# Consistent information criteria for regularized regression and loss-based learning problems

Qingyuan Zhang[1] [iD], Hien Duy Nguyen[2,3] [iD]

[1]School of Mathematics and Physics, University of Queensland, St Lucia, QLD 4072, Australia
[2]School of Computing, Engineering and Mathematical Sciences, La Trobe University, Bundoora, VIC 3086, Australia
[3]Institute of Mathematics for Industry, Kyushu University, Nishi Ward, Fukuoka 819-0395, Japan
Email: h.nguyen5@latrobe.edu.au

## Abstract

Many problems in statistics and machine learning can be formulated as model selection problems, where the goal is to choose an optimal parsimonious model among a set of candidate models. It is typical to conduct model selection by penalizing the objective function via information criteria (IC), as with the pioneering work by Akaike and Schwarz. Via recent work, we propose a generalized IC framework to consistently estimate general loss-based learning problems. In this work, we propose a consistent estimation method for Generalized Linear Model (GLM) regressions by utilizing the recent IC developments. We advance the generalized IC framework by proposing model selection problems, where the model set consists of a potentially uncountable set of models. In addition to theoretical expositions, our proposal introduces a computational procedure for the implementation of our methods in the finite sample setting, which we demonstrate via an extensive simulation study.

**Keywords**: Information criteria; model selection; generalized linear models; least absolute shrinkage and selection operator; regularization; asymptotic theory

## 1 Introduction

Model selection is a common problem in many statistical learning and machine learning applications. The typical objective is to select an appropriate model from a set of candidate models on the basis of sample data. Model selection criteria provide a methodology for achieving such a goal, where a good criterion defines a trade-off between the quality of a model's fit to the sample data and the complexity of the model.

Archetypal examples of such criteria are the Akaike Information Criteria (AIC; Akaike 1974) and the Bayesian Information Criteria (BIC; Schwarz 1978), which are among the earliest proposed methods and remain among the most popular approaches. Detailed expositions on the subject of model selection can be found in the volumes of McQuarrie and Tsai [1998], Vapnik [1998], Vapnik [2000], Claeskens and Hjort [2001], and Bühlmann and Van de Geer [2011].

Our study focuses on establishing model selection schemes that provide consistency guarantees. The seminal work of Sin and White [1996] provides consistency results for a broad class of IC for generic loss-based learning problems. However, their theorems impose strong conditions on the objective function and the candidate models that are challenging to verify in practice. In Nguyen [2023], the IC of Sin and White [1996] for generic risk minimization problems are further studied, yielding a generic framework for model selection, which we refer to as PanIC. The PanIC framework, together with the works of Sin and White [1996] and Baudry [2015], consider IC in generic loss-based learning problems. This contrasts with other works in the literature, which typically assume the learning objective to be a negative log-likelihood function [Leroux, 1992, Hui et al., 2015], or a negative composite, pseudo, or quasi log-likelihood function [Varin and Vidoni, 2005, Gao and Song, 2010, Ng and Joe, 2014, Hui, 2021].

In this work, we seek to specialize and refine the PanIC framework towards consistent estimation methods for general regression problems. Consistency conditions for parameter estimators and minimal risks of regression problems under the Empirical Risk Minimization (ERM) approach are well established and can be found, for example, in Vapnik [1998], Vapnik [2000], and Shapiro et al. [2021]. However, it is observed that learning model complexity according to the ERM principle tends to produce models that overfit the data. This fact has been widely recognized in the literature and has motivated the broad study of regularization techniques. Examples of such methods include $l_2$ regularization, leading to the Ridge regression problem [Hoerl and Kennard, 1970]; $l_1$ regularization, leading to the Lasso problem [Tibshirani, 1996]; and a convex combination of of $l_2$ and $l_1$ penalties, leading to the elastic-net problem [Zou and Hastie, 2005]. Results regarding the consistency of estimators obtained from solving such regularization problems, and others, can be found, for example, in Zhao and Yu [2006], Yuan and Lin [2007], Jia and Yu [2010], and Hastie et al. [2015].

To address the consistency of regularized regression problems, we extend the PanIC framework from the finite set of models setting of Nguyen [2023], to be able to handle sets of models that may be uncountably large. As we will argue in Section 2, in many learning problems, considering infinitely many models is both natural and advantageous. It is therefore useful to derive conditions under which consistency may be proved in such settings.

A summary of our contributions is as follows:

1. We outline sufficient conditions for consistently estimating various regression problems, including linear, logistic, Poisson, and Gamma regressions. Notably, we provide details of our derivations that can be followed by the reader to derive conditions for regression problems not addressed in this work.

2. We expand upon the theory of PanIC by demonstrating its consistency in certain model selection problems involving infinitely many models. Although not expansive enough to cover all model selection problems on uncountable spaces of models, we argue that our formulation is nevertheless natural for a ubiquitous class of learning problems and is thus of broad practical interest.

3. We propose a computational method for implementing PanIC in regression problems, and analyze the performance of our routine against popular benchmarks in finite sample studies. We note that our method maintains the asymptotic consistency as a PanIC estimator and performs effectively in simulated regression problems.

We note in particular that our work differs from previous works in that the consistency condition we derive concerns the norm of the optimal model, rather than consistency in terms of parameter estimation or model selection, as per Fan and Li [2001] and Zhao and Yu [2006]. A related problem appears in the literature in the work of Massart and Meynet [2011]. However, the emphasis of Massart and Meynet [2011] is on the oracle inequalities related to the Lasso problem, while our focus lies on the asymptotic consistency of a more general set of regularized regression problems without first obtaining finite sample oracles.

The rest of this article proceeds as follows: in Section 2, the PanIC framework is reviewed and our proposed estimation method is given with theoretical justification. In Section 3, a simulated study of PanIC in linear regression problems is reported with some discussions. Finally, concluding remarks and an outline of future research directions are given in Section 4.

# 2   Method

We review the PanIC framework in Section 2.1. The connection between regularized regression estimation and PanIC is established in Section 2.2. A formal justification of our proposed estimation method is presented in Section 2.3. Lastly, a computational method is given in Section 2.4.

## 2.1   PanIC

Consider the following setting of a general model selection problem. Let $(\Omega, \mathcal{F}, \mathbb{P})$ be the underlying probability space and $X : \Omega \to \mathbb{X} \subset \mathbb{R}^d$ be a random vector on the pushforward probability space $(\mathbb{X}, \mathcal{B}(\mathbb{X}), \mathbb{P}_X)$, where $\mathcal{B}(\mathbb{X})$ denotes the Borel $\sigma$-algebra of $\mathbb{X}$. We observe a sequence of independent and

identically distributed (i.i.d) random variables $(X_i)_{i\in[n]}$ on the probability space $(\mathbb{X}, \mathcal{B}(\mathbb{X}), \mathbb{P}_X)$, where $[n] = \{1, \ldots, n\}$. The set of candidate models is given by a sequence of hypotheses $(\mathcal{H}_k)_{k\in[m]}$. Each $\mathcal{H}_k$ defines a functional space identified by some parameter vector $\beta_k \in \mathbb{T}_k \subset \mathbb{R}^{q_k}$, where $(q_k)_{k\in[m]} \subset \mathbb{N}$. That is, for each $k \in [m]$,

$$\mathcal{H}_k = \{h_k(\cdot; \beta_k) : \mathbb{X} \to \mathbb{R} : \beta_k \in \mathbb{T}_k\}.$$

This construction admits two sources of flexibility. For each hypothesis $\mathcal{H}_k$ with index $k \in [m]$, the functional form of $h_k$ and the parameter space $\mathbb{T}_k$ are both allowed to vary. We further define a loss function $\ell : \mathbb{R} \to \mathbb{R}$, and for each $k \in [m]$, denote,

$$R_{k,n}(\beta_k) = \frac{1}{n} \sum_{i=1}^{n} \ell(h_k(X_i; \beta_k))$$

and

$$r_k(\beta_k) = \mathbb{E}\{\ell(h_k(X; \beta_k))\}$$

as the empirical risk and the expected risk, respectively. We shall abbreviate $\ell(h_k(X; \beta_k))$ by $\ell_k(X; \beta_k)$, which implicitly defines $\ell_k(\cdot; \beta_k) : \mathbb{X} \to \mathbb{R}$ and $\ell_k(x; \cdot) : \mathbb{T}_k \to \mathbb{R}$ for each hypothesis $\mathcal{H}_k$ with index $k \in [m]$. Examples of some classic model selection problems under this formulation are presented by Nguyen [2023].

Further, we call $k^*$ the optimal hypothesis index,

$$k^* = \min \arg\min_{k\in[m]} \left\{ \min_{\beta_k \in \mathbb{T}_k} r_k(\beta_k) \right\}.$$

The corresponding hypothesis space $\mathcal{H}_{k^*}$ is known as the class of parsimonious models. Let $\hat{K}_n$ denote an estimator of $k^*$. A model selection scheme is said to be consistent if

$$\plim_{n\to\infty} \hat{K}_n = k^*, \tag{1}$$

where plim denotes convergence in probability (as per Amemiya, 1985, Sec. 3.2). In the case where the set of candidate models is finite, this is equivalent to

$$\lim_{n\to\infty} \mathbb{P}\left(\hat{K}_n = k^*\right) = 1.$$

Now, suppose that the model selection problem satisfies conditions A1–A3, for each $k \in [m]$:

A1 : $\ell_k(x; \beta_k)$ is Caratheodory in the sense that $\ell_k(x; \cdot) : \mathbb{T}_k \to \mathbb{R}$ is continuous for each $x \in \mathbb{X}$, and $\ell_k(\cdot; \beta_k) : \mathbb{X} \to \mathbb{R}$ is $\mathcal{B}(\mathbb{X})$-measurable for each $\beta_k \in \mathbb{T}_k$.

A2 : $\mathbb{T}_k$ is compact and there exists some $\tau_k \in \mathbb{T}_k$ such that $\ell_k(X; \tau_k)^2$ is square integrable. That is,

$$\mathbb{E}\{\ell_k(X; \tau_k)^2\} < \infty.$$

A3 : There exists some measurable function $\mathcal{G}_k : \mathbb{X} \to \mathbb{R}_{\geq 0}$, such that $\mathbb{E}\{\mathcal{G}_k(X)^2\} < \infty$ and,

$$|\ell_k(x; \beta_k) - \ell_k(x; \tau_k)| \leq \mathcal{G}_k(x)\|\beta_k - \tau_k\|, \quad (\forall \beta_k, \tau_k \in \mathbb{T}_k, \text{ a.e. } x \in \mathbb{X}).$$

A PanIC estimator $\hat{K}_n$ is defined as

$$\hat{K}_n = \min \arg\min_{k \in [m]} \left\{ \min_{\beta_k \in \mathbb{T}_k} R_{k,n}(\beta_k) + P_{k,n} \right\}, \tag{2}$$

where $P_{k,n} : \Omega \to \mathbb{R}_{\geq 0}$ is allowed to be a stochastic, non-negative function that satisfies conditions B1 and B2, for each $k \in [m]$:

B1 : $P_{k,n} > 0$, for $n \in \mathbb{N}$, and $P_{k,n} = o_{\mathbb{P}}(1)$ as $n \to \infty$.

B2 : If $k < l$, then $\sqrt{n}\{P_{l,n} - P_{k,n}\} \xrightarrow{\mathbb{P}} \infty$, as $n \to \infty$.

Under these assumptions, Theorem 1 of Nguyen [2023] shows that PanIC estimators are consistent. Equivalently, an estimator taking the form of (2) consistently solves the model selection problem, provided that conditions A1–A3, B1, and B2 are satisfied.

Nguyen [2023] also introduced the notion of a BIC-like criterion, one that satisfies the condition B2*instead of B2,

B2* : If $k < l$, then $n\{P_{l,n} - P_{k,n}\} \xrightarrow{\mathbb{P}} \infty$, as $n \to \infty$,

Theorem 2 of Nguyen [2023] shows the consistency of BIC-like criteria under additional assumptions C1–C5. For $k \in [m]$,

C1 : $r_k$ is Lipschitz continuous on $\mathbb{T}_k$, and it is twice differentiable and uniquely minimized at some $\beta_k^* \in \mathbb{T}_k$.

C2 : $\ell_k(x, y; )$ is Lipschitz continuous and differentiable at $\beta_k^*$, for almost all $(x, y) \in \mathbb{X}$.

C3 : The set $\mathbb{T}_k$ is second order regular at $\beta_k^*$.

C4 : The quadratic growth condition holds at $\beta_k^*$.

C5 : If $r_k(\beta_k^*) = r_{k^*}(\beta_{k^*}^*)$, then $n(R_{k,n}(\beta_k^*) - R_{k,n}(\beta_{k^*}^*)) = O_{\mathbb{P}}(1)$.

As an important point to be revisited in Section 3, we note that conditions B1, B2, and B2* remain satisfied when the penalty function $P_{k,n}$ is multiplied by any positive constant. We will express a valid PanIC penalty in the form

$$P_{k,n} = \kappa \times \text{pen}_{\text{shape}}(k, n) \tag{3}$$

and call $\text{pen}_{\text{shape}} : \mathbb{N} \times \mathbb{N} \to \mathbb{R}_+$ the penalty shape of $P_{k,n}$. From this expression, it is clear that any valid PanIC penalty is defined up to a multiplicative constant. This observation holds practical significance, as it implies that the constant $\kappa$ should be treated as a hyperparameter, requiring calibration in finite sample settings.

## 2.2   Regression Problems

This section connects the regression estimation problem with a model selection problem in the PanIC framework. In regression problems, we observe some response $y \in \mathcal{Y} \subset \mathbb{R}$, and covariate vector $x \in \mathcal{X} \subset \mathbb{R}^d$, for some positive integer $d$. Each $x \in \mathcal{X}$ is assumed to be related to a $y \in \mathcal{Y}$ through some stochastic dependencies. The goal is to infer from an observed sample $(x_i, y_i)_{i \in [n]}$ this underlying relationship. In the Generalized Linear Model (GLM) framework, this relationship is assumed to be defined by a single functional form $h(\cdot; \beta_0, \beta) : \mathcal{X} \to \mathcal{Y}$ that is parametrized by a coefficient vector $\beta$ taking values in $\mathbb{R}^d$ and a bias term $\beta_0$ taking values in $\mathbb{R}$. Further, we assume the existence of a loss function $\ell : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$. Our focus is on linear, logistic, Poisson, and Gamma regressions. These regressions collectively model a comprehensive range of common response types, including continuous, binary, count, and positive continuous data. Table 1 reviews these regression problems,

Table 1: Mean and loss functions for generalized linear regression models. Here, $h = h(x; \beta_0, \beta)$.

| Regression | Mean function | Loss function |
|---|---|---|
| Linear | $h = \beta_0 + \beta^\top x$ | $\ell(x, y; \beta_0, \beta) = (y - h)^2$ |
| Logistic | $h = (1 + \exp(-(\beta_0 + \beta^\top x)))^{-1}$ | $\ell(x, y; \beta_0, \beta) = -\{y \log h + (1 - y) \log(1 - h)\}$ |
| Poisson | $h = \exp(\beta_0 + \beta^\top x)$ | $\ell(x, y; \beta_0, \beta) = -\{y \log h - h - \log(y!)\}$ |
| Gamma | $h = \exp(\beta_0 + \beta^\top x)$ | $\ell(x, y; \beta_0, \beta, \nu) = -\{-\nu \log(h) + (\nu - 1) \log(y) - y\nu/h\}$ |

Following the GLM convention, the loss function of each regression is defined by the negative log-likelihood of the corresponding distribution. Specifically, for the Gamma regression, we use the mean parameterization

of the Gamma density function

$$p(y; \mu, \nu) = \frac{1}{\Gamma(\nu)} \left(\frac{\nu}{\mu}\right)^{\nu} y^{\nu-1} \exp\left(-y\frac{\nu}{\mu}\right),$$

where $\mathbb{E}(Y) = \mu$ and $\text{Var}(Y) = \mu^2/\nu$.

Under the Empirical Risk Minimization principle (ERM), a regression problem is estimated by minimizing the empirical risk over the parameter space. Assuming a unique solution exists, learning according to the ERM principle yields the estimates $\hat{\beta}_0$ and $\hat{\beta}$, where

$$\left(\hat{\beta}_0, \hat{\beta}\right) = \operatorname*{arg\,min}_{\beta_0 \in \mathbb{R}, \beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^{n} \ell\left(x_i, y_i; \beta_0, \beta\right).$$

To correct the overfitting tendency of the ERM, the Lasso regression, the Ridge regression, or the Elastic-net regression are often solved instead. We note that both Lasso and Ridge regression problems can be represented in a unified manner with the following formulation. Define

$$\begin{aligned} \mathcal{H} &= \{h(\cdot; \beta_0, \beta) : \mathcal{X} \to \mathcal{Y}, (\beta_0, \beta) \in \mathbb{T}\}, \\ \mathbb{T} &= \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \, \beta \in \mathbb{R}^d : \|\beta\|_p \leq C, \, C \in \mathbb{R}_+\}, \end{aligned} \tag{4}$$

where $\mathbb{U}$ is some sufficiently large compact subset of $\mathbb{R}$ that contains the optimal bias term $\beta_0^*$ in its interior. Setting $p = 1$ in (4) corresponds to the Lasso problem [Tibshirani, 1996], and setting $p = 2$ corresponds to the Ridge regression problem [Hoerl and Kennard, 1970]. Slightly modifying (4), we also retrieve the Elastic-net problem [Zou and Hastie, 2005]:

$$\begin{aligned} \mathcal{H} &= \{h(\cdot; \beta_0, \beta) : \mathcal{X} \to \mathcal{Y}, (\beta_0, \beta) \in \mathbb{T}\}, \\ \mathbb{T} &= \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \, \beta \in \mathbb{R}^d : \alpha\|\beta\|_1 + (1-\alpha)\|\beta\|_2^2 \leq C, \, C \in \mathbb{R}_+\}, \end{aligned} \tag{5}$$

for some $\alpha \in (0, 1)$. Notably, both (4) and (5) are equivalent to a single hypothesis within the PanIC framework, dependent on a hyperparameter $C$. If we consider a sequence of positive real numbers denoted by $(C_k)_{k \in [m]}$, we obtain a sequence of hypotheses:

$$\mathcal{H}_k = \{h(\cdot; \beta_0, \beta) : \mathcal{X} \to \mathcal{Y}, (\beta_0, \beta) \in \mathbb{T}_k\}$$

$$\mathbb{T}_k = \begin{cases} \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \beta \in \mathbb{R}^d : \|\beta\|_1 \leq C_k\}, & \text{for Lasso regression,} \\ \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \beta \in \mathbb{R}^d : \|\beta\|_2 \leq C_k\}, & \text{for Ridge regression,} \\ \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \beta \in \mathbb{R}^d : \alpha\|\beta\|_1 + (1-\alpha)\|\beta\|_2^2 \leq C_k\}, & \text{for Elastic-net regression.} \end{cases} \tag{6}$$

Then, utilizing the theory of PanIC, we can construct a consistent estimator of the optimal hypothesis index $k^*$. Thus, the PanIC framework provides a unifying approach to consistently solve constrained regression problems.

## 2.3 Consistency

We outline sufficient conditions for consistent model selection in linear, logistic, Poisson, and Gamma regression problems. These conditions are derived from fulfilling conditions A1–A3 of PanIC, with complete derivations given in the Appendix.

**Proposition 1.** *Assume that $(X_i, Y_i)_{i \in [n]}$ is an i.i.d sequence and specific conditions are met for each regression problem. Then, the PanIC estimator satisfies the consistency property (Eq. 1) in the corresponding problem given the specific assumptions.*

1. *Linear regression: Both the covariate vector $X$ and the response $Y$ have a finite fourth moment.*

2. *Logistic regression: The covariate vector $X$ has a finite second moment.*

3. *Poisson and Gamma regressions: The covariate vector $X$ is Gaussian or sub-Gaussian distributed, and the response $Y$ has a finite fourth moment.*

*Remark*:

i) We note that a random vector $X$ is said to have a finite $p$-th moment if $\mathbb{E}(\|X\|^p) < \infty$, where $\|\cdot\|^p$ is the $L^p$ norm.

ii) Based on our derivations, we further observe that consistency is ensured if both the covariate vector and the response are supported on a compact set. This assumption is reasonable for many real-world problems, where the data typically falls within a reasonable range.

Furthermore, we have the following result regarding a BIC-like criterion for regularized linear regression. Recall that the usual BIC criterion for linear regression is defined as

$$\hat{K}_n = \min \underset{k \in [m]}{\arg\min} \left\{ R_n(\hat{\beta}_k) + \frac{\log(n)}{n} \hat{\mathrm{df}}(\hat{\beta}_{k,n}) \right\}, \tag{7}$$

where

$$\hat{\beta}_{k,n} = \underset{\beta \in \mathbb{T}_k}{\arg\min}\ R_n(\beta)$$

and $\hat{\mathrm{df}}(\hat{\beta}_{k,n})$ is the number of non-zero coefficients in $\hat{\beta}_{k,n}$, the unbiased estimator for the degrees of freedom established by Zou et al. [2007]. Details concerning this estimator can be found in Chapter 2.11 of Bühlmann and Van de Geer [2011]. To the best of our knowledge, the consistency of the BIC has not been justified in

the setting of regularized linear regression problems. Using the results from Nguyen [2023], we can establish the consistency of the following modifications of the BIC in (7):

$$\hat{K}_n = \min \arg\min_{k \in [m]} \left\{ R_n(\hat{\beta}_{k,n}) + \frac{\log(n)}{n} \left( \kappa \tilde{\mathrm{df}}(\hat{\beta}_{k,n}) + \varepsilon C_k \right) \right\}. \tag{8}$$

Here, $\kappa$ is a non-negative constant and $\varepsilon$ is a small positive number. The term $\tilde{\mathrm{df}}(\hat{\beta}_{k,n})$ is defined as follows: $\tilde{\mathrm{df}}(\hat{\beta}_1) = \hat{\mathrm{df}}(\hat{\beta}_1)$, and for $1 < k \in [m]$,

$$\tilde{\mathrm{df}}(\hat{\beta}_{k,n}) = \begin{cases} \hat{\mathrm{df}}(\hat{\beta}_{k,n}), & \text{if } \hat{\mathrm{df}}(\hat{\beta}_{k,n}) \geq \hat{\mathrm{df}}(\hat{\beta}_{j,n}), \forall j < k, \\ \hat{\mathrm{df}}(\hat{\beta}_{k-1,n}), & \text{otherwise.} \end{cases}$$

To motivate these modifications, we note that the penalty shape of the BIC, denoted by $P_{k,n}^{\mathrm{BIC}} = \hat{\mathrm{df}}(\hat{\beta}_{k,n})$, is not strictly monotone and does not satisfy condition B2*. In contrast, the modified penalty shape presented in (8), given by $P_{k,n} = \kappa \tilde{\mathrm{df}}(\hat{\beta}_{k,n}) + \varepsilon C_k$, ensures strict monotonicity in accordance with B2*, for any $\kappa \geq 0$ and $\varepsilon > 0$.

**Proposition 2.** *In the context of linear regression, let the covariate vector and response $(X, Y)$ be distributed on the measurable space $(\mathbb{R}^d \times \mathbb{R}, \mathcal{B}(\mathbb{R}^d) \otimes \mathcal{B}(\mathbb{R}))$ with probability measure $\Pi$. Assume both $X$ and $Y$ have a finite fourth moment, and the covariance matrix $\Sigma$ of $X$ is positive definite. That is,*

$$\Sigma = \mathbb{E}\{(X - \mathbb{E}X)(X - \mathbb{E}X)^\top\}$$

*satisfies*

$$x^\top \Sigma x > 0$$

*for any $x \in \mathbb{R}^d \backslash \{0\}$. Then, the criterion specified by (8) is consistent.*

So far, we have formulated regression problems as model selection problems involving sets of finitely many models. A more natural formulation arises when we allow for an infinite number of candidate models, which occurs when the hypotheses are indexed by the entire positive real line or some compact interval within it, instead of an increasing sequence of numbers. This motivates the following result, which is given for certain sets of continuous and compact correspondences. Recall that a correspondence is a set-valued function. Specifically, a correspondence $\varphi$ from a set $\mathbb{X}$ to a set $\mathbb{Y}$ assigns to each $x \in \mathbb{X}$ a subset $\varphi(x)$ of $\mathbb{Y}$ [Aliprantis and Border, 2006, Definition 17.1].

**Theorem 1.** *Let $(X_i)_{i \in [n]}$ be an i.i.d sequence taking values in a set $\mathbb{X} \subset \mathbb{R}^d$ and $(\mathcal{H}_k)_{k \in [a,b]}$ be a set of*

*hypothesis spaces of the form*

$$\mathcal{H}_k = \{h(\cdot; \beta) : \mathbb{X} \to \mathbb{R} : \beta \in \mathbb{T}_k \subset \mathbb{R}^d\},$$

*where $a < b$. Let $\mathcal{L} : k \mapsto \mathbb{T}_k$ be a continuous and compact-valued correspondence, and assume that the parameter spaces $(\mathbb{T}_k)_{k \in [a,b]}$ are nested. Assume that A1–A3, B1 and B2 hold for each $k$. Additionally, assume $P_{k,n}$ is a continuous and strictly increasing function in $k$ for all $n \in \mathbb{N}$. Define*

$$\mathcal{K} = \operatorname*{arg\,min}_{k \in [a,b]} \left\{ \min_{\beta \in \mathbb{T}_k} r(\beta) \right\}, \quad k^* = \min_{k \in \mathcal{K}} k.$$

*Then, the PanIC estimator satisfies,*

$$\operatorname*{plim}_{n \to \infty} \hat{K}_n = k^*. \tag{9}$$

In the context of regression problems, we shall define the parameter space $\mathbb{T}_k$ of the correspondence $\mathcal{L} : k \mapsto \mathbb{T}_k$ as

$$\mathbb{T}_k = \begin{cases} \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \beta \in \mathbb{R}^d : \|\beta\|_1 \leq k\}, & \text{for Lasso regression} \\ \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \beta \in \mathbb{R}^d : \|\beta\|_2 \leq k\}, & \text{for Ridge regression} \\ \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{U}, \beta \in \mathbb{R}^d : \alpha\|\beta\|_1 + (1-\alpha)\|\beta\|_2^2 \leq k\}, & \text{for Elastic-net regression.} \end{cases} \tag{10}$$

When thusly defined, $\mathcal{L} : k \mapsto \mathbb{T}_k$ is compact-valued and continuous. Therefore, we have the following result.

**Corollary 1.** *Under the Assumptions of Theorem 1, let $(\mathcal{H}_k)_{k \in [a,b]}$ be the set of hypothesis spaces for a regression problem. Specifically, the function $h$ is the mean function in linear, logistic, Poisson, or Gamma regression, as outlined in Table 1, with the parameter spaces defined in (10). Suppose the corresponding conditions in Proposition 1 hold and $P_{k,n}$ is appropriately defined as in Theorem 1. Then, the PanIC estimator is consistent in the sense of (9).*

## 2.4   Computation

In this section, we propose a computational method designed for implementing PanIC in regression problems within a finite sample setting. This is achieved via an application of duality in optimization theory, as covered, for example, in Boyd and Vandenberghe [2004], Kloft et al. [2011], and Oneto et al. [2016]. First, we review a useful result on duality. Suppose $f : \mathbb{R}^d \to \mathbb{R}$ is a convex function and $g : \mathbb{R}^d \to \mathbb{R}_{\geq 0}$ is a

non-negative convex function. Consider the following constrained optimization problem:

$$\min \quad f(\beta)$$
$$\text{s.t.} \quad \beta \in \mathbb{R}^d : g(\beta) \leq C,$$

(11)

for some value $C \in \mathbb{R}_+$. Assume that a solution exists and a constraint qualification holds, the solutions of (11) are equivalent to that of its Lagrange dual problem [Kloft et al., 2011]:

$$\min \quad f(\beta) + \lambda g(\beta)$$
$$\text{s.t.} \quad \beta \in \mathbb{R}^d,$$

(12)

for some $\lambda \in [0, \infty)$, which is known as a regularization constant. Hence, for each $C$, there exists a $\lambda$ such that the optimal solutions of both problems coincide. The converse statement also holds, and if $f$ is strictly convex, the corresponding constrained problem is identified by the constraint $g(\beta) \leq g(\beta_\lambda)$, where $\beta_\lambda$ is the unique solution of (12) [Kloft et al., 2011]. This result allows us to define the function $\mathcal{L} : \lambda \mapsto g(\beta_\lambda)$, for all $\lambda \in [0, \infty)$. Under this setting, we can further establish the monotonicity and the continuity of the function $\mathcal{L}$.

**Lemma 1.** *Suppose $f : \mathbb{R}^d \to \mathbb{R}_{\geq 0}$ is a strictly convex function and $g : \mathbb{R}^d \to \mathbb{R}_{\geq 0}$ is a non-negative convex function with a non-empty and bounded sub-level set $L_c = \{\beta \in \mathbb{R}^d \,|\, g(\beta) \leq c\}$, for some $c \in \mathbb{R}_{\geq 0}$. Then, the mapping $\mathcal{L} : \mathbb{R}_+ \to \mathbb{R}$, given by*

$$\mathcal{L}(\lambda) = g(\beta_\lambda),$$

*is a continuous function, where*

$$\beta_\lambda = \arg\min_{\beta \in \mathbb{R}^d} f(\beta) + \lambda g(\beta).$$

**Lemma 2.** *Let $\lambda_1$ and $\lambda_2$ be non-negative real numbers, with $\lambda_2 > \lambda_1$. Under the same assumptions as those of Lemma 1, it holds that $g(\beta_{\lambda_2}) \leq g(\beta_{\lambda_1})$. Further, if $g(\beta_{\lambda_2}) < g(\beta_{\lambda_1})$, then $f(\beta_{\lambda_2}) > f(\beta_{\lambda_1})$. Else, if $g(\beta_{\lambda_2}) = g(\beta_{\lambda_1})$, then $\beta_{\lambda_2} = \beta_{\lambda_1}$.*

Our formulation of regression problems in the PanIC framework is associated with a sequence of constrained optimization problems, which can be expressed in the general form

$$\min \quad R_n(\beta_0, \beta) = \frac{1}{n} \sum_{i=1}^{n} \ell(x_i, y_i; \beta_0, \beta)$$
$$\text{s.t.} \quad \beta_0 \in \mathbb{U}, \, \beta \in \mathbb{R}^d : g(\beta) \leq C, \, C \in \mathbb{R}_+,$$

(13)

where $g$ is defined as the $l_1$ norm, the $l_2$ norm, or a combination of the two. Assuming the empirical risk function $R_n(\beta_0, \beta)$ is strictly convex, Lemma 1 and 2 allow us to use any root-finding algorithm, such as

the bisection method, to find an (approximate) Lagrange dual problem of (14). Compared to directly solving the constrained optimization problem, this approach greatly simplifies the computations. Algorithm 1 presents the computational steps for solving the Lasso problem. Generalizing this algorithm to both the Ridge regression and the Elastic-net regression is straightforward.

---

**Algorithm 1** Computing the Lasso solution using PanIC

---

**Require:** a sample of data $(x_i, y_i)_{i \in [n]}$, an increasing sequence $(C_k)_{k \in [m]}$, a function $h(\cdot; \beta_0, \beta) : \mathcal{X} \to \mathcal{Y}$, a loss function $\ell : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}$, a valid PanIC penalty function $P_{k,n}$.

1: Formulate a sequence of hypotheses $(\mathcal{H}_k)_{k \in [m]}$ with

$$\mathbb{T}_k \leftarrow \{(\beta_0, \beta) \,|\, \beta_0 \in \mathbb{R}, \ \beta \in \mathbb{R}^d : \|\beta\|_1 \leq C_k\}$$

and

$$\mathcal{H}_k \leftarrow \{h(\cdot; \beta_0, \beta) \,|\, (\beta_0, \beta) \in \mathbb{T}_k\}.$$

Associated with each hypothesis, define the optimization problem

$$
\begin{aligned}
\arg\min \quad & R_n(\beta_0, \beta) \\
\text{s.t.} \quad & (\beta_0, \beta) \in \mathbb{T}_k.
\end{aligned}
\tag{14}
$$

2: **for** $k \in [m]$ **do**
3:     Determine the Lagrange dual problem of (14) using a root-finding algorithm.
4:     Denote $\hat{\vartheta}_k$ and $\hat{\beta}_k$ as the optimal value and the optimal solution of the Lagrange dual problem.
5: **end for**
6: Compute the set of optimal hypotheses $\hat{\mathcal{K}}$, with

$$\hat{\mathcal{K}} \leftarrow \underset{k \in [m]}{\arg\min} \left\{ \hat{\vartheta}_k + P_{k,n} \right\}.$$

7: Compute the PanIC estimate

$$\hat{K}_n \leftarrow \min_{k \in \hat{\mathcal{K}}} k.$$

**return** $\hat{K}_n, \hat{\beta}_{\hat{K}_n}$

---

# 3    Simulation Study

In this section, we present a set of simulated regression problems to assess the numerical performance of PanIC. We particularly focus on linear and logistic regression.

## 3.1 Simulation Setup

For each type of regression problem, we investigate three sample sizes: $n \in \{500, 1000, 2000\}$. The performance of PanIC on each sample size is simulated $N = 500$ times to obtain summary statistics. For each simulated regression problem, we employ Algorithm 1 to compute the PanIC, and then compare its performance to that of the 5-fold Cross-Validation (CV) and the modified BIC scheme, as described in Stone [1974] and (8), respectively.

For the $j$-th regression problem, where $j \in [N]$, we simulate $n$ i.i.d random covariate vectors $(X_i^{(j)})_{i \in [n]}$ from the 20-dimensional Gaussian distribution with zero mean and the identity covariance matrix, $i.e.$ $N(\mathbf{0}, \boldsymbol{I}_{20})$. Additionally, within the sparse learning framework, we generate a sparse true model vector $\beta^{*(j)}$ from $N(\mathbf{0}, \boldsymbol{I}_{20}^*)$, where

$$\boldsymbol{I}_{20}^* = \left( \begin{array}{c|c} \boldsymbol{I}_{10} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right).$$

This distribution ensures $\beta^{*(j)}$ is a 20-dimensional vector with 10 active variables. For each observation $i$, where $i \in [n]$, the response $Y_i^{(j)}$ is generated as a random variable following a normal distribution $N(\beta^{*(j)^\top} X_i^{(j)}, \sigma^2)$, in the context of linear regressions, with $\sigma^2 > 0$. We further investigate three noise levels, given by $\sigma \in \{1, 2, 5\}$, to assess the performance of PanIC across different noise settings. For logistic regressions, the response is generated as a random variable following a Bernoulli distribution $\text{Bernoulli}(p_i^{(j)})$, where

$$p_i^{(j)} = f(\beta^{*(j)^\top} X_i^{(j)})$$

and

$$f(x) = \frac{1}{1 + \exp(-x)}.$$

To explore a range of model complexities, we construct an evenly spaced sequence $(C_k^{(j)})_{k \in [m]} \subset \mathbb{R}_+$, where $C_1^{(j)} = 0$, $C_m^{(j)} = \|\hat{\beta}^{(j)}\|_1$, and $\hat{\beta}^{(j)}$ is the solution obtained by minimizing the empirical risk function.

As the final input for Algorithm 1, we seek a valid penalty shape and a value of the hyperparameter $\kappa$, as identified in the form of (3). In our context, a natural choice for the penalty shape is

$$\text{pen}_{\text{shape}}(k, n) = C_k \sqrt{\frac{\log n}{n}}. \tag{15}$$

Verifying that (15) qualifies as a valid PanIC penalty is straightforward, as it satisfies conditions B1 and B2.

## 3.2 Hyperparameter Calibration for PanIC

The optimal calibration of $\kappa$ remains an open-ended task. To this end, we introduce several statistics, which help us evaluate a learning method's performance in finite sample settings. For each of the $N = 500$ simulation problems, we first employ an error function defined by

$$\text{error}(\hat{\beta}_{\hat{k}}^{(j)}) = \|\hat{\beta}_{\hat{k}}^{(j)}\|_1 - \|\beta^{*(j)}\|_1,$$

where $\hat{\beta}_{\hat{k}}^{(j)}$ denotes the chosen model of the $j$-th simulated problem, for $j \in [N]$. We use the signed value of this function to indicate the position of the selected model relative to the truth, and its absolute value to indicate the model's proximity to that truth. We note that this error function is of primary interest. We have established its asymptotic properties via the PanIC theory; however, its behavior in finite samples remains unknown. As complementary statistics, we report the number of active variables and the number of wrongly selected variables[1] of the chosen model $\hat{\beta}_{\hat{k}}^{(j)}$, denoted as #var and as #w.var, respectively. Including these metrics enables a further understanding of PanIC's finite sample behavior, especially regarding its model selection capacity. However, we note that the theory for PanIC does not provide insights into the behavior of these metrics, whether in finite samples or asymptotically.

We experiment with different values of $\kappa$ and record the preceding statistics of their learned model. The performance of each $\kappa$ in ten randomly selected regression problems is visualized in Figure 1.



Figure 1: Performance of different $\kappa$ values in ten randomly selected linear regression problems of the $N = 500$ simulations.

---

[1]A variable is considered wrongly selected in two situations: (i) its coefficient in $\beta^{*(j)}$ is non-zero, but its coefficient in $\hat{\beta}_{\hat{k}}^{(j)}$ is zero; (ii) its coefficient in $\hat{\beta}_{\hat{k}}^{(j)}$ is non-zero, but its coefficient in $\beta^{*(j)}$ is zero.

## 3.3 Results and Discussion

Concerning the calibration of $\kappa$, we make the following observations. First, for our primary objective of estimating the norm of the optimal model, small values of $\kappa$ appear to be effective, as indicated by our simulation results. This is clear from the plots in the second row of Figure 1, where it can be seen that small $\kappa$ values lead to a near-optimal value of the absolute error function. However, for model selection, this range of $\kappa$ values is inadequate. The plots in the first row of Figure 1 illustrate that models learned with small $\kappa$ values are dense, which is in direct contrast to the sparsity of the true model.

In general, calibrating $\kappa$ for reasonable performance on #w.var remains a challenging task. As shown in Figure 1, the #w.var statistic's sensitivity to $\kappa$ values varies significantly. For problems characterized by low noise, or a high signal-to-noise ratio, the learned model performs well across a broad range of $\kappa$ values. However, this robustness diminishes in high noise problems. For $\sigma = 5$, the range of $\kappa$ values that yield reasonable performance is significantly narrowed. Moreover, it appears impossible to identify a universal range of $\kappa$ values that performs well across different problem settings.

Tables 2–5 report summary statistics for models learned under PanIC, the modified BIC, and the CV schemes, in both linear and logistic regression scenarios. For PanIC, $\kappa$ is set to the value that yields optimal performance in #w.var for each respective simulation setting, as determined by our calibration results. In the case of the modified BIC, $\kappa$ is fixed at 1, on account of the similarity to the BIC (Eq. 7) as shown in (8). According to the theory for PanIC, both the PanIC and the modified BIC are asymptotically consistent in the linear setting, with PanIC also known to be consistent in the logistic setting. This aligns with the results reported in Tables 2 and 4. More interesting are the results in Tables 3 and 5, where we observe notable differences between the two schemes in their model selection capacity. It is evident that the modified BIC, denoted as $\text{BIC}_\text{m}$, is effective in selecting the correct model across all simulation settings. In comparison, while PanIC has the potential to be effective in model selection, this efficacy is conditional on the understanding of the problem, which enables prior calibration of $\kappa$.

We note that both $\text{BIC}_\text{m}$ and PanIC are better suited for learning under sparsity than 5-fold CV. Tables 3 and 5 indicate that specific $\kappa$ values for the modified BIC and PanIC lead to sparse solutions that closely approximate the true model. In contrast, CV consistently selects denser models. CV's ineffectiveness for learning under sparsity is briefly discussed by Bühlmann and Van de Geer [2011], where it is noted that the CV scheme operates on optimizing prediction accuracy, which is often in conflict with variable selection. The goal of the latter is to recover the set of active variables, which often requires a sufficiently large regularization constant $\lambda$ that is not optimal for prediction. As a final remark on comparing PanIC with CV, our profiling demonstrates that PanIC is computationally more efficient than CV for sample sizes ranging from 500 to 2000, as visualized in Figure 2.

Table 2: error and |error| statistics of simulated linear regression problems

| | $n = 500$ | | | | $n = 1,000$ | | | | $n = 2,000$ | | | |
| | error | | |error| | | error | | |error| | | error | | |error| | |
| | mean | se | mean | se | mean | se | mean | se | mean | se | mean | se |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | $\sigma = 1$ | | | | | | |
| PanIC | -1.1907 | 0.0463 | 1.1907 | 0.0463 | -0.7581 | 0.0050 | 0.7581 | 0.0050 | -0.5598 | 0.0035 | 0.5598 | 0.0035 |
| BIC$_m$ | -0.4117 | 0.0104 | 0.4218 | -0.2884 | -0.2884 | 0.0079 | 0.3014 | 0.0068 | -0.2233 | 0.0053 | 0.2259 | 0.0051 |
| CV | 0.3012 | 0.0081 | 0.3061 | 0.0077 | 0.2308 | 0.0057 | 0.2332 | 0.0055 | 0.1647 | 0.0039 | 0.1674 | 0.0036 |
| | | | | | | $\sigma = 2$ | | | | | | |
| PanIC | -1.8612 | 0.0211 | 1.8612 | 0.0211 | -1.3461 | 0.0149 | 1.3461 | 0.0149 | -0.9931 | 0.0072 | 0.9931 | 0.0072 |
| BIC$_m$ | -0.8315 | 0.0224 | 0.8578 | 0.0203 | -0.5906 | 0.0159 | 0.6095 | 0.0144 | -0.4431 | 0.0103 | 0.4533 | 0.0094 |
| CV | 0.6648 | 0.0158 | 0.6707 | 0.0153 | 0.4607 | 0.0110 | 0.4673 | 0.0104 | 0.3221 | 0.0081 | 0.3283 | 0.0076 |
| | | | | | | $\sigma = 5$ | | | | | | |
| PanIC | -3.3285 | 0.0731 | 3.3285 | 0.0731 | -2.3476 | 0.0285 | 2.3476 | 0.0285 | -1.7429 | 0.0181 | 1.7429 | 0.0181 |
| BIC$_m$ | -2.3934 | 0.0544 | 2.4285 | 0.0512 | -1.6987 | 0.0391 | 1.7285 | 0.0364 | -1.1329 | 0.0268 | 1.1547 | 0.0248 |
| CV | 1.7180 | 0.0372 | 1.7246 | 0.0366 | 1.1857 | 0.0284 | 1.2049 | 0.0268 | 0.8681 | 0.0193 | 0.8743 | 0.0187 |

Table 3: #var and #w.var statistics of simulated linear regression problems

| | $n = 500$ | | | | $n = 1,000$ | | | | $n = 2,000$ | | | |
| | #var | | #w.var | | #var | | #w.var | | #var | | #w.var | |
| | mean | se | mean | se | mean | se | mean | se | mean | se | mean | se |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | $\sigma = 1$ | | | | | | |
| PanIC | 9.1580 | 0.0447 | 1.1340 | 0.0454 | 9.5140 | 0.0370 | 0.7900 | 0.0374 | 9.6380 | 0.0315 | 0.5500 | 0.0314 |
| BIC$_m$ | 11.3580 | 0.0832 | 2.4580 | 0.0665 | 11.3940 | 0.0799 | 2.1460 | 0.0676 | 11.3460 | 0.0690 | 1.8700 | 0.0631 |
| CV | 19.4640 | 0.0347 | 9.4960 | 0.0345 | 19.4820 | 0.0354 | 9.5260 | 0.0353 | 19.5620 | 0.0314 | 9.5900 | 0.0299 |
| | | | | | | $\sigma = 2$ | | | | | | |
| PanIC | 8.7480 | 0.0564 | 1.8720 | 0.0588 | 9.1280 | 0.0511 | 1.4800 | 0.0500 | 9.4840 | 0.0429 | 0.9560 | 0.0420 |
| BIC$_m$ | 10.6000 | 0.0938 | 2.7360 | 0.0707 | 10.8860 | 0.0877 | 2.3980 | 0.0672 | 11.0060 | 0.0725 | 1.9380 | 0.0631 |
| CV | 19.5020 | 0.0364 | 9.5820 | 0.0353 | 19.4880 | 0.0348 | 9.5360 | 0.0334 | 19.5020 | 0.0356 | 9.5460 | 0.0345 |
| | | | | | | $\sigma = 5$ | | | | | | |
| PanIC | 7.9860 | 0.0891 | 3.9460 | 0.0792 | 8.6940 | 0.0776 | 3.1620 | 0.0730 | 9.1360 | 0.0657 | 2.4080 | 0.0621 |
| BIC$_m$ | 7.8860 | 0.1183 | 4.0660 | 0.0790 | 8.7940 | 0.1138 | 3.4700 | 0.0732 | 9.6640 | 0.0941 | 2.7680 | 0.0632 |
| CV | 19.3720 | 0.0415 | 9.6000 | 0.0388 | 19.4540 | 0.0371 | 9.6260 | 0.0324 | 19.5260 | 0.0327 | 9.6300 | 0.0304 |

Table 4: error and |error| statistics of simulated logistic regression problems

| | $n = 500$ | | | | $n = 1,000$ | | | | $n = 2,000$ | | | |
| | error | | |error| | | error | | |error| | | error | | |error| | |
| | mean | se | mean | se | mean | se | mean | se | mean | se | mean | se |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PanIC | -4.6068 | 0.0610 | 4.6068 | 0.0610 | -3.9045 | 0.0534 | 3.9045 | 0.0534 | -3.2755 | 0.0501 | 3.2755 | 0.0501 |
| BIC$_m$ | -3.5468 | 0.0560 | 3.5560 | 0.0548 | -2.5519 | 0.0436 | 2.5540 | 0.0433 | -1.9428 | 0.0341 | 1.9428 | 0.0341 |
| CV | 1.5344 | 0.0478 | 1.5677 | 0.0456 | 0.9107 | 0.0298 | 0.9514 | 0.0272 | 0.5996 | 0.0197 | 0.6288 | 0.0177 |

Table 5: #var and #w.var statistics of simulated logistic regression problems

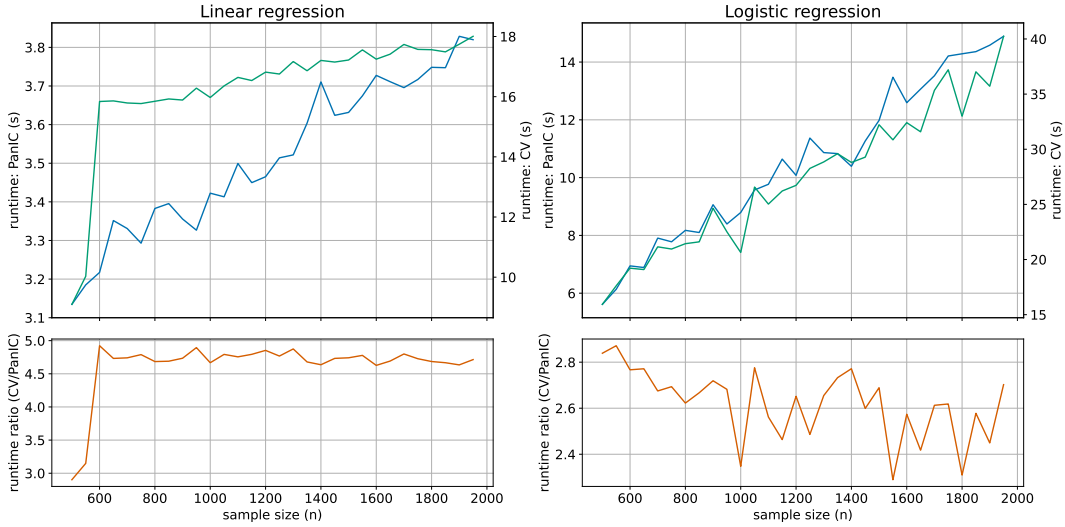| | $n = 500$ | | | | $n = 1,000$ | | | | $n = 2,000$ | | | |
| | #var | | #w.var | | #var | | #w.var | | #var | | #w.var | |
| | mean | se | mean | se | mean | se | mean | se | mean | se | mean | se |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PanIC | 8.3020 | 0.0603 | 2.6260 | 0.0613 | 8.6620 | 0.0543 | 1.9420 | 0.0555 | 9.0440 | 0.0469 | 1.4080 | 0.0474 |
| BIC$_m$ | 8.6020 | 0.0985 | 2.9940 | 0.0661 | 9.6320 | 0.0852 | 2.3760 | 0.0602 | 10.2840 | 0.0793 | 2.0840 | 0.0602 |
| CV | 19.5940 | 0.0283 | 9.7060 | 0.0288 | 19.6300 | 0.0292 | 9.6980 | 0.0277 | 19.6640 | 0.0263 | 9.6920 | 0.0258 |

Figure 2: Total runtime of PanIC (in blue) and CV (in green) for linear and logistic regressions

# 4 Conclusion

In this work, we focus on consistently estimating regression problems using a class of information criteria known as PanIC for regression problems, including linear, logistic, Poisson, and Gamma regressions. In addition to theoretical discussions, we present a simulation study of PanIC in the finite sample setting. Specifically, the results of our simulations indicate that PanIC's performance is comparable to that of cross-validation and BIC. Extending from our consistency result on sets of finitely many models, we introduce a new result (Theorem 1) to consistently estimate over certain sets of uncountably infinitely many models. The setting adopted in Theorem 1 is natural for regressions, but can be extended to a broader scope.

We observe two interesting extensions of this research topic. First, we are interested in establishing consistency when the models are identified by random parameter spaces. For example, in a regression setting it is often useful to formulate a sequence of model spaces $(\mathcal{H}_k)$ given by $\mathcal{H}_k = \{h(\cdot; \beta) : \mathbb{X} \to \mathbb{R} : \beta \in \mathbb{T}_k\}$ and $\mathbb{T}_k = \{\beta \in \mathbb{R}^d : \|\beta\|_1 \le C_k(\omega)\}$, where $C_k(\omega)$ is some random variable that depends on the data. This formulation corresponds to solving regression problems as a sequence of regularization problems, where the sequence of regularization constants $(\lambda_k)$ induces a data-dependent restriction on the parameter spaces $(\mathbb{T}_k)$. For general learning problems that involve regularization, consistency under this setting allows us to directly estimate these problems in their unconstrained form, which is often useful. Secondly, we are interested in exploring the finite sample property of a PanIC estimator in greater detail. Further investigation may focus on the calibration of the unknown multiplicative factor $\kappa$ of a PanIC penalty term.

# Acknowledgements

# A    Proofs

## A.1    Proof of Theorem 1

*Proof.* Since $\mathcal{L}$ is a compact-valued continuous correspondence, by the Berge Maximum Theorem [Aliprantis and Border, 2006, Theorem 17.31], the value function defined by $m(k) = \min_{\beta \in \mathbb{T}_k} r(\beta)$ is continuous. A continuous function attains its minimum on a compact set, so $\mathcal{K}$ is non-empty. As the expected risk function $r$ is continuous, the set $\mathcal{K}$ contains a minimum value and $k^*$ is well-defined. By a similar argument, for each $n \in \mathbb{N}$, $\hat{K}_n$ is well-defined. To show convergence in probability, we will show for all $\sigma > 0$, $\mathbb{P}\left( \left| \hat{K}_n - k^* \right| \geq \sigma \right) \to 0$ as $n \to \infty$. Fix $\sigma > 0$, consider the case $\hat{K}_n \leq k^* - \sigma$. If $\sigma > k^* - a$, the statement is trivial. So assume $\sigma \leq k^* - a$. Let $\varepsilon$ be

$$\varepsilon = \frac{1}{2} \left\{ \min_{\beta \in \mathbb{T}_{k^* - \sigma}} r(\beta) - \min_{\beta \in \mathbb{T}_{k^*}} r(\beta) \right\}.$$

Conditions A1–A3 and the i.i.d assumption are sufficient for the uniform strong law of large numbers Shapiro et al. [2021, Thm. 9.60], which gives $R_n(\beta) \to r(\beta)$ uniformly on any compact set $\mathbb{T}_k$ with probability one as $n \to \infty$. Using this result, and assumption in B1, for every $\delta > 0$, there exists some $N_\delta \in \mathbb{N}$ such that for all $n > N_\delta$, the events

$$\left| \min_{\beta \in \mathbb{T}_{k^* - \sigma}} R_n(\beta) - \min_{\beta \in \mathbb{T}_{k^* - \sigma}} r(\beta) \right| \leq \frac{\varepsilon}{3},$$

$$\left| \min_{\beta \in \mathbb{T}_{k^*}} R_n(\beta) - \min_{\beta \in \mathbb{T}_{k^*}} r(\beta) \right| \leq \frac{\varepsilon}{3},$$

and $P_{k^*, n} \leq \varepsilon/3$ occur with probability at least $1 - \delta$. These events imply the following

$$\min_{\beta \in \mathbb{T}_{k^* - \sigma}} R_n(\beta) \geq \min_{\beta \in \mathbb{T}_{k^* - \sigma}} r(\beta) - \frac{\varepsilon}{3}$$

$$= \min_{\beta \in \mathbb{T}_{k^*}} r(\beta) + \frac{5\varepsilon}{3}$$

$$\geq \min_{\beta \in \mathbb{T}_{k^*}} R_n(\beta) + P_{k^*, n} + \varepsilon.$$

So for all $k \in [a, k^* - \sigma]$,

$$\min_{\beta \in \mathbb{T}_{k^*}} R_n(\beta) + P_{k^*,n} < \min_{\beta \in \mathbb{T}_{k^*-\sigma}} R_n(\beta) + P_{a,n}$$

$$\leq \min_{\beta \in \mathbb{T}_k} R_n(\beta) + P_{a,n}$$

$$\leq \min_{\beta \in \mathbb{T}_k} R_n(\beta) + P_{k,n}$$

occurs with probability at least $1 - \delta$. Take $n \to \infty$, we have $\mathbb{P}\left(\hat{K}_n \leq k^* - \sigma\right) \to 0$. Now consider the case $\hat{K}_n \geq k^* + \sigma$. If $\sigma > b - k^*$, the statement is trivial. So assume $\sigma \leq b - k^*$. Since $\mathbb{T}_b$ is compact, [Shapiro et al., 2021, Thm. 5.7] implies that

$$\sqrt{n}\left(\min_{\beta \in \mathbb{T}_b} R_n(\beta) - \min_{\beta \in \mathbb{T}_b} r(\beta)\right) = O_{\mathbb{P}}(1).$$

Since

$$\min_{\beta \in \mathbb{T}_{k^*}} r(\beta) = \min_{\beta \in \mathbb{T}_b} r(\beta),$$

we have

$$\sqrt{n}\left(\min_{\beta \in \mathbb{T}_b} R_n(\beta) - \min_{\beta \in \mathbb{T}_{k^*}} R_n(\beta)\right) = O_{\mathbb{P}}(1).$$

Fix $\delta > 0$, there exists $M > 0$ such that,

$$\sqrt{n}\left|\min_{\beta \in \mathbb{T}_b} R_n(\beta) - \min_{\beta \in \mathbb{T}_{k^*}} R_n(\beta)\right| \leq M$$

with probability at least $1 - \delta$ for large $n$. By B2, for each $\delta > 0$ and $M > 0$,

$$\sqrt{n}\{P_{k^*+\sigma,n} - P_{k^*,n}\} > M$$

with probability at least $1 - \delta$ for large $n$. So with probability $1 - 2\delta$, we have

$$\min_{\beta \in \mathbb{T}_{k^*}} R_n(\beta) - \min_{\beta \in \mathbb{T}_b} R_n(\beta) \leq \frac{M}{\sqrt{n}} < P_{k^*+\sigma,n} - P_{k^*,n}$$

and thus for all $k \in [k^* + \sigma, b]$,

$$\min_{\beta \in \mathbb{T}_{k^*}} R_n(\beta) + P_{k^*,n} < \min_{\beta \in \mathbb{T}_b} R_n(\beta) + P_{k^*+\sigma,n}$$

$$\leq \min_{\beta \in \mathbb{T}_k} R_n(\beta) + P_{k^*+\sigma,n}$$

$$\leq \min_{\beta \in \mathbb{T}_k} R_n(\beta) + P_{k,n}.$$

19

So $\lim_{n\to\infty} \mathbb{P}(\hat{K}_n - k^* \geq \sigma) = 0$ for all $\sigma > 0$. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ □

## A.2  Proof of Proposition 1

For simplicity, we show the results under the setting of Lasso regression, with a simplified mean function $h(x; \beta) = \beta^\top x$ (without bias term). Upon examining the argument, it is clear that these results are also applicable to both Ridge and Elastic-net regressions, and with the complete mean function $h(x; \beta_0, \beta) = \beta_0 + \beta^\top x$.

### A.2.1  Linear Regression

*Proof.* It suffices to verify the PanIC conditions,

A1 : $\ell_k(x, y; \cdot) : \mathbb{T}_k \to \mathbb{R}$ is continuous for each $(x, y) \in \mathbb{X}$ as it is differentiable with respect to $\beta$. Fix $\beta \in \mathbb{T}_k$, $\ell_k(\cdot; \beta) : \mathbb{X} \to \mathbb{R}$ is given by

$$\begin{aligned}
\ell_k(x, y; \beta) &= (\tilde{f}(x, y; \beta) - \tilde{g}(x, y; \beta))^2 \\
&= (f(y; \beta) - g(x; \beta))^2,
\end{aligned}$$

where $\tilde{f}(x, y; \beta) = f(y; \beta) = y$ and $\tilde{g}(x, y; \beta) = g(x; \beta) = \beta^\top x$. The functions $f$ and $g$ are measurable on their respective domain, so $\tilde{f}$ and $\tilde{g}$ are measurable on the product space $\mathbb{X}$. It follows that $\ell_k(\cdot; \beta)$ is measurable on $\mathbb{X}$.

A2 : $\mathbb{T}_k = \{\beta \in \mathbb{R}^d : \|\beta\|_1 \leq C_k\}$ is compact. Let $\tau_k = 0 \in \mathbb{R}^d$, as $\|0\|_1 \leq C_k$ for all $k \in [m]$, $\tau_k \in \mathbb{T}_k$. This condition is satisfied if $Y$ has a finite fourth moment:

$$\mathbb{E}\left(\ell_k(X, Y; \tau_k)^2\right) = \mathbb{E}(Y^4) < \infty.$$

A3 : Fix $(x, y) \in \mathbb{X}$, $k \in [m]$ and $\beta, \beta' \in \mathbb{T}_k$. By the multivariate Mean Value Theorem [Davidson, 2021],

$$|\ell_k(x, y; \beta) - \ell_k(x, y; \beta')| \leq \sup_{\beta^*} \left\| \frac{\partial \ell_k}{\partial \beta}\Big|_{\beta = \beta^*} \right\|_2 \|\beta - \beta'\|_2,$$

where $\beta^*$ is any point on the line segment joining $\beta$ and $\beta'$. We can write $\beta^* = \lambda\beta + (1 - \lambda)\beta'$, for

$\lambda \in [0, 1]$. Note that,

$$\left\| \frac{\partial \ell_k}{\partial \beta} \bigg|_{\beta = \beta} \right\|_2 = \|-2(y - \beta^\top x)x\|_2$$

$$= 2 \left| y - \beta^\top x \right| \|x\|_2.$$

This gives us

$$\sup_{\beta^*} \left\| \frac{\partial \ell_k}{\partial \beta} \bigg|_{\beta = \beta^*} \right\|_2 = \sup_{\lambda \in [0,1]} 2 \left| y - (\lambda\beta + (1 - \lambda)\beta')^\top x \right| \|x\|_2.$$

Using triangle inequality, we achieve the bound

$$\sup_{\beta^*} \left\| \frac{\partial \ell_k}{\partial \beta} \bigg|_{\beta = \beta^*} \right\|_2 \leq \sup_{\lambda \in [0,1]} 2 \left( |y| + \lambda \left| \beta^\top x \right| + (1 - \lambda) \left| \beta'^\top x \right| \right) \|x\|_2$$

$$\leq \sup_{\lambda \in [0,1]} 2 \left( |y| + \lambda \|x\|_1 \|\beta\|_1 + (1 - \lambda) \|x\|_1 \|\beta'\|_1 \right) \|x\|_2$$

$$\leq 2 \left( |y| + \|x\|_1 \|\beta\|_1 + \|x\|_1 \|\beta'\|_1 \right) \|x\|_2$$

$$\leq 2 \left( |y| + \|x\|_1 \|\beta\|_1 + \|x\|_1 \|\beta'\|_1 \right) \|x\|_1$$

$$\leq 2 \left( |y| + 2 \|x\|_1 C_k \right) \|x\|_1.$$

Let $\mathcal{G}_k(x, y) = 2 \left( |y| + 2 \|x\|_1 C_k \right) \|x\|_1$, then the condition $\mathbb{E}(\mathcal{G}_k(X, Y)^2) < \infty$ is satisfied if both $X$ and $Y$ have a finite fourth moment.

$\square$

### A.2.2   Logistic Regression

*Proof.*

A1 : $\ell_k(x, y; \cdot) : \mathbb{T}_k \to \mathbb{R}$ is continuous for each $(x, y) \in \mathbb{X}$ as it is differentiable with respect to $\beta$. Fix $\beta \in \mathbb{T}_k$, $\ell_k(\cdot; \beta) : \mathbb{X} \to \mathbb{R}$ is given by

$$\ell_k(x, y; \beta) = \tilde{f}(x, y; \beta)\tilde{g}(x, y; \beta) + \tilde{h}(x, y; \beta)$$

$$= f(y; \beta)g(x; \beta) + h(x; \beta),$$

where $\tilde{f}(x, y; \beta) = f(y; \beta) = y$,

$$\tilde{g}(x, y; \beta) = g(x; \beta) = -\log\left(\frac{1}{1 + e^{-\beta^\top x}}\right) + \log\left(1 - \frac{1}{1 + e^{-\beta^\top x}}\right),$$

$$\tilde{h}(x, y; \beta) = h(x; \beta) = -\log\left(1 - \frac{1}{1 + e^{-\beta^\top x}}\right).$$

Both $g$ and $h$ are continuous functions, so they are measurable with respect to the Borel $\sigma$-algebra on $\mathbb{R}^d$. The function $f$ is measurable with respect to the discrete $\sigma$-algebra on $\{0, 1\}$. So $\tilde{f}$, $\tilde{g}$ and $\tilde{h}$ are measurable on the product space $\mathbb{X}$. It follows that $\ell_k(\cdot; \beta)$ is measurable on $\mathbb{X}$.

A2 : $\mathbb{T}_k$ is compact. Let $\tau_k = 0 \in \mathbb{R}^d$, then $\tau_k \in \mathbb{T}_k$ for all $k \in [m]$. This condition is satisfied as

$$\mathbb{E}\left(\ell_k(X, Y; \tau_k)^2\right) = \mathbb{E}\left((-Y\log(1/2) - (1 - Y)\log(1 - 1/2))^2\right)$$

$$= \mathbb{E}\left((-\log(1/2))^2\right)$$

$$= (\log(1/2))^2 < \infty.$$

A3 : Fix $(x, y) \in \mathbb{X}$, $k \in [m]$ and $\beta, \beta' \in \mathbb{T}_k$. By the multivariate Mean Value Theorem, we have

$$|\ell_k(x, y; \beta) - \ell_k(x, y; \beta')| \le \sup_{\beta^*} \left\|\frac{\partial \ell_k}{\partial \beta}\bigg|_{\beta=\beta^*}\right\|_2 \|\beta - \beta'\|_2,$$

where $\beta^*$ is any point on the line segment joining $\beta$ and $\beta'$. We again write $\beta^* = \lambda\beta + (1 - \lambda)\beta'$, for $\lambda \in [0, 1]$. Note that,

$$\frac{\partial \ell_k}{\partial \beta}\bigg|_{\beta=\beta} = -\left(\frac{e^{-\beta^\top x}}{1 + e^{-\beta^\top x}}y - \frac{e^{-\beta^\top x}(1 - y)}{(1 + e^{-\beta^\top x})^2(1 - (1 + e^{-\beta^\top x})^{-1})}\right)x.$$

Then, we have

$$\left\|\frac{\partial \ell_k}{\partial \beta}\bigg|_{\beta=\beta}\right\|_2 = \left\|-\left(\frac{e^{-\beta^\top x}}{1 + e^{-\beta^\top x}}y - \frac{e^{-\beta^\top x}(1 - y)}{(1 + e^{-\beta^\top x})^2(1 - (1 + e^{-\beta^\top x})^{-1})}\right)x\right\|_2$$

$$= \left\|\frac{e^{-\beta^\top x}}{1 + e^{-\beta^\top x}}y - \frac{e^{-\beta^\top x}(1 - y)}{(1 + e^{-\beta^\top x})^2(1 - (1 + e^{-\beta^\top x})^{-1})}\right\|_2 \|x\|_2$$

$$\le \left(\left\|\frac{e^{-\beta^\top x}}{1 + e^{-\beta^\top x}}\right\|_2 + \left\|\frac{e^{-\beta^\top x}}{(1 + e^{-\beta^\top x})^2 - (1 + e^{-\beta^\top x})}\right\|_2\right) \|x\|_2$$

$$\le 2 \|x\|_2$$

$$\le 2 \|x\|_1.$$

It follows that

$$\sup_{\beta^*} \left\| \left. \frac{\partial \ell_k}{\partial \beta} \right|_{\beta=\beta^*} \right\|_2 \leq 2 \left\| x \right\|_1.$$

Let $\mathcal{G}_k(x, y) = 2\|x\|_1$, then the condition $\mathbb{E}(\mathcal{G}_k(X, Y)^2) < \infty$ is satisfied if $X$ has a finite second moment.

$\square$

### A.2.3 Poisson Regression

*Proof.*

A1 : $\ell_k(x, y; \cdot) : \mathbb{T}_k \to \mathbb{R}$ is continuous for each $(x, y) \in \mathbb{X}$ as it is differentiable with respect to $\beta$. Fix $\beta \in \mathbb{T}_k$, $\ell_k(\cdot; \beta) : \mathbb{X} \to \mathbb{R}$ is given by

$$\ell_k(x, y; \beta) = \tilde{f}(x, y; \beta)\tilde{g}(x, y; \beta) + \tilde{h}_1(x, y; \beta) + \tilde{h}_2(x, y; \beta)$$
$$= f(y; \beta)g(x; \beta) + h_1(x; \beta) + h_2(y; \beta),$$

where $\tilde{f}(x, y; \beta) = f(y; \beta) = y$,

$$\tilde{g}(x, y; \beta) = g(x; \beta) = -\beta^\top x,$$
$$\tilde{h}_1(x, y; \beta) = h_1(x; \beta) = \exp(\beta^\top x),$$
$$\tilde{h}_2(x, y; \beta) = h_2(y; \beta) = \log(y!).$$

Both $g$ and $h_1$ are continuous functions, so they are measurable with respect to the Borel $\sigma$-algebra on $\mathbb{R}^d$. Both $f$ and $h_2$ are measurable with respect to the discrete $\sigma$-algebra on $\mathbb{Z}_{\geq 0}$. So $\tilde{f}$, $\tilde{g}$, $\tilde{h}_1$ and $\tilde{h}_2$ are measurable on the product space $\mathbb{X}$. It follows that $\ell_k(\cdot; \beta)$ is measurable on $\mathbb{X}$.

A2 : $\mathbb{T}_k$ is compact. Let $\tau_k = 0 \in \mathbb{R}^d$, then $\tau_k \in \mathbb{T}_k$ for all $k \in [m]$. This condition is satisfied if $Y$ has a finite fourth moment

$$\mathbb{E}\left( \ell_k(X, Y; \tau_k)^2 \right) = \mathbb{E}\left( (1 + \log(Y!))^2 \right)$$
$$\leq \mathbb{E}\left( (1 + Y\log(Y))^2 \right)$$
$$\leq \mathbb{E}\left( (1 + Y^2)^2 \right)$$
$$= \mathbb{E}\left( Y^4 + 2Y^2 + 1 \right) < \infty.$$

A3 : Fix $(x, y) \in \mathbb{X}$, $k \in [m]$ and $\beta, \beta' \in \mathbb{T}_k$. By the multivariate Mean Value Theorem,

$$\left|\ell_k(x, y; \beta) - \ell_k(x, y; \beta')\right| \leq \sup_{\beta^*} \left\| \left.\frac{\partial \ell_k}{\partial \beta}\right|_{\beta = \beta^*} \right\|_2 \|\beta - \beta'\|_2,$$

where $\beta^*$ is any point on the line segment joining $\beta$ and $\beta'$. We write $\beta^* = \lambda\beta + (1 - \lambda)\beta'$, for $\lambda \in [0, 1]$. Note that,

$$\left.\frac{\partial \ell_k}{\partial \beta}\right|_{\beta = \beta} = \left(\exp(\beta^\top x) - y\right) x.$$

Using triangle inequality,

$$
\begin{aligned}
\sup_{\beta^*} \left\| \left.\frac{\partial \ell_k}{\partial \beta}\right|_{\beta = \beta^*} \right\|_2 &= \sup_{\lambda \in [0,1]} \left|\exp\left((\lambda\beta + (1 - \lambda)\beta')^\top x\right) - y\right| \|x\|_2 \\
&\leq \left|\exp\left(|\beta^\top x| + |\beta'^\top x|\right) - y\right| \|x\|_2 \\
&\leq \left(\exp\left(|\beta^\top x| + |\beta'^\top x|\right) + |y|\right) \|x\|_2 \\
&\leq \left(\exp(\|\beta\|_1 \|x\|_1 + \|\beta'\|_1 \|x\|_1) + |y|\right) \|x\|_2 \\
&\leq \left(\exp(2C_k \|x\|_1) + y\right) \|x\|_1.
\end{aligned}
$$

Let $\mathcal{G}_k(x, y) = \left(\exp(2C_k \|x\|_1) + y\right) \|x\|_1$, then the condition $\mathbb{E}(\mathcal{G}_k(X, Y)^2) < \infty$ is satisfied if

$$\mathbb{E}\left(\mathcal{G}_k(X, Y)^2\right) = \mathbb{E}\left(\left(\exp(2C_k \|X\|_1) + Y\right)^2 \|X\|_1^2\right) < \infty. \tag{16}$$

(16) is satisfied when the covariate vector $X$ is Gaussian or sub-Gaussian distributed and the response $Y$ has a finite fourth moment.

$\square$

### A.2.4 Gamma Regression

*Proof.*

A1 : $\ell_k(x, y; \cdot) : \mathbb{T}_k \to \mathbb{R}$ is continuous for each $(x, y) \in \mathbb{X}$ as it is differentiable with respect to $\beta$. Fix $\beta \in \mathbb{T}_k$, $\ell_k(\cdot; \beta) : \mathbb{X} \to \mathbb{R}$ is given by

$$
\begin{aligned}
\ell_k(x, y; \beta) &= \tilde{f}(x, y; \beta)\tilde{g}(x, y; \beta) + \tilde{h}_1(x, y; \beta) + \tilde{h}_2(x, y; \beta) \\
&= f(y; \beta)g(x; \beta) + h_1(x; \beta) + h_2(y; \beta),
\end{aligned}
$$

where $\tilde{f}(x, y; \beta) = f(y; \beta) = y$,

$$\tilde{g}(x, y; \beta) = g(x; \beta) = \nu \exp(-\beta^\top x),$$

$$\tilde{h}_1(x, y; \beta) = h_1(x; \beta) = \nu\beta^\top x + \log(\Gamma(\nu)) - \nu \log(\nu),$$

$$\tilde{h}_2(x, y; \beta) = h_2(y; \beta) = -(\nu - 1)\log(y).$$

Both $g$ and $h_1$ are continuous functions, so they are measurable with respect to the Borel $\sigma$-algebra on $\mathbb{R}^d$. Both $f$ and $h_2$ are measurable with respect to the Borel $\sigma$-algebra on $\mathbb{R}_+$. So $\tilde{f}$, $\tilde{g}$, $\tilde{h}_1$ and $\tilde{h}_2$ are measurable on the product space $\mathbb{X}$. It follows that $\ell_k(\cdot; \beta)$ is measurable on $\mathbb{X}$.

A2 : $\mathbb{T}_k$ is compact. Let $\tau_k = 0 \in \mathbb{R}^d$, then $\tau_k \in \mathbb{T}_k$ for all $k \in [m]$. This condition is satisfied if $Y$ has a finite second moment,

$$\mathbb{E}\left(\ell_k(X, Y; \tau_k)^2\right) = \mathbb{E}\left((-\log(\Gamma(\nu)) + \nu\log(\nu) + (\nu - 1)\log(Y) - Y\nu)^2\right)$$

$$< \infty.$$

A3 : Fix $(x, y) \in \mathbb{X}$, $k \in [m]$ and $\beta, \beta' \in \mathbb{T}_k$. By the multivariate Mean Value Theorem,

$$|\ell_k(x, y; \beta) - \ell_k(x, y; \beta')| \le \sup_{\beta^*} \left\| \left.\frac{\partial \ell_k}{\partial\beta}\right|_{\beta=\beta^*} \right\|_2 \|\beta - \beta'\|_2,$$

where $\beta^*$ is any point on the line segment joining $\beta$ and $\beta'$. We write $\beta^* = \lambda\beta + (1 - \lambda)\beta'$, for $\lambda \in [0, 1]$.

$$\left.\frac{\partial \ell_k}{\partial\beta}\right|_{\beta=\beta} = \left(1 - y\exp(-\beta^\top x)\right)\nu x.$$

Using triangle inequality,

$$\sup_{\beta^*} \left\| \left.\frac{\partial \ell_k}{\partial\beta}\right|_{\beta=\beta^*} \right\|_2 = \sup_{\lambda \in [0,1]} \left\|\left(1 - y\exp\left(-(\lambda\beta + (1-\lambda)\beta')^\top x\right)\right)\nu x\right\|_2$$

$$= \sup_{\lambda \in [0,1]} \nu\left|1 - y\exp\left(-(\lambda\beta + (1-\lambda)\beta')^\top x\right)\right|\|x\|_2$$

$$\le \nu\left(1 + y\exp\left(|\beta^\top x| + |\beta'^\top x|\right)\right)\|x\|_1$$

$$\le \nu\left(1 + y\exp(\|\beta\|_1\|x\|_1 + \|\beta'\|_1\|x\|_1)\right)\|x\|_1$$

$$\le \nu\left(1 + y\exp(2C_k\|x\|_1)\right)\|x\|_1.$$

Let $\mathcal{G}_k(x, y) = \nu\left(1 + y\exp(2C_k\|x\|_1)\right)\|x\|_1$, we require the pair of covariates and response $(X, Y)$ to

satisfy

$$\mathbb{E}\left(\mathcal{G}_k(X,Y)^2\right) = \mathbb{E}\left(\nu^2\left(1 + Y\exp(2C_k\|X\|_1)\right)^2\|X\|_1^2\right) < \infty.$$

Similar to the Poisson regression, this is satisfied when the covariate vector $X$ is Gaussian or sub-Gaussian distributed and the response $Y$ has a finite fourth moment.

$\square$

## A.3   Proof of Proposition 2

*Proof.* By Theorem 2 of Nguyen [2023], it suffices to show conditions C1–C5. We start with C1. Using the functional form of $\ell_k$ from Table 1, we have

$$
\begin{aligned}
r_k(\beta) &= \mathbb{E}(\ell_k(\beta; X, Y)) \\
&= \mathbb{E}\left(Y - \beta^\top X\right)^2 \\
&= \mathbb{E}(Y^\top Y) - 2\beta^\top \mathbb{E}(XY) + \operatorname{tr}(\beta\beta^\top \Sigma) + \mu^\top \beta\beta^\top \mu,
\end{aligned}
$$

where $\mu$ is the mean and $\Sigma$ be the covariance matrix of $X$, respectively. Differentiating $r_k$ twice, we get

$$\frac{\partial^2 r_k}{\partial\beta\partial\beta^\top}(\beta) = 2\Sigma + 2\mu\mu^\top.$$

By our assumption that $\Sigma$ is a positive definite matrix, we see that the Hessian of $r_k$ is also positive definite. Hence, $r_k$ is strictly convex. As $r_k$ is strictly convex, it has, on a convex set, at most one minimizer, which we assume exists. To show Lipschitz continuity, first, we note that our assumptions on $X$ and $Y$ satisfy the conditions in Proposition 1. This gives us the Lipschitz continuity of $\ell_k(\cdot; x, y)$, for all $(x, y) \in \mathbb{R}^d \times \mathbb{R}$. Then, for $\beta, \beta' \in \mathbb{T}_k$,

$$
\begin{aligned}
|r_k(\beta) - r_k(\beta')| &\leq \int_{\mathbb{R}^d \times \mathbb{R}} |\ell_k(\beta; x, y) - \ell_k(\beta'; x, y)| \, \Pi(dx, dy) \\
&\leq \int_{\mathbb{R}^d \times \mathbb{R}} \mathcal{G}(x, y)\|\beta - \beta'\|_1 \, \Pi(dx, dy) \\
&= c\|\beta - \beta'\|_1,
\end{aligned}
$$

for some $c < \infty$. This shows C1. Next, we note that, in our current setting, C2 is implied by condition A3. To show C5, let's first assume $\mathbb{T}_{k^*} \subset \mathbb{T}_k$. Since $r_k$ and $r_{k^*}$ are the same strictly convex function, $\beta_k^*$ and $\beta_{k^*}^*$ must be identical. The same argument can be applied to the case where $\mathbb{T}_{k^*} \supset \mathbb{T}_k$. Hence, condition C5 is satisfied. In the case of Lasso regression, condition C3 follows from the fact that $\mathbb{T}_k$ is polyhedral, for

all $k \in \mathbb{R}_+$. In the case of the Ridge regression, this condition is justified using the Mangasarian–Fromovitz constraint qualification which holds for $\mathbb{T}_k$, for all $k \in \mathbb{R}_+$. Finally, in the case of Elastic-net regression, note the set $\mathbb{T}_k$ can be constructed by intersecting sets that satisfy the Mangasarian–Fromovitz constraint qualification, which leads to C3 by Proposition 3.90 from Bonnans and Shapiro [2000]. Condition C4 is satisfied if the Hessian of $r_k$ is positive definite, see Nguyen [2023, p.25]. Hence, C4 is satisfied. $\square$

*Remark*: The condition that the covariance matrix $\Sigma$ of $X$ is positive definite can be elaborated upon further. Suppose $\Sigma$ takes the form

$$\Sigma = \begin{bmatrix} 0 & 0 \\ 0 & \Omega \end{bmatrix},$$

which is a representation of the covariance matrix of $\bar{X} = (1, X^\top)$, the covariate vector augmented with a bias term. Suppose that $\Omega$ is a positive definite matrix in $\mathbb{R}^{d \times d}$, one can show that $\Sigma$ is positive definite. Indeed, let $M$ be the rank one matrix defined by

$$M = \mu \mu^\top$$

$$= \begin{bmatrix} a \\ b \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}^\top$$

$$= \begin{bmatrix} a^2 & ab^\top \\ ab & bb^\top \end{bmatrix},$$

where $a \in \mathbb{R}$ and $b \in \mathbb{R}^d$. Let $z = (x, y)^\top \neq 0$, where $x \in \mathbb{R}$ and $y \in \mathbb{R}^d$. Observe that

$$z^\top (\Sigma + M) z = \begin{bmatrix} x & y^\top \end{bmatrix} \left( \begin{bmatrix} 0 & 0 \\ 0 & \Omega \end{bmatrix} + \begin{bmatrix} a^2 & ab^\top \\ ab & bb^\top \end{bmatrix} \right) \begin{bmatrix} x \\ y \end{bmatrix}$$

$$= \begin{bmatrix} x & y^\top \end{bmatrix} \begin{bmatrix} a^2 & ab^\top \\ ab & \Omega + bb^\top \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$= \begin{bmatrix} a^2 x + a y^\top b & a x b^\top + y^\top \Omega + y^\top bb^\top \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

$$= a^2 x^2 + 2 a x y^\top b + y^\top \Omega y + y^\top bb^\top y.$$

For any value of $x$ and $y$, we have

$$a^2 x^2 + 2 a x y^\top b + y^\top \Omega y + y^\top bb^\top y \geq 0.$$

If $y \neq 0$, the equality is strict as $\Omega$ is positive definite. So it suffices to check that the quadratic form is positive for $x \neq 0$ and $y = 0$. Note that,

$$a^2 x^2 + 2axy^\top b + y^\top \Omega y + y^\top bb^\top y = a^2 x^2 > 0$$

for all $a \neq 0$. In particular, it is true when $a = 1$.

## A.4 Proof of Lemma 1

*Proof.* A finite-valued convex function is continuous on $\mathbb{R}^d$ [Rockafellar, 1997, Cor. 10.1.1], so it suffices to show the arg min function is continuous. The continuity of the arg min function results from the Berge Maximum Theorem [Aliprantis and Border, 2006, Thm. 17.31]. Let $\phi : \mathbb{R}_+ \twoheadrightarrow \mathbb{R}^d$ be the constant correspondence

$$\phi(\lambda) = \{\beta \in \mathbb{R}^d : g(\beta) \leq g(\hat{\beta})\},$$

where $\hat{\beta}$ is the unique minimizer of $f$ on $\mathbb{R}^d$. Since $\phi$ is a constant correspondence, it is continuous. Moreover, the set $\phi(\lambda) \subset \mathbb{R}^d$ is compact for every $\lambda > 0$. To verify compactness, it suffices to show that $\phi(\lambda)$ is bounded. First, if $g(\hat{\beta}) \leq c$, then $\phi(\lambda) \subseteq L_c$ and it is necessarily bounded. So let $c < g(\hat{\beta})$, note that this implies $L_c \subseteq \phi(\lambda)$. Let's suppose for the sake of contradiction that $\phi(\lambda)$ is unbounded. Then, by Theorem 8.4 of Rockafellar [1997], there exists some non-zero recession direction $y \in \mathbb{R}^d$ in the recession cone of $\phi(\lambda)$ such that, for all $\beta \in \phi(\lambda)$ and $t \geq 0$, $\beta + ty \in \phi(\lambda)$. This implies $g(\beta + ty)$ must be non-increasing in $t$, for all $\beta \in \phi(\lambda)$. So if we choose a $\beta$ that is also in $L_c$, then it holds that $\beta + ty \in L_c$ for all $t \geq 0$. This contradicts the boundedness of $L_c$, so it must be true that $\phi(\lambda)$ is bounded.

Now, let $L : \mathrm{Gr}_\phi \to \mathbb{R}$ be defined as

$$L(\lambda, \beta) = f(\beta) + \lambda g(\beta).$$

Note that $L$ is continuous and strictly convex in $\beta$. Define the value function $m : \mathbb{R}_+ \to \mathbb{R}$ as

$$m(\lambda) = \min_{\beta \in \phi(\lambda)} L(\lambda, \beta).$$

Fix any $\lambda > 0$, by Lemma 2, the unique optimal solution $\beta_\lambda$ of $L(\lambda, \beta)$ on $\mathbb{R}^d$ satisfies $g(\beta_\lambda) \leq g(\hat{\beta})$. Hence, $\beta_\lambda \in \phi(\lambda)$ and the argmin correspondence $\mu$ of the value function $m$ is single-valued. By the Berge Maximum Theorem, $\mu$ is an upper hemicontinuous correspondence. Then, by Lemma. 17.6 of Aliprantis and Border [2006], $\mu$ is continuous as a function. $\qquad\square$

## A.5 Proof of Lemma 2

We use the same proof technique as Oneto et al. [2016], but for the general case of $f$ and $g$.

*Proof.* Define $K^{\lambda_1}$ as the minimum value of the following problem

$$\min \quad f(\beta) + \lambda_1 g(\beta)$$
$$\text{s.t.} \quad \beta \in \mathbb{R}^d,$$

so that

$$K^{\lambda_1} = f(\beta_{\lambda_1}) + \lambda_1 g(\beta_{\lambda_1}).$$

Similarly, define

$$K^{\lambda_2} = f(\beta_{\lambda_2}) + \lambda_2 g(\beta_{\lambda_2}).$$

Let us show the desired result by eliminating other outcomes.

1. $g(\beta_{\lambda_1}) < g(\beta_{\lambda_2})$: Suppose $f(\beta_{\lambda_1}) \leq f(\beta_{\lambda_2})$, this implies

$$f(\beta_{\lambda_1}) + \lambda_2 g(\beta_{\lambda_1}) < f(\beta_{\lambda_2}) + \lambda_2 g(\beta_{\lambda_2}) = K^{\lambda_2},$$

which contradicts the definition of $K^{\lambda_2}$. Now, suppose $f(\beta_{\lambda_1}) > f(\beta_{\lambda_2})$. Using the definition of $K^{\lambda_1}$, we have

$$K^{\lambda_1} = f(\beta_{\lambda_1}) + \lambda_1 g(\beta_{\lambda_1}) \leq f(\beta_{\lambda_2}) + \lambda_1 g(\beta_{\lambda_2}). \tag{17}$$

The above inequality implies

$$\lambda_1 \geq \frac{f(\beta_{\lambda_2}) - f(\beta_{\lambda_1})}{g(\beta_{\lambda_1}) - g(\beta_{\lambda_2})}.$$

However, following the definition of $K^{\lambda_2}$, we also know that

$$K^{\lambda_2} = f(\beta_{\lambda_2}) + \lambda_2 g(\beta_{\lambda_2}) \leq f(\beta_{\lambda_1}) + \lambda_2 g(\beta_{\lambda_1}), \tag{18}$$

which implies

$$\lambda_2 \leq \frac{f(\beta_{\lambda_1}) - f(\beta_{\lambda_2})}{g(\beta_{\lambda_2}) - g(\beta_{\lambda_1})}$$
$$\leq \lambda_1.$$

This contradicts the assumption that $\lambda_1 < \lambda_2$. Hence, $g(\beta_{\lambda_1}) < g(\beta_{\lambda_2})$ is impossible.

2. $g(\beta_{\lambda_1}) > g(\beta_{\lambda_2})$: A similar argument to that used for the case $g(\beta_{\lambda_1}) < g(\beta_{\lambda_2})$ shows that $f(\beta_{\lambda_1}) \geq f(\beta_{\lambda_2})$ is impossible. Let us show $f(\beta_{\lambda_1}) < f(\beta_{\lambda_2})$ is plausible. Using (17), here we have

$$\lambda_1 \leq \frac{f(\beta_{\lambda_2}) - f(\beta_{\lambda_1})}{g(\beta_{\lambda_1}) - g(\beta_{\lambda_2})}.$$

Similarly, (18) implies that
$$\frac{f(\beta_{\lambda_1}) - f(\beta_{\lambda_2})}{g(\beta_{\lambda_2}) - g(\beta_{\lambda_1})} \leq \lambda_2.$$

This observation is consistent with our original assumption that $\lambda_1 \leq \lambda_2$.

3. $g(\beta_{\lambda_1}) = g(\beta_{\lambda_2})$: The case $f(\beta_{\lambda_1}) \neq f(\beta_{\lambda_2})$ is impossible, as it would lead to a contradiction with the definitions of $K^{\lambda_1}$ or $K^{\lambda_2}$. As for the case $f(\beta_{\lambda_1}) = f(\beta_{\lambda_2})$, first let us suppose that $\beta_{\lambda_1} \neq \beta_{\lambda_2}$. This is also impossible, as the number of optimal solutions is at most one. The remaining outcome $\beta_{\lambda_1} = \beta_{\lambda_2}$ is plausible. For example, let $f(\beta) = \|\beta\|_2^2$, $g(\beta) = \|\beta\|_1$. Then, for any $\lambda_1, \lambda_2 > 0$,

$$0 = \underset{\beta \in \mathbb{R}^d}{\arg\min} \, f(\beta) + \lambda_1 g(\beta) = \underset{\beta \in \mathbb{R}^d}{\arg\min} \, f(\beta) + \lambda_2 g(\beta).$$

$\square$

# References

H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19 (6):716–723, December 1974. ISSN 0018-9286.

Charalambos D. Aliprantis and Kim C. Border. *Infinite dimensional analysis: a hitchhiker's guide*. Springer, Berlin ; New York, 3rd [rev. and enl.] ed edition, 2006. ISBN 9783540295860. OCLC: ocm69983226.

Takeshi Amemiya. *Advanced Econometrics*. Harvard University Press, Cambridge, 1985.

Jean-Patrick Baudry. Estimation and model selection for model-based clustering with the conditional classification likelihood. *Electronic Journal of Statistics*, 9:1041–1077, 2015.

J. Frédéric Bonnans and Alexander Shapiro. *Perturbation analysis of optimization problems*. Springer New York, New York, NY, 2000. ISBN 9781461271291 9781461213949.

Stephen P. Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, Cambridge, UK ; New York, 2004. ISBN 9780521833783.

Peter Bühlmann and Sara Van de Geer. *Statistics for high-dimensional data: methods, theory and applications.* Springer Series in Statistics. Springer Berlin Heidelberg, Berlin, Heidelberg, 2011. ISBN 9783642201912 9783642201929.

Gerda Claeskens and Nils Lid Hjort. *Model selection and model averaging.* Cambridge University Press, 1 edition, January 2001. ISBN 9780521852258 9780511790485.

James Davidson. *Stochastic limit theory: an introduction for econometricians.* Advanced texts in econometrics. Oxford University Press, Oxford, second edition edition, 2021. ISBN 9780192844507. OCLC: on1272885940.

Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001. ISSN 01621459.

Xin Gao and Peter X.-K. Song. Composite likelihood bayesian information criteria for model selection in high-dimensional data. *Journal of the American Statistical Association*, 105(492):1531–1540, 2010. ISSN 0162-1459.

Trevor Hastie, Robert Tibshirani, and Martin Wainwright. *Statistical learning with sparsity: the lasso and generalizations.* Number 143 in Monographs on statistics and applied probability. CRC Press, Taylor & Francis Group, Boca Raton, 2015. ISBN 9781498712163.

Arthur E. Hoerl and Robert W. Kennard. Ridge regression: biased estimation for nonorthogonal problems. *Technometrics*, 12(1):55–67, February 1970. ISSN 0040-1706, 1537-2723.

Francis K C Hui. On the use of a penalized quasilikelihood information criterion for generalized linear mixed models. *Biometrika*, 108(2):353–365, May 2021. ISSN 0006-3444, 1464-3510.

Francis K.C. Hui, David I. Warton, and Scott D. Foster. Order selection in finite mixture models: complete or observed likelihood information criteria? *Biometrika*, 102(3):724–730, September 2015. ISSN 0006-3444, 1464-3510.

Jinzhu Jia and Bin Yu. On model selection consistency of the elastic net when p ¿¿ n. *Statistica Sinica*, 20 (2):595–611, 2010. ISSN 1017-0405.

Marius Kloft, Ulf Brefeld, Sören Sonnenburg, and Alexander Zien. $l_p$-norm multiple kernel learning. *Journal of Machine Learning Research*, 12(26):953–997, 2011. ISSN 1533-7928.

Brian G. Leroux. Consistent estimation of a mixing distribution. *The Annals of Statistics*, 20(3):1350–1360, 1992. ISSN 0090-5364.

Pascal Massart and Caroline Meynet. The Lasso as an l1-ball model selection procedure. *Electronic Journal of Statistics*, 5:669 – 687, 2011.

Allan D R McQuarrie and Chih-Ling Tsai. *Regression and time series model selection*. WORLD SCIENTIFIC, May 1998. ISBN 9789810232429 9789812385451.

Chi Tim Ng and Harry Joe. Model comparison with composite likelihood information criteria. *Bernoulli*, 20(4):1738–1764, 2014. ISSN 1350-7265.

Hien Duy Nguyen. PanIC: consistent information criteria for general model selection problems. 2023.

Luca Oneto, Sandro Ridella, and Davide Anguita. Tikhonov, Ivanov and Morozov regularization for support vector machine learning. *Machine Learning*, 103(1):103–136, April 2016. ISSN 0885-6125, 1573-0565.

Ralph Tyrrell Rockafellar. *Convex analysis*. Princeton Landmarks in mathematics and physics. Princeton Univ. Press, Princeton, NJ, 10. print. and 1. paperb. print edition, 1997. ISBN 9780691015866 9780691080697.

Gideon Schwarz. Estimating the dimension of a model. *The Annals of Statistics*, 6(2), March 1978. ISSN 0090-5364.

Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczynski. *Lectures on stochastic programming: modeling and theory, third edition*. Society for Industrial and Applied Mathematics, Philadelphia, PA, July 2021. ISBN 9781611976588 9781611976595.

Chor-Yiu Sin and Halbert White. Information criteria for selecting possibly misspecified parametric models. *Journal of Econometrics*, 71(1-2):207–225, March 1996. ISSN 03044076.

M. Stone. Cross-validatory choice and assessment of statistical predictions. *Journal of the Royal Statistical Society. Series B (Methodological)*, 36(2):111–147, 1974. ISSN 0035-9246.

Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, 58(1):267–288, 1996. ISSN 0035-9246.

Vladimir N. Vapnik. *The nature of statistical learning theory*. Springer New York, New York, NY, 2000. ISBN 9781441931603 9781475732641.

Vladimir Naumovich Vapnik. *Statistical learning theory*. Adaptive and learning systems for signal processing, communications, and control. Wiley, New York, 1998. ISBN 9780471030034.

Cristiano Varin and Paolo Vidoni. A note on composite likelihood inference and model selection. *Biometrika*, 92(3):519–528, 2005. ISSN 0006-3444.

Ming Yuan and Yi Lin. On the non-negative garrotte estimator. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 69(2):143–161, 2007. ISSN 1369-7412.

Peng Zhao and Bin Yu. On model selection consistency of lasso. *Journal of Machine Learning Research*, 7 (90):2541–2563, 2006. ISSN 1533-7928.

Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 67(2):301–320, 2005. ISSN 1369-7412.

Hui Zou, Trevor Hastie, and Robert Tibshirani. On the "degrees of freedom" of the lasso. *The Annals of Statistics*, 35(5), October 2007. ISSN 0090-5364.