Pathikrit Basu

## 2 Model

Let  $X \subseteq \mathbb{R}^K$  be a compact set of characteristics and  $A = \{a_1, a_2, \dots, a_m\}$  denote a finite set indexing finitely many alternatives. In certain contexts,  $x = (x_1, x_A) \in \mathbb{R}^K$  shall be a vector including both individual characteristics ( $x_1$ ) and product characteristics ( $x_A = (x_{a_1}, \dots, x_{a_m})$ ). In other contexts, involving choice from menus,  $X \subseteq \mathcal{A} = \mathcal{F}(A)$ . Hence, at each  $x \in X$  which is a menu, a choice of an alternative is made,  $a \in x$ . We define  $Z = X \times A$  and will denote a typical element of  $Z$  as  $z$ . A stochastic choice function is a map  $\sigma : X \rightarrow \Delta(A)$ . The interpretation of  $\sigma$  is that at each  $x$ ,

$\mu_a(x) = \int_{a \in A} \mu_a(x) g_{a \in A}$  is a probability vector where  $\mu_a(x)$  denotes the probability that alternative  $a \in A$  will be chosen when the underlying characteristics are given by  $x$ .

## 2.1 Data Generating Process

The analyst has access to a finite data set of choices from a population. Each data point consists of an individual's characteristic and the alternative chosen i.e.  $(x_i; a_i) \in Z$ . Hence, a data set is a finite sequence

$$z_n = (x_1; a_1); (x_2; a_2); \dots; (x_n; a_n) \in Z^n \quad (2.1)$$

Hence, a data set is an element  $z_n = (z_1; \dots; z_n)$  of  $Z^n$ .

We now describe the data generating process. There is a probability measure  $\mu \in \mathcal{P}(X)$  which defines the distribution of characteristics in the population. Given  $\mu$  and a stochastic choice map  $\sigma$ , the data is generated as follows. Independent across  $i$ ,  $x_i$  is drawn according to  $\mu$  and then  $a_i$  is drawn according to the choice probabilities given by  $\sigma(x_i)$ . Note here that the analyst only observes the characteristic and the alternative chosen i.e.  $(x_i; a_i)$  through the data in 2.1. The analyst knows neither the distribution  $\mu$  nor the stochastic choice function  $\sigma$ , but assumes that it satisfies certain properties. We denote as  $\mu \otimes \sigma$ , the probability measure induced by  $\mu; \sigma$  together on the set  $X \times A$ . This represents the joint distribution from which the data is generated, by taking  $n$  i.i.d. samples from  $\mu \otimes \sigma$ . We shall denote as  $\mu_n \otimes \sigma_n$ , the  $n$ -fold product measure induced by  $\mu \otimes \sigma$  on  $Z^n$ . This is essentially the distribution of the data  $z_n$ .

## 2.2 Prediction

A model is any family of stochastic choice functions  $\sigma$ . The objective of the analyst is to learn the true choice probabilities based on the data and a model  $\mu$  represents the analyst's hypothesis. Formally, a prediction method is a map

$$\hat{\sigma} : \left[ \prod_{i=1}^n (X \times A) \right] \rightarrow \mathcal{P}(A) \quad (2.2)$$

Suppose now that the true choice probabilities are given by  $\sigma_0$  and suppose  $\mu_0$  governs the distribution over characteristics. We shall require that a learning map  $\hat{\sigma}$  be so that with enough data,  $z_n = (x_i; a_i)_{i=1}^n$ , the estimate of the prediction method  $\hat{\sigma}(z_n)$  would be close to the true choice probabilities  $\sigma_0$ . Here, our notion of closeness between two choice probability maps  $\sigma; \sigma_0$  will be given by

$$d_{\mu_0}(\sigma; \sigma_0) = \int_X \sum_{a \in A} d(\sigma(x); \sigma_0(x)) d_{\mu_0}(x); \quad (2.3)$$

Throughout the paper, for any metric space  $Y$ , we will denote as  $\mathcal{P}(Y)$ , the set of all Borel probability measures on  $Y$ . For any  $\mu \in \mathcal{P}(Y)$ , we shall denote as  $\mu_n$ , the  $n$ -fold product measure on  $Y^n$  defined by  $\mu_n := \prod_{i=1}^n \mu$ .

5

Prediction and Stochastic Choice, Pathikrit Basu

where  $d : \mathcal{P}(A) \times \mathcal{P}(A) \rightarrow \mathbb{R}$  denotes a divergence function or metric on the space of all choice probability vectors on  $A$  i.e.  $\mathcal{P}(A)$ . For example,  $d$  could be squared Euclidean distance, KL divergence or the total variation distance. This leads us to the following definition of consistency for prediction methods.

**Definition 2.1.** A prediction method  $\hat{\sigma}$  is consistent (with respect to  $d$  and  $\mu_0$ ) if for all  $0 < \epsilon < 1$ , there exists  $N(\epsilon; \mu_0)$  such that for all  $n \geq N(\epsilon; \mu_0)$ ,  $(\int_{Z^n} d_{\mu_0}(\hat{\sigma}(z_n); \sigma_0) d_{\mu_0^n}(z_n)) < \epsilon$

$$\int_{Z^n} d_{\mu_0}(\hat{\sigma}(z_n); \sigma_0) d_{\mu_0^n}(z_n) < \epsilon$$

$$\int_{Z^n} d_{\mu_0}(\hat{\sigma}(z_n); \sigma_0) d_{\mu_0^n}(z_n) < \epsilon$$

$$\int_{Z^n} d_{\mu_0}(\hat{\sigma}(z_n); \sigma_0) d_{\mu_0^n}(z_n) < \epsilon$$

$$\int_{Z^n} d_{\mu_0}(\hat{\sigma}(z_n); \sigma_0) d_{\mu_0^n}(z_n) < \epsilon$$

We say that a model  $\mu$  is predictable with respect to  $d$  if there exists a prediction method  $\hat{\sigma}$ , which is consistent with respect to  $d$  and  $\mu$ . Finally, for a given  $\epsilon > 0$ , we denote as  $N(\epsilon; \mu)$ , the smallest  $n$  for which 2.4 holds. The function  $N : (0; 1) \rightarrow \mathbb{N}$  is called the sample complexity of  $\mu$  (with respect to  $\hat{\sigma}$ ).

The above definition of predictability is based on PAC predictability (see, for example, Valiant

[1984], Yamanishi [1992]). In what follows, we shall discuss consistent prediction methods for various stochastic choice models.

### 3 Consistent prediction methods

#### 3.1 Empirical Risk Minimization

We shall construct consistent prediction methods based on the principle of empirical risk minimization. For a detailed treatment, see Vapnik [1998]. Much of this section contains standard ideas from machine learning. However, one should acknowledge Proposition 3.1, which provides a novel way to solve the problem of PAC learning stochastic functions w.r.t divergences via the use of scoring rules (see Yamanishi [1992], Kearns and Schapire [1994], Abe et al. [1991] Abe et al. [2001]).

For the learning problem, we first define a loss function  $V : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ . Suppose the true distribution and choice probabilities are given by  $\mu, \nu$ . Then, the expected risk corresponding to  $\mu, \nu$  is defined as

$$R(\mu, \nu; V) = \int_{\mathcal{X} \times \mathcal{A}} V(x; a) d_{\mu, \nu}(x; a) \quad (3.1)$$

A minimizer of expected risk  $R(\mu, \nu; V)$  is defined as

$$R(\mu, \nu; V) = \arg \min_{\mu, \nu} R(\mu, \nu; V) \quad (3.2)$$

We shall consider loss functions  $V$  for which it will hold that  $\mu = \nu$  i.e the true choice probabilities would minimise the expected risk (a property also known as Fisher consistency). In the next section, we will introduce risk in a specific manner via the application of incentive compatible scoring rules.

The principle of empirical risk minimization involves estimating the expected risk by the empirical risk on the sample  $Z_n = \{(x_i, a_i)\}_{i=1}^n$

For each  $\mu, \nu$ , the empirical risk is given by

$$\hat{V}(\mu, \nu; Z_n) = \frac{1}{n} \sum_{i=1}^n V(x_i; a_i) \quad (3.3)$$

The prediction method that corresponds to empirical risk minimization, denoted by  $\hat{\mu}_\epsilon$ , is defined as

$$\hat{\mu}_\epsilon(Z_n) = \arg \min_{\mu} \hat{V}(\mu, \nu; Z_n) \quad (3.4)$$

6 Prediction and Stochastic Choice, Pathikrit Basu

The minimum in 3.4 need not always exist. However, if it holds for some  $M$  that  $V(x; a) \leq M$  for all  $\mu, \nu$  and  $(x; a) \in \mathcal{X} \times \mathcal{A}$ , then the infimum exists and we can define an almost-ERM prediction method as follows. We have  $f_n$  such that  $f_n > 0$  and  $\lim_{n \rightarrow \infty} f_n = 0$ . The almost-ERM rule selects

$$\hat{\mu}_\epsilon(Z_n) \text{ such that } \hat{V}(\hat{\mu}_\epsilon(Z_n), \nu) \leq \inf_{\mu} \hat{V}(\mu, \nu; Z_n) + f_n \quad (3.5)$$

In this context, consistency relies on  $\hat{\mu}_\epsilon(Z_n)$  being close to the true choice probability function  $\mu$  as the function  $\hat{V}$  approximates  $R$ , for large enough  $n$ . Consistency here is defined as follows, in terms of  $V$ .

**Definition 3.1.** A prediction method  $\hat{\mu}_\epsilon$  is said to be consistent (with respect to  $V$  and  $\mu, \nu$ ) if for all  $0 < \epsilon; \delta < 1$ , there exist an  $N(\epsilon; \delta)$  such that for all  $n \geq N(\epsilon; \delta)$ , it holds that

$$\frac{1}{n} \sum_{i=1}^n V(x_i; a_i) < \inf_{\mu} \int_{\mathcal{X} \times \mathcal{A}} V(x; a) d_{\mu, \nu}(x; a) + \epsilon \quad (3.6)$$

We say that a model  $\mu$  is predictable with respect to  $V$  if there exists a prediction method  $\hat{\mu}_\epsilon$ , which is consistent with respect to  $V$  and  $\mu, \nu$ .

It turns out that a model  $\mathcal{F}$  is predictable if the family of real-valued functions  $\mathcal{V}_{\mathcal{F}} = \{fV_{\mathcal{F}}(\cdot; \cdot; \cdot) : \mathcal{F} \in \mathcal{F}\}$  is a uniform Glivenko-Cantelli class of functions (see, for example, Vapnik [1998]; Shalev-Shwartz and Ben-David [2014]). Each  $f \in \mathcal{V}_{\mathcal{F}}$  is a function of the form  $f : X \rightarrow \mathbb{R}$ . We provide a definition below (see also Dudley [2014]).

**Definition 3.2.** A class of real valued functions  $\mathcal{F}$  on  $Z$  is said to be a uniform Glivenko-Cantelli class of functions if for all  $\epsilon > 0$ , there exists  $N(\epsilon)$  such that for all  $n \geq N(\epsilon)$ , for all  $\mathcal{Z}_n \subset Z$ ,

$$\sup_{\mathcal{Z}_n} \sum_{j=1}^n |f_j(z) - \mathbb{E}_n(f_j)| < \epsilon$$

$$\epsilon \leq \frac{1}{n} \sum_{j=1}^n \mathbb{E} f_j^2 \quad (3.7)$$

where  $\mathbb{E}_n(f)$  denotes the expectation of  $f$  under  $\mathbb{P}_n$ .

A necessary and sufficient condition for a class of real-valued functions to be a uniform Glivenko-Cantelli class of functions is that it have finite VC-dimension for each  $\epsilon > 0$  (see, for example, Alon et al. [1997]). Other combinatorial measures and notions of dimensions and capacity also guarantee the uniform Glivenko-Cantelli property for a class of functions for. eg. VC-dimension, Pollard's pseudodimension, P-dimension, Covering numbers and Rademacher Complexity. We shall define and apply these notions and their study implications for sample complexity bounds in the subsequent sections.

### 3.1.1 Scoring Rules

In this section, we will define loss functions based on scoring rules. A scoring rule is a function  $S : \mathcal{A} \rightarrow \mathbb{R}_+$  (see, for example, Savage [1971], Selten [1998]). The interpretation of  $S$  is that it is a mechanism for eliciting subjective probability judgements. If  $p$  is a probabilistic prediction about the alternative to be chosen in  $A$  and suppose  $a$  is the alternative chosen, then  $S_a(p)$  is the reward obtained. For each  $p; q \in \mathcal{P}(A)$ , we can define  $S(p; q) := \int_A S_a(p) q_a$  to be the expected score

when  $q$  is the true distribution over  $A$ . A scoring rule is said to be incentive compatible if

$$S(q; q) \geq S(p; q) \text{ for all } p; q \in \mathcal{P}(A) \quad (3.8)$$

i.e  $p = q$  maximises the function  $S(\cdot; q)$ . We say that  $S$  is strongly incentive compatible if additionally,  $p = q$  is the unique maximizer for each  $q \in \mathcal{P}(A)$ .

We use the terminology of Selten [1998]. Incentive compatible (strongly) scoring rules are also referred to as proper (strictly) scoring rules. See, for example, Gneiting and Raftery [2007]

7

Prediction and Stochastic Choice, Pathikrit Basu

We now define a loss function based on  $S$ .

$$V_s(\cdot; x; a) := \int_A S_a(\cdot) d_{x,a} \quad (3.9)$$

We can now prove the following lemma.

**Lemma 3.1.** Let  $\mathbb{P}_0 \in \mathcal{P}(X)$  be the true distribution over characteristics and let  $\mathbb{P}_0 \in \mathcal{P}$  be the true choice probability function. Suppose  $S$  is incentive compatible. Then,  $\mathbb{P}_0$  minimises  $\int V_s(\cdot)$ . Furthermore, if  $S$  is strongly incentive compatible and  $\mathbb{P}_0$  minimizes  $\int V_s(\cdot)$ , then  $\mathbb{P}_0(x) = \mathbb{P}_0(x)$  with probability one according to  $\mathbb{P}_0$ .

Proof. For any  $\mathbb{P}_0 \in \mathcal{P}$ , we have

$$\begin{aligned} \int V_s(\cdot) d_{\mathbb{P}_0} &= \int_{X \times A} V_s(\cdot; x; a) d_{\mathbb{P}_0}(x; a) \\ &= \int_{X \times A} S_a(\cdot) d_{\mathbb{P}_0}(x; a) \\ &= \int_{X \times A} S(\cdot; x; a) d_{\mathbb{P}_0}(x; a) \end{aligned}$$

□

Z

$$S(\rho(x); \rho_0(x))d_{\rho}(x) \quad (3.10)$$

$$= \int_X V_s(\rho_0);$$

where the inequality 3.10 follows from the fact that  $S$  is an incentive compatible scoring rule. Now, suppose  $\rho_0$  minimizes  $\int V_s(\cdot)$  and let  $E = \{x : \rho(x) \neq \rho_0(x)\}$ . Since  $S$  is strongly incentive compatible, note that  $S(\rho_0(x); \rho_0(x)) > S(\rho(x); \rho_0(x))$  for all  $x \in E$ . We already have that  $S(\rho_0(x); \rho_0(x)) \leq S(\rho(x); \rho_0(x))$  for all  $x \in X$ . Hence, if  $\rho_0(E) > 0$ , then we will have that  $\int V_s(\rho_0) < \int V_s(\rho)$ . This contradicts that  $\rho_0$  minimizes  $\int V_s(\cdot)$ .

The above result shows that  $\rho_0$  is the minimiser of expected risk. Each incentive compatible scoring rule leads to a divergence function

$$d_s(p; q) = S(q; q) - S(p; q) \quad \text{for all } p; q \in \mathcal{A} \quad (3.11)$$

Hence, from Lemma 1, it follows that

$$\int V_s(\rho) \leq \inf_{\rho \in \mathcal{P}_M} \int V_s(\rho)$$

$$\int V_s(\rho) = \int V_s(\rho) \leq \int V_s(\rho_0)$$

Z

$$d_s(\rho(x); \rho_0(x))d_{\rho}(x):$$

Hence, if we can construct consistent prediction methods with respect to  $\int V_s$ , then we can construct consistent prediction methods with respect to the divergence function  $d_s$ .

The prediction method based on empirical risk minimisation corresponds to maximising the empirical score based on the data  $D = \{(x_i; a_i)\}_{i=1}^n$

$\hat{\rho}_n$ . Hence,  $\hat{\rho}_n$  maximises

$$\int \rho(x) V_s(\rho(x)) =$$

1

n

$X_n$

$i=1$

$$S_{a_i}(\rho(x_i)); \quad (3.12)$$

which can be thought of as the empirical score. For a scoring rule  $S$  and stochastic choice function  $\rho$ , we define the function  $S_{\rho}(x; a) := S_a(\rho(x))$ . Consistency of the prediction method  $\hat{\rho}_n$  with respect to  $d_s$  relies on the nature of the real-valued function class

$$S_{\rho} = \{S_{\rho}(x; a) : \rho \in \mathcal{P}\}$$

The following is a key result.

8

Prediction and Stochastic Choice, Pathikrit Basu

**Proposition 3.1.** Let  $\rho$  be a model of stochastic choice and let  $S$  be an incentive compatible scoring rule. Suppose the class of functions  $S_{\rho}$  is bounded above i.e. there exists an  $M$  such that  $f(z) \leq M$  for all  $f \in S_{\rho}$  and  $z \in Z$ . If  $S_{\rho}$  is a uniform Glivenko-Cantelli class, then the model  $\rho$  is predictable with respect to the divergence function  $d_s$ .

Proof. The proof is in the appendix. It follows from the fact that the ERM rule yields consistency (see Vapnik [1998]) in the sense of Definition 3.1 and from Lemma 3.1.

We now give some examples of incentive compatible scoring rules and their associated divergence functions.

1. (Log Rule) :  $S(p; a) = \ln(p_a)$ . The log scoring rule is incentive compatible and its divergence function corresponds to KL divergence  $S(p; q) = d_{KL}(p||q) =$

P

$\frac{a}{2A}$

$\ln(p_a)$

$\ln(q_a) p_a$ . Note that

maximisation of the empirical score with respect to the log rule corresponds to choosing choice probability functions according to the conditional maximum likelihood procedure .

2. (Brier/Quadratic Scoring Rule) :  $S(p; a) = 2p_a -$

P

$b p^2$

b . The scoring rule is called the

Brier or quadratic scoring rule and is strongly incentive compatible. The divergence function associated with  $S$  is the square of the euclidean distance between  $p$  and  $q$  i.e.  $d_s(p; q) =$

$$\sum_j (p_j - q_j)^2$$

2.

Since all  $L_p$  metrics on  $R^d$  are equivalent (as all norms on a finite dimensional vector space

are equivalent), it follows that consistency of a prediction method with respect to  $d_s$  implies consistency with respect to any other  $L_p$  metric. Further, from Pinsker's inequality, we have that

$$d_{TV}(p, q) \leq \frac{1}{2} d_{KL}(p, q)$$

1

2

$d_{KL}(p, q)$ :

Recall the total variational norm is equal to half times the  $L_1$  distance i.e.  $d_{TV}(p, q) = \frac{1}{2} \|p - q\|_1$ .

Hence, the above inequality implies that empirical score maximisation using the log scoring rule, which corresponds to conditional maximum likelihood, leads to consistency with respect to all  $L_p$  metrics.

3. (Manski's score with tie breaking) : Let  $\prec$  be a complete strict order on  $A$  and let  $p \in \mathcal{P}(A)$ . Now, define another strict order  $\prec_p$  on  $A$  as follows :  $a \prec_p b$  if either  $p_a > p_b$  or it is the case that  $p_a = p_b$  and  $a \prec b$ . Finally, we let  $W(1) \prec W(2) \prec \dots \prec W(j) \prec \dots \prec W(|A|) \prec 1$ .

The Manski scoring rule with tie breaking is defined as  $S(p; a) = W(\text{rank}_{\prec_p}(a))$ . Note that  $a \prec_p b$  implies both  $p_a > p_b$  and  $S(p; a) > S(p; b)$ . Now, by the rearrangement inequality, this means that  $S(p; q) \geq S(p; a)$  for all  $p, q \in \mathcal{P}(A)$  i.e.  $S$  is an incentive compatible scoring rule.

4. (Subgradient of a convex function) : Suppose we have a convex function  $f : \mathcal{P}(A) \rightarrow \mathbb{R}$ . Let  $r_f(x) : \mathcal{P}(A) \rightarrow \mathbb{R}_+$  be a subgradient for the function  $f$ . We can now define a scoring rule based on  $f$  as follows.

$$S_a(p) = f(p) + r_f(p)(a - p).$$

$S$  is incentive compatible and its associated divergence function corresponds to the Bregman divergence of  $f$ . Moreover, any incentive compatible scoring rule corresponds to some convex function. For example,  $f(p) = \sum p_j^2$  leads to the Brier score.

9

Prediction and Stochastic Choice, Pathikrit Basu

### 3.2 Dimension and Sample Complexity

In this section, we discuss sufficient conditions for a class of real valued functions to be uniform Glivenko-Cantelli. These place bounds on the capacity or complexity of the function class. There exist several alternative measures to quantify such complexity and we will discuss the central notions here. In particular, we discuss two notions we will use extensively : Rademacher complexity and Pollard Dimension.

**Rademacher Complexity.** Consider a real-valued function class  $F$  defined on the set  $Z$ . For the sample  $z_n$ , the empirical Rademacher complexity is defined as

$$R_{z_n}(F) = E \left[ \sup_{f \in F} \sum_{i=1}^n \epsilon_i f(z_i) \right]$$

h

sup

$f \in F$

1

n

$X_n$

$i=1$

$\epsilon_i f(z_i)$

i

:

Here, each  $\epsilon_i$  is an independent bernoulli random variable which takes value 1 with 0.5 probability and  $-1$  with 0.5 probability. Rademacher complexity of the function class  $F$  (with respect to a distribution  $\mu \in \mathcal{P}(Z)$ ) is defined as

$$R_n(F) = E_n[R_{z_n}(F)]$$

Now, it turns out that the Rademacher complexity of a function class is closely related to its Glivenko-Cantelli properties. Suppose  $F$  is uniformly bounded i.e. there exists a real number  $M > 0$  such that  $|f(z)| \leq M$  for all  $z \in Z$  and  $f \in F$ . Then, we have that (see, for example, Shalev-Shwartz and Ben-David [2014]), for each  $\epsilon > 0$

$\frac{1}{n}$

$$z_n : \sup_{f \in F} \sum_{i=1}^n \epsilon_i f(z_i) \leq \epsilon$$

$f \in F$

$j=1$

$X_n$

$i=1$

$$f(z_i) - E_{f_j} < 2R_{z_n}(F) + 6M$$

$$\log \frac{4}{2^n} \leq \dots$$

Hence, if the empirical Rademacher complexity  $R_{z_n}(F)$  converges to zero (say irrespective of the sample  $z_n$ ), then we obtain the uniform Glivenko-Cantelli property for the function class  $F$ . For example, if one shows that  $R_{z_n}(F) \leq C n^{-p}$ , where  $C$  is an absolute constant, then we get that the above inequality is satisfied for sample size  $n$  above

$$\frac{1}{2C + 6M} \log \frac{4}{2^n} \leq \dots$$

This has implications for the sample complexity of stochastic choice models. Indeed, the proof of Proposition 3.1 shows that for a stochastic choice model, the sample complexity  $N(\epsilon, \gamma)$ , is at most  $\max\{3, N_0(\epsilon, \gamma)\}$ , where  $N_0(\epsilon, \gamma)$  is the threshold sample size that guarantees inequality 3.7 for values  $(\epsilon, \gamma)$  in the Glivenko-Cantelli definition. We next discuss Pollard dimension, which is related to Rademacher complexity.

**Pollard Dimension.** We say that a real-valued function class  $F$  shatters a set of points  $(z_1, \dots, z_n)$  if there exists a vector  $(t_1, \dots, t_n) \in \mathbb{R}^n$  such that for all  $B \subseteq \{1, \dots, n\}$ , there exists  $f_B \in F$  such that  $f_B(z_i) > t_i$  for all  $i \in B$  and  $f_B(z_i) \leq t_i$  for all  $i \notin B$ .

The Pollard dimension of  $F$ , denoted as  $P(F)$ , is defined as  $P(F) = \max\{n \mid \exists (z_1, \dots, z_n) \text{ which can be shattered by } F\}$ .

This complexity measure was introduced by Pollard [1984]. It is a generalisation of the VC dimension (see Vapnik and Chervonenkis [1968], Vapnik [1998]), which is defined for sets. Indeed, for a class

Prediction and Stochastic Choice, Pathikrit Basu  
of 0-1 functions, the Pollard and VC dimensions coincide. Hence, for a collection of sets  $P$ , we have that  $VC(P) = P(P)$ . The more standard (equivalent) definition would involve the notion of shattering for sets. A class of sets  $P$  shatters a set of points  $(z_1, \dots, z_n)$  if for each  $B \subseteq \{1, \dots, n\}$ , there exists  $P \in \mathcal{P}$  such that  $z_i \in P$  if and only if  $i \in B$ . Then, the VC dimension of  $P$  is defined as  $VC(P) = \max\{n \mid \exists (z_i)_n \text{ which can be shattered by } P\}$ .

For the function class  $F$  with Pollard dimension  $d$ , the empirical Rademacher complexity is  $R_{z_n}(F) \leq C \frac{d}{n}$ .

where  $C > 0$  is an absolute constant (see Pollard [1984]). Hence, an immediate implication of this is that finite Pollard dimension implies the uniform Glivenko-Cantelli property. In view of Proposition 3.1, the sample complexity of a stochastic choice model  $\sigma$ , to be learned with respect to  $d_s$ , is upper bounded by a linear function of the Pollard dimension of  $S_{\sigma}$ .

## 4 Stochastic Choice Models

In this section, we will consider a variety of stochastic choice models. In particular, we are interested in the recoverability of heterogeneity. By heterogeneity, we mean the distribution over preferences or parameters of the model, that leads to the stochasticity in choice. The main result of the paper is that when we consider stochastic preferences, we can recover preference heterogeneity if the support of the distribution has finite VC dimension. We discuss several familiar examples for which the result applies. We also establish a converse result on the non-predictability of such models. Lastly, we also consider choice with consideration sets and the recoverability of cognitive heterogeneity.

### 4.1 Stochastic Preferences

Let  $X \subseteq \mathbb{R}^k$  be a set of alternatives. Suppose  $\mathcal{P}$  is a class of continuous preference relations (complete and transitive) on  $X$ . We are interested in defining stochastic choice functions where preferences are random in  $\mathcal{P}$ . Formally, we have a probability measure  $\mu$  on  $\mathcal{P}$  i.e.  $\mu \in \mathcal{P}(\mathcal{P})$ . In this setting, the characteristic vector shall be a vector of alternatives. Hence,  $X \subseteq \mathbb{R}^k$ , as in earlier. We shall assume that for each  $x = (x_a)_{a \in A} \in X$  we have  $x_a \geq x_b$  for all  $a \geq b$ . Further,

$$\int_{\mathcal{P}} \mathbb{1}_{\{x_a \geq x_b \text{ for all } b \geq a\}} d\mu = 1 \quad (4.1)$$

These assumptions ensure that in the random choice process, there are no ties with probability one. This is natural and simplifies the analysis. We now define the stochastic choice function corresponding to  $\mu$  as

$$c_a(x; \mu) = \int_{\mathcal{P}} \mathbb{1}_{\{x_a \geq x_b \text{ for all } b \geq a\}} d\mu;$$

which is the probability of choosing a given characteristics  $x$  and the distribution  $\mu$  over preferences.

The stochastic choice model corresponding to a random preference model on the class of preferences  $\mathcal{P}$  is defined as follows.

$$c_{\mathcal{P}} = \int_{\mathcal{P}} c_a(x; \mu) d\mu \quad (4.1) \text{ for all } x \in X;$$

We show the following result.

The set  $\mathcal{P}(\mathcal{P})$  denotes the set of all Borel probability measures on  $\mathcal{P}$ , which has the subspace topology generated by the closed convergence topology on the closed subsets of  $X \times X$ . See also the appendix.

11

Prediction and Stochastic Choice, Pathikrit Basu

**Proposition 4.1.** Suppose  $\mathcal{P}$  is a class of continuous preferences with finite VC dimension  $d$ . Then,  $c_{\mathcal{P}}$  is predictable with respect to euclidean distance. Moreover, for the Quadratic scoring rule  $S_{br}$ , the Rademacher complexity of  $S_{br} \circ c_{\mathcal{P}}$  is at most  $O(\sqrt{\frac{1}{n} \sum_{j=1}^n \log(jA) d \log(d)})$ .

$\sqrt{\frac{1}{n} \sum_{j=1}^n \log(jA) d \log(d)}$ .

**Proof.** Firstly, we shall define a class of deterministic choice rules based on  $\mathcal{P}$ . Then, we will show that  $c_{\mathcal{P}}$  is derived by taking continuous convex combinations of the deterministic class. We will then bound the Rademacher complexity corresponding to  $S_{br} \circ c_{\mathcal{P}}$ .

For each  $\mu \in \mathcal{P}(\mathcal{P})$  and  $a \in A$ , define the 0-1 valued function

$$c_a(x; \mu) = \begin{cases} 1 & \text{if } x_a \geq x_b \text{ for all } b \geq a \\ 0 & \text{otherwise} \end{cases}$$

By assumption (4.1), we get that for all  $\mu$ , and  $a \in A$ ,

$$\int_{\mathcal{P}} c_a(x; \mu) d\mu = \int_{\mathcal{P}} c_a(x; \mu) d\mu;$$

We will now show that for each  $a$ , the class of functions  $S_{\mathcal{P}}^a$

$S_{\mathcal{P}}^a = \int_{\mathcal{P}} c_a(x; \mu) d\mu$  also has finite VC dimension.

Consider a finite set of  $n$  points  $D = \{x_1, \dots, x_n\} \subseteq X$ . Suppose that  $D$  can be shattered by  $S_{\mathcal{P}}^a$ .

Then, all  $2^n$  possible 0-1 valued labellings of  $D$  can be generated by  $\mathcal{P}$ . Now, consider the following set of points in  $X \times X$ .

$$D_0 = \{(x_i, x_j) \mid i < j\}$$

Note that  $D_0$  is a set of  $n(n-1)/2$  points. Let  $d$  be the VC dimension of  $\mathcal{P}$ . On the one hand, from Sauer's Lemma (see Kearns and Vazirani [1994]), we have that at most  $\sum_{i=0}^d \binom{n(n-1)/2}{i}$  many labellings

in  $D_0$  can be generated by  $\mathcal{P}$ . On the other hand, we have that since  $S_{\mathcal{P}}^a$  shatters  $D$ , at least  $2^n$

labellings are obtained by  $\mathcal{P}$  in  $D_0$ . For large  $n$ , we obtain a contradiction as

$$\sum_{i=0}^d \binom{n(n-1)/2}{i} < 2^n \quad (4.2)$$

Hence, we establish that  $S_{\mathcal{P}}^a$

has finite VC dimension, say  $d_a$ . Now, from a well known result (see

Pollard [1984]), it follows that there is an absolute constant  $C_a$  such that the Rademacher complexity of  $S_P$

$a$ , say  $R_{X_n}(S_P)$ , is at most  $C_a$

$d_a = n$  for all  $X_n$ . Now, consider the probability of  $a$  in the stochastic choice model i.e. the function class  $\mathcal{F}_P$

$a = f_a(\cdot; \cdot)_j$  satisfies (4.1)g. Since,  $\mathcal{F}_P$

$a$  is essentially

derived by taking convex combinations (including continuous combinations) from  $S_P$

$a$  (see Lemma

6.3 in the appendix), we obtain that the Rademacher complexity of  $\mathcal{F}_P$

$a$  is at most equal to that of  $S_P$

$a$ .

Hence, we have  $R_{X_n}(\mathcal{F}_P)$

$a) \leq C_a$

$P$

$d_a = n$ .

We now complete the proof by applying the Quadratic scoring rule and bounding the Rademacher complexity of the class  $F = S_{br} \mathcal{F}_P$ . Firstly, consider the function class  $F_a = \{f(\cdot; a) \mid f \in S_{br} \mathcal{F}_P\}$ .

Note that

$f(x; a) = 2_a(x; \cdot) \square$

$X$

$b \in \mathcal{A}$

$\mathcal{F}_b(x; \cdot)^2$

for some  $\cdot$ . Since Rademacher complexity is additive (applied to  $f_P$

$a$   $g_a$ ) and has nice Lipschitz

composition properties (in this case  $y_2$ ), we get that  $R_{Z_n}(F_a) \leq 2C_a$

$P$

$d_a = n +$

$P$

$b \leq 2C_b$

$P$

$d_b = n$  (see

Bartlett and Mendelson [2002] and also Shalev-Shwartz and Ben-David [2014]). Define  $C_0 = \max_{a \in \mathcal{A}} 2C_a$

$P$

$d_a +$

12

Prediction and Stochastic Choice, Pathikrit Basu

$P$

$b \leq 2C_b$

$P$

$d_b$ . Now, for a given  $Z_n = (X_i; a_i)_{i=1}^n$

$i=1, \dots, n$ , define  $X_{n_a} = (X_i)_{i: |j| a_i = a}$  and  $n_a = \sum_j |j| a_i = a$ .

Using Lemma 6.4 in the appendix, we get

$R_{Z_n}(F) \leq$

$X$

$a \in \mathcal{A}$

$n_a$

$n$

$R_{X_{n_a}}(F_a)$

$\leq$

$X$

$a \in \mathcal{A}$

$n_a$

$n$

$(C_0 =$

$P$

$n_a)$

$\leq$

$X$

$a \in \mathcal{A}$

$r$

$n_a$

$$\frac{1}{n} \sum_{i=1}^n \left( C_0 + \frac{p}{n} \sum_{j=1}^n |A_j| \right) \log \left( \frac{1}{n} \sum_{j=1}^n |A_j| + 1 \right) \quad (4.3)$$

The last inequality follows from the fact that  $\log$  is concave. Hence, we establish that  $\mathcal{P}$  is predictable with respect to any  $L_p$  metric.

Finally to obtain the sample complexity bound, consider the inequality 4.2. Based on the inequality, we can in fact say that the following holds.

$$\frac{1}{n} \sum_{i=1}^n \log \left( \frac{1}{n} \sum_{j=1}^n |A_j| + 1 \right) \leq \frac{1}{n} \sum_{i=1}^n \log \left( \frac{1}{n} \sum_{j=1}^n |A_j| \right) + 2 \log \left( \frac{1}{n} \sum_{j=1}^n |A_j| + 1 \right) \quad (4.4)$$

The above follows from a fact of arithmetic that  $x \geq a \log(2a) + 2b$  implies  $x \geq a \log(x) + b$  (see, for example, Shalev-Shwartz and Ben-David [2014]) whenever  $a > 1$  and  $b > 0$ . Hence, the derived inequality in 4.4, together with the conclusion in 4.3 implies that the Rademacher complexity of  $F$  is at most  $O$

$$\frac{1}{n} \sum_{i=1}^n |A_j| \log \left( \frac{1}{n} \sum_{j=1}^n |A_j| \right) \log(d).$$

The above result also allows us to recover the distribution over preferences. For this, we will define the following notion of distance between distributions.

$$d(\mu; \mu_0) := \int_{\mathcal{X}} d(\mu(x); \mu_0(x)) d\mu_0(x)$$

Then, the following result holds. (Recovering heterogeneity) Suppose  $\mathcal{P}$  is a class of continuous preferences with finite VC dimension. Then, there exists an estimator  $\hat{\mu}$  (based on a prediction method  $\hat{\mu}$  derived from almost-ERM with the Quadratic scoring rule) of the true distribution  $\mu_0$  over preferences such that

$d(\hat{\mu}; \mu_0) \rightarrow_p 0$ ; uniformly over  $\mu_0$  and  $\mu$ . Moreover, the rate of convergence is given by the sample complexity bound derived from Theorem 4.1.

In this setting, the class of preferences  $\mathcal{P}$  would depend on the context. Stochastic preferences supported on linear utility would, for example, lead to the random coefficients model. Several canonical utility based models such as Cobb-Douglas or CES preferences may be considered. If preferences are over risky prospects, then we may have stochastic preferences over expected utility models and models exhibiting ambiguity aversion. We would bound the VC dimension of the preference class (see also Basu and Echenique [2020]) and then apply the above result to obtain predictability of the corresponding stochastic choice model supported on the preferences.

**Applications** We discuss some applications of the above result below.

13

Prediction and Stochastic Choice, Pathikrit Basu

1. (Linear Preferences) Consider the model of preferences  $\mathcal{P}_{LIN}$  defined as follows. For each

$\%2 P_{LIN}$  there exists a vector  $v \in \mathbb{R}^k$  such that

$x \succ y$  if and only if  $v \cdot x > v \cdot y$ :

Such preferences satisfy the well-studied independence axiom from decision theory. This means that for all  $x, y, z \in X$  and  $\alpha \in [0, 1]$ ,

$x \succ y$  if and only if  $\alpha x + (1 - \alpha)z \succ \alpha y + (1 - \alpha)z$ :

From a result in Basu and Echenique [2020], it follows that the VC dimension of the set of preferences that satisfy the independence axiom is at most  $k + 1$ . Hence, it follows that  $VC(P_{LIN}) \leq k + 1$ . From Theorem 4.1, it follows that the corresponding stochastic choice model, in which the vector of coefficients is random, would be predictable.

Consider the following widely used econometric specification. The utility from an alternative  $a$  is

$$u(x; \alpha) = \alpha \cdot x_a + \epsilon_a;$$

where  $\alpha \in \mathbb{R}^k$  and  $\epsilon = (\epsilon_a)_{a \in A} \in \mathbb{R}^A$  are vectors that are randomly drawn according to the probability measure  $\mu$ . For example, for the Logit model,  $\alpha$  is deterministically fixed and each  $\epsilon_a$  independently follows a logistic distribution (similarly for Probit). For Mixed Logit, we would have that  $\alpha$  would follow an unknown distribution for eg. multivariate normal. Here we have  $v = (\alpha; \epsilon)$  and one may rewrite alternatives as  $x_0$

$$a = (x_a; e_a), \text{ where } e_a \text{ is the unit vector in } \mathbb{R}^A$$

where the value corresponding to  $a$  is one. This also ensures that  $x_0$

$$a \succ b \iff x_0 \cdot a > x_0 \cdot b$$

for all  $a \neq b$ . Note

that our formulation here essentially pertains to the case where  $x$  is independent of  $\alpha$ .

2. (Random Utility) Suppose now that preferences are defined by an underlying utility function.

Hence, suppose we have a class of utility functions  $U = \{u_j : X \rightarrow \mathbb{R}\}$ . For an given utility

function  $u \in U$ , the corresponding preference relation  $\succ_u$  is defined as

$x \succ_u y$  if and only if  $u(x) > u(y)$ :

The model of preferences generated by  $U$  is then defined as  $P_U := \{ \succ_u : u \in U \}$ . It turns out that if the class of utility functions has finite Pollard dimension, say  $d$ , then the VC dimension of  $P_U$  is also finite, and is of the order  $O(d)$ . It can be shown that utility classes corresponding to canonical models such as Cobb-Douglas and CES preferences have finite Pollard dimension.

3. (Decisions under Risk/Uncertainty) Suppose that the set of alternatives  $X$  corresponds to a set

of monetary acts over a finite state space (uncertainty) i.e.  $X \subseteq \mathbb{R}^Z$  or the lotteries over a finite

set of consequences  $Z$  i.e.  $X \subseteq \Delta(Z)$  (risk). Such contexts have been studied extensively

in decision theory (see, for example, Kreps [1988] or Gilboa [2009]). Several studies have

investigated different models of preferences based on different notions of expected utility. As was shown in Basu and Echenique [2020], the standard expected utility and choquet expected utility models had finite VC dimension, where the latter has VC dimension exponential in  $|Z|$ .

In the context of risk, one may show that expected utility and Quadratic expected utility (see

Chew et al. [1991]) preferences have finite VC dimension. Indeed, the study of Gul and

Pesendorfer [2006] consider random preferences with expected utility. The related model of

random choice with private information due to Lu [2016] can be studied within our framework of linear preferences as well.

Finally, one should note that the nature of axioms of the model and VC dimension bounds are

closely related. As was noted earlier, in Basu and Echenique [2020], we show the preferences

satisfying the independence axiom have finite VC dimension, linear in  $|Z|$ . In that paper, it

14

Prediction and Stochastic Choice, Pathikrit Basu

was also argued that if preferences satisfy comonotonic independence and admit existence of certainty equivalents, then the VC dimension would be finite and upper bounded by a function exponential in  $|Z|$ . Hence, for Choquet Expected Utility, denoted by  $P_{CEU}$ , the bounds

$$\frac{1}{|Z|} \leq VC(P_{CEU}) \leq |Z|$$

follow.

4. (Euclidean Preferences) Consider the context of spatial voting (see Downs et al. [1957]) where

$X \subseteq \mathbb{R}^k$  is an ideological space and each candidate in an election is represented by a point in

$X$  (alternative). In this voting setting, preferences are defined through an ideal point  $x_*$  such

that candidates closer to the ideal point are preferred by the voter. Hence, we have that

$x \succ_{x_*} y$  if and only if  $\|x - x_*\| < \|y - x_*\|$ :

Hence  $P_{EUC} = \{ \succ_{x_*} : x_* \in X \}$ . In this case, the probability measure  $\mu$  essentially gives us a distribution over ideal points, which is also referred to as voter heterogeneity. Since euclidean preferences can be defined by finitely many arithmetic computations (see Proposition 6.2 in the appendix), it follows that the VC dimension is finite and of the order  $O(k)$ . Our estimation

results provide us a way to recover voter heterogeneity in these settings.

## 4.2 Non-predictability

We now establish a result on the non-predictability of  $\mathcal{P}$ . With some additional conditions, we may prove that a result converse to Theorem 4.1 follows. This requires defining a notion of shattering based on strict preferences. We say that a model of preferences  $P$  strictly shatters a set of pairs of alternatives  $((x_i; y_i))_{i=1}^n$  if for any labelling  $(a_i)_{i=1}^n$  of the pairs, there exists a preference relation  $\succsim_P$  such that

$x_i \succsim_P y_i$  whenever  $a_i = 1$  and  $y_i \succsim_P x_i$  whenever  $a_i = 0$

We now consider the following modified definition of VC dimension, which we will call  $VC_+(P)$ . It is defined as follows.

$VC_+(P) = \max_{((x_i; y_i))_{i=1}^n}$  which can be strictly shattered by  $P$ :

We will assume here for simplicity that  $X = \mathbb{R}^k$  and that  $x = (x_1, \dots, x_k)$  for all  $x \in X$ . A preference relation  $\succsim$  is said to be monotonic if  $x \succ y$  implies that  $x \succ z$ . We now state and prove the result.

**Proposition 4.2.** Let  $P$  be a class of monotone continuous preferences such that  $VC_+(P) = +1$ . Furthermore, suppose that there exists a probability measure satisfying 4.1 that has full support in  $P$ . Then, the following hold.

1. The stochastic choice model  $\mathcal{P}$  is not predictable with respect to any  $L_p$  metric.
2. Suppose we fix the euclidean distance on  $(A)$ . Then, for any prediction method  $\hat{\cdot}$  and sample size  $n$ , the worst case expected distance (error) between true and estimated choice probabilities is at least  $\frac{1}{2}$

2, which is half of the maximum possible distance. For  $|A| = 2$ , this lower bound for error is above what can be achieved by predicting that each alternative will be chosen with probability  $\frac{1}{2}$ .

**Proof.** The proof applies the probabilistic method where we randomly choose a stochastic choice function from  $\mathcal{P}$  and show that expected error is bounded away from zero. Then, we argue that there exists some function in  $\mathcal{P}$  that generates error above a threshold level.

This assumption is not crucial for the result.

15

Prediction and Stochastic Choice, Pathikrit Basu

Fix any prediction method  $\hat{\cdot}$  and suppose that  $0 < \epsilon < \frac{1}{2}$

$\frac{1}{2} \leq \epsilon < 1$  and  $0 < \epsilon < 1 - \epsilon$

$\frac{1}{2}$

2. Also,

let  $n \geq N$  be a sample size. Now, further fix  $k \geq 2$  such that  $k \geq 2$ . Since  $VC_+(P) = +1$ , there exists a strictly shattered set  $((x_i; y_i))_{i=1}^n$ .

Finally, let  $t \in (0, 1)$  such that  $t \geq 0.5$ . We will show then, that there exists a set of points  $(x_i)_{i=1}^n$

$\{x_i\}_{i=1}^n \subset X$  such that for all  $B \subset \{1, 2, \dots, n\}$ , there exists  $\{a_i\}_{i \in B} \in \mathcal{P}$  such that

$a_i(x_i) \geq t$  for all  $i \in B$  and

$a_i(x_i) \leq 1 - \epsilon$  for all  $i \in B$ ; (4.5)

where  $a \in \mathcal{A}$  is fixed beforehand. We now show how to construct the set of points  $(x_i)_{i=1}^n$ .

$\{x_i\}_{i=1}^n \subset X$ . Let

$b \in \mathcal{A}$ . We first define  $x_i$

$a := x_i$  and  $x_i$

$b := y_i$ . Further, for each  $i$  we let  $f_{x_i}$

$g_{c \in \mathcal{A}} : \mathbb{R}^k \rightarrow \mathbb{R}$

be such that all the vectors in  $f_{x_i}$

$g_{c \in \mathcal{A}}$  are distinct and also that  $x_i$

$a \succ x_i$

$c$  and  $x_i$

$b \succ x_i$

$c$  for all

c)  $2 \text{ Anfa; bg.}$  Due to monotonicity, this ensures that from the set of alternatives in  $x_i$  either  $x_i$

a)  $\text{or } x_i$

b)

is maximal.

Since  $P$  strictly shatters  $((x_i$

a);  $x_i$

b)) $_{k \times n}$

$i=1, \dots, k$ , for each  $B \subseteq \{1, 2, \dots, k\}$ , there exists a preference relation

$\succsim_B$  such that

$x_i$

a)  $\succsim_B x_i$

b) whenever  $i \in B$  and  $x_i$

$x_i$

b)  $\succsim_B x_i$

a) whenever  $i \in B$ :

Now, consider the set of preferences  $Q = \{f \mid x_i$

a)  $\succsim_B x_i$

b) for all  $i \in B$  and  $x_i$

b)  $\succsim_B x_i$

a) for all  $i \in B$ . It

is open in  $P$  under the closed convergence topology. Hence, it is measurable in the corresponding

Borel sigma algebra. Now, consider a probability measure  $\mu_2(P)$  that satisfies the regularity

condition in 4.1 and has full support in  $P$ . Then, since  $Q$  is open, we must have that  $\mu_2(Q) > 0$ .

We will now construct another probability measure  $\mu_B$  satisfying 4.1 such that its corresponding

stochastic choice function  $\mu_B$  satisfies 4.5. Firstly, let  $\epsilon > 0$  be such that  $\mu_2(Q_c) \leq \epsilon$ . Then, let

$\mu := (1 - \epsilon)\mu_2(Q_c)$

$\mu_2(Q)$ . This implies that  $\mu(Q) + \mu(Q_c) = 1$  and  $\mu(Q) \geq 1 - \epsilon$ . Now, define the probability

measure  $\mu_B(E) := \mu(Q \setminus E) + \mu(Q_c \setminus E)$ . Note that  $\mu_B(E) = 0$  if and only if  $\mu(E) = 0$ , which

implies that  $\mu_B$  also satisfies 4.1. By definition, it now follows that the corresponding stochastic choice

function  $\mu_B$  satisfies 4.5.

We now proceed with the proof of non-predictability. Since all  $L_p$  metrics are equivalent on  $\mathcal{A}$ ,

it suffices to show the claim for euclidean distance. Let  $\mu$  be the uniform distribution over the

set of points  $x_i \in \mathcal{X}_{k \times n}$

$i=1, \dots, k$ . Also, let  $\mu$  be the uniform distribution over all stochastic choice functions in

$\mathcal{F}_{\text{BG}, B, \{1, 2, \dots, k\}, n}$ . We will lower bound the expected error after  $n$  samples.

$E \sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

Now, from Fubini's theorem, we will change the order of expectations. The above expectation then

becomes

$E \sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

Since  $\mu$  is uniform, it follows that  $P_{x_i} \|x_i - \hat{x}_i\| = \int_{\mathcal{X}_{k \times n}} \|x_i - \hat{x}_i\| d\mu(x_i)$

$\mu$ . We will now simply lower bound the

conditional expectation  $E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

$\sum_{i=1}^k$

$\sum_{i=1}^k E_{x_i} \|x_i - \hat{x}_i\|$

$\sum_{i=1}^k$

. Consider some fixed  $x = x_n$ . Now, let us write the expectation  $E_{\mathcal{H}}$

$$\mathbb{E}_{\mathcal{H}} \sum_{i=1}^n \mathbb{1}_{\{z_i(x) \leq t\}}$$

more explicitly as

$$\mathbb{E}_{\mathcal{H}} \sum_{i=1}^n \mathbb{1}_{\{x_i \in C_n\}} \mathbb{1}_{\{x_i \leq t\}}$$

:

16 Prediction and Stochastic Choice, Pathikrit Basu

where  $\mu(x_n)$  is the joint distribution induced on  $\mathcal{F}_n$  by the product of the choice probabilities in  $(\mu(y))_{y \in \mathcal{F}_n}$ . Consider some fixed  $B$  for some  $B \subseteq \{1, 2, \dots, n\}$ . Now, let  $d(B; x_n) \in \mathcal{F}_n$  be such that for each  $j \in \{1, \dots, n\}$  we define  $d_j(B; x_n) = a$  if and only if  $x_j \in \{x_1, \dots, x_n\} = x_n$  is such that  $x_j = x_i$  for some  $i \in B$ . Now, recall that for  $i \in B$ , we have  $\mu(x_i) \geq t$  and for  $i \in \{1, 2, \dots, n\} \setminus B$ , we have  $\mu(x_i) < t$ .

Hence, by definition, it follows that  $\mu(x_n)$  draws the vector  $d(B; x_n)$  with probability at least  $t^n$ . This gives us

$$\mathbb{E}_{\mathcal{H}} \sum_{i=1}^n \mathbb{1}_{\{x_i \in C_n\}} \mathbb{1}_{\{x_i \leq t\}}$$

$\geq t^n$

$$\mathbb{E}_{\mathcal{H}} \sum_{i=1}^n \mathbb{1}_{\{x_i \in d_n(B; x_n)\}} \mathbb{1}_{\{x_i \leq t\}}$$

;

which implies

$$\mathbb{E}_{\mathcal{H}} \sum_{i=1}^n \mathbb{1}_{\{x_i \in C_n\}} \mathbb{1}_{\{x_i \leq t\}}$$

$\geq t^n \mathbb{E}_{\mathcal{H}}$

$$\sum_{i=1}^n \mathbb{1}_{\{x_i \in d_n(B; x_n)\}} \mathbb{1}_{\{x_i \leq t\}}$$

:

Now, note in the latter expectation, since  $\mu$  is uniform and  $x = x_n$ , it follows that the random variable  $\sum_{i=1}^n \mathbb{1}_{\{x_i \in d_n(B; x_n)\}} \mathbb{1}_{\{x_i \leq t\}}$  is independent of the event " $\mu(x) \geq t^n$ ", with the latter event having probability  $0.5$ . This is because  $d_n(B; x_n)$  only depends on the set of indexes  $B(x_n) = \{i \mid x_i \in x_n\}$ . Hence, conditional on  $\sum_{i=1}^n \mathbb{1}_{\{x_i \in d_n(B; x_n)\}} \mathbb{1}_{\{x_i \leq t\}}$ , the expectation

$$\mathbb{E}_{\mathcal{H}} \sum_{i=1}^n \mathbb{1}_{\{x_i \in d_n(B; x_n)\}} \mathbb{1}_{\{x_i \leq t\}}$$

is such that with probability  $0.5$  we have  $\sum_{i=1}^n \mathbb{1}_{\{x_i \in d_n(B; x_n)\}} \mathbb{1}_{\{x_i \leq t\}} \geq t^n$  and with

probability  $0.5$  we have  $\sum_{i=1}^n \mathbb{1}_{\{x_i \in d_n(B; x_n)\}} \mathbb{1}_{\{x_i \leq t\}} < t^n$  such that  $\mathbb{E}_{\mathcal{H}}$

value of the above conditional expectation is at least  $0.5 \sum_{i=1}^n \mathbb{1}_{\{x_i \leq t\}} \geq (1 - t) \sum_{i=1}^n \mathbb{1}_{\{x_i \leq t\}} = 0.5$

$\geq 2 \sum_{i=1}^n \mathbb{1}_{\{x_i \leq t\}}$ .

Combining these, we get the lower bound

$$\mathbb{E}_{\mathcal{H}} \sum_{i=1}^n \mathbb{1}_{\{z_i(x) \leq t\}} \geq 0.5 t^n$$

$\geq 0.5 t^n$

$\geq 0.5 t^n$



at most 1=

p  
2.

From Markov's inequality, one can show that for a random variable  $X \in [0; 1]$  and  $\epsilon \in (0; 1)$ , we have that if  $E[X] \leq \epsilon$ , then  $P[X > 1 - \epsilon] \leq \epsilon$ .

17

Prediction and Stochastic Choice, Pathikrit Basu

One can argue that the class of all monotone, continuous and convex preferences  $\mathcal{P}$  has  $VC_+(\mathcal{P}) = +1$ . In the context of uncertainty from a result in Basu and Echenique [2018], it follows that Max-min preferences (MEU) have infinite VC dimension whenever  $j \geq 3$ . However, the arguments can be modified to show that  $VC_+(\mathcal{P}_{MEU}) = +1$  as well. Of course, non-monotonic preferences can be non-predictable as well. Consider again the following econometric specification: an alternative has utility  $u(x_a; \theta) = \theta_a$ , where each  $\theta_a$  independently follows a normal distribution with  $(\theta_a; \sigma_a^2)$

a) as mean

and variance. Suppose that for the utility class  $\mathcal{U} = \{u(\cdot; \theta) : \theta \in \mathbb{R}^j\}$ , the corresponding preference class has  $VC_+(\mathcal{P}_{\mathcal{U}}) = +1$ . We define now preferences of the form  $x_0$

$$u(x_a; \theta_a) = \theta_a - \frac{1}{2} \theta_a^2$$

b

if and only if  $u(x_a; \theta_a) + \theta_a \geq u(x_b; \theta_b) + \theta_b$ . Then, it follows that  $VC_+(\mathcal{P}_{\mathcal{U}}) = +1$ .

Preference heterogeneity could be specified by a product probability measure  $\mu = \prod_{a \in A} \mu_a$  on  $\mathbb{R}^A$ , where  $\mu_a$  follows a multivariate normal distribution as a product of independent normal distributions with parameters  $(\theta_a; \sigma_a^2)$

a)  $\mu_{a2A}$ . Then, without any restrictions on  $\mu$  and not restricting  $f(\theta_a; \sigma_a^2)$

a)  $\mu_{a2A}$ ,

it obtains that the implied stochastic choice model would be non-predictable. This follows from the same arguments as in the proof of Theorem 4.2.

For these models, the above result establishes non-predictability or irrecoverability of heterogeneity in a certain sense (uniform). The dimension being infinite has the consequence that the preference class is complex enough to significantly affect the predictability of the model. The variability in choice behaviour permitted by the model is large enough to deter good worst case guarantees on out-of-sample estimation error.

### 4.3 Cognitive Heterogeneity

In this section, we will study stochastic choice in the setting where there is data on choice from menus of alternatives. As before,  $A$  is a finite set of alternatives. Hence, we will have  $X = 2^A$ . For a choice probability function  $\mu : X \rightarrow [0; 1]$ , for each menu of alternatives  $x = A \in X$ , the probability vector  $\mu(x)$  has full support in  $A$ . For each  $a \in A$ , the quantity  $\mu_a(x)$  denotes the probability of alternative  $a$  being chosen from the menu  $x$ . We will assume the following holds for  $\mu$ .

**Definition 4.1.** (Positivity) A stochastic choice function  $\mu$  satisfies positivity if for each menu  $A$  and  $a \in A$ , we have  $\mu_a(A) > 0$ .

In this context, the choice probabilities from the Logit or the Luce model are as follows. There exists a function  $w : A \rightarrow \mathbb{R}_+$ , which assigns weights to the various alternatives. For each menu  $A$ , the alternative  $a \in A$  is chosen with probability

$$\mu_a(A) = \frac{w(a)}{\sum_{b \in A} w(b)}$$

$$w(a)$$

$$w(b)$$

:

It is well-known that under positivity, a stochastic choice map  $\mu$  has the Luce representation if and only if  $\mu$  satisfies the Independence of Irrelevant Alternatives axiom (IIA). The IIA axiom says that if  $A$  and  $B$  are two menus and we have two alternatives  $a; b \in A \setminus B$ . Then,

$$\frac{\mu_a(A)}{\mu_b(A)}$$

$$= \frac{\mu_a(B)}{\mu_b(B)}$$

$$=$$

$$\frac{\mu_a(B)}{\mu_b(B)}$$

$$\frac{\mu_a(B)}{\mu_b(B)}$$

:

Another example is the stochastic choice model involving consideration sets (Manzini and Mariotti [2007]). The model is defined as follows. There is a function  $\alpha : A \rightarrow (0; 1)$  and a strict order  $\succ$  on  $A$ . The choice probabilities are given as follows.

When  $j = 2$ , the Max-min model has VC dimension equal to 2.

18

Prediction and Stochastic Choice, Pathikrit Basu

---

$$p_a(A) := (1 - \alpha) \prod_{a \in A} p_a$$

Y  
b2A.b\_a  
(b).

The interpretation is that  $1 - \alpha$  is the probability with which alternative  $a$  will be considered in choosing from a given menu  $A$ . Once the considered alternatives are determined, the agent chooses according to the strict order  $\succ$ . The following result obtains in this setting.

**Proposition 4.3.** Suppose  $\succ$  is either the Logit/Luce model or the model with consideration sets.

Then, under the log rule  $S_{\log}$ , the Pollard dimension of  $S_{\log, \succ}$  is at most  $O(jAj_2)$ .

**Proof.** The proof is in the appendix and applies a result on the VC dimension of neural networks (see Proposition 6.2). The bound is derived by computing the total number of arithmetic operations needed to express the choice probabilities in the respective models. Given that the expressions can be written via addition, multiplication and division of real numbers, the number of operations, is at most linear in  $jAj$ . Together with the dimensionality of the parameter space, which is also linear in  $jAj$  (every strict order has a one-to-one numeric representation), we get the quadratic bound.

The techniques from the previous section can be extended to allow recoverability of cognitive heterogeneity. This type of heterogeneity corresponds to a distribution  $\mu$  over  $(; \succ)$ . By arguments similar to the proof of Theorem 4.1, one shows that finite pollard dimension of  $S_{\log, \succ}$  allows for consistent estimation of the true distribution  $\mu_0$ , if we place a restriction that the choice probabilities are bounded below. Since we are applying the log scoring rule, the estimator is essentially based on the conditional maximum likelihood procedure.

#### 4.4 Non-uniform Prediction

The notion of consistency considered above requires that convergence should take place uniformly over both the distribution over characteristics  $\mu_0$  and the true stochastic choice function  $\mu_0$ . The main advantage of such a requirement is that it leads to estimates that are robust with respect to the data generating process. Suppose, we instead consider non-uniform consistency which only requires that the divergence between true and estimated choice probabilities converges to zero. We provide a definition.

**Definition 4.2.** A prediction method  $\hat{\mu}$  is pointwise consistent (with respect to  $d$  and  $\succ$ ) if for all  $0 < \epsilon < 1$ , and for all  $\mu_0 \in \mathcal{M}(X)$ ,  $\mu_0 \succ \mu$ , there exists  $N(\epsilon; \mu; \mu_0)$  such that for all  $n \geq N(\epsilon; \mu; \mu_0)$ ,

$$d_0(\hat{\mu}_n; \mu_0) < \epsilon$$

□  
— 1 □ :

We say that a model  $\succ$  is pointwise predictable with respect to  $d$  if there exists a prediction method  $\hat{\mu}$ , which is pointwise consistent with respect to  $d$  and  $\succ$ .

The above is also a kind of PAC criterion but only requires that the random variable  $d_0(\hat{\mu}_n; \mu_0)$  converge in probability to zero. Note that we had a stronger notion of consistency before. Hence, if a stochastic choice model is predictable according to Definition 2.1, then it is also pointwise predictable.

It turns out that one can study pointwise predictability via a weaker definition of Glivenko-Cantelli classes called universal Glivenko-Cantelli classes (see, for example, Dudley et al. [1991], van Handel [2013]). This requires that for a real-valued function class  $F$  (defined on a set  $Z$ ) we have that

$$\sup_{f \in F} \frac{1}{n} \sum_{i=1}^n f(z_i) \xrightarrow{P_n} \int f(z) d\mu(z)$$

converges to zero for each probability measure  $\mu \in \mathcal{M}(Z)$ . In contrast to uniform Glivenko-Cantelli classes, the convergence need not be uniform over  $\mu$ . Note, however, we still have uniformity over  $F$  as we have convergence of the supremum. This would also be closer to the uniformity criterion typically required in extremum estimation (see Amemiya [1985]). Similar to the proof of Proposition 3.1, one can argue that for a stochastic choice model  $\succ$  and scoring rule  $S$ , if  $S_{\succ}$  is universal Glivenko-Cantelli, then  $\succ$  is pointwise predictable w.r.t  $d_S$ .

19

Prediction and Stochastic Choice, Pathikrit Basu

A useful characterisation of universal Glivenko-Cantelli classes is given by van Handel [2013] in terms of Boolean  $\mu$ -independence. A sequence of functions  $f_j: Z \rightarrow \mathbb{R}$  is Boolean  $\mu$ -independent at levels  $\epsilon > 0$  if for every subset  $F \subseteq \mathbb{N}$ ,

$$\prod_{j \in F} \int f_j(z) d\mu(z) < \epsilon \quad \text{and} \quad \prod_{j \in F} \int f_j(z) d\mu(z) > \epsilon$$

It then holds that a function class is universal Glivenko-Cantelli if and only if there exists no Boolean  $\mathcal{C}$ -independent sequence in  $F$  for any  $\mathcal{C}$ . From this, it follows that a parametrised function class  $F = \{f_{\mathbf{z}}(\cdot) : \mathbf{z} \in \mathcal{Z}\}$ , where  $\mathcal{Z}$  is a compact metric space and  $f_{\mathbf{z}}(\cdot)$  is continuous in  $\mathcal{X}$  for each  $\mathbf{z}$ , is a universal Glivenko-Cantelli class. We can then show that under some conditions, a stochastic preference model supported on continuous preferences would be pointwise predictable. As before, it will be convenient to work with the closed convergence topology on continuous preferences (see appendix). We state the result below.

**Proposition 4.4.** Let  $\mathcal{P}_{\text{CONT}}$  be the set of all continuous preferences over the set of alternatives  $X$ , equipped with the closed convergence topology. Further, let  $\mathcal{M}_{\text{CONT}}(\mathcal{P}_{\text{CONT}})$  be any set of probability measures satisfying the condition 4.1, that is closed in the topology of weak convergence. Then,  $\mathcal{M}_{\text{CONT}}(\mathcal{P}_{\text{CONT}})$  is pointwise predictable with respect to any  $L_p$  metric. Proof. Can be found in the appendix.

It is interesting to compare the above result with the results in Theorem 4.1 and Theorem 4.2. Firstly, there may be continuous preferences with infinite VC dimension (for eg. all continuous convex preferences) and there may be discontinuous preferences with finite VC dimension (for eg. lexicographic preferences). Secondly, conditional on considering only continuous preferences, which allows pointwise predictability, robust estimation is characterised by the dimension of the preference classes ( $VC$  and  $VC^+$ ). If the dimension is finite, we get uniform consistency but for infinite dimension, the estimates are non-robust, to an extent that the same estimation procedure that would guarantee pointwise learning, would lead to erroneous out-of-sample predictions in the worst case. Lastly, the techniques for obtaining rates of convergence for pointwise predictability, would involve deriving bounds on the bracketing or covering numbers of the function class  $\mathcal{S}_{\text{CONT}}$  (see Dudley et al. [1991]). This would also be similar to PAC learning problems in the deterministic case with a non-uniform predictability criterion (see Blumer et al. [1987], Benedek and Itai [1994], Shalev-Shwartz and Ben-David [2014]) and also structural risk minimisation (see Vapnik and Chervonenkis [1974] and Vapnik [1998]).

## 5 Conclusion

This paper presents new learning rules for estimation of choice probabilities in various stochastic choice models. The main feature of preferences that is taken in account is its model complexity, for example, the VC dimension of the model. Perhaps interestingly, given that VC dimension is a characteristic of the collection of sets that defines the model, it also characterises uniform learning involving distributions over the collection, leading to random discrete choice. We also present a wide class of models and learning rules, allowing for applicability to discrete choice models used in statistics.

## References

- Naoki Abe, Manfred K Warmuth, and Jun-ichi Takeuchi. Polynomial learnability of probabilistic concepts with respect to the kullback-leibler divergence. In Proceedings of the fourth annual Prediction and Stochastic Choice, Pathikrit Basu workshop on Computational learning theory, pages 277–289. Morgan Kaufmann Publishers Inc., 1991.
- Naoki Abe, Jun-ichi Takeuchi, and Manfred K Warmuth. Polynomial learnability of stochastic rules with respect to the kl-divergence and quadratic distance. IEICE Transactions on Information and Systems, 84(3):299–316, 2001.
- Charalambos D Aliprantis and Kim C Border. Infinite Dimensional Analysis: A Hitchhiker’s Guide. Springer Science & Business Media, 2006.
- Noga Alon, Shai Ben-David, Nicolo Cesa-Bianchi, and David Haussler. Scale-sensitive dimensions, uniform convergence, and learnability. Journal of the ACM (JACM), 44(4):615–631, 1997.
- Takeshi Amemiya. Advanced econometrics. Harvard university press, 1985.
- Martin Anthony and Peter L Bartlett. Neural network learning: Theoretical foundations. cambridge university press, 2009.
- Susan Athey and Guido W Imbens. Machine learning methods that economists should know about. Annual Review of Economics, 11, 2019.
- Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and structural results. Journal of Machine Learning Research, 3(Nov):463–482, 2002.
- Pathikrit Basu and Federico Echenique. Learnability and models of decision making under uncertainty. In Proceedings of the 2018 ACM Conference on Economics and Computation, pages 53–53. ACM, 2018.
- Pathikrit Basu and Federico Echenique. On the falsifiability and learnability of decision theories. Theoretical Economics, 15(4):1279–1305, 2020.
- Eyal Beigman and Rakesh Vohra. Learning from revealed preference. In Proceedings of the 7th ACM

Conference on Electronic Commerce, pages 36–42. ACM, 2006.

Gyora M Benedek and Alon Itai. Nonuniform learnability. *Journal of Computer and System Sciences*, 48(2):311–323, 1994.

Steven Berry, James Levinsohn, and Ariel Pakes. Automobile prices in market equilibrium. *Econometrica: Journal of the Econometric Society*, pages 841–890, 1995.

Anselm Blumer, Andrzej Ehrenfeucht, David Haussler, and Manfred K Warmuth. Occam's razor. *Information processing letters*, 24(6):377–380, 1987.

Kim C Border and Uzi Segal. Dynamic consistency implies approximately expected utility preferences. *Journal of Economic Theory*, 63(2):170–188, 1994.

Richard A Briesch, Pradeep K Chintagunta, and Rosa L Matzkin. Nonparametric discrete choice models with unobserved heterogeneity. *Journal of Business & Economic Statistics*, 28(2):291–307, 2010.

Andreas Buja, Werner Stuetzle, and Yi Shen. Loss functions for binary class probability estimation and classification: Structure and applications. Working draft, November, 3, 2005.

Christopher P Chambers, Federico Echenique, and Nicolas S Lambert. Recovering preferences from finite data. *Econometrica*, 89(4):1633–1664, 2021.

21

Prediction and Stochastic Choice, Pathikrit Basu

Zachary Chase and Siddharth Prasad. Learning Time Dependent Choice. In Avrim Blum, editor, 10th Innovations in Theoretical Computer Science Conference (ITCS 2019), volume 124 of Leibniz International Proceedings in Informatics (LIPIcs), pages 62:1–62:19. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, 2018. ISBN 978-3-95977-095-8.

Soo Hong Chew, Larry G Epstein, Uzi Segal, et al. Mixture symmetry and quadratic utility. *Econometrica*, 59(1):139–63, 1991.

Pierpaolo De Blasi, Lancelot F James, John W Lau, et al. Bayesian nonparametric estimation and consistency of mixed multinomial logit choice models. *Bernoulli*, 16(3):679–704, 2010.

Anthony Downs et al. An economic theory of democracy. 1957.

Richard M Dudley. Uniform central limit theorems, volume 142. Cambridge university press, 2014.

Richard M Dudley, Evarist Giné, and Joel Zinn. Uniform and universal glivenko-cantelli classes. *Journal of Theoretical Probability*, 4(3):485–510, 1991.

Drew Fudenberg, Ryota Iijima, and Tomasz Strzalecki. Stochastic choice and revealed perturbed utility. *Econometrica*, 83(6):2371–2409, 2015.

Itzhak Gilboa. Theory of decision under uncertainty, volume 45. Cambridge university press, 2009.

Tilmann Gneiting and Adrian E Raftery. Strictly proper scoring rules, prediction, and estimation. *Journal of the American Statistical Association*, 102(477):359–378, 2007.

Faruk Gul and Wolfgang Pesendorfer. Random expected utility. *Econometrica*, 74(1):121–146, 2006.

Werner Hildenbrand. On economies with many agents. *Journal of economic theory*, 2(2):161–188, 1970.

Werner Hildenbrand. Core and Equilibria of a Large Economy.(PSME-5). Princeton university press, 2015.

Gil Kalai. Learnability and rationality of choice. *Journal of Economic theory*, 113(1):104–117, 2003.

Yakar Kannai. Continuity properties of the core of a market. *Econometrica* (pre-1986), 38(6):791, 1970.

Michael J Kearns and Robert E Schapire. Efficient distribution-free learning of probabilistic concepts. *Journal of Computer and System Sciences*, 48(3):464–497, 1994.

Michael J Kearns and Umesh Vazirani. An introduction to computational learning theory. MIT press, 1994.

Alexander Kechris. Classical descriptive set theory, volume 156. Springer Science & Business Media, 2012.

David Kreps. Notes on the Theory of Choice. Routledge, 1988.

Jay Lu. Random choice and private information. *Econometrica*, 84(6):1983–2027, 2016.

R Duncan Luce. Individual choice behavior: A theoretical analysis. John Wiley, 1959.

R Duncan Luce. The choice axiom after twenty years. *Journal of mathematical psychology*, 15(3):215–233, 1977.

22

Prediction and Stochastic Choice, Pathikrit Basu

Charles F Manski. Maximum score estimation of the stochastic utility model of choice. *Journal of econometrics*, 3(3):205–228, 1975.

Charles F Manski. Semiparametric analysis of discrete response: Asymptotic properties of the maximum score estimator. *Journal of econometrics*, 27(3):313–333, 1985.

Paola Manzini and Marco Mariotti. Sequentially rationalizable choice. *American Economic Review*, 97(5):1824–1839, 2007.

Andreu Mas-Colell. On the continuous representation of preorders. *International Economic Review*, pages 509–513, 1977.

Rosa L Matzkin. Nonparametric and distribution-free estimation of the binary threshold crossing and the binary choice models. *Econometrica: Journal of the Econometric Society*, pages 239–270, 1992.

Daniel McFadden. Modeling the choice of residential location. *Transportation Research Record* (673), 1978.

Daniel McFadden and Kenneth Train. Mixed mnl models for discrete response. *Journal of applied Econometrics*, 15(5):447–470, 2000.

Daniel L McFadden. Econometric analysis of qualitative response models. *Handbook of econometrics*, 2:1395–1457, 1984.

Ilya Molchanov. *Theory of random sets*, volume 19. Springer, 2005.

Partha Niyogi. *The computational nature of language learning and evolution*, volume 43. MIT press Cambridge, MA, 2006.

Partha Niyogi. *The informational complexity of learning: perspectives on neural networks and generative grammar*. Springer Science & Business Media, 2012.

Partha Niyogi, Stephen Smale, and Shmuel Weinberger. A topological view of unsupervised learning from noisy data. *SIAM Journal on Computing*, 40(3):646–663, 2011.

Matthew Parry et al. Linear scoring rules for probabilistic binary classification. *Electronic Journal of Statistics*, 10(1):1596–1607, 2016.

David Pollard. *Convergence of stochastic processes*. Springer Science & Business Media, 1984.

Leonard J Savage. Elicitation of personal probabilities and expectations. *Journal of the American Statistical Association*, 66(336):783–801, 1971.

Reinhard Selten. Axiomatic characterization of the quadratic scoring rule. *Experimental Economics*, 1(1):43–61, 1998.

Shai Shalev-Shwartz and Shai Ben-David. *Understanding machine learning: From theory to algorithms*. Cambridge university press, 2014.

Michel Talagrand et al. The glivenko-cantelli problem. *The Annals of Probability*, 15(3):837–870, 1987.

Leslie G Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984. 23

Prediction and Stochastic Choice, Pathikrit Basu

Aad W Van Der Vaart and Jon A Wellner. Weak convergence. In *Weak convergence and empirical processes*, pages 16–28. Springer, 1996.

Ramon van Handel. The universal glivenko–cantelli property. *Probability Theory and Related Fields*, 155(3-4):911–934, 2013.

Vladimir Vapnik. *Statistical learning theory*. 1998, volume 3. Wiley, New York, 1998.

Vladimir Vapnik and Alexey Chervonenkis. *Theory of pattern recognition*, 1974.

Vladimir N Vapnik and Aleksei Yakovlevich Chervonenkis. The uniform convergence of frequencies of the appearance of events to their probabilities. In *Doklady Akademii Nauk*, volume 181, pages 781–783. Russian Academy of Sciences, 1968.

Kenji Yamanishi. A learning criterion for stochastic rules. *Machine Learning*, 9(2-3):165–203, 1992.

Morteza Zadimoghaddam and Aaron Roth. Efficiently learning from revealed preference. In *International Workshop on Internet and Network Economics*, pages 114–127. Springer, 2012.

## 6 Appendix

### 6.1 Technical

When dealing with preferences in economic contexts, a popular topology on the space of preference relations is that given by the topology of closed convergence (Kannai [1970], Hildenbrand [1970]). We provide some general definitions in this section and note some important results from the literature that will be useful in our setting. Let  $(Y; \tau_Y)$  be a topological space. Let  $C$  denote the set of all closed subsets of  $Y$ . The closed convergence topology on  $C$  is defined as follows (see also Aliprantis and Border [2006]). For any compact  $K \subseteq Y$  and finitely many open sets  $\{U_i\}_{i=1}^n$  in  $Y$ , we define the following subset of  $C$ .

$$O(K; \{U_i\}) = \{F \subseteq Y \mid K \subseteq F \text{ and } F \setminus U_i \text{ is compact for each } i = 1, 2, \dots, n\}$$

All sets of the form  $O(K; \{U_i\})$  constitute a basis for the closed convergence topology on  $C$ . We will denote this topology as  $\tau_c$ . When  $Y$  is a locally compact and separable metric space, then it turns out that  $(C; \tau_c)$  is a compact metrizable space (see Hildenbrand [2015]). Further, it also holds that  $(C; \tau_c)$  is second countable (see Molchanov [2005]).

We now note how the closed convergence topology provides us a way to topologize the space of continuous preferences. Let  $(X; \tau_X)$  be a topological space of alternatives. A preference relation on  $X$  is a binary relation  $\succsim_X$  on  $X$  which is complete and transitive. A continuous preference relation

$\succsim$  is a preference relation that is closed in  $X \times X$  (under the product topology). Hence, a class of continuous preferences  $\mathcal{P} \subseteq 2^{X \times X}$  can be endowed with the subspace topology, say  $\tau_{\mathcal{P}}$ , generated by the topology of closed convergence on closed subsets of  $2^{X \times X}$ .

The following theorem, due to Mas-Colell [1977], provides conditions under which a utility representation exists that is jointly continuous in the alternatives and the preferences.

**Proposition 6.1.** (Mas-Colell [1977]) Let  $X$  be a locally compact and second countable space of alternatives. Further, let  $\mathcal{P}_{\text{CONT}}$  denote the set of all continuous preference relations on  $X$ , equipped with the subspace topology induced by the closed convergence topology. Then, there exists a utility function  $U : X \times \mathcal{P}_{\text{CONT}} \rightarrow \mathbb{R}$ , which is jointly continuous and each  $\succsim \in \mathcal{P}_{\text{CONT}}$  is represented by  $U(\cdot, \succsim)$ .

24

Prediction and Stochastic Choice, Pathikrit Basu

An immediate corollary of the above result is the following. For any finite set of alternatives  $(x_a)_{a \in A}$ , the set of preferences  $\{\succsim \in \mathcal{P}_{\text{CONT}} \mid x_a \succsim x_b \text{ for all } b \in A\}$  is closed and  $\{\succsim \in \mathcal{P}_{\text{CONT}} \mid x_a \succ x_b \text{ for all } b \in A\}$  is open with respect to the closed convergence topology.

Now, consider a subspace of preferences  $\mathcal{P} \subseteq \mathcal{P}_{\text{CONT}}$ . We denote as  $\mathcal{M}(\mathcal{P})$ , the set of all Borel probability measures defined on  $\mathcal{P}$ . We will endow  $\mathcal{M}(\mathcal{P})$  with the topology of weak convergence (see Aliprantis and Border [2006]). Hence, for a sequence  $\mu_n$  and probability measure  $\mu$  in  $\mathcal{M}(\mathcal{P})$  we have that  $\mu_n \rightarrow \mu$  in the weak convergence topology if and only if for every bounded continuous function  $f : \mathcal{P} \rightarrow \mathbb{R}$ , the convergence  $\int f(\succsim) d\mu_n \rightarrow \int f(\succsim) d\mu$  holds. Now, if the space of alternatives  $X$  is a locally compact and separable metric space, then it is also the case that  $X \times X$  is a locally compact and separable metric space. Hence, the subspace  $\mathcal{P}_{\text{CONT}} \subseteq 2^{X \times X}$  is a second countable compact metric space. This means that for any subspace  $\mathcal{P} \subseteq \mathcal{P}_{\text{CONT}}$ , the set of probability measures  $\mathcal{M}(\mathcal{P})$  with the weak convergence topology is a separable metrizable space (see Aliprantis and Border [2006]). We now show some measurability and continuity properties of the stochastic choice probabilities  $\sigma_a(x; \cdot)$ .

We will first define the regularity condition that we assume on  $\mathcal{M}(\mathcal{P})$ . Consider the set of preferences where  $x_a$  is strictly preferred from a menu of alternatives  $x = (x_b)_{b \in A}$  i.e.  $Q(x; a) = \{\succsim \in \mathcal{P} \mid x_a \succ x_b \text{ for all } b \in A\}$ . The regularity assumption imposes on  $\mathcal{M}(\mathcal{P})$  that we have  $\int \mathbb{1}_{Q(x; a)} d\mu = 1$ . Also, by the definition of choice probabilities,  $\sigma_a(x; \cdot) = \int \mathbb{1}_{Q(x; a)} d\mu$ . The following lemma is useful.

**Lemma 6.1.** Let  $X$  be a locally compact, separable metric space of alternatives and let  $X \times X$  be such that  $x_a \succ x_b$  for all  $x = (x_c)_{c \in A} \in X$ . Suppose  $\mathcal{P} \subseteq \mathcal{P}_{\text{CONT}}$  and let  $\mathcal{M}(\mathcal{P}) \subseteq \mathcal{M}(X \times X)$  be the set of all probability measures on the preference class  $\mathcal{P}$  satisfying the regularity condition. The stochastic choice function  $\sigma_a(x; \cdot)$  is continuous in  $\cdot$  for each  $x \in X$ ; continuous in  $x$  for each  $\succsim \in \mathcal{M}(\mathcal{P})$ . Hence,  $\sigma_a(x; \cdot)$  is jointly measurable under the Borel sigma-algebra induced by the product topology on  $X \times \mathcal{M}(\mathcal{P})$ .

*Proof.* Consider a fixed  $x \in X$  and let  $\mu_n$  be a sequence and let  $\mu$  be a probability measure in  $\mathcal{M}(\mathcal{P})$  such that  $\mu_n \rightarrow \mu$  in the weak convergence topology. It suffices to argue that  $Q(x; a)$  is a continuity set according to the measure  $\mu$  i.e.  $\mu(\partial Q(x; a)) = 0$ , where the set  $\partial Q(x; a)$  is the boundary of the set  $Q(x; a)$ . Suppose we have a convergent sequence  $\mu_n$  in  $\mathcal{M}(\mathcal{P})$  with limit  $\mu$ . As  $\mu$  satisfies the regularity condition  $\int \mathbb{1}_{Q(x; a)} d\mu = 1$ , it suffices to show that  $\int \mathbb{1}_{Q(x; a)} d\mu_n \rightarrow \int \mathbb{1}_{Q(x; a)} d\mu$  for all  $\mu_n \in \mathcal{M}(\mathcal{P})$ . But this is true since the latter set is closed and contains the set  $Q(x; a)$  as noted above. It now follows that  $\sigma_a(x; \cdot)$  is continuous.

Now, let  $\mu_n \in \mathcal{M}(\mathcal{P})$  and suppose  $x_n \in X$  is a convergent sequence with limit  $x \in X$ . Note that since preferences are continuous, we have that  $Q(x; a) = \liminf_n Q(x_n; a)$ . This implies that  $\mathbb{1}_{Q(x; a)} = \liminf_n \mathbb{1}_{Q(x_n; a)}$ . Hence,  $\sigma_a(\cdot; \cdot)$  is lower semi-continuous. We now show upper-semicontinuity. Define the set  $Q(x; \text{Anfag}) := \{\succsim \in \mathcal{P} \mid x \succ y \text{ for all } y \in A\}$ . Note that since  $\mu$  satisfies regularity, we have  $\int \mathbb{1}_{Q(x; \text{Anfag})} d\mu = 1$  for all  $x \in X$ . Now, let  $x_n \rightarrow x$ . Again, as preferences are continuous, we obtain that  $Q(x; \text{Anfag}) = \liminf_n Q(x_n; \text{Anfag})$  which means  $\mathbb{1}_{Q(x; \text{Anfag})} = \liminf_n \mathbb{1}_{Q(x_n; \text{Anfag})}$ . But this implies that  $\mathbb{1}_{Q(x; a)} = \limsup_n \mathbb{1}_{Q(x_n; a)}$ . Hence, the function  $\sigma_a(\cdot; \cdot)$  is continuous as it is both lower and upper semi-continuous.

The two facts established above imply that  $\sigma_a(x; \cdot)$  is a Caratheodory function (see, for example, Aliprantis and Border [2006], Kechris [2012]). Hence,  $\sigma_a(x; \cdot)$  is jointly measurable in  $(x; \cdot)$ . Now it follows that under the conditions satisfied by the above lemma, we obtain appropriate measurability of the function class  $F = \{\sigma_a(\cdot; \cdot) \mid a \in A\}$  for the Glivenko-Cantelli theorems to apply. From Van Der Vaart and Wellner [1996], a sufficient condition is a separability condition called pointwise

25

Prediction and Stochastic Choice, Pathikrit Basu

measurability. It requires that there exist a countable subset  $F_0 \subseteq F$  such that for any  $f \in F$ , there exists a sequence in  $F_0$  which converges to  $f$  pointwise. Since  $\mathcal{M}(\mathcal{P})$  is a separable metric space, it follows that  $\mathcal{M}(\mathcal{P})$  is a separable metric space as well. Let  $\mathcal{M}_0 \subseteq \mathcal{M}(\mathcal{P})$  be a countable dense subset.

Since from Lemma 6.1 above, the choice probability  $\mu_a(x; \cdot)$  is continuous in  $\cdot$  and from the definition of the Quadratic scoring rule  $S_{br}$ , it follows that the countable set  $F_0 = \{f_{S_{br}}(\cdot; \cdot) \mid \cdot \in \mathcal{G}\}$  makes  $F$  pointwise measurable. Alternatively, Talagrand et al. [1987] requires that  $F$  belong to  $L_1(\cdot)$ , for each  $\cdot \in \mathcal{X} \setminus A$ .

We will now discuss the existence of probability measures  $\mu$  which satisfy the regularity condition.

The following result applies when preferences are defined through a parametrised utility function  $u : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}$ . For any such function  $u$  and  $\cdot \in \mathcal{G}$ , we will denote as  $\succ_{\cdot}$ , the preference relation  $u$  induces on  $\mathcal{X}$ . Hence, we define  $x \succ_{\cdot} y$  if and only if  $u(x; \cdot) > u(y; \cdot)$ . We will say that a probability measure  $\mu \in \mathcal{P}(\mathcal{X})$  satisfies the regularity condition for  $u$  if for all  $x = (x_b)_{b \in A} \in \mathcal{X}_A$  with  $x_a = x_b$  for all  $a = b$ , we have  $\int [a \cdot 2^a \int u(x_a; \cdot) > \max_{x_b = a} u(x_b; \cdot)] d\mu = 1$ .

The following result holds. We will assume that preferences are locally strict (see Border and Segal [1994]).

**Lemma 6.2.** Let  $\mathcal{X}$  be a locally compact and separable metric space of alternatives,  $\mathcal{P}_{CONT}$  be the set of all continuous preference relations on  $\mathcal{X}$  (endowed with the topology induced by the closed convergence topology) and  $\mathcal{G}$  be a metric space of parameters. Let  $u : \mathcal{X} \times \mathcal{G} \rightarrow \mathbb{R}$  be a jointly continuous utility function. Further, suppose that  $u$  is locally strict i.e. for all  $\cdot$ , if  $u(x; \cdot) = u(y; \cdot)$ , then every neighbourhood of  $(x; y)$  contains a pair  $(x_0; y_0)$  such that  $u(x_0; \cdot) > u(y_0; \cdot)$ . Then, the function  $\mu : \mathcal{P}_{CONT} \rightarrow \mathbb{R}$  defined as

$$\mu(\cdot) = \mu_{\cdot};$$

is a continuous function. Now, suppose  $\mu \in \mathcal{P}(\mathcal{P}_{CONT})$  is a Borel probability measure on the space  $\mathcal{P}_{CONT}$  which satisfies the regularity condition for  $u$ . Also suppose that the class of preferences  $\mathcal{P}_u = \{f_{\cdot, g} \mid \cdot \in \mathcal{G}\}$  is Borel. Then, the induced probability measure  $\nu$  on the Borel subsets of  $\mathcal{P}_{CONT}$ , defined as

$$\nu(B) = \mu(\mu^{-1}(B));$$

satisfies the regularity condition on  $\mathcal{P}_u$ . If  $\mathcal{P}_u$  is also closed in  $\mathcal{P}_{CONT}$  and  $\mu$  has full support on  $\mathcal{P}_{CONT}$ , then  $\nu$  has full support on  $\mathcal{P}_u$ .

**Proof.** To show that  $\nu$  is continuous, it will suffice to show that the pre-image of any basic open set is open in  $\mathcal{P}_{CONT}$ . Let  $O(K; f_{U_i, g_i})$  be a basic open set. Now, each  $U_i$  is open in  $\mathcal{X} \times \mathcal{X}$ . Hence, from local strictness of  $u$  it follows that for any  $\cdot$ , if  $u(x; \cdot) = u(y; \cdot)$ , where  $(x; y) \in U_i$ , there exists another pair  $(x_0; y_0) \in U_i$  such that  $u(x_0; \cdot) > u(y_0; \cdot)$ . For convenience, define the function  $\mu(x; y; \cdot) = u(x; \cdot) - u(y; \cdot)$ . Then, we get that the pre-image of  $O(K; f_{U_i, g_i})$  can be written as

$$\mu^{-1}(O(K; f_{U_i, g_i})) = \{ \cdot \in \mathcal{P}_{CONT} \mid \max_{(x,y) \in 2K} \mu(x; y; \cdot) > K \}$$

$$\mu(x; y; \cdot) < 0 \text{ and } \sup_{(x,y) \in 2U_i} \mu(x; y; \cdot) > 0 \text{ for each } i \in I_g;$$

$$\mu(x; y; \cdot) > 0 \text{ for each } i \in I_g;$$

Now, as  $K$  is compact, from the theorem of the maximum, it follows that  $\max_{(x,y) \in 2K} \mu(x; y; \cdot)$  is continuous in  $\cdot$ . Hence, the set  $\mu^{-1}(K) = \{ \cdot \in \mathcal{P}_{CONT} \mid \max_{(x,y) \in 2K} \mu(x; y; \cdot) < 0 \}$  is open in  $\mathcal{P}_{CONT}$ . Now, note that  $\sup_{(x,y) \in 2U_i} \mu(x; y; \cdot)$

$\mu(x; y; \cdot)$  is a supremum of lower semi-continuous functions and is hence, lower semi-continuous in  $\cdot$ . Hence, the set  $\mu^{-1}(i) = \{ \cdot \in \mathcal{P}_{CONT} \mid \sup_{(x,y) \in 2U_i} \mu(x; y; \cdot) > 0 \}$  is open, for each  $i \in I$ .

To conclude, observe that  $\mu^{-1}(O(K; f_{U_i, g_i})) = \mu^{-1}(K) \cap \bigcap_{i \in I} \mu^{-1}(i)$  is open as an intersection of finitely many open sets. Hence,  $\nu$  is continuous.

We now argue that  $\nu$  satisfies the regularity condition if  $\mu$  satisfies the regularity condition for  $u$ . Let  $B = \{ [a \cdot 2^a \int u(x_a; \cdot) > \max_{x_b = a} u(x_b; \cdot)] \}$ . Now, further

$$\nu(B) = \mu(\mu^{-1}(B)) = \int [a \cdot 2^a \int u(x_a; \cdot) > \max_{x_b = a} u(x_b; \cdot)] d\mu = 1. \text{ Now, further}$$

26 Prediction and Stochastic Choice, Pathikrit Basu

assume that  $\mu$  has full support on  $\mathcal{P}_{CONT}$  and  $\mathcal{P}_u$  is closed. Suppose for contradiction that there exists a closed strict subset  $P \subset \mathcal{P}_u$  such that  $\nu(P) = 1$ . But this means that  $\mu(\mu^{-1}(P)) = 1$  where the set  $\mu^{-1}(P)$  is closed as  $\mu$  is continuous and  $P$  is closed in  $\mathcal{P}_{CONT}$ . But  $\mu^{-1}(P) \cap \mu^{-1}(\mathcal{P}_u) = \emptyset$ , where both sets are closed but  $\mu(\mu^{-1}(P)) = 1$  meaning that  $\mu$  does not have full support, which is a contradiction.

## 6.2 Proofs from Section 3

The following is the proof for Proposition 3.1.

**Proof.** We show that  $\mu$  is predictable with respect to  $d_s$ . We shall construct an almost-ERM learning  $\hat{\mu}_{\epsilon}$ , which would be consistent with respect to  $d$  and  $\mu$ . For each  $n \in \mathbb{N}$ , define  $\mu_n := \mu$ . Now, for

$$Z_n = ((x_i; a_i))_{i=1}^n$$

$i=1 \dots n$ , let  $\hat{\mu}_{\epsilon}(Z_n)$  be such that

$$\hat{\mu}_{\epsilon}(Z_n) = \inf_{\mu \in \mathcal{P}} \sum_{i=1}^n \mu(x_i; a_i)$$

$$\sum_{i=1}^n \mu(x_i; a_i) + \epsilon;$$

where recall that for any  $\mu \in \mathcal{P}$ ,  $\sum_{i=1}^n \mu(x_i; a_i) := \sum_{i=1}^n \mu(x_i; a_i)$

$\mu_n$

$\sum_{i=1}^n S_{a_i}(\cdot(x_i))$ . Note that the infimum exists in the above definition since  $S_{\cdot}$  is bounded above by  $M$ . Hence, the  $\hat{S}_{\epsilon}$  is well-defined. We now show it is consistent.

Let  $(\epsilon, \delta) \in (0, 1)^2$ . Since  $S_{\cdot}$  is a Glivenko-Cantelli class of functions, for  $(\epsilon, \delta) \in (0, 1)^2$ , there exists  $N_0(\epsilon, \delta)$  such that for all  $n \geq N_0(\epsilon, \delta)$ , for all  $z \in Z$

$$\sum_{j=1}^n \sup_{f \in S_{\epsilon}} \sum_{i=1}^n f(z_i) - E_{\cdot}(f) < \epsilon \quad (6.1)$$

Now, let  $N(\epsilon, \delta) := \max\{N_0(\epsilon, \delta)\}$ . Now, suppose  $n \geq N(\epsilon, \delta)$ . Let  $\mu_{\epsilon, \delta}(X)$  be a distribution over characteristics and suppose the true choice probabilities are given by  $\mu_{\epsilon, \delta}$ . For convenience, denote  $E_{z_n}(f) = \int f(z) d\mu_{\epsilon, \delta}(z)$

Now suppose  $z_n$  is such that  $\sum_{j=1}^n E_{z_n}(f_j) - E_{\mu_{\epsilon, \delta}}(f_j) < \epsilon$ . Note that this happens with probability at least  $1 - \delta$  under  $\mu_{\epsilon, \delta}$

Consider the almost-ERM rule  $\hat{S}_{\epsilon}(z_n; a)$ . Define also  $f_{\epsilon, z_n}(x; a) := S(\hat{S}_{\epsilon}(z_n)(x); a)$ ,  $f_0(x; a) := S(\mu_{\epsilon, \delta}(x); a)$  and  $\mu_0 = \mu_{\epsilon, \delta}$ . It follows that

$$\sum_{j=1}^n |E_{\mu_0}(f_{\epsilon, z_n}) - E_{\mu_0}(f_0)| = E_{\mu_0}(f_0) - E_{\mu_0}(f_{\epsilon, z_n}) \quad (6.2)$$

$$= E_{\mu_0}(f_0) - E_{z_n}(f_{\epsilon, z_n}) + E_{z_n}(f_{\epsilon, z_n}) - E_{\mu_0}(f_{\epsilon, z_n}) \quad (6.3)$$

$$= E_{\mu_0}(f_0) - E_{z_n}(f_0) + 1/n + \epsilon \quad (6.4)$$

$$= \epsilon + \epsilon + \epsilon \quad (6.5)$$

Here, 6.2 follows since  $\mu_0$  minimises expected risk as  $S$  is incentive compatible (Lemma 3.1). The inequality 6.3 follows by assumption that  $\sum_{j=1}^n E_{z_n}(f_j) - E_{\mu_0}(f_j) < \epsilon$ . This yields  $E_{z_n}(f_{\epsilon, z_n}) - E_{\mu_0}(f_{\epsilon, z_n}) < \epsilon$ . Also, 6.4 follows since  $\hat{S}_{\epsilon}$  is almost-ERM. Lastly, 6.5 follows since  $n \geq 3/\epsilon$  and again from our assumption that  $\sum_{j=1}^n E_{z_n}(f_j) - E_{\mu_0}(f_j) < \epsilon$ , which yields  $E_{z_n}(f_0) - E_{\mu_0}(f_0) < \epsilon$ .

Now, since we have

$$E_{\mu_0}(f_0) - E_{\mu_0}(f_{\epsilon, z_n}) =$$

$$\int_X dS(\hat{S}_{\epsilon}(z_n)(x); \mu_0(x)) d\mu_0(x);$$

it follows that  $\hat{S}_{\epsilon}$  is consistent with respect to  $dS$  and  $\mu_0$ .

27

Prediction and Stochastic Choice, Pathikrit Basu

### 6.3 Proofs from Section 4

The following two lemmas will be useful.

**Lemma 6.3.** Suppose  $F$  is a uniformly bounded class of real-valued functions defined on the set  $Z$  and let  $M_{\cdot}(F)$ . For  $\mu \in \mathcal{P}(Z)$ . Define the function  $f_{\cdot}$  as

$$f_{\cdot}(z) =$$

$$\int_F$$

$$f(z) d_{\cdot}(f);$$

Consider the function class  $F_M = \{f_{\cdot} - \mu\}$ . Then,

$$R_{z_n}(F_M) \leq R_{z_n}(F);$$

for all  $z_n \in Z^n$ .

**Proof.** We first develop some notation. For  $z_n$  and a real-valued function  $f$ , define  $f(z_n) = (f(z_1); \dots; f(z_n)) \in \mathbb{R}^n$ . From the definition of Rademacher Complexity, it suffices to show that for any real vector  $\mu \in \mathbb{R}^n$ ,

$$\sup_{\mu \in \mathbb{R}^n} \int_F f(z_n) d_{\cdot}(f) \leq \sup_{\mu \in \mathbb{R}^n} \int_F f(z_n); \quad (6.6)$$

Suppose for contradiction that 6.6 does not hold. Then,  $\sup_{\mu \in \mathbb{R}^n}$

R

$$F_{\infty} f(z_n) d_{\infty}(f) > \sup_{f \in F} f(z_n).$$

Then, there exists  $\epsilon > 0$  such that

R

$$F_{\infty} f(z_n) d_{\infty}(f) > \sup_{f \in F} f(z_n).$$

But this means that there exists  $f_{\infty}$  such that  $f_{\infty}(z_n) > \sup_{f \in F} f(z_n)$ . This is a contradiction and hence, the result follows.

The following is another useful lemma.

**Lemma 6.4.** Suppose, as in the present context,  $Z = X \times A$ , where  $\|A_j\| < 1$ . Further, let  $F$  be a real-valued function class defined on  $Z$ . For each  $a \in A$ , define the following function class on  $X$ :

$$F_a = \{f(\cdot; a) : f \in F\}$$

Then, for all  $z_n = (x_i; a_i)_{i=1}^n$

$$R_{z_n}(F) \leq$$

$R_{x_n}$

$(F_a)$

where

$x_n = (x_i)_{i=1}^n$

and

$a_i = a_j$

for

$j \in N_a$ .

Proof.

Let

$z_n \in Z_n$ .

Define

$N_a = \{j : a_j = a\}$

and then

$z_n = (x_i; a)_{i \in N_a}$ .

Then,

$$R_{z_n}(F) = E_{z_n}$$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$\leq$

$E_{z_n}$

$i$   
 $=$   
 $X$   
 $a \in A$   
 $n_a$   
 $n$   
 $E$   
 $h$   
 $\sup_{f \in F}$   
 $1$   
 $n_a$   
 $X$   
 $i \in N_a$   
 $\prod_{i \in N_a} f(x_i; a)$   
 $i$   
 $=$   
 $X$   
 $a \in A$   
 $n_a$   
 $n$   
 $R^{n_a} (F_a)$ :

This is the set of all probability measures on  $F$ .

28

Prediction and Stochastic Choice, Pathikrit Basu

### 6.4 Proofs from Section 4.3

Proposition 4.3 follows straightforwardly from the following result on the VC dimension of neural networks.

**Proposition 6.2.** (Anthony and Bartlett [2009]) Let  $\mathcal{W} \subseteq \mathbb{R}^p$ . Suppose  $F = \{f(w; \cdot) : w \in \mathcal{W}\}$  is a parametrized class of 0-1 valued functions defined on a set  $W \subseteq \mathbb{R}^k$ . Further, suppose, for each  $f \in F$  and each  $(w; \cdot) \in \mathcal{W} \times W$ , computing the value of  $f(w; \cdot)$  takes no more than  $t$  many operations from the following

1. The arithmetic operations  $+$ ,  $-$ ,  $\cdot$ , and  $/$  defined on real numbers.
2. Pairwise comparisons of two real numbers involving the relations  $>$ ,  $<$ , and  $=$ .
3. Output 0 or 1.

Then, the VC dimension of  $F$  is at most  $4p(t + 2)$ .

We now prove Proposition 4.4.

Proof. We will first show that if we have a parametrised function class  $F = \{f(z; \cdot) : z \in \mathcal{Z}\}$ , where  $\mathcal{Z}$  is a compact metric space and  $f(z; \cdot)$  is continuous in  $\mathcal{Z}$  for  $z$ , then  $F$  contains no Boolean  $\epsilon$ -independent sequence for any  $\epsilon > 0$ . Suppose not. Then, there exists a sequence of parameters  $\{z_n\}_{n=1}^{\infty}$

such that  $\{f(z_n; \cdot)\}_{n=1}^{\infty}$  is a Boolean  $\epsilon$ -independent sequence in  $F$  for some  $0 < \epsilon < 1$ .

Since  $\mathcal{Z}$  is compact, there exists a convergent subsequence  $\{z_{n_k}\}_{k=1}^{\infty}$

It then follows that  $\{f(z_{n_k}; \cdot)\}_{k=1}^{\infty}$

is also a Boolean  $\epsilon$ -independent sequence at  $\mathcal{Z}$ . Now consider the sets  $K_{\text{odd}} = \{k : f(z_{n_k}; \cdot) < \epsilon\}$  and  $K_{\text{even}} = \{k : f(z_{n_k}; \cdot) > \epsilon\}$ . The following set is non-empty.

$$\bigcap_{k \in K_{\text{odd}}} \{z : f(z; z_{n_k}) < \epsilon\}$$

$$\bigcap_{k \in K_{\text{even}}} \{z : f(z; z_{n_k}) > \epsilon\}$$

Consider an element  $z$  in the above intersection. Further, let  $z_{n_k}$  be the limit of  $\{z_{n_k}\}_{k=1}^{\infty}$ .

Then,

$\lim_{k \rightarrow \infty} f(z; z_{n_k}) = f(z; z_{n_k})$ . However, from the above intersection, it follows that  $f(z; z_{n_k}) < \epsilon$

for  $k$  odd and  $f(z; z_{n_k}) > \epsilon$  for  $k$  even. Hence, we have two subsequences  $\{f(z; z_{n_k})\}_{k \in K_{\text{odd}}}$  and

$\{f(z; z_{n_k})\}_{k \in K_{\text{even}}}$  with different limits. But this is a contradiction as all subsequences of a convergent sequence must have the same limit.

We now show the pointwise predictability of continuous preferences. Note that  $P_{\text{CONT}} = \{f : \mathcal{X} \rightarrow \mathcal{Y}\}$

$\mathcal{X}$  is closed in  $\mathcal{X}$ . With the topology of closed convergence on  $P_{\text{CONT}}$ , which is also

the topology induced by the Hausdorff metric on  $P_{\text{CONT}}$  (see, for example, Aliprantis and Border

[2006]), it follows that  $P_{\text{CONT}}$  is compact. Hence,  $\mathcal{P}(P_{\text{CONT}})$  is compact in the weak convergence

topology which can be metrized by the Prokhorov metric. Since  $M$  is now a closed subset of a compact

metric space, it is also compact. Now, we will define  $\mathcal{M} := M$ .

Consider the Quadratic scoring rule  $S_{br}$  and the stochastic choice model  $\mu_{PCONT}(M)$ . Then, for any  $f \in S_{br, \mu_{PCONT}}(M)$ , we have

$$f(x; a; \mu) = 2_{-a}(x; \mu) \quad \square$$

X

b2A

$\mu_b(x; \mu)^2$ :

Now we need only show that if  $\mu_n \neq \mu$ , then  $f(x; a; \mu_n) \neq f(x; a; \mu)$ . This follows from Lemma 6.1.

---

© 2024 Pathikrit Basu; This is an Open Access article distributed under the terms of the Creative Commons Attribution License <http://creativecommons.org/licenses/by/2.0>, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.