

Limit of Learning and Efficient Strategies

Jie Ding, Enmao Diao, Jiawei Zhou, and Vahid Tarokh

Abstract

Machine learning systems learn from and make predictions by building models from observed data. Because large models tend to inflate the variance (e.g. overfitting) while small models tend to cause biases (e.g. underfitting) for a given fixed dataset, a critical challenge is to select an appropriate model (e.g. set of variables/features). Model selection aims to strike a balance between the goodness of fit and model complexity, and thus to gain reliable predictive power. In the first part of this paper, we study a penalized model selection technique that achieves the optimal possible prediction loss (referred to as the limit of learning) offered by a set of candidate models (model class). We prove that the proposed procedure asymptotically approaches the limit of learning for a wide variety of model classes, loss functions, and high dimensions (i.e. the models' complexity can grow with data size). As we will show, allowing high dimensionality is necessary if “all models are wrong” and data are sequentially observed (which often happens in practice). A reliable learning procedure naturally demands “bigger models” to accommodate “bigger data”. In other words, a model class of increasing complexity levels helps to search/track the best model as data size increases. In the second part of this paper, we propose a computationally efficient algorithm that sequentially expands the model complexity, in order to enhance selection stability and reduce computation cost. We demonstrate the proposed methodology by applying it to generalized linear models and neural networks, where the variables/neurons are reliably selected “on the fly”. In particular, the proposed procedure has more predictive power, and less computational cost compared to various types of cross validations.

Index Terms

Cross-validation; Computational efficiency; Ensemble learning; Feature selection; Graph-based expert learning; High dimensionality; Limit of learning; Model expansion.

I. INTRODUCTION

How much knowledge can we learn from a given set of data? Statistical modeling is a simplification of real world complexity. It can be used to learn the key *patterns* or *relationships* from available data and to predict the future data. In order to model the data, typically the first step in data analysts is to narrow the scope by specifying a set of candidate parametric models (referred to as model class). The model class can be determined by exploratory studies or scientific reasoning. For data with specific types and sizes, each postulated model may have its own advantages. In the second step, data analysts estimate the parameters and “goodness of fit” of each candidate model. An illustration of a typical learning procedure is plotted in Fig. 1, where the true data generating model may or may not be included in the model class. Simply selecting the model with the best fitting performance usually leads to suboptimal results. For example, the largest model always fits the best in a nested model class. But too large a model can lead to inflated variance and thus severe overfitting. Therefore, the third step is to apply a model selection procedure. State-of-art selection procedure can be roughly categorized into two classes, the penalized selection and cross-validation. We shall elaborate on those in the next section.

Example 1 (Generalized linear models). *In a generalized linear model (GLM), each response variable Y is assumed to be generated from a particular distribution (e.g. Gaussian, Binomial, Poisson, Gamma), with its mean μ linked with potential covariates X_1, X_2, \dots through $E_*(Y) = \mu = g^{-1}(\beta_1 X_1 + \beta_2 X_2 + \dots)$. In this example, data $Z = [Y, X_1, X_2, \dots]^T$, unknown parameters are $\theta = [\beta_1, \beta_2, \dots]^T$, and models are subsets of $\{\beta_1, \beta_2, \dots\}$. We may be interested in the most appropriate distribution form (pattern) as well as the most significant variables X_j 's (relationships).*

Example 2 (Neural networks). *In establishing a neural network (NN) model, we need to choose the number of neurons and hidden layers, activation function (relationships), and the configuration of their connectivity (pattern).*

J. Ding, E. Diao, J. Zhou, and V. Tarokh are with the John A. Paulson School of Engineering and Applied Sciences, Harvard University, Cambridge, MA, 02138 USA e-mail: (jieding@fas.harvard.edu).

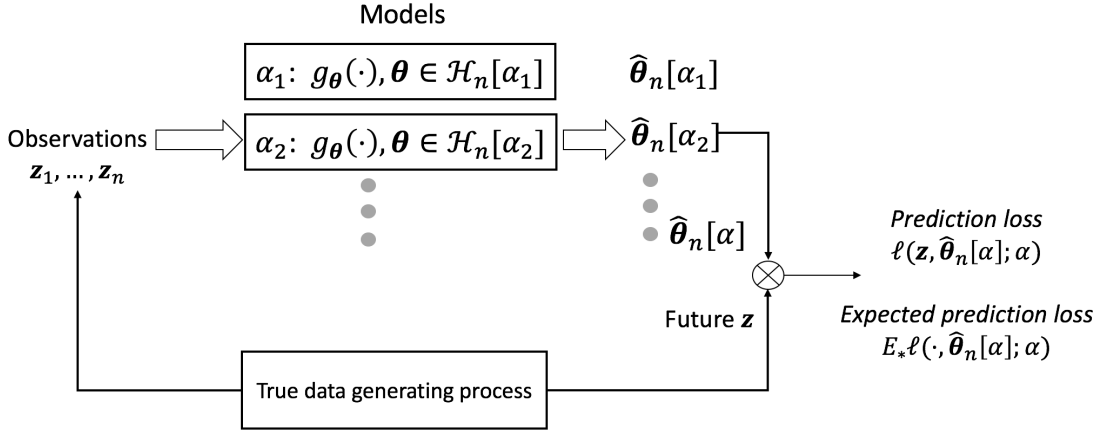


Fig. 1: Illustration of a typical learning procedure, where each candidate model α_j is trained in terms of $\hat{\theta}_n[\alpha_j]$ in its parameter space $\mathcal{H}_n[\alpha_j]$, and then evaluated for future data under some loss function ℓ .

In this example, data are similar to that of the above example, and unknown parameters are the weights on connected edges. Clearly, with larger number of neurons and connections, more complex functional relationships can be modeled. But selecting models with too large of dimensions may result in overfitting and more computational complexity.

How can we quantify the theoretical limits of learning procedures? We first introduce the expected prediction loss that quantifies the predictive power of each candidate model.

Definition 1 (Expected prediction loss). The loss function for each data size n and $\alpha \in \mathcal{A}_n$ (model class) is a map $l_n : \mathcal{Z} \times \mathcal{H}_n[\alpha] \rightarrow \mathbb{R}$, usually written as $l_n(\mathbf{z}, \boldsymbol{\theta}; \alpha)$, where \mathcal{Z} is the data domain, $\mathcal{H}_n[\alpha]$ is the parameter space associated with model α , and α is included to emphasize the model under consideration. As Fig. 1 shows, for a loss function and a given dataset $\mathbf{z}_1, \dots, \mathbf{z}_n$, each candidate model α produces an estimator $\hat{\theta}_n[\alpha]$ (referred to as the minimum loss estimator) defined by

$$\hat{\theta}_n[\alpha] \triangleq \arg \min_{\boldsymbol{\theta} \in \mathcal{H}_n[\alpha]} \frac{1}{n} \sum_{i=1}^n l_n(\mathbf{z}_i, \boldsymbol{\theta}; \alpha). \quad (1)$$

Moreover, the expected prediction loss given by candidate model α , denoted by $\mathcal{L}_n(\alpha)$, is defined by

$$\mathcal{L}_n(\alpha) \triangleq E_* l_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) = \int_{\mathcal{Z}} p(\mathbf{z}) l_n(\mathbf{z}, \hat{\theta}_n[\alpha]; \alpha) d\mathbf{z}. \quad (2)$$

Here, E_* denotes the expectation with respect to the distribution of a future (unseen) random variable \mathbf{z} . We also define the risk by $\mathcal{R}_n[\alpha] = E_* \mathcal{L}_n[\alpha]$, where the expectation in $\mathcal{R}_n[\alpha]$ is taken with respect to the observed data.

Typically, \mathbf{z} consists of response \mathbf{y} and covariates \mathbf{x} , and only the entries of \mathbf{x} associated with α are involved in the evaluation of l_n . Throughout the paper, we consider loss functions $l_n(\dots)$ such that $\mathcal{L}_n[\alpha]$ is always nonnegative. A common choice is to use negative log-likelihood of model α minus that of the true data generating model. Table I lists some other loss functions widely used in machine learning. Based on Definition 1, a natural way to define the limit of learning is through the optimal prediction loss that one could possibly achieve.

TABLE I: Some common loss functions in addition to negative log-likelihood

Name	quadratic	exponential	hinge	perceptron	logistic
Formula	$(y - \boldsymbol{\theta}^T \mathbf{x})^2$	$e^{-y \boldsymbol{\theta}^T \mathbf{x}}$	$\max\{0, 1 - y \boldsymbol{\theta}^T \mathbf{x}\}$	$\max\{0, -y \boldsymbol{\theta}^T \mathbf{x}\}$	$\log(1 + e^{-y \boldsymbol{\theta}^T \mathbf{x}})$
Domain	$y \in \mathbb{R}$	$y \in \mathbb{R}$	$y \in \mathbb{R}$	$y \in \mathbb{R}$	$y \in \{0, 1\}$

Definition 2 (Limit of learning). For a given data (of size n) and model class \mathcal{A}_n , the limit of learning (LoL) is defined as $\arg \min_{\alpha \in \mathcal{A}_n} \mathcal{L}_n(\alpha)$, the optimal expected prediction loss offered by candidate models.

We note that the LoL is associated with three key elements: data, loss function, and model class. Motivated by the original derivation of Akaike information criterion (AIC) [1], [2] and Takeuchi's information criterion (TIC) [3], we propose a penalized selection procedure and prove that it can approach the LoL under reasonable assumptions. Those assumptions allow a wide variety of loss functions, model classes (i.e. nested, non-overlapping or partially-overlapping), and high dimensionality (i.e. the models' complexity can grow with data size). Our theoretical results extend the classical statistical theory on AIC for linear (fixed-design) regression models. Moreover, we also review the conceptual and technical connections between cross validation and information theoretical criteria. In particular, we show that the proposed procedure can be much more computationally efficient than cross validation (with the same level of predictive power).

Why is it necessary to consider a high dimensional model class, e.g., to let the number of candidate models or each model's complexity grow with data size? In the context of fixed design regression, technical discussions that address the question have been elaborated in [4], [5]. Here, we give an intuitive explanation for a general setting. We let $\theta_n^*[\alpha]$ denote the minimum loss parameter defined by

$$\theta_n^*[\alpha] \triangleq \arg \min_{\theta \in \mathcal{H}_n[\alpha]} E_* l_n(\cdot, \theta; \alpha). \quad (3)$$

We shall show in the appendix that $\mathcal{L}_n[\alpha]$ may be expressed as

$$\mathcal{L}_n[\alpha] = E_* l_n(\mathbf{z}, \theta_n^*[\alpha]; \alpha) + \frac{1}{2} \left\| \hat{\theta}_n[\alpha] - \theta_n^*[\alpha] \right\|_{V_n(\theta_n^*[\alpha])}^2 \times \{1 + o_p(1)\} \quad (4)$$

under some regularity conditions. The expected prediction loss consists of two terms: the first being the bias term, and the second being the variance term. Large models tend to reduce the bias but inflate the variance (*overfitting*), while small models tend to reduce the variance but increase the bias (*underfitting*) for given a fixed dataset. Suppose that “all models are wrong”, meaning that the data generating model is not included in the model class. Usually, the bias is non-vanishing (with n) for a fixed model complexity (say d), and it is approximately a decreasing function of d ; while on the other hand, the variance vanishes at rate n^{-1} for a fixed d , and it is an increasing function of d . Suppose for example that the bias and variance terms are approximately $c_1 \gamma^{-d}$ and $c_2 d/n$, respectively, for some positive constants c_1, c_2, γ . Then the optimal d is at the order of $\log(n)$.

In view of the above arguments, as more data become available, the model complexity need to be enlarged in order to strike a balance between bias and variance (or *approach the LoL*). We refer to this phenomena as “bigger models for bigger data” (BMBD). To illustrate, we generated $n = 100, 200$ data from a logistic regression model, where coefficients are $\beta_i = 10/i$ and covariates x_i 's are independent standard Gaussian (for $i = 1, \dots, 100$). We consider the nested model class $\mathcal{A}_n = \{\{1\}, \{1, 2\}, \dots, \{1, 2, \dots, 50\}\}$. We summarize the results in Fig. 2. As model complexity increases, the “goodness of fit” (measured by in-sample loss) improves (Fig. 2a), while the “predictive power” (measured by expected prediction loss) first improves and then deteriorates after some “optimal dimension” (Fig. 2b). Moreover, the optimal dimension becomes larger as sample size increases. It means that better fitting does not mean better predictive power, and large data sizes requires the search over a larger model class.

The outline of the paper is given as follows. In Section II, we propose a computational efficient method that determines the most appropriate learning model as more data become available. We prove that the LoL can be asymptotically approached under some regularity assumptions. Based on the proposed method, we further propose a sequential model expansion strategy that aims to facilitate interpretability of learning. As data sequentially arrives, the selected model from our proposed method (and many other existing method such as cross validation) suffer from fluctuations (due to randomness). A philosophically appealing and computationally efficient way is to move from small model to larger models sequentially. In Section III, we propose a model expansion techniques inspired by “graph-based” online learning. In Section IV, we demonstrate the applications of the proposed methodology to generalized linear models and neural networks, in order to select the variables/neurons with optimal predictive power and low computational cost.

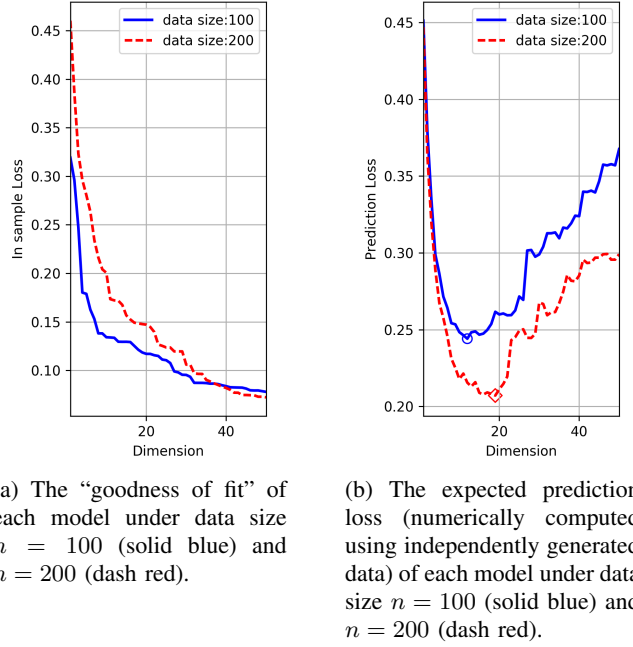


Fig. 2: Experiment showing BMBD phenomena.

II. LIMIT OF LEARNING

A. Notation

Let \mathcal{A}_n , α , $d_n[\alpha]$, $\mathcal{H}_n[\alpha] \subset \mathbb{R}^{d_n[\alpha]}$ denote respectively a set of candidate models, a candidate model, its dimension, its associated parameter space. Let $d_n \triangleq \max_{\alpha \in \mathcal{A}_n} d_n[\alpha]$ denote the dimension of the largest candidate model. We shall frequently use subscript n to emphasize the dependency on n , and include an α in the arguments of many variables or functions in order to emphasize their dependency on the model (and parameters space) under consideration. For a measurable function $f(\cdot)$, we define $E_n f(\cdot) = n^{-1} \sum_{i=1}^n f(z_i)$. For example, $E_n l_n(\cdot, \theta; \alpha) = n^{-1} \sum_{i=1}^n l_n(z_i, \theta; \alpha)$. We let $\psi_n(z, \theta; \alpha) \triangleq \nabla_{\theta} l_n(z, \theta; \alpha)$, and $\nabla_{\theta} \psi_n(\cdot, \theta; \alpha) \triangleq \nabla_{\theta}^2 l_n(\cdot, \theta; \alpha)$, which are respectively measurable vector-valued and matrix-valued functions of θ . We define the matrices

$$\begin{aligned} V_n(\theta; \alpha) &\triangleq E_* \nabla_{\theta} \psi_n(\cdot, \theta; \alpha) \\ J_n(\theta; \alpha) &\triangleq E_* \{ \psi_n(\cdot, \theta; \alpha) \times \psi_n(\cdot, \theta; \alpha)^T \} \end{aligned}$$

Recall the definition of $\mathcal{L}_n[\alpha]$. Its sample analog (also referred to as the *in-sample loss*) is defined by $\hat{\mathcal{L}}_n[\alpha] \triangleq E_n l_n(\cdot, \hat{\theta}_n[\alpha]; \alpha)$. Similarly, we define

$$\begin{aligned} \hat{V}_n(\theta; \alpha) &\triangleq E_n \nabla_{\theta} \psi_n(\cdot, \theta; \alpha) \\ \hat{J}_n(\theta; \alpha) &\triangleq E_n \{ \psi_n(\cdot, \theta; \alpha) \times \psi_n(\cdot, \theta; \alpha)^T \} \end{aligned}$$

For a sequence j_1, \dots, j_T , we define its size, $\text{size}(j_1, \dots, j_T)$, as the number of t 's such that $j_t \neq j_{t+1}$. Let $\text{int}(S)$ denote the interior of a set S .

Throughout the paper, the vectors are arranged in column and marked in bold. Let $\|\cdot\|$ denote Euclidean norm of a vector or spectral norm of matrix. For any vector $c \in \mathbb{R}^d$ ($d \in \mathbb{N}$) and scalar $r > 0$, let $B(c, r) \triangleq \{x \in \mathbb{R}^d : \|x - c\| \leq r\}$. For a positive semidefinite matrix V and a vector x of the same dimension, we shall abbreviate $x^T V x$ as $\|x\|_V^2$. For a given probability measure P_* and a measurable function m , let $\|m\|_{P_*} \triangleq (E_* m^2)^{1/2}$ denote the $L_2(P_*)$ -norm. Let $\text{eig}_{\min}(V)$ (resp. $\text{eig}_{\max}(V)$) denote the smallest (resp. largest) eigenvalue of a symmetric matrix V . For a sequence of scalar random variables f_n , we write $f_n = o_p(1)$ if $\lim_{n \rightarrow \infty} f_n = 0$ in probability, and $f_n = O_p(1)$, if it is stochastically bounded. For a fixed measurable vector-valued function \mathbf{f} , we define

$$\mathbb{G}_n \mathbf{f} \triangleq \sqrt{n} (E_n - E_*) \mathbf{f},$$

the empirical process evaluated at \mathbf{f} . For $a, b \in \mathbb{R}$, we write $a \lesssim b$ if $a \leq cb$ for a universal constant c . For a vector \mathbf{a} or a vector-valued function \mathbf{f} , we let a_i or f_i denote the i th component.

We use \rightarrow and \rightarrow_p to respectively denote the deterministic and in probability convergences. Unless stated explicitly, all the limits throughout the paper are with respect to $n \rightarrow \infty$ where n is the sample size.

B. Approaching the LoL – Selection Procedure

To obtain the optimal predictive power, an appropriate model selection procedure is necessary to strike a balance between the *goodness of fit*, and *model complexity* based on the observed data. The basic idea of penalized selection is to impose an additive penalty term to the in-sample loss (i.e. goodness of fit), so that larger models are more penalized. In this paper, we follow the aphorism that “all models are wrong”, and assume that the model class under consideration is mis-specified.

Definition 3 (Efficient learning). Our goal is to select $\hat{\alpha}_n \in \mathcal{A}_n$ that is asymptotically efficient, in the sense that

$$\frac{\mathcal{L}_n[\hat{\alpha}_n]}{\min_{\alpha \in \mathcal{A}_n} \mathcal{L}_n[\alpha]} \rightarrow_p 1 \quad (5)$$

as $n \rightarrow \infty$.

Note that this requirement is weaker than selecting the exact optimal model $\arg \min_{\alpha \in \mathcal{A}_n} \mathcal{L}_n[\alpha]$. Similar definition has been adopted in the study of the optimality of AIC in the context of autoregressive order selection [6] and variable selection in linear regression models [7]. We propose to use the following penalized model selection procedure, which generalizes TIC from negative log-likelihood to general loss functions.

Generalized TIC (GTIC) procedure: Given data $\mathbf{z}_1, \dots, \mathbf{z}_n$ and a specified model class \mathcal{A}_n . We select a model $\hat{\alpha} \in \mathcal{A}_n$ in the following way: 1) for each $\alpha \in \mathcal{A}_n$, find the minimal loss estimator $\hat{\theta}_n[\alpha]$ defined in (1), and record the minimum as $\hat{\mathcal{L}}_n[\alpha]$; 2) select $\hat{\alpha} = \arg \min_{\alpha \in \mathcal{A}_n} \mathcal{L}_t^c[\alpha]$, where

$$\mathcal{L}_t^c[\alpha] \triangleq \hat{\mathcal{L}}_n[\alpha] + n^{-1} \text{tr} \left\{ \hat{V}_n(\hat{\theta}_n[\alpha]; \alpha)^{-1} \hat{J}_n(\hat{\theta}_n[\alpha]; \alpha) \right\}. \quad (6)$$

Remark 1 (TIC and GTIC). The quantity $\mathcal{L}_t^c[\alpha]$, also referred to as the corrected prediction loss, can be calculated from data, and it serves as a surrogate of the expected prediction loss $\mathcal{L}_n[\alpha]$ that is usually not analytically computable. The in-sample loss $\hat{\mathcal{L}}_n[\alpha]$ cannot be directly used as an approximation for $\mathcal{L}_n[\alpha]$, because it uses the sample approximation twice: once in the estimation of θ_n^* , and then in the approximation of $E_* l_n(\cdot, \theta; \alpha)$ using $E_n l_n(\cdot, \theta; \alpha)$ (the law of large numbers). For example, in a nested model class, the largest model always has the least $\hat{\mathcal{L}}_n[\alpha]$ (i.e. fits data the best). But as we discussed in the introduction, $\mathcal{L}_n[\alpha]$ is typically decreasing first and then increasing as the dimension increases.

TIC was heuristically derived as an alternative of AIC, also from an information theoretical point of view (using Kullback-Leibler divergence), assuming model mis-specification. It does not appear to be widely appreciated nor used [8] compared with other information criteria such as AIC or Bayesian information criterion (BIC) [9]. In terms of provable asymptotic performance, only AIC is known to be asymptotically efficient for variable selection in linear regression models [7] and autoregressive order selection in linear time series models [6], [10], in mis-specified settings. It can be shown that TIC is close to AIC in linear models (and exactly the same for fixed-design regressions). Conceptually, TIC was proposed as a surrogate for AIC in general mis-specified settings. However, the theoretical optimality of AIC and TIC in the general context remains unknown.

Why should TIC be preferred over AIC in nonlinear models in general? Intuitively speaking, TIC has the potential of exploiting the nonlinearity while AIC does not. Recall our Example 2 in the introduction, with loss being the negative log-likelihood. It is well known from machine learning practices that neural network structures play a key role in effective prediction. However, information criteria such as AIC impose the same amount of penalty as long as the number of neurons remains the same, regardless of how neurons are configured.

In this paper, we extend the scope of allowable loss functions, and theoretically justify the use of GTIC (and thus TIC). Under some regularity conditions (elaborated in the Appendix), we shall prove that the $\hat{\alpha}$ selected by the GTIC procedure is efficient (in the sense of Definition 3). It is formally stated as a theorem in Subsection II-D. Our theoretical results extend some existing statistical theories on AIC for linear models. We note that the technical analysis of high dimensional (non)linear model classes is highly nontrivial. We will develop some new technical tools in the Appendix, which may be interesting on their own rights.

C. Related Works

A wide variety of model selection techniques have been proposed in the last fifty years, motivated by different viewpoints and justified under various circumstances. State-of-art methods can be roughly categorized into two classes, the penalized selection and cross-validation. Examples are final prediction error criterion [11], AIC [1], [2], BIC [9] (and its finite sample counterpart Bayes factor [12]), minimum description length criterion [13], Hannan and Quinn criterion [14], predictive minimum description length criterion [15], [16], C_p method [17], generalized information criterion (GIC_{λ_n}) with $\lambda_n \rightarrow \infty$ [4], [18], [19], generalized cross-validation method (GCV) [20], and the bridge criterion (BC) [5], [21].

Is Cross-Validation Really The Best Choice?

The basic idea of cross-validation (CV) is to split the data into two parts, one for training and one for testing. The model with the best testing performance is selected, in the hope that it will perform well for future data as well. It is a common practice to apply 10-fold CV, 5-fold CV, 2-fold CV, or 30%-for-testing. In general, the advantages of CV method are its stability and easy implementation. However, it has been shown that only the delete- d CV method with $\lim_{n \rightarrow \infty} d/n = 1$ [22]–[25], or the delete-1 CV method [26] (or leave-one-out, LOO) can exhibit asymptotic (large sample) optimality. In fact, the former CV exhibits the same asymptotic behavior as BIC, which is typically consistent in a well-specified model class (i.e. it contains the true data generating model), but is suboptimal in a mis-specified model class. The latter CV is shown to be asymptotically equivalent to AIC and GCV if $d_n[\alpha] = o(n)$ [26], which is asymptotically efficient in a mis-specified model class, but usually overfits in a well-specified model class. We refer to [4], [5], [21], [27], [28] for more detailed discussions on the discrepancy and reconciliation of the two types of selection criteria.

In other words, common folklores that advocate the use of k -fold or 30%-for-testing CV are all asymptotically suboptimal (in the sense of Definition 3), even in linear regression models [4]. Since the only optimal CV is LOO-type (in mis-specified settings), it is more appealing to apply AIC or TIC that gives the same asymptotic performance and *significantly reduces the computational complexity* by n times. For general (mis-specified) nonlinear model class, we shall prove that GTIC procedure approaches the LoL. While the asymptotic performance of LOO is not clear in that case, it is typically more complex to implement. To demonstrate that, we shall provide some experimental studies in the Appendix. As a result, the GTIC procedure can be a promising competitor of various types of standard CVs adopted in practice.

D. Asymptotic Analysis of the GTIC Procedure

We need the following assumptions for asymptotic analysis.

Assumption 1. Data $\mathbf{Z}_i, i = 1, \dots, n$ are independent and identically distributed (i.i.d.).

Assumption 2. For each model $\alpha \in \mathcal{A}_n$, $\boldsymbol{\theta}_n^*[\alpha]$ (as was defined in (3)) is in the interior of the compact parameter space $\mathcal{H}_n[\alpha]$, and for all $\varepsilon > 0$ we have

$$\liminf_{n \rightarrow \infty} \inf_{\alpha \in \mathcal{A}_n} \left(\inf_{\boldsymbol{\theta} \in \mathcal{H}_n[\alpha]: \|\boldsymbol{\theta} - \boldsymbol{\theta}_n^*[\alpha]\| \geq \varepsilon} E_* \ell_n(\cdot, \boldsymbol{\theta}; \alpha) - E_* \ell_n(\cdot, \boldsymbol{\theta}_n^*[\alpha]; \alpha) \right) \geq \eta_\varepsilon$$

for some constant $\eta_\varepsilon > 0$ that depends only on ε . Moreover, we have

$$\sup_{\alpha \in \mathcal{A}_n} \sup_{\boldsymbol{\theta} \in \mathcal{H}_n[\alpha]} \left| E_n \ell_n(\cdot, \boldsymbol{\theta}; \alpha) - E_* \ell_n(\cdot, \boldsymbol{\theta}; \alpha) \right| \rightarrow_p 0,$$

as $n \rightarrow \infty$, and $\ell_n(\cdot, \boldsymbol{\theta}_n^*[\alpha]; \alpha)$ is twice differentiable in $\text{int}(\mathcal{Z})$ for all $n, \alpha \in \mathcal{A}_n$.

Assumption 3. There exist constants $\tau \in (0, 0.5)$ and $\delta > 0$ such that for all $\alpha \in \mathcal{A}_n$,

$$\begin{aligned} & \sup_{\alpha \in \mathcal{A}_n} \sup_{\boldsymbol{\theta} \in \mathcal{H}_n[\alpha] \cap B(\boldsymbol{\theta}_n^*[\alpha], \delta)} n^\tau \left\| E_n \psi_n(\cdot, \boldsymbol{\theta}; \alpha) - E_* \psi_n(\cdot, \boldsymbol{\theta}; \alpha) \right\| \\ &= O_p(1). \end{aligned}$$

Additionally, the map $\theta \mapsto E_*\psi_n(\cdot, \theta; \alpha)$ is differentiable at $\theta \in \text{int}(\mathcal{H}_n[\alpha])$ for all n and $\alpha \in \mathcal{A}_n$.

Assumption 4. There exist constants $c_1, c_2 > 0$ such that

$$\liminf_{n \rightarrow \infty} \min_{\alpha \in \mathcal{A}_n} \text{eig}_{\min}(V_n(\theta_n^*; \alpha)) \geq c_1, \quad \limsup_{n \rightarrow \infty} \max_{\alpha \in \mathcal{A}_n} \text{eig}_{\max}(V_n(\theta_n^*; \alpha)) \leq c_2.$$

Assumption 5. There exist constants $r > 0, \gamma > 1$, and measurable functions $m[\alpha] : \mathcal{Z} \rightarrow \mathbb{R}^+ \cup \{0\}$, $z \mapsto m[\alpha](z)$ for each $\alpha \in \mathcal{A}_n$, such that for all n and $\theta_1, \theta_2 \in B(\theta_n^*[\alpha], r)$,

$$\|\psi_n(z, \theta_1; \alpha) - \psi_n(z, \theta_2; \alpha)\| \leq m_n[\alpha](z) \|\theta_1 - \theta_2\|, \quad (7)$$

$$E_* m_n[\alpha] < \infty. \quad (8)$$

Moreover, we have

$$\begin{aligned} & \max \left\{ d_n^\gamma \text{card}(\mathcal{A}_n)^{\gamma/2}, d_n \sqrt{\log\{d_n \text{card}(\mathcal{A}_n)\}} \right\} \\ & \times n^{-\tau} \left\| \sup_{\alpha \in \mathcal{A}_n} m_n[\alpha] \right\|_{P_*} \rightarrow 0, \end{aligned} \quad (9)$$

and for all n ,

$$E_* \sup_{\theta \in \mathcal{H}_n[\alpha]} \|\psi_n(z, \theta_1; \alpha)\| < \infty. \quad (10)$$

Assumption 6. There exists a constant $\delta > 0$ such that for all $\alpha \in \mathcal{A}_n$,

$$\sup_{\alpha \in \mathcal{A}_n} \sup_{\theta \in \mathcal{H}_n[\alpha] \cap B(\theta_n^*[\alpha], \delta)} \|\hat{J}_n(\theta; \alpha) - J_n(\theta; \alpha)\| \rightarrow_p 0, \quad (11)$$

$$\sup_{\alpha \in \mathcal{A}_n} \sup_{\theta \in \mathcal{H}_n[\alpha] \cap B(\theta_n^*[\alpha], \delta)} \|\hat{V}_n(\theta; \alpha) - V_n(\theta; \alpha)\| \rightarrow_p 0, \quad (12)$$

$$\sup_{\alpha \in \mathcal{A}_n} \sup_{\theta, \theta' \in \mathcal{H}_n[\alpha] \cap B(\theta_n^*[\alpha], \delta)} \|V_n(\theta; \alpha) - V_n(\theta'; \alpha)\| \rightarrow_p 0. \quad (13)$$

We define

$$\mathbf{w}_n[\alpha] = \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_n(z_i, \theta_n^*[\alpha]; \alpha).$$

Clearly, $\mathbf{w}_n[\alpha]$ has zero mean and variance matrix $J_n(\theta_n^*[\alpha]; \alpha)$, and thus

$$E_* \|\mathbf{w}_n[\alpha]\|_{V_n(\theta_n^*[\alpha]; \alpha)^{-1}}^2 = \text{tr}\{V_n(\theta_n^*[\alpha]; \alpha)^{-1} J_n(\theta_n^*[\alpha]; \alpha)\}.$$

Assumption 7. Suppose that the following regularity conditions are satisfied.

$$\inf_{\alpha \in \mathcal{A}_n} n^{2\tau} \mathcal{R}_n[\alpha] \rightarrow \infty, \quad (14)$$

$$\sup_{\alpha \in \mathcal{A}_n} \frac{d_n[\alpha]}{n \mathcal{R}_n[\alpha]} \rightarrow 0. \quad (15)$$

Moreover, there exists a fixed constant $m_1 > 0$ such that

$$\begin{aligned} & \sum_{\alpha \in \mathcal{A}_n} (n \mathcal{R}_n[\alpha])^{-2m_1} \sum_{\alpha \in \mathcal{A}_n} E_* \{l_n(z_i, \theta_n^*[\alpha]; \alpha) - E_* l_n(z, \theta_n^*[\alpha]; \alpha)\}^{2m_1} \\ & \rightarrow 0, \end{aligned} \quad (16)$$

there exists a fixed constant $m_2 > 0$ such that

$$\begin{aligned} & \sum_{\alpha \in \mathcal{A}_n} (n \mathcal{R}_n[\alpha])^{-2m_2} \sum_{\alpha \in \mathcal{A}_n} E_* \left[\|\mathbf{w}_n[\alpha]\|_{V_n(\theta_n^*[\alpha]; \alpha)^{-1}}^2 \right. \\ & \left. - \text{tr}\{V_n(\theta_n^*[\alpha]; \alpha)^{-1} J_n(\theta_n^*[\alpha]; \alpha)\} \right]^{2m_2} \rightarrow 0, \end{aligned} \quad (17)$$

and there exists a fixed constant $m_3 > 0$ such that

$$\limsup_{n \rightarrow \infty} \sum_{\alpha \in \mathcal{A}_n} (n\mathcal{R}_n[\alpha])^{-m_3} \{E_* \|\mathbf{w}_n[\alpha]\|^{m_3} + E_* \|\mathbf{w}_n[\alpha]\|^{2m_3}\} < \infty. \quad (18)$$

Remark 2 (Intuitive explanations of each assumption). Assumption 1 is *reasonable for many practical applications*. In the context of regression analysis, it corresponds to the random design. In our technical proofs, it is possible to extend the assumption of i.i.d. to strong mixing which is more commonly assumed for time series data.

Assumption 2 is the counterpart of the “separated mode” and “uniform law of large numbers” conditions that have been commonly required in proving the consistency of maximum likelihood estimator for classical statistical models [29, Theorem 5.7]. The $\theta_n^*[\alpha]$ can be interpreted as the oracle optimum under model α , or the “projection” of true data generating model onto α .

Assumption 3 is a weaker statement to the central limit theorem (and its extension to Donsker classes) in a classical (non-high dimensional) setting. In our high dimensional setting, the assumption ensures that each projected model $\theta_n^*[\alpha]$ behaves regularly. It implicitly builds a relation between d_n (the dimension of the largest candidate models) and n .

Assumption 4 assumes that the second derivative of the expected prediction loss has bounded eigenvalues at the optimum $\theta_n^*[\alpha]$. It is useful because our asymptotic analysis requires “well-behaved” Taylor expansion up to the second order.

Assumption 5 is a Lipschitz-type condition. Similar but simpler *form* has been used in classical analysis of asymptotic normality [29, Theorem 5.21]. We note that the condition (9) explicitly requires that the largest dimension d_n and the candidate size $\text{card}(\mathcal{A}_n)$ *cannot* grow too fast.

Assumption 6 requires that the sample analogs of matrices $J_n(\theta; \alpha)$ and $V_n(\theta; \alpha)$ are asymptotically close to the truth (in spectral norm) in a neighborhood of $\theta_n^*[\alpha]$. In the classical setting, it is guaranteed by the law of large numbers (applied to each matrix element). Assumption 6 also requires the continuity of $V_n(\theta; \alpha)$ in a neighborhood of $\theta_n^*[\alpha]$.

In Assumption 7, the assumptions in (14), (15), and (18) mean that the model class is “quite mis-specified”, since they require that the risks $\mathcal{R}_n[\alpha]$ (which are deterministic numbers) cannot be too small. The assumptions in equalities (16) and (17) are central moment constraints that control both the regularity of loss functions and the speed of growth of the model class.

Theorem 1. Suppose that Assumptions 1-7 hold. Then the $\tilde{\alpha}$ selected by GTIC procedure is *efficient* (in the sense of Definition 3).

Remark 3 (Sketch of Technical Ideas). Classical asymptotic analysis typically relies on a type of uniform convergence of empirical process around $\theta_n^*[\alpha]$. Because our functions are vector valued with dimension depending on data size, we cannot directly use state-of-art technical tools such as [29, Theorem 19.28]. The classical proof by White [30] (in proving asymptotic normality in mis-specified class) cannot be directly adapted, either, for parameter spaces that depend on n . Some new technical tools are needed in our proof. Our ideas of proof are sketched below.

We first prove that $\hat{\theta}_n[\alpha]$ is n^τ -consistent (instead of the classical \sqrt{n} -consistency). We then prove the first key *lemma*, Lemma 6, that states a type of local uniform convergence. Note that its proof is nontrivial as both the empirical process and $\hat{\theta}_n$ depend on the same observed data. Our technical tools resemble those for proving a Donsker class, but the major difference is that our model dimensions depend on n . We then prove the second key lemma, *Lemma* 7. It directly leads to the asymptotic normality of maximum likelihood estimators in the classical setting. It is somewhat interesting to see that the proof of Lemma 7 does not require the \sqrt{n} -consistency of $\hat{\theta}_n[\alpha]$ (which may not hold in high dimensional settings).

III. SEQUENTIAL MODEL EXPANSION

As explained in the introduction, in terms of predictive power, a model in a mis-specified model class could be determined to be unnecessarily large, ~~to~~ suitable, and finally to inadequately small, depending on specific data size (see Fig. 2). A realistic learning procedure thus requires models of different complexity levels as more data become available. For example, new-born children can only do simple tasks, but gradually they can succeed at doing more complicated tasks. Biological evidence has been found in a recent study on how children reason about

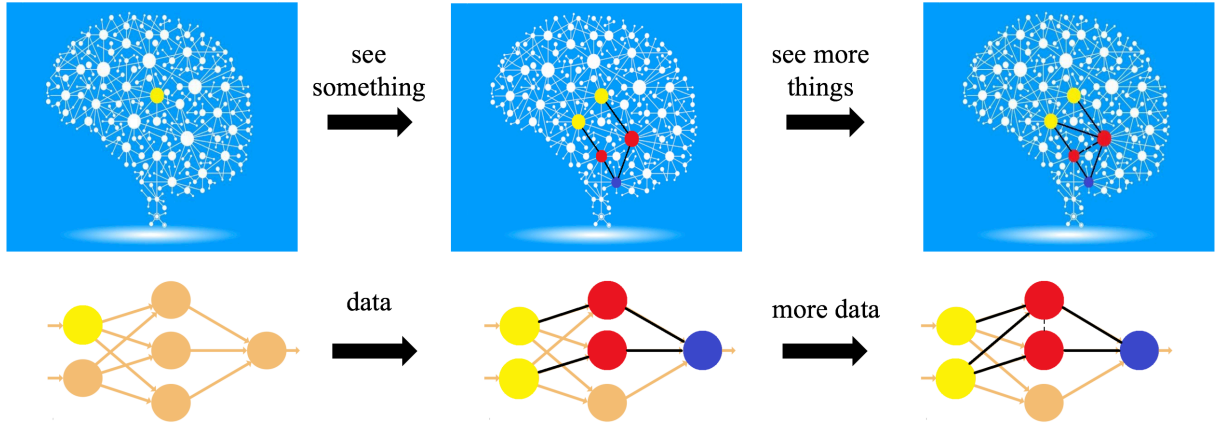


Fig. 3: Illustration of neurons being activated to adapt to the currently observed data in an online fashion, similar to a human brains.

other people’s mind, where cognitive scientists found that the same brain regions of children get more functionally specified as they get older [31]. Using the example of neural networks, the idea is illustrated in Fig. 3.

We have addressed the challenge of selecting an efficient model for a given number of observations. In many practical situations, data are sequentially observed. A straightforward model selection is to repeatedly apply GTIC procedure upon each newly arrived data. However, in a sequential setting, the following issue naturally arises:

Suppose that we successively select a model and use it to predict at each time step. The path of the historically selected models may fluctuate a lot (which will be illustrated in our numerical experiments). Instead, it is more appealing (either philosophically or computationally) to force the selected models to evolve continuously: from the smaller to the larger (in dimension).

To address the above challenge, we further propose an algorithm (Algo. 2) for **on-the-fly** prediction, which is **“efficient”** not only statistically (Definition 3), but also computationally. Specifically, we propose a predictor that has cumulative expected prediction loss (over time) close to the following optimum benchmark:

$$\min_{\text{size}(j_1, \dots, j_T) \leq k, j_{t-1} \leq j_t \leq j_{t-1} + 1} \sum_{t=1}^T \mathcal{L}_n[\alpha_{j_t}].$$

In other words, the minimization is taken over all tuples (j_1, \dots, j_T) that have at most k switches and that are restricted to the chain $1 \rightarrow 2 \rightarrow \dots$. For example, $(j_1, \dots, j_5) = (1, 2, 2, 3, 3)$. Throughout this section, we shall use T (instead of the previously used n) to denote sample size, and subscript t as the data index, in order to emphasize the sequential setting.

A. **Tracking the Best Expert with Graphical Constraints**

In this subsection, we propose a novel graph-based expert tracking technique that motivates our algorithm in the following subsection. The discussion may be interesting on its own right, as it includes the state-of-art expert tracking framework as a special case (when the underlying graph is complete).

Suppose there are N experts. At each discrete time step $t = 1, 2, \dots, T$, each expert gives its prediction, after which the environment reveals the truth $z_t \in \mathcal{Y}$. The performance of each prediction is measured by a loss $l : \{1, 2, \dots, N\} \times \mathcal{Y} \rightarrow \mathbb{R}$. In this subsection, with a slight abuse of notation, we shall also use l to denote loss functions in the context of online learning. The smaller the loss is, the better the prediction is. In order to aggregate all the predictions that the experts make, we maintain a weight for each expert, and update them upon the arrival of each new data point based on the qualities of the predictions. We denote the weight for expert $i \in \{1, \dots, N\}$ at time t as $w_{i,t}$, and the normalized version as $W_{i,t}$. The goal is to optimally update the weights for better prediction, which is measured by the cumulative loss minus the best achievable (benchmark) loss. This measure is often called “regret” in the online learning literature.

Algorithm 1 Tracking the best expert with graphical transitional constraints**input** Parameters $\eta > 0$, $0 < \kappa < 1/D$ **output** $\mathbf{p}_t = [p_{t,1}, \dots, p_{t,K}]^T$ (predictive distribution over the active models) for each $t = 1, \dots, T$ 1: Initialize $w_{1,0} = 1$ $w_{i,0} = 0$ for all $i \in \{2, \dots, N\}$ 2: **for** $t = 1 \rightarrow T$ **do**3: Calculate the predictive distribution $p_{i,t} = w_{i,t-1} / \sum_{j=1}^N w_{j,t-1}$, for each $i \in \{1, \dots, N\}$ 4: Read \mathbf{z}_t , and compute $v_{i,t} = w_{i,t-1} \exp(-\eta \cdot l(i, \mathbf{z}_t))$, for each $i \in \{1, \dots, N\}$ 5: Let $w_{i,t} = \kappa \sum_{j=1}^N \beta_{ji} v_{j,t} + (1 - \kappa \beta_i) v_{i,t}$, for each $i \in \{1, \dots, N\}$ 6: **end for**

If the benchmark in the regret is defined as the minimum cumulative loss achieved by a single expert, then the exponential re-weighting algorithm can be applied easily. In many cases the best performing expert can be different from one time segment to another, for which the fixed share algorithm can be a good solution with guaranteed regret bound. We consider the following problem setting. The best performing expert is restricted to switch according to a *directed graph*, $G = (V, E)$ (without self-loops), with $V = \{1, \dots, N\}$ denoting the set of nodes (representing experts) and E denoting the set of directed edges. At each time point, the best performing expert can either stay the same or jump to another node which is directly connected from the current node. Let $\beta_{ij} = \mathbf{1}\{\exists(i, j) \in E\}$, which is 1 if there is a directed edge (i, j) on the graph, and 0 otherwise. Let $\beta_i = \sum_{j=1}^N \beta_{ij}$, which is the out-degree of the node i . In addition, we assume that $\max_{i \in \{1, \dots, N\}} \beta_i \leq D$, where $0 < D < N$.

We propose the following algorithm to follow the best expert with the graphical transitional constraints. We use a special prior $w_{i,0}$ here to motivate content in the next subsection. It is not difficult to extend our discussion to more general priors here. The classical fixed-share algorithm [32, Chapter 5] can be seen as a special case when the graph is complete.

Theorem 2. Suppose the loss function range is in $[0, 1]$. For all $T \geq 1$, the output of the algorithm in Algo. 1 satisfies

$$\begin{aligned} & \sum_{t=1}^T \left(\sum_{i=1}^N l(i, \mathbf{z}_t) p_{i,t} - l(i_t, \mathbf{z}_t) \right) \\ & \leq \frac{1}{\eta} (T - k - 1) \log \frac{1}{1 - \kappa D} + \frac{1}{\eta} k \log \frac{1}{\kappa} + \eta \frac{T}{8} \end{aligned}$$

for all expert sequence (i_1, i_2, \dots, i_T) and all observation sequence $(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T)$, given that (i_1, i_2, \dots, i_T) form directed paths or self loops in graph G and $\text{size}(i_1, i_2, \dots, i_T) \leq k$. Here, the left hand side of the above inequality is referred to as regret.

In particular, the above regret bound has minimum $\sqrt{TS/2}$ which is achieved when $\kappa = k/\{(T-1)D\}$ and $\eta = \sqrt{8S/T}$, where

$$S = (T-1)H(k/(T-1)) + k \log D,$$

and $H(\cdot)$ is the binary entropy function defined by $H(x) \triangleq -x \log x - (1-x) \log(1-x)$ for $x \in (0, 1)$, $H(0) = H(1) = 0$.

It is interesting to see that with graphical constraint, the regret bound does not depend on N , but on the maximum out-degree D instead. Thus, the bound can be tight even when N grows exponentially in T , as long as $D \ll N$ (i.e. sparse graph).

B. Algorithm for Sequential Model Expansion

The new online learning theory proposed in the last subsection is motivated by graph-based expert tracking. Intuitively speaking, instead of using the exponentially updated weights directly, a “tax” of rate $\kappa \beta_i$ is imposed on each expert i , and then the total tax revenue is redistributed equally among experts according to a directed graph. In that way, the experts are encouraged to rejuvenate their past performance and “start a new life”, so that we can track the best expert in different time epochs. The classical fixed-share algorithm [32, Chapter 5] is a special case

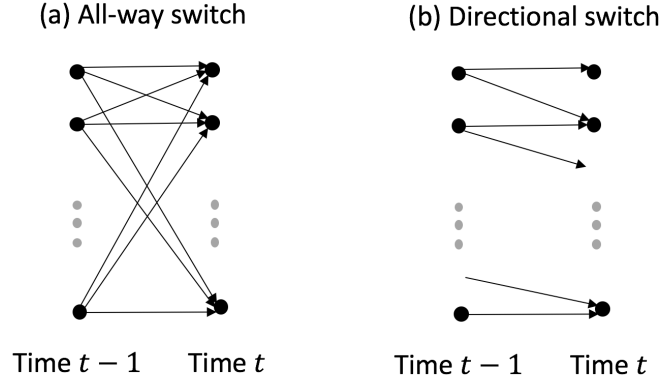


Fig. 4: Illustration of the state-of-art and our new way of redistributing the share of weights

when $\beta_{ij} = 1$ for all i, j , illustrated in Fig. 4(a). Our algorithm in this subsection is motivated by the particular *path graph* $1 \rightarrow 2 \rightarrow \dots \rightarrow N$, where $1, 2, \dots, N$ index the model class. In other words, we redistribute the revenue in a directional way, thus encouraging the experts to switch in a chain. The update rule is illustrated by Fig. 4(b).

Our algorithm for sequential model expansion is summarized in Algo. 2, where each candidate model and its corrected prediction loss can be regarded respectively as an expert and loss at each time. The labeling of models $\alpha_1, \alpha_2, \dots$ is generally in the ascending order of their dimensions. To further reduce the computational cost, we maintain only an active subset (of size K) instead of all the candidate models at each time. The active subset starts from $\{\alpha_1, \dots, \alpha_K\}$; it switches to $\{\alpha_2, \dots, \alpha_{K+1}\}$ when the weight of the smallest model α_1 becomes small and that of the largest model α_K becomes large; it continues to switch upon the aggregation of data.

The output of Algo. 1 is a predictive distribution over the active models. It can be used in the following two ways in practice: 1) we randomly draw a model according to the predictive distribution and use the predictor of that model, or 2) we weighted average the predictors of each model according to the predictive distribution. This can be regarded as a specific ensemble learning (or model averaging) method. The following Proposition 1 shows that with appropriate learning parameters, the average predictive performance of our algorithm is asymptotically close to the average of a series of truly optimal models (i.e. optimal model expansion) allowing moderately many switches.

Algorithm 2 Sequential model expansion using GTIC-corrected loss (GTIC-sequential)

input $\{z_t : t = 1, \dots, T\}$, $\eta > 0$, $\zeta \in [0, 1]$, $w_{0,1} = 1, w_{0,2} = \dots = w_{0,K} = 0$, candidate models $\mathcal{A}_T = \{\alpha_1, \alpha_2, \dots, \alpha_{\text{card}(\mathcal{A}_T)}\}$, $s = 0$ ($\alpha_{s+1}, \dots, \alpha_{s+K}$ are the maintained active subsets of models), $K \in \mathbb{N}$, threshold $\rho \in [0, 1]$

output $p_t = [p_{t,1}, \dots, p_{t,K}]^T$ (predictive distribution over the active models) for each $t = 1, \dots, T$

1: **for** $t = 1 \rightarrow n$ **do**

2: Obtain z_t and compute $v_{t,k} = w_{t-1,k} \exp\{-\eta \mathcal{L}_t^c[\alpha_{s+k}]\}$ for each $k = 1, \dots, K$, where $\mathcal{L}_t^c[\alpha]$ was defined using the data z_1, \dots, z_t .

3: Let

$$w_{t,k} = \begin{cases} (1 - \zeta) v_{t,k} & \text{if } k = 1 \\ (1 - \zeta) v_{t,k} + \zeta v_{t,k-1} & \text{if } 1 < k < K \\ v_{t,k} + \zeta v_{t,k-1} & \text{if } k = K \end{cases}$$

4: Let $p_{t,k} = (\sum_{k=1}^K w_{t,k})^{-1} w_{t,k}$, $k = 1, \dots, K$

5: **if** $p_{t,1} \leq \rho$ and $p_{t,K} \geq 1 - \rho$ and $s + K \leq \text{card}(\mathcal{A}_T)$ **then**

6: Let $s = s + 1$

7: Let $w_{t,k} = w_{t,k'}$, where $k = 1, \dots, K$ and $k' = (k + 1 \bmod K)$ (relabeling the active models)

8: **end if**

9: **end for**

Proposition 1. Suppose that Assumptions 1-7 hold, and that $\sup_{1 \leq t \leq T} \sup_{\alpha \in \mathcal{A}_n} \mathcal{L}_t^c[\alpha_{s+k}] < c$ for some fixed

constant $c > 0$ **almost surely**. Suppose that the lines 5-8 are removed from **Algo. 2**, then its output satisfies

$$\frac{1}{T} \left(\sum_{t=1}^T \sum_{i=1}^{\text{card}(\mathcal{A}_T)} p_{i,t} \mathcal{L}_t^c[\alpha_i] - \min_{\text{size}(i_1, i_2, \dots, i_T) \leq k} \sum_{t=1}^T \mathcal{L}_t^c[\alpha_{i_t}] \right) \leq \frac{c}{\sqrt{2}} \sqrt{H\left(\frac{k}{T-1}\right)} \quad (19)$$

for all $T \geq 1$, given that

$$\kappa = \frac{k}{T-1}, \quad \eta = \sqrt{8 \frac{T-1}{T} H\left(\frac{k}{T-1}\right)}.$$

In particular, if $k = o(T)$, we have

$$\limsup_{T \rightarrow \infty} \frac{1}{T} \left(\sum_{t=1}^T \sum_{i=1}^{\text{card}(\mathcal{A}_T)} p_{i,t} \mathcal{L}_n[\alpha_i] - \min_{\text{size}(i_1, i_2, \dots, i_T) \leq k} \sum_{t=1}^T \mathcal{L}_t^c[\alpha_{i_t}] \right) \leq 0 \quad (20)$$

almost surely.

Remark 4 (Intuitive explanations of Algorithm 2 and Proposition 1). In addition to the (sequential) data and model class, other inputs to **Algo. 1** are two learning parameters η, ζ , the number of active models K , and the threshold ρ . The parameters η and ζ control the rate of learning and the rate of model expansion, respectively. The number of active models K is set to reduce the computation cost when data size is small compared with model dimensions, and the threshold ρ is used to update our active models under consideration.

In particular, upon the arrival of a new data point or a set of data points, denoted by \mathbf{z}_t , at each time step t (line 1), we update the weight of each candidate model by a Bayes-type procedure (line 2). The loss employed in the update is the corrected prediction loss, which is directly computable from the data and which serves as an approximation of the expected prediction loss (as was discussed in Subsection II-B). The weights of each model are then updated following the path graph (line 3). When the weight of the smallest model becomes small and that of the largest model becomes large, it means the current active models are inadequately small. So we drop the smallest model and include the next large model into the active set, and adjust their weights accordingly (lines 5-8). In line 7, the weight of the removed model is assigned to the newly included one, so that the sum of the weights remains the same.

Proposition 1 states that the average predictive performance of our algorithm is asymptotically close to that of the optimal model expansion allowing $k = o(T)$ switches. For example, if only one point arrives at each time step, and the dimension of the optimal model is at the order of T^δ for $\delta \in (0, 1)$, then the condition is trivially satisfied. The proof of Proposition 1 follows directly from Theorem 1 and Theorem 2, by using simple manipulations. For technical convenience, Proposition 1 is only proved by removing the part of maintaining an active subset (lines 5-8). We maintain an active subset only for computational **concerns**, and we notice that it does not deteriorate the performance from various numerical experiments. **It is** because the weights of the removed models are negligible.

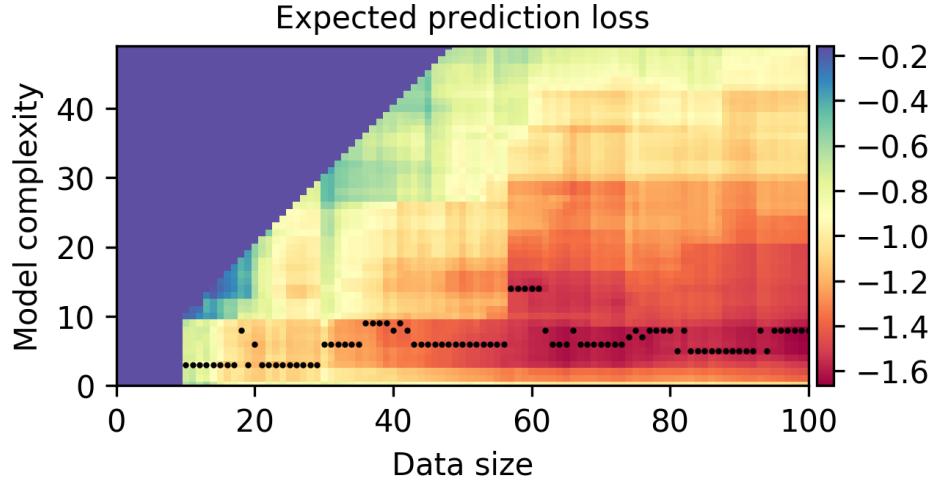
IV. NUMERICAL EXPERIMENTS

In view of Examples 1 and 2, the model classes under consideration are logistic regression and single-layer feed-forward neural networks.

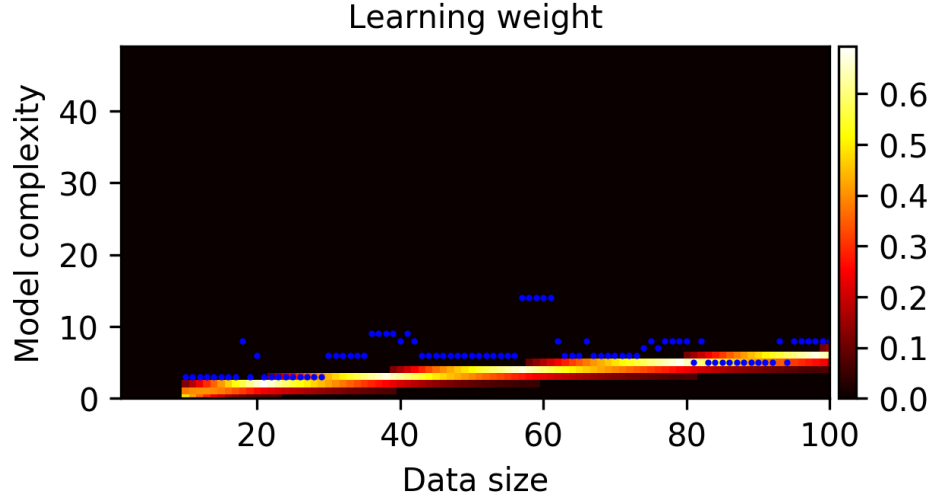
Logistic regression models

We consider the model class to be logistic regression. We generate data from a logistic regression model, where the coefficient vector is $\beta = 10 \times [1^{-1.5}, \dots, 100^{-1.5}]^T$, and covariates x_1, \dots, x_{100} are independent standard Gaussian. Suppose that we sequentially obtain and learn the data, starting from $n = 10$, and then $n = 11, \dots, 100$. We choose the candidate models to be logistic regression models of dimension $1, \dots, [n] = 10$. Here, a model of dimension d means that the first d covariates are nonzero. The model class is nested because a small model is a special case of a large model. We restrict the maximum dimension of candidate models to be \sqrt{t} at time t (see our theoretical assumptions). We summarize the results in Fig. 5 and 6.

To illustrate the efficiency of GTIC, we first simulate model selection results with batch data. We numerically compute the true prediction loss of each trained model (obtained by testing on a large dataset), and then identify the optimal model (with the least loss). In Fig. 6a, we compare the performance of GTIC to different types of CV.



(a) Heat-map showing the true prediction loss of estimated candidate models of each dimension (y-axis) at each data size (x-axis), where the black dots indicate the model of optimal loss at each data size. The true loss is numerically computed from independently generated test data.

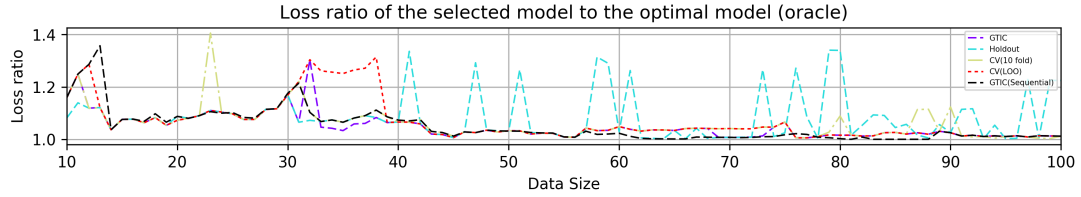


(b) Heat-map showing our predictive weights over the candidate models (y-axis) at each data size (x-axis), using sequential model expansion.

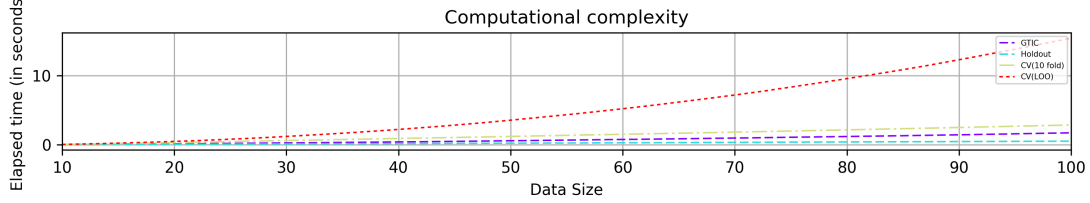
Fig. 5: Experiment 1

Holdout takes 70% data for training and tests on 30% data. It fluctuates throughout the experiment, and most of time it yields the worst performance. GTIC, 10-fold CV and LOO perform well in this experiment. However, both GTIC and 10-fold CV fluctuate a bit. Our proposed sequential model expansion algorithm smoothly expands the model and yields the best performance compared to all the other approaches. As shown in Fig. 5a and 5b, although the optimal model of each data size is not always identical to the selected model from our model expansion algorithm, the loss of our selected model is almost the same as the optimal model. This result is consistent with our definition of efficient Learning.

The computation cost of all approaches is provided in Fig. 6b. As shown in the figures, under logistic regression, GTIC is slightly better than 10-fold CV but worse than Holdout. Indeed, we need to compute the penalty term in GTIC. However, depending on the problem and data, we may need different number of folds for CV in order to have a satisfactory result. Since GTIC performs almost as well as LOO and 10-fold CV, we suggest using GTIC instead of guessing the optimal number of fold for CV. With GTIC, we do not need to sacrifice much on computation cost, but can still achieve theoretically justifiable result which is as good as LOO.

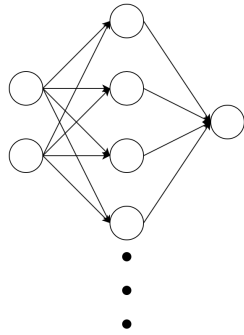


(a) Plot showing the loss of our predictor (GTIC) and cross validations at each data size



(b) Plot showing the computational costs.

Fig. 6: Experiment 1



(a) An illustration of the single-layer Feed-Forward Neural Network

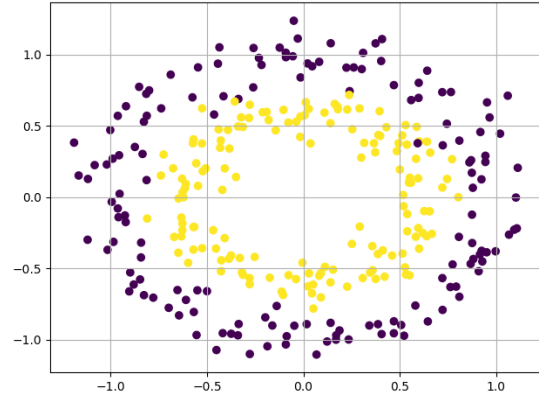

 (b) A set of 300 data uniformly sampled from two circles corrupted by Gaussian noise ($\mu = 0$, $\sigma^2 = 0.1$, radius ratio = 0.6)

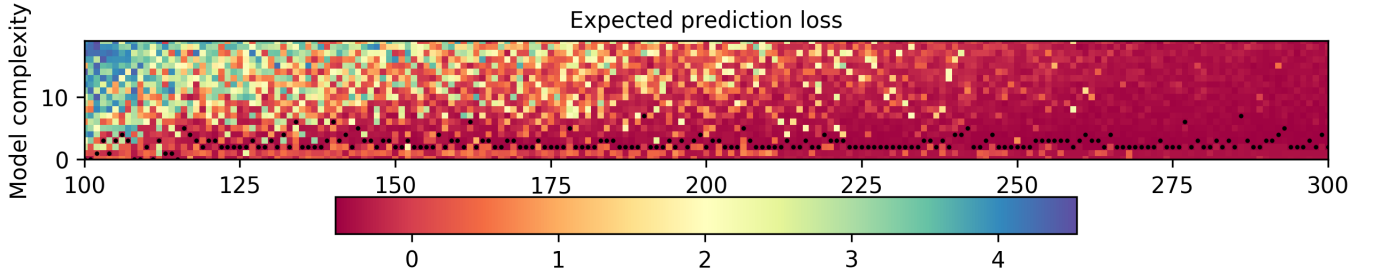
Fig. 7: Experiment 2

Neural networks

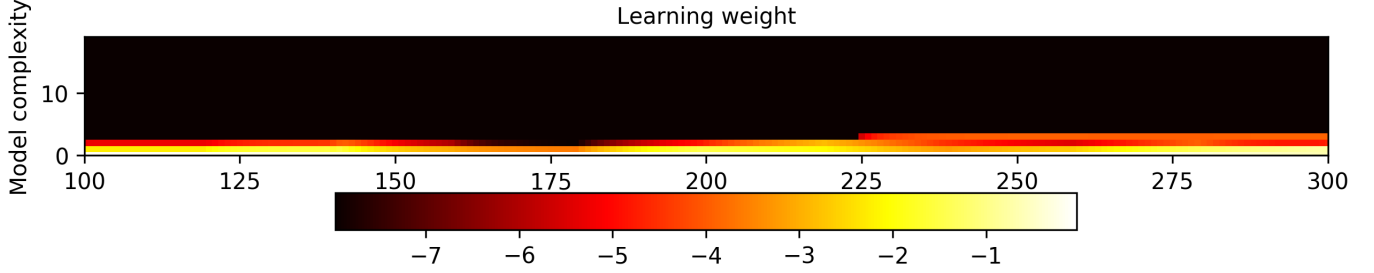
We consider the model class to be single-layer feed-forward neural networks (see Fig. 7a). Neural networks are inherently miss-specified models.

Data are generated from the following way. A set of two-dimensional data are uniformly sampled from two circles (with radius ratio 0.6), corrupted by independent Gaussian noise with mean 0 and variance 0.1 (generated from python package sklearn dataset “make_circle”). The goal is to correctly classify the data into two groups, the larger and smaller rings. Since we have two-dimensional data, our input dimension for the model is two. And because we want to classify into two groups, the output dimension is one. In this experiment, the model complexity of our model is the number of hidden nodes in the single hidden layer.

We sequentially obtain and learn the data, starting from $t = 100$, then $t = 101, \dots, 300$. We start from 100 samples because Neural Network likely converges to a local optimal for small sample size. The path of expansion in this case is the number of hidden nodes in the single hidden layer. Since the true model is not linearly-separable, we do need more than one hidden node to accurately classify the data. We restrict the maximum number of hidden nodes to be $\sqrt{t}/(\text{input dimension})$ due to our assumption. The path of expansion is in increasing order of the number of hidden nodes, since having a small number of hidden nodes is a special case of having more number of hidden



(a) Heat-map showing the prediction loss of estimated candidate models of each dimension (y-axis) at each data size (x-axis), where the black dots indicate the model of optimal loss at each data size.



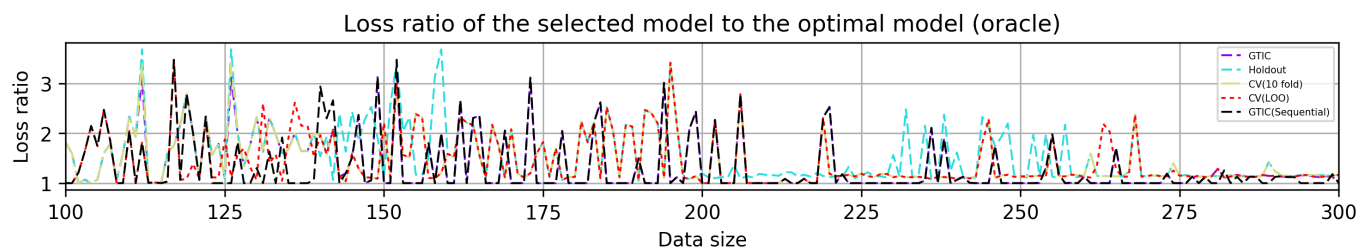
(b) Heat-map showing our predictive weights over the candidate models (y-axis) at each data size (x-axis).

Fig. 8: Experiment 2

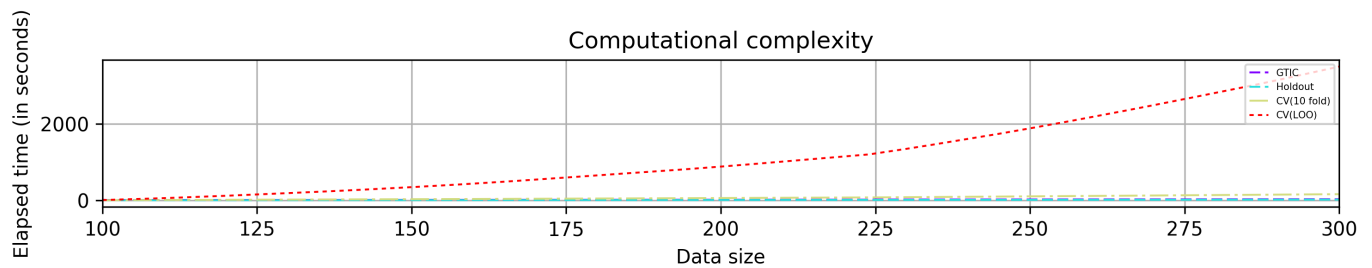
nodes.

Similarly, the optimal model (oracle) is obtained by testing the trained model on a large dataset. The oracle loss of different models at different data size is shown in Fig. 8a. With a small sample size, the cost of overfitting is considerably high. When we have enough samples for training, the cost of overfitting decreases. This effect may also depend on the dimension of input data and labels. In Fig. 9a, the loss ratio varies quite a lot when the sample size is small, but gradually converges. This is partially because the influence of overfitting on the predictive power decreases as sample size increases. In other words, even if we choose a model that is slightly overfitting, the loss ratio is still close to one. Our proposed sequential algorithm is superior to other approaches as shown in Fig. 8b, because the weight of smaller models in the active set is large enough to prevent the model to expand. As a result, we alleviate the tendency to choose the overfitting models even when their loss is relatively small.

The computation cost is shown in Fig. 9b. As expected, the computation of 10-fold CV and LOO increases significantly. However, since we can analytically compute the gradient and hessian involved in the GTIC penalty term, using symbolic expression computation software and saving them on the disk in advance, our computation cost is almost constant at each time step. Therefore, our overall computation cost is almost identical to Holdout. Furthermore, we can utilize warm-start in our implementation, which is a benefit that CV cannot enjoy in naive sequential model selection framework. Therefore, we encourage the use of GTIC in sequential model expansion scheme.



(a) Plot showing the loss of our predictor (GTIC) and cross validations at each data size.



(b) Plot showing the computational costs.

Fig. 9: Experiment 2

V. CONCLUSION

In a framework of high-dimensional parametric models with possibly expanding model dimensions to reflect the challenge of high dimension and small sample size, we aim to investigate strategies that approach the best predictive power in a sequential setting. ~~In the first part of this paper, we pointed out the “bigger models for bigger data” phenomena that is almost ubiquitous in statistical prediction and machine learning tasks.~~ We then studied a method to approach the limit of learning, in the sense that the expected prediction loss of the selected model is asymptotically close to the best offered from a model class (that can depend on data size). The theoretical analysis of GTIC extends classical analysis on AIC for rather general mis-specified model classes. In the second part of the paper, we propose a sequential model expansion strategy that is motivated by our graph-based expert tracking techniques. In summary, the proposed methodology is theoretically optimal and practically useful, and may be a promising competitor of various types of cross-validations.

~~As future work, we may improve on the current estimators of $V(\theta, \alpha)$ and $J(\theta, \alpha)$ by techniques such as bootstrap aggregating, in terms of stability and accuracy. We will apply the proposed methods to more complex neural networks, where computing estimators of the above two quantities could be a nontrivial task. We are investigating alternative quantities that are more easily computable.~~

~~Moreover, sequentially observed data may come from dramatically different data generating processes in different time epochs. If that is the case, any single prescribed model class (instead of model) can be inadequate in predicting future data. In the framework of neural networks, for example, it is therefore interesting to apply our methods not only to select the number of neurons and connections in a prescribed type of networks, but also to adaptively determine the types of networks as well, in a sequential adaptive manner.~~

ACKNOWLEDGEMENT

This work is supported by Defense Advanced Research Projects Agency (DARPA) grant numbers W911NF-14-1-0508, N66001-15-C-4028, and W911NF-16-1-0561.

The authors thank Yang Wu for pointing out some interesting developments in cognitive science, Kathryn Heal for helps in improving the presentation, and Bingfeng Xia from Harvard for suggesting some implementation tools for neural networks.

APPENDIX

We first provide some experimental studies, showing that the GTIC procedure can be a promising competitor for various types of cross-validations. The advantages of GTIC are both superior predictive performance and lower computational cost (especially when the derivatives of the loss function are easy to calculate). We then prove the technical results.

APPENDIX A

PROOF OF THEOREM 1

We start with the following technical lemmas and additional definitions.

Lemma 1. *Suppose that Assumptions 1, 2, 3, 4, 6 hold. Then $\hat{\theta}_n$ is n^τ -consistent uniformly over \mathcal{A}_n , namely $\sup_{\alpha \in \mathcal{A}_n} n^\tau \|\hat{\theta}_n[\alpha] - \theta_n^*[\alpha]\| = O_p(1)$.*

Proof: Using Assumption 2 and a direct adaptation of the techniques in [29, Theorem 5.7], we can prove that $\hat{\theta}_n[\alpha]$ is consistent in the sense that

$$\sup_{\alpha \in \mathcal{A}_n} \|\hat{\theta}_n[\alpha] - \theta_n^*[\alpha]\| = o_p(1) \quad (21)$$

as $n \rightarrow \infty$.

By (10) in Assumption 5, and the definition of $\theta_n^*[\alpha]$, we have

$$E_*\{\psi_n(\cdot, \theta_n^*[\alpha]; \alpha) = \nabla_{\theta} E_* \ell_n(\cdot, \theta_n^*[\alpha]; \alpha) = 0.$$

From the definitions of $\hat{\theta}_n$ and θ_n^* , we have for each $\alpha \in \mathcal{A}_n$

$$\begin{aligned} & n^\tau E_* \{ \psi_n(\cdot, \theta_n^*[\alpha]; \alpha) - \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) \} \\ &= n^\tau \{ 0 - E_* \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) \} \\ &= n^\tau \{ E_n \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) - E_* \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) \} \end{aligned} \quad (22)$$

From the differentiability of the map $\theta \mapsto E_* \psi_n(\cdot, \theta; \alpha)$, there exists $\tilde{\theta}[\alpha]$ such that $\|\tilde{\theta}[\alpha] - \theta_n^*[\alpha]\| \leq \|\hat{\theta}_n[\alpha] - \theta_n^*[\alpha]\|$, and

$$\begin{aligned} & E_* \{ \psi_n(\cdot, \theta_n^*[\alpha]; \alpha) - \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) \} \\ &= \nabla_\theta E_* \{ \psi_n(\cdot, \tilde{\theta}[\alpha]; \alpha) \} (\theta_n^*[\alpha] - \hat{\theta}_n[\alpha]) \\ &= V_n(\tilde{\theta}[\alpha]; \alpha) (\theta_n^*[\alpha] - \hat{\theta}_n[\alpha]), \end{aligned} \quad (23)$$

where the exchangeability of integral and differentiation (in the second identity) is guaranteed by (7) and (8) in Assumption 5.

Therefore, with probability tending to one, we have

$$\begin{aligned} & \sup_{\alpha \in \mathcal{A}_n} n^\tau \|V_n(\tilde{\theta}[\alpha]; \alpha) (\theta_n^*[\alpha] - \hat{\theta}_n[\alpha])\| \\ &= \sup_{\alpha \in \mathcal{A}_n} n^\tau \|E_* \{ \psi_n(\cdot, \theta; \alpha) - \psi_n(\cdot, \hat{\theta}_n; \alpha) \}\| \\ &= \sup_{\alpha \in \mathcal{A}_n} n^\tau \|E_n \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) - E_* \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha)\| \\ &= O_p(1) \end{aligned}$$

where the first equality is due to (23), the second equality is due to (22), and the third equality comes from Assumption 3. By the (13) in Assumption 6 and Assumption 4, $V_n(\tilde{\theta}[\alpha]; \alpha)$ is invertible for each $\alpha \in \mathcal{A}_n$, and

$$\sup_{\alpha \in \mathcal{A}_n} \|V_n(\tilde{\theta}[\alpha]; \alpha)^{-1}\| < 1/(2c_1)$$

with probability tending to one. It follows that

$$\begin{aligned} & \sup_{\alpha \in \mathcal{A}_n} n^\tau \|\theta_n^*[\alpha] - \hat{\theta}_n[\alpha]\| \\ &\leq \sup_{\alpha \in \mathcal{A}_n} \left\{ \|V_n(\tilde{\theta}[\alpha]; \alpha)^{-1}\| \times \|n^\tau V_n(\tilde{\theta}[\alpha]; \alpha) (\theta_n^*[\alpha] - \hat{\theta}_n[\alpha])\| \right\} \\ &= O_p(1), \end{aligned} \quad (24)$$

which concludes the proof. ■

Before we proceed, we need the following definition.

Definition 4 (Bracketing number). Given two scalar functions f_1 and f_2 , the bracket $[f_1, f_2]$ is the set of all functions f such that $f_1 \leq f \leq f_2$. An ε -bracket in $L_2(P_*)$ is a bracket $[f_1, f_2]$ with $E_*(f_2 - f_1)^2 < \varepsilon^2$. The bracketing number $N_{[\cdot]}(\varepsilon, \mathcal{F}, L_2(P_*))$ is the minimum number of ε -brackets needed to cover a set \mathcal{F} . Moreover, the bracketing integral is defined by

$$I_{[\cdot]}(\delta, \mathcal{F}, L_2(P_*)) = \int_0^\delta \sqrt{\log N_{[\cdot]}(\varepsilon, \mathcal{F}, L_2(P_*))} d\varepsilon \quad (25)$$

for $\delta > 0$.

Remark 5. The logarithm of the above bracketing number is also referred to as bracketing entropy relative to the $L_2(P_*)$ -norm. It is commonly used to describe the size of a class of functions. We will use the above definition in order to prove uniform convergence results.

We have the following lemma whose proof follows directly from Definition 4 and Assumption 5.

Lemma 2. Suppose that Assumption 5 holds, and $r_n \leq r$ for all n (where r has been defined in Assumption 5). Let $\mathcal{F}_n[\alpha] = \{\psi_n(\cdot, \theta; \alpha) : \theta \in B(\theta_n^*[\alpha], r_n) \subset \mathbb{R}^{d_n[\alpha]}\}$ be a collection of (vector-valued) measurable functions. Then

$$N_{[\cdot]}(\varepsilon, \mathcal{F}_n[\alpha], L_2(P_*)) \leq (\varepsilon^{-1} r_n \|m_n\|_{P_*})^{d_n[\alpha]} \quad \text{for all } 0 < \varepsilon < r_n.$$

We prove the following technical lemmas.

Lemma 3. For any sets of functions $\mathcal{F}_j, j = 1, \dots, k$, we have

$$I_{[\cdot]}(\delta, \cup_{1 \leq j \leq k} \mathcal{F}_j, L_2(P_*)) \leq \sqrt{2 \log k} \delta + \sqrt{k} \sup_{1 \leq j \leq k} I_{[\cdot]}(\delta, \mathcal{F}_j, L_2(P_*))$$

Proof: The case $k = 1$ is straightforward. We only need to prove for $k \geq 2$. By Definition 4, we have

$$N_{[\cdot]}(\varepsilon, \cup_{1 \leq j \leq k} \mathcal{F}_j, L_2(P_*)) \leq \sum_{1 \leq j \leq k} N_{[\cdot]}(\varepsilon, \mathcal{F}_j, L_2(P_*)).$$

From (25), it suffices to prove the following result, and then let v_j 's be $N_{[\cdot]}(\varepsilon, \mathcal{F}_j, L_2(P_*))$'s. For any numbers $v_1 \geq \dots \geq v_k \geq 1$ ($k \geq 2$), we have

$$\sqrt{\log \sum_{j=1}^k v_j} \leq \max\{\sqrt{2 \log k}, \sqrt{k \log v_1}\} \leq \sqrt{2 \log k} + \sqrt{k \log v_1}.$$

Furthermore, it suffices to prove that

$$\log(kv_1) \leq \max\{2 \log k, k \log v_1\}. \quad (26)$$

In fact, if $v_1 \leq k^{1/(k-1)}$, then

$$\log(kv_1) \leq \log(k \cdot k^{1/(k-1)}) = \frac{k}{k-1} \log k \leq 2 \log k.$$

Otherwise, $\log(kv_1) \leq k \log v_1$, because $g : v \mapsto k \log v - \log(kv)$ is increasing on $v \geq 1$ and it equals zero when $v = k^{1/(k-1)}$. ■

Definition 5. For any class \mathcal{F} of functions $f : \mathcal{Z} \rightarrow \mathbb{R}$, a function $F : \mathcal{Z} \rightarrow \mathbb{R}$ is called an envelope function of \mathcal{F} , if $\sup_{f \in \mathcal{F}} |f(z)| \leq F(z) < \infty$ for every $z \in \mathcal{Z}$.

Lemma 4. ([29, Lemma 19.34]) For any class \mathcal{F} of measurable functions $f : \mathcal{Z} \rightarrow \mathbb{R}$ such that $E_* f^2 < \delta^2$ for all f , with

$$a(\delta) = \delta / \sqrt{\max\{1, \log N_{[\cdot]}(\delta, \mathcal{F}, L_2(P_*))\}}$$

and F an envelope function, that

$$E_* \sup_{f \in \mathcal{F}} |\mathbb{G}_n f| \lesssim I_{[\cdot]}(\delta, \mathcal{F}, L_2(P_*)) + \sqrt{n} E_* \{F \cdot 1_{F > \sqrt{n} a(\delta)}\}.$$

Here, 1_A is the indicator function of event A .

Lemma 5. Let $\mathcal{F}_n = \cup_{\alpha \in \mathcal{A}_n} \mathcal{F}_n[\alpha]$ where $\mathcal{F}_n[\alpha] = \{\mathbf{f}_{n,u} : u \in U[\alpha]\}$ be a class of measurable vector-valued functions. In other words, for each $\alpha \in \mathcal{A}_n$ and $u \in U[\alpha]$, $\mathbf{f}_{n,u} = [f_{n,u,1}, \dots, f_{n,u,d_n[\alpha]}]^T$ with $f_{n,u,i} : \mathcal{Z} \rightarrow \mathbb{R}$ being a scalar-valued function. The dimension $d_n[\alpha]$ can be different for $\alpha \in \mathcal{A}_n$, and we let $d_n = \max_{\alpha \in \mathcal{A}_n} d_n[\alpha]$. Assume that the following conditions hold.

(i) There is an envelope function F_n that satisfies

$$\sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha] \subset \mathbb{R}^{d_n[\alpha]}, 1 \leq i \leq d_n[\alpha]} |f_{n,u,i}(\mathbf{z})| \leq F_n(\mathbf{z}) < \infty$$

for every $\mathbf{z} \in \mathcal{Z}$;

(ii) There exists a deterministic sequence $\{\delta_n\}$ such that

$$d_n \sqrt{\log\{d_n \text{card}(\mathcal{A}_n)\}} \delta_n \rightarrow 0, \quad (27)$$

and

(iii) *The bounded moment condition:*

$$\delta_n^{-2} E_* F_n^2 \rightarrow 0;$$

(iv) *The bounded class condition:*

$$\sqrt{d_n^{3/2} \text{card}(\mathcal{A}_n)} \times \sup_{\alpha \in \mathcal{A}_n, 1 \leq i \leq d_n[\alpha]} I_{[\cdot]}(\delta_n, \mathcal{F}_{n,i}[\alpha], L_2(P_*)) \rightarrow 0,$$

where we let $\mathcal{F}_{n,i}[\alpha] = \{f_{n,u,i} : u \in U[\alpha]\}$.

Then we have $\sup_{\mathbf{f} \in \mathcal{F}_n} \|\mathbb{G}_n \mathbf{f}\| \rightarrow_p 0$ as $n \rightarrow \infty$.

Proof: By Markov's inequality, it suffices to prove that $E_* \sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha]} \|\mathbb{G}_n f_{n,u}\| \rightarrow 0$ as $n \rightarrow \infty$. Condition (iii) implies that for all sufficiently large n ,

$$\sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha], i=1, \dots, d_n} E_* f_{n,u,i}^2 \leq E_* \sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha], i=1, \dots, d_n} f_{n,u,i}^2 < \delta_n^2. \quad (28)$$

Let $\delta_n, a_n(\delta_n)$ be the constants given in Lemma 4 corresponding to $\delta = \delta_n$ and

$$\tilde{\mathcal{F}}_n = \bigcup_{\alpha \in \mathcal{A}_n, 1 \leq i \leq d_n[\alpha]} \mathcal{F}_{n,i}[\alpha].$$

From inequality (28) and Lemma 4, we have

$$\begin{aligned} & E_* \sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha], 1 \leq i \leq d_n[\alpha]} |\mathbb{G}_n f_{n,u,i}| \\ & \leq I_{[\cdot]}(\delta_n, \tilde{\mathcal{F}}_n, L_2(P_*)) + \sqrt{n} E_* \{F_n \cdot 1_{F_n > \sqrt{n} a_n(\delta_n)}\} \\ & \leq I_{[\cdot]}(\delta_n, \tilde{\mathcal{F}}_n, L_2(P_*)) + \frac{1}{a_n(\delta_n)} E_* F_n^2, \end{aligned} \quad (29)$$

where the second inequality comes from the fact that

$$1_{F_n > \sqrt{n} a_n(\delta_n)} \leq \frac{F_n}{\sqrt{n} a_n(\delta_n)} 1_{F_n > \sqrt{n} a_n(\delta_n)} \leq \frac{F_n}{\sqrt{n} a_n(\delta_n)}.$$

By the definition of $a_n(\cdot)$, $I_{[\cdot]}(\delta, \tilde{\mathcal{F}}_n, L_2(P_*))$, and the fact that $N_{[\cdot]}(\delta, \tilde{\mathcal{F}}_n, L_2(P_*))$ is non-increasing in δ , we have

$$\begin{aligned} \frac{1}{a_n(\delta_n)} &= \frac{1}{\delta_n} \sqrt{\max\{1, \log N_{[\cdot]}(\delta_n, \tilde{\mathcal{F}}_n, L_2(P_*))\}} \\ &\leq \frac{1}{\delta_n^2} I_{[\cdot]}(\delta_n, \tilde{\mathcal{F}}_n, L_2(P_*)). \end{aligned}$$

It follows that the right hand side of (29) is upper bounded by

$$I_{[\cdot]}(\delta_n, \tilde{\mathcal{F}}_n, L_2(P_*)) (1 + \delta_n^{-2} E_* F_n^2).$$

Therefore, by Lemma 3 and simple manipulations, we have

$$\begin{aligned} & E_* \sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha]} \|\mathbb{G}_n \mathbf{f}_{n,u}\| \\ & \leq E_* \sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha]} \sum_{i=1}^{d_n} |\mathbb{G}_n f_{n,u,i}| \\ & \leq d_n E_* \sup_{\alpha \in \mathcal{A}_n, u \in U[\alpha], 1 \leq i \leq d_n[\alpha]} |\mathbb{G}_n f_{n,u,i}| \\ & \leq (A_1 + A_2) (1 + \delta_n^{-2} E_* F_n^2), \end{aligned} \quad (30)$$

where

$$\begin{aligned} A_1 &= d_n \sqrt{2 \log\{d_n \text{card}(\mathcal{A}_n)\}} \delta_n, \\ A_2 &= d_n^{3/2} \sqrt{\text{card}(\mathcal{A}_n)} \sup_{\alpha \in \mathcal{A}_n, 1 \leq i \leq d_n[\alpha]} I_{[\cdot]}(\delta_n, \mathcal{F}_{n,i}[\alpha], L_2(P_*)), \end{aligned}$$

Assumptions (ii), (iii), and (iv) guarantee that the right hand side of (30) goes to zero as $n \rightarrow \infty$, which concludes the proof. ■

Using the above results, we can prove the following key technical lemma.

Lemma 6. *Suppose that Assumptions 1-6 hold. Then*

$$\sup_{\alpha \in \mathcal{A}_n} \left\| \mathbb{G}_n \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) - \mathbb{G}_n \psi_n(\cdot, \theta_n^*[\alpha]; \alpha) \right\| = o_p(1). \quad (31)$$

Proof: For a constant c , consider the class $\mathcal{F}_n = \cup_{\alpha \in \mathcal{A}_n} \mathcal{F}_n[\alpha]$, with $\mathcal{F}_n[\alpha] \triangleq \{\mathbf{f}_{n,\mathbf{u}} : \mathbf{u} \in U[\alpha]\}$, $U[\alpha] = \{[u_1, \dots, u_{d_n[\alpha]}]^T : \sum_{i=1}^{d_n[\alpha]} u_i^2 = c\}$, and

$$\mathbf{f}_{n,\mathbf{u}}(\cdot) = \psi_n(\cdot, \theta_n^*[\alpha] + n^{-\tau} \mathbf{u}; \alpha) - \psi_n(\cdot, \theta_n^*[\alpha]; \alpha).$$

Suppose that $\varepsilon, \delta > 0$ are fixed constants. It suffices to prove that the left hand side of (31) is less than δ with probability at least $1 - \varepsilon$ for all sufficiently large n . By Lemma 1, there exists a constant $c > 0$ such that $\mathbb{G}_n \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) - \mathbb{G}_n \psi_n(\cdot, \theta_n^*[\alpha]; \alpha)$ falls into the class \mathcal{F}_n with probability at least $1 - \varepsilon/2$ for all sufficiently large n . Therefore, we only need to prove that for any given constant $c > 0$, $\sup_{\mathbf{f} \in \mathcal{F}_n} \|\mathbb{G}_n \mathbf{f}\| \rightarrow_p 0$. It remains to prove that there are δ_n 's that satisfy Conditions (i)-(iv) of Lemma 5.

We define $\mathcal{F}_{n,i}[\alpha]$ as was in Lemma 5, and define $m_n(\cdot) = \sup_{\alpha \in \mathcal{A}_n} m_n[\alpha](\cdot)$. By Assumption 5, we can use $F_n(\cdot) \triangleq cn^{-\tau} \sup_{\alpha \in \mathcal{A}_n} m_n[\alpha](\cdot)$ as the envelop function for each $\mathbf{f}_{n,u,i}(\cdot)$, and we have

$$E_* F_n^2 \leq C_1 \triangleq c^2 n^{-2\tau} E_* m_n^2$$

Let

$$C_2 = d_n \sqrt{\log\{d_n \text{card}(\mathcal{A}_n)\}}.$$

Because of (9) in Assumption 5, we have

$$C_2^2 C_1 = c^2 n^{-2\tau} d_n^2 \log\{d_n \text{card}(\mathcal{A}_n)\} E_* m_n^2 \rightarrow 0 \quad (32)$$

This implies the existence of a sequence δ_n (e.g. $\delta_n = C_1^{1/4} C_2^{-1/2}$) such that

$$\delta_n C_2 \rightarrow 0, \quad \delta_n^{-2} C_1 \rightarrow 0,$$

which further implies Conditions (ii) and (iii) in Lemma 5.

To conclude the proof, we prove that Condition (iv) in Lemma 5 holds for any $\delta_n \rightarrow 0$. From Lemma 2, we have for each $\alpha \in \mathcal{A}_n, 1 \leq i \leq d_n[\alpha]$ that

$$\begin{aligned} & I_{[\cdot]}(\delta_n, \mathcal{F}_{n,i}[\alpha], L_2(P_*)) \\ & \leq \int_0^{\delta_n} \left[\max \left\{ 0, d_n \log(\varepsilon^{-1} cn^{-\tau} \|m_n\|_{P_*}) \right\} \right]^{1/2} d\varepsilon \\ & = \int_0^{\min\{\delta_n, cn^{-\tau} \|m_n\|_{P_*}\}} \left[d_n \log(\varepsilon^{-1} cn^{-\tau} \|m_n\|_{P_*}) \right]^{1/2} d\varepsilon \end{aligned} \quad (33)$$

Because condition (9) implies that $n^{-\tau} \|m_n\|_{P_*} \rightarrow 0$, the value of ε in the integral is close to one. This implies that for all sufficiently large n , the integrand in (33) is upper bounded by $d_n^{1/2} \varepsilon^{-\rho}$, where $1/(1 - \rho) = \gamma$ and γ is given in Assumption 5. Therefore, for all sufficiently large n , the right hand side of (33) is upper bounded by

$$\int_0^{cn^{-\tau} \|m_n\|_{P_*}} d_n^{1/2} \varepsilon^{-\rho} d\varepsilon = (1 - \rho)^{-1} d_n^{1/2} (cn^{-\tau} \|m_n\|_{P_*})^{1-\rho},$$

which does not depend on α, i . This further implies

$$\begin{aligned}
& \sqrt{d_n^{3/2} \text{card}(\mathcal{A}_n)} \times \sup_{\alpha \in \mathcal{A}_n, 1 \leq i \leq d_n[\alpha]} I_{[\cdot]}(\delta_n, \mathcal{F}_{n,i}[\alpha], L_2(P_*)) \\
& \leq (1 - \rho)^{-1} c^{1-\rho} d_n \sqrt{\text{card}(\mathcal{A}_n)} (n^{-\tau} \|m_n\|_{P_*})^{1-\rho} \\
& = (1 - \rho)^{-1} c^{1-\rho} \left(d_n^\gamma \text{card}(\mathcal{A}_n)^{\gamma/2} n^{-\tau} \|m_n\|_{P_*} \right)^{1-\rho} \\
& \rightarrow 0
\end{aligned} \tag{34}$$

where the last limit is due to (9) in Assumption 5. ■

We next prove the second key technical lemma.

Lemma 7. *Suppose that Assumptions 1-6 hold. Assume that the map $\theta \mapsto E_*\psi_n(\cdot, \theta; \alpha)$ is differentiable at a θ_n^* for all n . Then we have*

$$\sqrt{n}(\hat{\theta}_n[\alpha] - \theta_n^*[\alpha]) = -\{V_n(\theta_n^*[\alpha]; \alpha)^{-1} + \nu_{1,n}[\alpha]\} \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_n(z_i, \theta_n^*[\alpha]; \alpha) + \nu_{2,n}[\alpha]$$

where $\nu_{1,n}[\alpha]$ is a positive semidefinite matrix and $\nu_{2,n}[\alpha]$ is a vector such that $\sup_{\alpha \in \mathcal{A}_n} \|\nu_{1,n}[\alpha]\| \rightarrow_p 0$ and $\sup_{\alpha \in \mathcal{A}_n} \|\nu_{2,n}[\alpha]\| \rightarrow_p 0$.

Proof: By the definitions of θ_n^* and $\hat{\theta}_n$, we have

$$\begin{aligned}
& \sqrt{n} E_* \{ \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) - \psi_n(\cdot, \theta_n^*[\alpha]; \alpha) \} \\
& = \sqrt{n} \{ E_* \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) - 0 \} \\
& = \sqrt{n} \{ E_* \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) - E_n \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) \} \\
& = -\mathbb{G}_n \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) \\
& = -\mathbb{G}_n \psi_n(\cdot, \theta_n^*[\alpha]; \alpha) + \nu_n
\end{aligned} \tag{35}$$

where the last equality is due to Lemma 6, and $\|\nu_n\| = o_p(1)$.

From the differentiability of the map $\theta \mapsto E_*\psi_n(\cdot, \theta; \alpha)$, there exists $\tilde{\theta}[\alpha]$ such that $\|\tilde{\theta}[\alpha] - \theta_n^*[\alpha]\| \leq \|\hat{\theta}_n[\alpha] - \theta_n^*[\alpha]\|$, and

$$\begin{aligned}
& E_* \{ \psi_n(\cdot, \theta_n^*[\alpha]; \alpha) - \psi_n(\cdot, \hat{\theta}_n[\alpha]; \alpha) \} \\
& = \nabla_{\theta} E_* \{ \psi_n(\cdot, \tilde{\theta}[\alpha]; \alpha) \} (\theta_n^*[\alpha] - \hat{\theta}_n[\alpha]) \\
& = V_n(\tilde{\theta}[\alpha]; \alpha) (\theta_n^*[\alpha] - \hat{\theta}_n[\alpha]),
\end{aligned} \tag{36}$$

where the exchangeability of integral and differentiation (in the second identity) is guaranteed by (7) and (8) in Assumption 5. Multiplying the matrix $\sqrt{n}V_n(\tilde{\theta}[\alpha]; \alpha)^{-1}$ to both sides of (36) and using equality (35), we have

$$\sqrt{n}(\hat{\theta}_n[\alpha] - \theta_n^*[\alpha]) = -V_n(\tilde{\theta}[\alpha]; \alpha)^{-1} \mathbb{G}_n \psi_n(\cdot, \theta_n^*[\alpha]; \alpha) + V_n(\tilde{\theta}[\alpha]; \alpha)^{-1} \nu_n. \tag{37}$$

We conclude the proof by applying Assumption 4 (with the constant c_2) and (12) in Assumption 6 to equality (37). ■

Proof of Theorem 1

In order to prove that the minimum of $\mathcal{L}_t^c[\alpha]$ asymptotically approaches the minimum of $\mathcal{L}_n[\alpha]$ (in the sense of Definition 3), we only need to prove that $\mathcal{L}_t^c[\alpha]/\mathcal{L}_n[\alpha] = 1 + o_p(1)$ where $o_p(1)$ is uniform in $\alpha \in \mathcal{A}_n$. In other words,

$$\sup_{\alpha \in \mathcal{A}_n} \left| \frac{\mathcal{L}_t^c[\alpha] - \mathcal{L}_n[\alpha]}{\mathcal{L}_n[\alpha]} \right| \rightarrow_p 0.$$

Recall the definition of $\mathcal{R}_n[\alpha]$. It further suffices to prove that

$$\sup_{\alpha \in \mathcal{A}_n} \left| \frac{\mathcal{L}_t^c[\alpha] - \mathcal{L}_n[\alpha]}{\mathcal{R}_n[\alpha]} \right| \rightarrow_p 0, \quad (38)$$

and

$$\sup_{\alpha \in \mathcal{A}_n} \frac{\mathcal{L}_n[\alpha]}{\mathcal{R}_n[\alpha]} \rightarrow_p 1. \quad (39)$$

By the definition of loss $\mathcal{L}_n[\alpha]$ and Taylor expansion, we have for each $\alpha \in \mathcal{A}_n$

$$\begin{aligned} \mathcal{L}_n[\alpha] &= E_* l_n(\mathbf{z}, \hat{\boldsymbol{\theta}}_n[\alpha]; \alpha) \\ &= E_* l_n(\mathbf{z}, \boldsymbol{\theta}_n^*[\alpha]; \alpha) + (\hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha])^\top \frac{\partial}{\partial \boldsymbol{\theta}} E_* l_n(\mathbf{z}, \boldsymbol{\theta}_n^*[\alpha]; \alpha) + \frac{1}{2} \left\| \hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha] \right\|_{\nabla_{\boldsymbol{\theta}}^2 E_* l_n(\mathbf{z}, \tilde{\boldsymbol{\theta}}[\alpha]; \alpha)}^2 \\ &= E_* l_n(\mathbf{z}, \boldsymbol{\theta}_n^*[\alpha]; \alpha) + \frac{1}{2} \left\| \hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha] \right\|_{V_n(\tilde{\boldsymbol{\theta}}[\alpha]; \alpha)}^2 \end{aligned} \quad (40)$$

where $\tilde{\boldsymbol{\theta}}[\alpha]$ in the second equality is a vector satisfying $\|\tilde{\boldsymbol{\theta}}[\alpha] - \boldsymbol{\theta}_n^*[\alpha]\| \leq \|\hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha]\|$, and the exchangeability of expectation and differentiation in the third equality is guaranteed by (8) in Assumption 5, and the consistency of $\hat{\boldsymbol{\theta}}_n[\alpha]$. We note that by Assumption 4, the equality (40) further implies (4) presented in our introduction.

Similarly, we have

$$\begin{aligned} \hat{\mathcal{L}}_n[\alpha] &= \frac{1}{n} \sum_{i=1}^n l_n(\mathbf{z}_i, \hat{\boldsymbol{\theta}}_n[\alpha]; \alpha) \\ &= \frac{1}{n} \sum_{i=1}^n l_n(\mathbf{z}_i, \boldsymbol{\theta}_n^*[\alpha]; \alpha) + (\hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha])^\top \frac{1}{n} \sum_{i=1}^n \psi_n(\mathbf{z}_i, \boldsymbol{\theta}_n^*[\alpha]; \alpha) + \frac{1}{2} \left\| \hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha] \right\|_{\hat{V}_n(\tilde{\boldsymbol{\theta}}[\alpha])}^2 \end{aligned} \quad (41)$$

From identities (40) and (41), we may write

$$\begin{aligned} \mathcal{L}_n[\alpha] - \hat{\mathcal{L}}_n[\alpha] &= \frac{1}{n} \text{tr} \left\{ \hat{V}_n(\hat{\boldsymbol{\theta}}_n[\alpha]; \alpha)^{-1} \hat{J}_n(\hat{\boldsymbol{\theta}}_n[\alpha]; \alpha) \right\} \\ &= A_3[\alpha] + A_4[\alpha] + A_5[\alpha] + A_6[\alpha] \end{aligned}$$

where we define

$$\begin{aligned} A_3[\alpha] &= \frac{1}{2} \left\| \hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha] \right\|_{V_n(\tilde{\boldsymbol{\theta}}[\alpha]; \alpha) - \hat{V}_n(\tilde{\boldsymbol{\theta}}[\alpha])}^2 \\ A_4[\alpha] &= -\frac{1}{n} \sum_{i=1}^n \{ l_n(\mathbf{z}_i, \boldsymbol{\theta}_n^*[\alpha]; \alpha) - E_* l_n(\mathbf{z}, \boldsymbol{\theta}_n^*[\alpha]; \alpha) \} \\ A_5[\alpha] &= \frac{1}{n} \left\{ \text{tr} \{ V_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha)^{-1} J_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha) \} \right. \\ &\quad \left. - \text{tr} \{ \hat{V}_n(\hat{\boldsymbol{\theta}}_n[\alpha]; \alpha)^{-1} \hat{J}_n(\hat{\boldsymbol{\theta}}_n[\alpha]; \alpha) \} \right\} \\ A_6[\alpha] &= -(\hat{\boldsymbol{\theta}}_n[\alpha] - \boldsymbol{\theta}_n^*[\alpha])^\top \frac{1}{n} \sum_{i=1}^n \psi_n(\mathbf{z}_i, \boldsymbol{\theta}_n^*[\alpha]; \alpha) \\ &\quad - \frac{1}{n} \text{tr} \{ V_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha)^{-1} J_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha) \}. \end{aligned}$$

In view of (38), it suffices to prove that

$$\sup_{\alpha \in \mathcal{A}_n} \frac{|A_k[\alpha]|}{\mathcal{R}_n[\alpha]} \rightarrow_p 0 \quad (42)$$

as $n \rightarrow \infty$ for $k = 3, 4, 5, 6$, and the limit (39).

By the n^τ -consistency of $\hat{\theta}_n[\alpha]$ uniformly over \mathcal{A}_n (Lemma 1) and Assumption 6,

$$\sup_{\alpha \in \mathcal{A}_n} \frac{|A_3[\alpha]|}{\mathcal{R}_n[\alpha]} = \sup_{\alpha \in \mathcal{A}_n} \frac{1}{2} \frac{n^{-2\tau}}{\mathcal{R}_n[\alpha]} \|\nu_n\|_{V_n(\hat{\theta}_n[\alpha]; \alpha) - \hat{V}_n(\hat{\theta}_n[\alpha])}^2$$

where $\sup_{\alpha \in \mathcal{A}_n} \|\nu_n\| = O_p(1)$. Thus, given assumption (14), (42) with $k = 3$ can be proved.

By Chebyshev's inequality, for any positive constant $\delta > 0$, we have

$$\begin{aligned} & P_* \left(\sup_{\alpha \in \mathcal{A}_n} \frac{|A_4[\alpha]|}{\mathcal{R}_n[\alpha]} > \delta \right) \\ & \leq \sum_{\alpha \in \mathcal{A}_n} P_* \left(\frac{|A_4[\alpha]|}{\mathcal{R}_n[\alpha]} > \delta \right) \\ & \leq \sum_{\alpha \in \mathcal{A}_n} \frac{E_* \{ l_n(\mathbf{z}_1, \boldsymbol{\theta}_n^*; \alpha) - E_* l_n(\mathbf{z}, \boldsymbol{\theta}_n^*[\alpha]; \alpha) \}^{2m_1}}{\delta^{2m_1} n^{2m_1} \mathcal{R}_n[\alpha]^{2m_1}}. \end{aligned} \quad (43)$$

Thus, given assumption (16), (42) with $k = 4$ can be proved.

For brevity, we temporarily denote

$$V_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha), \hat{V}_n(\hat{\boldsymbol{\theta}}_n[\alpha]; \alpha), J_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha), \text{ and } \hat{J}_n(\hat{\boldsymbol{\theta}}_n[\alpha]; \alpha)$$

respectively by

$$V[\alpha], \hat{V}[\alpha], J[\alpha], \text{ and } \hat{J}[\alpha].$$

Then

$$\begin{aligned} & tr\{V[\alpha]^{-1}J[\alpha]\} - tr\{\hat{V}[\alpha]^{-1}\hat{J}[\alpha]\} \\ & = tr\{V[\alpha]^{-1}(J[\alpha] - \hat{J}[\alpha])\} + tr\{(V[\alpha]^{-1} - \hat{V}[\alpha]^{-1})\hat{J}[\alpha]\}. \end{aligned}$$

To prove (42) with $k = 5$, we only need to show that

$$\sup_{\alpha \in \mathcal{A}_n} \frac{1}{n\mathcal{R}_n[\alpha]} tr\{V[\alpha]^{-1}(J[\alpha] - \hat{J}[\alpha])\} \rightarrow_p 0, \quad (44)$$

$$\sup_{\alpha \in \mathcal{A}_n} \frac{1}{n\mathcal{R}_n[\alpha]} tr\{(V[\alpha]^{-1} - \hat{V}[\alpha]^{-1})\hat{J}[\alpha]\} \rightarrow_p 0. \quad (45)$$

We only prove (44), and then (45) follows similar arguments. Suppose that \mathbf{z} is a $\mathcal{N}(0, I)$ random variable of dimension $d_n[\alpha]$, and $V[\alpha]^{-1/2}$ is a positive semidefinite matrix whose square equals $V[\alpha]^{-1}$. Because of Assumption 4 and 5, (44) could be rewritten as

$$\begin{aligned} & \sup_{\alpha \in \mathcal{A}_n} \frac{1}{n\mathcal{R}_n[\alpha]} E \left\{ \mathbf{z}^T V[\alpha]^{-1/2} (J[\alpha] - \hat{J}[\alpha]) V[\alpha]^{-1/2} \mathbf{z} \right\} \\ & = o_p(1) \sup_{\alpha \in \mathcal{A}_n} \frac{1}{n\mathcal{R}_n[\alpha]} E \left\| V[\alpha]^{-1/2} \mathbf{z} \right\|^2 \\ & = o_p(1) \sup_{\alpha \in \mathcal{A}_n} \frac{1}{n\mathcal{R}_n[\alpha]} E \|\mathbf{z}\|^2 \\ & = o_p(1) \sup_{\alpha \in \mathcal{A}_n} \frac{d_n[\alpha]}{n\mathcal{R}_n[\alpha]} \rightarrow_p 0 \end{aligned}$$

where the first equality is due to (11) in Assumption 6, the second equality is due to Assumption 4, and the last equality is guaranteed by assumption (15).

Next, we prove (42) with $k = 6$. Applying Lemma 7, we could rewrite

$$\frac{|A_6[\alpha]|}{\mathcal{R}_n[\alpha]} = A_7[\alpha] + A_8[\alpha] + A_9[\alpha],$$

where we define

$$A_7[\alpha] = \frac{\|\mathbf{w}_n[\alpha]\|_{V_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha)^{-1}}^2 - \text{tr}\{V_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha)^{-1} J_n(\boldsymbol{\theta}_n^*[\alpha]; \alpha)\}}{n\mathcal{R}_n[\alpha]},$$

$$A_8[\alpha] = \frac{\|\mathbf{w}_n[\alpha]\|_{\nu_{1,n}[\alpha]}^2}{n\mathcal{R}_n[\alpha]}, \quad A_9[\alpha] = \frac{\boldsymbol{\nu}_{2,n}[\alpha]^\top \mathbf{w}_n[\alpha]}{n\mathcal{R}_n[\alpha]}.$$

Using assumption (17) and similar arguments as in (43), we can prove $\limsup_{\alpha \in \mathcal{A}_n} |A_7[\alpha]| \rightarrow_p 0$. Similarly, because

$$|A_8[\alpha]| = o_p(1) \frac{\|\mathbf{w}_n[\alpha]\|^2}{n\mathcal{R}_n[\alpha]}$$

where $o_p(1)$ is uniform in \mathcal{A}_n , assumption (18) guarantees that $\sup_{\alpha \in \mathcal{A}_n} A_8[\alpha] \rightarrow_p 0$. Cauchy inequality and assumption (18) also imply that

$$\sup_{\alpha \in \mathcal{A}_n} |A_9[\alpha]| \leq \sup_{\alpha \in \mathcal{A}_n} \frac{\|\boldsymbol{\nu}_{2,n}[\alpha]\| \times \|\mathbf{w}_n[\alpha]\|}{n\mathcal{R}_n[\alpha]} \rightarrow_p 0. \quad (46)$$

Finally, we prove (39). From (40) and τ -consistency of $\hat{\boldsymbol{\theta}}_n[\alpha]$, we have

$$\begin{aligned} \mathcal{L}_n[\alpha] &= E_* l_n(\cdot, \hat{\boldsymbol{\theta}}_n[\alpha]; \alpha) \\ &= E_* l_n(\cdot, \boldsymbol{\theta}_n^*[\alpha]; \alpha) + n^{-2\tau} O_p(1) \end{aligned}$$

where $O_p(1)$ is uniformly in \mathcal{A}_n . Therefore

$$\begin{aligned} \sup_{\alpha \in \mathcal{A}_n} \frac{\mathcal{L}_n[\alpha]}{\mathcal{R}_n[\alpha]} &= 1 + \sup_{\alpha \in \mathcal{A}_n} \frac{\mathcal{L}_n[\alpha] - E_* \mathcal{L}_n[\alpha]}{\mathcal{R}_n[\alpha]} \\ &= 1 + O_p(1) \sup_{\alpha \in \mathcal{A}_n} \frac{1}{n^{2\tau} \mathcal{R}_n[\alpha]} \rightarrow_p 1. \end{aligned}$$

APPENDIX B PROOF OF THEOREM 2

First, we introduce the concept of “compound experts”. A compound expert is defined as an expert sequence (i_1, i_2, \dots, i_T) whose size $\leq k$ with some prescribed $k > 0$. Then in order to tackle the problem of “tracking the best expert”, we could simply apply the exponentially re-weighting algorithm over all the possible compound experts, which can yield provable tight regret bounds. The reason why this simple strategy is not used in practice is that the number of compound experts is usually too large to manage, while the fixed share algorithm greatly reduces the computational complexity and has similar regret bounds.

For our extension of “tracking the best expert” with graphical transitional constraints, following a similar proving strategy used in [32, Chapter 5], we first prove an equivalence between the results of the exponentially re-weighting algorithm over compound experts and the algorithm that we propose, and then apply the regret bound for the former algorithm directly.

The exponentially re-weighting algorithm that we are considering here is as follows. At each time $t = 0, 1, \dots, T$, the distribution over the compound experts is maintained by $w'_t(i_1, i_2, \dots, i_T)$ (not necessarily normalized) for all the sequences (i_1, i_2, \dots, i_T) . The initial distribution is

$$\begin{aligned} w'_0(i_1, i_2, \dots, i_T) &= w'_0(i_1) w'_0(i_2|i_1) w'_0(i_3|i_1, i_2) \cdots w'_0(i_T|i_1, \dots, i_{T-1}) \\ &= w'_0(i_1) w'_0(i_2|i_1) w'_0(i_3|i_2) \cdots w'_0(i_T|i_{T-1}) \quad (\text{Assume Markovian property}) \\ &= w'_0(i_1) \prod_{t=1}^{T-1} w'_0(i_{t+1}|i_t) \\ &= \mathbf{1}\{i_1 = 1\} \prod_{t=1}^{T-1} \left[(1 - \kappa\beta_{i_t}) \mathbf{1}\{i_{t+1} = i_t\} + \kappa\beta_{i_t, i_{t+1}} \mathbf{1}\{i_{t+1} \neq i_t\} \right] \end{aligned}$$

This initial distribution over compound experts ensures that only the “valid” expert sequences (those follow graphical transitions) have positive probabilities.

Based on the exponentially re-weighting updating rule, the distribution at each time instant $t = 1, 2, \dots, T$ becomes

$$w'_t(i_1, i_2, \dots, i_T) = w'_0(i_1, i_2, \dots, i_T) \exp\left(-\eta \sum_{s=1}^t l(i_s, z_s)\right)$$

Marginally, at time t ,

$$w'_{i,t} = \sum_{i_1, \dots, i_t, i_{t+2}, \dots, i_T} w'_t(i_1, \dots, i_t, i_{t+2}, \dots, i_T)$$

And we then have

$$p'_{i,t} = \frac{w'_{i,t}}{W'_t} \quad \text{with } W'_t = \sum_{j=1}^N w'_{j,t}, \quad \text{and} \quad p'_{i,0} = w'_{i,0} = \mathbf{1}\{i = 1\}$$

The exponentially forecaster draws action according to expert i at time $t + 1$ with probability $p'_{i,t}$.

Lemma 8. *For all $\kappa \in (0, 1/D)$, for any sequence of T outcomes, and for all $t = 0, 1, \dots, T$, the predictive distribution $p_{i,t}$ for $i = 1, \dots, N$ generated by our proposed Algorithm 1 is the same as the predictive distribution $p'_{i,t}$ for $i = 1, \dots, N$ that is maintained by the special exponentially re-weighting algorithm described above.*

Proof: It is enough to show that for all i and t , $w_{i,t} = w'_{i,t}$. We proceed by induction on t . For $t = 0$, $w_{i,0} = w'_{i,0} = \mathbf{1}\{i = 1\}$ for all i . For the induction step, assume that $w_{i,s} = w'_{i,s}$ for all i and all $s < t$. We then have

$$\begin{aligned} w'_{i,t} &= \sum_{i_1, \dots, i_t, i_{t+2}, \dots, i_T} w'_t(i_1, \dots, i_t, i_{t+2}, \dots, i_T) \\ &= \sum_{i_1, \dots, i_t, i_{t+2}, \dots, i_T} \exp\left(-\eta \sum_{s=1}^t l(i_s, z_s)\right) w'_0(i_1, \dots, i_t, i_{t+2}, \dots, i_T) \\ &= \sum_{i_1, \dots, i_t} \exp\left(-\eta \sum_{s=1}^t l(i_s, z_s)\right) w'_0(i_1, \dots, i_t, i) \\ &= \sum_{i_1, \dots, i_t} \exp\left(-\eta \sum_{s=1}^t l(i_s, z_s)\right) w'_0(i_1, \dots, i_t) \frac{w'_0(i_1, \dots, i_t, i)}{w'_0(i_1, \dots, i_t)} \\ &= \sum_{i_1, \dots, i_t} \exp\left(-\eta \sum_{s=1}^t l(i_s, z_s)\right) w'_0(i_1, \dots, i_t) \left[(1 - \kappa \beta_{i_t}) \mathbf{1}\{i = i_t\} + \kappa \beta_{i_t, i} \mathbf{1}\{i \neq i_t\} \right] \\ &= \sum_{i_1, \dots, i_t} \exp\left(-\eta l(i_t, z_t)\right) \exp\left(-\eta \sum_{s=1}^{t-1} l(i_s, z_s)\right) w'_0(i_1, \dots, i_t) \left[(1 - \kappa \beta_{i_t}) \mathbf{1}\{i = i_t\} + \kappa \beta_{i_t, i} \mathbf{1}\{i \neq i_t\} \right] \\ &= \sum_{i_t} \exp\left(-\eta l(i_t, z_t)\right) w'_{i_t, t-1} \left[(1 - \kappa \beta_{i_t}) \mathbf{1}\{i = i_t\} + \kappa \beta_{i_t, i} \mathbf{1}\{i \neq i_t\} \right] \\ &= \sum_{i_t} \exp\left(-\eta l(i_t, z_t)\right) w_{i_t, t-1} \left[(1 - \kappa \beta_{i_t}) \mathbf{1}\{i = i_t\} + \kappa \beta_{i_t, i} \mathbf{1}\{i \neq i_t\} \right] \quad (\text{induction hypothesis}) \\ &= \sum_{i_t} v_{i_t, t-1} \left[(1 - \kappa \beta_{i_t}) \mathbf{1}\{i = i_t\} + \kappa \beta_{i_t, i} \mathbf{1}\{i \neq i_t\} \right] \\ &= (1 - \kappa \beta_i) v_{i,t} + \kappa \sum_{j=1}^N \beta_j v_{j,t} \quad (\text{notice that } \beta_{ii} = 0) \\ &= w_{i,t} \end{aligned}$$

■

Lemma 9. For all $T \geq 1$, if $l \in [0, 1]$ and we run the exponentially weighted forecaster over compound experts as described before, we will have

$$\sum_{t=1}^T \sum_{i=1}^N p'_{i,t} l(i, z_t) \leq \frac{1}{\eta} \ln \frac{1}{W'_T} + \frac{\eta}{8} T$$

Proof: First, notice that

$$\begin{aligned} W'_t &= \sum_{i=1}^N w'_{i,t} = \sum_{i=1}^N \sum_{i_1, \dots, i_t, i_{t+2}, \dots, i_T} w'_t(i_1, \dots, i_t, i, i_{t+2}, \dots, i_T) \\ &= \sum_{i_1, \dots, i_T} w'_t(i_1, \dots, i_T) \end{aligned}$$

Then, we also have

$$\begin{aligned} \sum_{i=1}^N p'_{i,t} l(i, z_t) &= \sum_{i_t} l(i_t, z_t) \frac{w'_{i_t,t}}{W'_{t-1}} \\ &= \sum_{i_t} l(i_t, z_t) \frac{\sum_{i_1, \dots, i_{t-1}, i_{t+1}, \dots, i_T} w'_{t-1}(i_1, \dots, i_T)}{W'_{t-1}} \\ &= \sum_{i_1, \dots, i_T} \frac{w'_t(i_1, \dots, i_T)}{W'_{t-1}} l(i_t, z_t) \end{aligned}$$

Then we can directly apply Lemma 5.1 in [32, Chapter 5] by noticing that $W'_0 = 1$. ■

Proof of Theorem 2

Proof: According to Lemma 8, it is equivalent to prove the bound for the equivalent exponentially weighted forecaster. There we have

$$\begin{aligned} w'_0(i_1, \dots, i_T) &= \mathbf{1}\{i_1 = 1\} \prod_{t=1}^{T-1} \left[(1 - \kappa \beta_{i_t}) \mathbf{1}\{i_{t+1} = i_t\} + \kappa \beta_{i_t, i_{t+1}} \mathbf{1}\{i_{t+1} \neq i_t\} \right] \\ &\geq (1 - \kappa d)^{T-k-1} \kappa^k \end{aligned}$$

for all the sequence (i_1, \dots, i_T) with size $\leq k$ and transitions restricted on the graph.

Also, we have

$$\ln w'_T(i_1, \dots, i_T) = \ln w'_0(i_1, \dots, i_T) - \eta \sum_{t=1}^T l(i_t, z_t)$$

And $W'_T \geq w'_T(i_1, \dots, i_T)$. Then by Lemma 9 and some simple manipulations, we will get

$$\sum_{t=1}^T \left(\sum_{i=1}^N l(i, z_t) p_{i,t} - l(i_t, z_t) \right) \leq \frac{1}{\eta} (T - k - 1) \log \frac{1}{1 - \kappa d} + \frac{1}{\eta} k \log \frac{1}{\kappa} + \eta \frac{T}{8}$$

In order to minimize the above bound with respect to (κ, η) , we first take derivative w.r.t κ and setting it to zero gives $\kappa = k / ((T - 1)D)$. Then the bound becomes $S/\eta + \eta T/8$. Minimizing w.r.t. η gives the minimal bound $\sqrt{TS/2}$ with $\eta = \sqrt{8S/T}$. This concludes the proof. ■

REFERENCES

- [1] H. Akaike, "Statistical predictor identification," *Ann. Inst. Statist. Math.*, vol. 22, no. 1, pp. 203–217, 1970.
- [2] —, "Information theory and an extension of the maximum likelihood principle," in *Selected Papers of Hirotugu Akaike*. Springer, 1998, pp. 199–213.
- [3] K. Takeuchi, "Distribution of informational statistics and a criterion of model fitting," pp. 12–18, 1976.
- [4] J. Shao, "An asymptotic theory for linear model selection," *Statist. Sinica*, vol. 7, no. 2, pp. 221–242, 1997.
- [5] J. Ding, V. Tarokh, and Y. Yang, "Optimal variable selection in regression models," <http://jding.org/jie-uploads/2017/03/variable-selection.pdf>, 2016.

- [6] R. Shibata, "Asymptotically efficient selection of the order of the model for estimating parameters of a linear process," *Ann. Statist.*, vol. 8, no. 1, pp. 147–164, 1980.
- [7] —, "An optimal selection of regression variables," *Biometrika*, vol. 68, no. 1, pp. 45–54, 1981.
- [8] K. P. Burnham and D. R. Anderson, *Model selection and multimodel inference: a practical information-theoretic approach*. Springer Science & Business Media, 2003.
- [9] G. Schwarz, "Estimating the dimension of a model," *Ann. Statist.*, vol. 6, no. 2, pp. 461–464, 1978.
- [10] C.-K. Ing and C.-Z. Wei, "Order selection for same-realization predictions in autoregressive processes," *Ann. Statist.*, vol. 33, no. 5, pp. 2423–2474, 2005.
- [11] H. Akaike, "Fitting autoregressive models for prediction," *Ann. Inst. Statist. Math.*, vol. 21, no. 1, pp. 243–247, 1969.
- [12] G. Casella, F. J. Girón, M. L. Martínez, and E. Moreno, "Consistency of bayesian procedures for variable selection," *Ann. Stat.*, pp. 1207–1228, 2009.
- [13] M. H. Hansen and B. Yu, "Model selection and the principle of minimum description length," *J. Amer. Statist. Assoc.*, vol. 96, no. 454, pp. 746–774, 2001.
- [14] E. J. Hannan and B. G. Quinn, "The determination of the order of an autoregression," *J. Roy. Statist. Soc. Ser. B*, vol. 41, no. 2, pp. 190–195, 1979.
- [15] J. Rissanen, "Stochastic complexity and modeling," *Ann. Statist.*, pp. 1080–1100, 1986.
- [16] C.-Z. Wei, "On predictive least squares principles," *Ann. Statist.*, pp. 1–42, 1992.
- [17] C. L. Mallows, "Some comments on c p," *Technometrics*, vol. 15, no. 4, pp. 661–675, 1973.
- [18] R. Nishii *et al.*, "Asymptotic properties of criteria for selection of variables in multiple regression," *Ann. Stat.*, vol. 12, no. 2, pp. 758–765, 1984.
- [19] R. Rao and Y. Wu, "A strongly consistent procedure for model selection in a regression problem," *Biometrika*, vol. 76, no. 2, pp. 369–374, 1989.
- [20] P. Craven and G. Wahba, "Smoothing noisy data with spline functions," *Numerische Mathematik*, vol. 31, no. 4, pp. 377–403, 1978.
- [21] J. Ding, V. Tarokh, and Y. Yang, "Bridging AIC and BIC: a new criterion for autoregression," *IEEE Trans. Inf. Theory*, 2017.
- [22] S. Geisser, "The predictive sample reuse method with applications," *J. Amer. Statist. Assoc.*, vol. 70, no. 350, pp. 320–328, 1975.
- [23] P. Burman, "A comparative study of ordinary cross-validation, v-fold cross-validation and the repeated learning-testing methods," *Biometrika*, vol. 76, no. 3, pp. 503–514, 1989.
- [24] J. Shao, "Linear model selection by cross-validation," *J. Amer. Statist. Assoc.*, vol. 88, no. 422, pp. 486–494, 1993.
- [25] P. Zhang, "Model selection via multifold cross validation," *Ann. Stat.*, pp. 299–313, 1993.
- [26] M. Stone, "An asymptotic equivalence of choice of model by cross-validation and akaike's criterion," *J. R. Stat. Soc. Ser. B*, pp. 44–47, 1977.
- [27] Y. Yang, "Can the strengths of AIC and BIC be shared? a conflict between model identification and regression estimation," *Biometrika*, vol. 92, no. 4, pp. 937–950, 2005.
- [28] Y. Zhang and Y. Yang, "Cross-validation for selecting a model selection procedure," *J. Econometrics*, vol. 187, no. 1, pp. 95–112, 2015.
- [29] A. W. Van der Vaart, *Asymptotic statistics*. Cambridge university press, 2000, vol. 3.
- [30] H. White, "Maximum likelihood estimation of misspecified models," *Econometrica*, pp. 1–25, 1982.
- [31] R. R. Saxe, S. Whitfield-Gabrieli, J. Scholz, and K. A. Pelphrey, "Brain regions for perceiving and reasoning about other people in school-aged children," *Child development*, vol. 80, no. 4, pp. 1197–1209, 2009.
- [32] N. Cesa-Bianchi and G. Lugosi, *Prediction, learning, and games*. Cambridge university press, 2006.