# Adaptive Noisy Data Augmentation for Regularized Estimation and Inference in Generalized Linear Models

Yinan Li[1] and Fang Liu[1*]

[1] Department of Applied and Computational Mathematics and Statistics

University of Notre Dame, Notre Dame, IN 46556, U.S.A.

## Abstract

We propose the AdaPtive Noise Augmentation (PANDA) procedure to regularize the estimation and inference of generalized linear models (GLMs). PANDA iteratively optimizes the objective function given noise augmented data until convergence to obtain the regularized model estimates. The augmented noises are designed to achieve various regularization effects, including $l_0$, bridge (lasso and ridge included), elastic net, adaptive lasso, and SCAD, as well as group lasso and fused ridge. We examine the tail bound of the noise-augmented loss function and establish the almost sure convergence of the noise-augmented loss function and its minimizer to the expected penalized loss function and its minimizer, respectively. We derive the asymptotic distributions for the regularized parameters, based on which, inferences can be obtained simultaneously with variable selection. PANDA exhibits ensemble learning behaviors that help further decrease the generalization error. Computationally, PANDA is easy to code, leveraging existing software for implementing GLMs, without resorting to complicated optimization techniques. We demonstrate the superior or similar performance of PANDA against the existing approaches of the same type of regularizers in simulated and real-life data. We show that the inferences through PANDA achieve nominal or near-nominal coverage and are far more efficient compared to a popular existing post-selection procedure.

**keywords**: $l_0$ penalty, augmented Fisher information, ensemble learning, noise injection and augmentation, regularization and penalization, inference

## 1 Introduction

Regularization of generalized linear models (GLMs) to mitigate overfitting and conduct variable selection is a well-studied topic. There exist a variety of regularizers, such as bridge (Frank and Friedman, 1993), ridge ($l_2$), lasso ($l_1$) (Tibshirani, 1996), elastic net (Zou and Hastie, 2005), SCAD (Fan and Li, 2001), adaptive lasso (Zou, 2006), group lasso (Yuan and Lin, 2014), fused lasso (Tibshirani et al., 2005), sparse group lasso (SGL) (Simon et al., 2013), among others. As for the inference for regression coefficients in penalized GLMs, many existing approaches are post-election procedures, meaning the inference is initiated after variable selection and oftentimes non-selected variables are assumed to be of no inferential interest and no uncertainty quantification are provided for the corresponding regression coefficients (Leeb and Pötscher, 2005; Leeb et al., 2006; Berk et al., 2013; Zhang and Zhang, 2013; Javanmard

---

*Corresponding author email: fang.liu.131@nd.edu

and Montanari, 2014; Lockhart et al., 2014; Efron, 2014; Lee et al., 2016; Tibshirani et al., 2016; Reid et al., 2017; Taylor and Tibshirani, 2017). Procedures for simultaneous variable selection and inference do exist. Fan and Li (2001) provide simultaneous variance estimates for the coefficients of selected estimated variables (so non-zero coefficient estimates). For linear regression, (Zhang and Zhang, 2013; Javanmard and Montanari, 2014) provide simultaneous variable selection with the lasso penalty and inference for both zero or non-zero coefficient estimates. Van de Geer et al. (2014) propose a procedure for constructing confidence intervals and hypothesis testing for a low-dimensional subset of a large parameter vector in the high-dimensional GLM setting with convex loss functions with the lasso penalty.

Despite the extensiveness of the work on regularized variable selection in GLMs, there is still room for improvement over the existing solutions. Two of these areas are the $l_0$ regularization and inference in regularized GLMs. Optimization with the $l_0$ penalty is NP-hard. Dicker et al. (2013) propose the seamless-$l_0$ (SELO) penalty to approximate the $l_0$ penalty and a coordinate descent algorithm to obtain solutions in the context of the least-squares optimization and $p < n$. SELO outperforms SCAD in model error and variable selection accuracy rate per the empirical studies. Liu and Li (2016) propose an EM algorithm that approximates the $l_0$ regularized regression by solving a sequence of $l_2$ optimizations. The method deals with $p > n$, but is examined only in the least-squares setting and does not provide inferential procedures. Regarding the inference for regularized GLMs, as mentioned above, the majority of existing methods operate in a post-selection matter and thus focus on the inference for selected variables only. Fan and Li (2001) and Dicker et al. (2013) provide standard errors for the parameter estimates in non-convex optimization, and again for selected variables only.

We propose a novel general regularization framework, AdaPtive Noisy Data Augmentation (PANDA), for GLMs that 1) achieves the $l_0$ penalty in addition to all the above mentioned existing penalty types, 2) obtains inference in regularized GLMs for both zero and non-zero coefficients, and 3) enjoys simple practical implementation that would greatly appeal to practitioners. In brief, PANDA augments the original $n$ observations with properly designed $n_e$ noise terms to achieve the desired regularization effects on model parameters. PANDA is iterative and the variance terms of the augmented noise are adaptive to the most updated parameter estimates until the algorithm converges. One requirement on $n_e$ is the augmented data size $n + n_e > p$ so to allow for the ordinary least squares (OLS) or maximum likelihood estimation (MLE) procedures to be applied without resorting to complicated optimization algorithms. As such, *PANDA is computationally straightforward and efficient. PANDA is also flexible and general.* By properly designing the variance of the augmented noise, PANDA can achieve various regularization effects, including $l_\gamma$, for $0 \leq \gamma \leq 2$ (including $l_0$, lasso, ridge as special cases), elastic net, SCAD, group lasso, and fused ridge. PANDA achieves close-to-exact $l_0$ regularization by promoting orthogonality between the coefficients and the augmented noise vector. When $n_e < p$, PANDA shrinks exactly $n_e$ parameters towards 0 upon convergence. *PANDA is more capable and more efficient inferentially compared to existing inferential approaches for regularized GLMs.* It conducts variable selection and provides inference for coefficients simultaneously, whether the coefficients are estimated to be zero or not. Our empirical results suggest the inference based on PANDA is valid and more efficient compared to some existing post-selection procedures. Finally, *PANDA is theoretically justified.* We establish the Gaussian tail of the noise-augmented loss function and the almost sure

2

convergence to its expectation under some regularity conditions, providing theoretical justification for PANDA as a regularization technique and that the noise-augmented loss function is trainable for practical implementation.

The optimizer calculated by PANDA from a GLM is similar to the local quadratic approximation (LQA) technique (Tibshirani, 1996; Fan and Li, 2001), but with several important differences. First, LQA cannot yield the $l_0$ penalty while PANDA can achieve close-to-exact $l_0$ regularization; second, LQA relies on analytical work to approximate penalized loss function with a quadratic form, followed by the optimization of the quadratic function, whereas PANDA only needs to augment the original data with noisy samples and then leverage existing software to compute OLS/MLE from GLMs.

The rest of the paper is organized as follows. Sec 2 presents the PANDA algorithm and the regularization effects it brings to GLMs. Sec 3 establishes the consistency on the noise-augmented loss function and the regularized parameter estimates, presents the Fisher information of the model parameters in augmented data, examines PANDA's ensemble learning behavior, and provides the asymptotic distributions for the parameter estimates via PANDA. Sec 4 demonstrates the $l_0$ penalty realized by PANDA, compares PANDA to a popular post-selection inferential approach in statistical inferences for GLMs, and implements PANDA in simulated and real-life studies to show its effectiveness in regularizing GLM estimation. Sec 5 provides some concluding remarks and offers future research directions on PANDA.

# 2  Methodology

## 2.1  Noise Augmentation Scheme and Regularization Effect in PANDA

Let $Y$ be the outcome variable and $\mathbf{X} = (X_1, \ldots, X_p)^T$ be the independent variables. GLM is based on the assumption that the conditional distribution of $Y$ given $\mathbf{X}$ comes from an exponential family
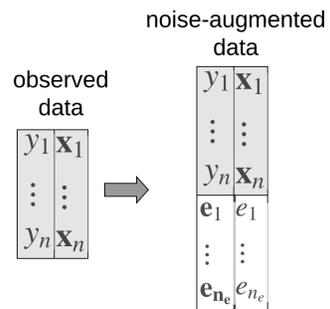
$$p(Y|\mathbf{X}) = \exp\left(Y\eta - B(\eta) + h(Y)\right), \tag{1}$$

where $\eta = \boldsymbol{\theta}_0 + \boldsymbol{\theta}\mathbf{X}$ if the canonical link is used (e.g., the identity link for Gaussian $Y$; the logit link for Bernoulli $Y$). When $p$ is large, regularization or penalty is often imposed on $\boldsymbol{\theta}$ when estimating $\boldsymbol{\theta}$.

PANDA regularizes the estimation of $\boldsymbol{\theta}$ by first augmenting the observed data with a noisy data matrix. Fig 1 depicts a schematic of data augmentation in PANDA, where the augmented noise $e_i$ to $\mathbf{y}$ is $\bar{y}$, the sample average of $\mathbf{y}$. For logistic regression, $e_i \sim \text{Bern}(\hat{p})$, where $\hat{p}$ is the sample proportion of an event. The augmented data $\mathbf{e}_x$ to $\mathbf{x}$ are drawn from the *Noise Generating Distributions* (NGD), the variance term of which is function of $\boldsymbol{\theta}$ and tuning parameters $\boldsymbol{\lambda}$.

$$\mathbf{e}_x \sim N\left(0, \text{V}(\boldsymbol{\theta}; \boldsymbol{\lambda})\right) \tag{2}$$

**Proposition 1 (regularization effects of PANDA for GLM).** Denote the loss function given the observed data



$$\mathbf{x}_i = (x_{i1}, \ldots, x_{ip}) \text{ for } i = 1, \ldots, n$$
$$e_i = \bar{y} \text{ and } \mathbf{e}_i = (e_{1,1}, \ldots, e_{1,p})$$
$$\text{for } i = 1, \ldots, n_e$$

Figure 1: A schematic of the data augmentation for GLM in PANDA (for logistic regression, $e_{y,i} \sim \text{Bern}(\hat{p})$, where $\hat{p}$ is the sample proportion of an event)

3

$(\mathbf{x}, \mathbf{y})$ by $l(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}) = -\sum_{i=1}^{n}\left(h(y_i) + \left(\theta_0 + \sum_j \theta_j x_{ij}\right) y_i - j(\theta_0 + \sum_j \theta_j x_{ij})\right)$ (the negative log-likelihood function), and that given the noise augmented data $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$) by

$$l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = -\left\{\sum_{i=1}^{n+n_e}\left(h(\tilde{y}_i) + \left(\theta_0 + \sum_j \theta_j \tilde{x}_{ij}\right) \tilde{y}_i\right) - B_j\left(\theta_0 + \sum_j \theta_j \tilde{x}_{ij}\right)\right\}, \quad (3)$$

where $\boldsymbol{\theta} = (\theta_0, \{\theta_j\}_{j=1,\dots,p})$. The expectation of the Taylor series of $l_p$ around $\sum_j \theta_j x_{ij} = 0$ over the distribution of $\mathbf{e}_x$ is

$$\mathrm{E}_{\mathbf{e}_x}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})) = l(\boldsymbol{\theta}|\mathbf{x}, \mathbf{y}) + P(\boldsymbol{\theta}), \text{ where}$$

$$P(\boldsymbol{\theta}) = n_e\left(C_1 \sum_j \theta_j^2 \mathrm{V}(e_j)\right) + O\left(n_e \sum_j \left(\theta_j^4 \mathrm{V}^2(e_j)\right)\right) + C, \quad (4)$$

where $C_1 = 2^{-1}B''(\theta_0)$ and $C = \sum_{i=1}^{n_e}(h(e_{y,i}) + e_{y,i}\theta_0) + B(\theta_0)$ are constants independent of $\boldsymbol{\theta}_{-0}$.

The proof is given in Sec S.1 of the supplementary materials. The regularization effect $P(\boldsymbol{\theta})$ in Eqn (4) depends on the variance term of the NGD in Eqn (2) from which the augmented noise to $\mathbf{x}$ is sampled. Eqns (5) to (8) list some examples of the NGD from which $e_j$, the noise term that augments $\mathbf{x}_j$ for $j = 1, \dots, p$, is drawn and their expected regularization effects.

$$\text{bridge: } e_j \sim N\left(0, \lambda|\theta_j|^{-\gamma}\right) \text{ for } \gamma \in [0, 2] \quad (5)$$

including $l_0$ when $\gamma = 2$, lasso when $\gamma = 1$, and ridge when $\gamma = 0$;

$$\text{elastic net: } e_j \sim N\left(0, \lambda|\theta_j|^{-1} + \sigma^2\right); \quad (6)$$

$$\text{adaptive lasso: } e_j \sim N\left(0, \lambda|\theta_j|^{-1}|\hat{\theta}_j|^{-\gamma}\right), \text{ where } \hat{\theta}_j \text{ is a consistent estimate for } \theta_j; \quad (7)$$

$$\text{SCAD: } e_j \begin{cases} = 0 \text{ if } |\theta_j| > an_e\lambda \\ \sim N\left(0, \left(\frac{\lambda}{|\theta_j|} - \frac{(a+1)}{2a^2 n_e}\right)1_{(0, n_e\lambda]}(|\theta_j|) + \frac{1}{(a-1)}\left(\frac{a\lambda}{|\theta_j|} - \frac{\lambda^2 n_e}{2\theta_j^2} - \frac{2a^2-1}{2a^2 n_e}\right)1_{(n_e\lambda, an_e\lambda]}(|\theta_j|)\right) \text{ o.w.} \end{cases},$$

$$\text{where } 1_{(l,u)}(|\theta_j|) = 1 \text{ if } l < |\theta_j| < u, 0 \text{ o.w.}; \quad (8)$$

For regularizing a group of $q$ parameters $\boldsymbol{\theta}_q = (\theta_1, \dots, \theta_q)$ simultaneously (e.g., genes on the same pathway, binary dummy variables created from the same categorical attribute), the NGDs in Eqns (9) and (11) can be used. Specifically, the group-lasso penalty sets all $q$ parameters in $\boldsymbol{\theta}_q$ either at zero or nonzero simultaneously; and the fused-ridge and fused-lasso penalties promote numerical similarity among $\boldsymbol{\theta}_q$.

$$\text{group lasso: } e_{j(l)} \sim N\left(0, \frac{\lambda\sqrt{q_l}}{||\boldsymbol{\theta}_l||_2}\right), \text{ where } \boldsymbol{\theta}^{(l)} = \{\theta_{j(l)}\} \text{ for } j = 1, \dots, q_l; l = 1, \dots, g \text{ groups} \quad (9)$$

$$\text{fused ridge: } \mathbf{e} = (e_1, \dots, e_q) \sim N_{(q)}\left(0, \lambda(\mathbf{TT}')\right), \text{ where entries in } \mathbf{T} \text{ are} \quad (10)$$

$$T_{k,k} = 1, T_{k+1-k\cdot 1(k=q),k} = -1; \text{ and } 0 \text{ o.w.};$$

$$\text{fused lasso: } \mathbf{e} = (e_1, \dots, e_q) \sim N_{(q)}\left(0, \lambda(\mathbf{TT}')\right), \text{ where } T_{kk'} = \lambda|\theta_j - \boldsymbol{\theta}_{k'}|^{-1} \text{ for } k \neq k'. \quad (11)$$

The tuning parameters $\sigma^2 \geq 0, \lambda > 0, 0 \leq \gamma < 2, a > 2$ in Eqn (5) and (11) can be user-specified or chosen by a model selection criterion such as cross-validation (CV), AIC, or BIC. The dispersion of the noise term varies by $X$ in general. $X_j$ associated with small $|\theta_j|$ is augmented with more spread-out noises, and $X_j$ with large $|\theta_j|$ is augmented with noises concentrated around 0. The exceptions are the ridge ($\gamma = 0$ in Eqn (5)) and fused ridge regularizations (Eqn (10)), where the variance term remains constant for $\theta_j$ for all $j$.

For linear regression, the noise-augmented loss function is $l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \sum_{i=1}^{n+n_e}\left(\tilde{y}_i - \sum_j \tilde{x}_{ij}\theta_j\right)^2$, and the penalty $P(\boldsymbol{\theta})$ realized by PANDA with different types of NGD can be obtained in closed form (Table 1) and are exact as suggested by the names. $P(\boldsymbol{\theta})$ in Table 1 can be

easily derived based on the results in Sec S.1 of the supplementary materials. Compared to the original SCAD in Fan and Li (2001) for linear regression, the SCAD penalty realized by PANDA for $|\theta_j| < \lambda n_e$ is not $l_1$ as in the original SCAD but closer to a $l_{0.5}$ penalty; the middle segment penalty $|\theta_j| \in (\lambda n_e, a n_e \lambda]$ is not exactly the same as the original SCAD either, but it has the same functionality by shrinking $|\theta_j| \in (\lambda n_e, |\theta_j| < a\lambda n_e]$ toward 0 and connecting the two end segments to form an overall smooth penalty for $\theta_j$; for $|\theta_j| > a\lambda n_e$, there is no penalty as in the original SCAD.

Table 1: Close-formed penalty term in regularized linear regression (Gaussian $Y$) via PANDA

| | NGD | $P(\boldsymbol{\theta})$ when $Y$ is Gaussian |
|---|---|---|
| $l_\gamma$ | Eqn (5) | $(\lambda n_e) \sum_{j=1}^p \sum_{j \neq k} |\theta_j|^{2-\gamma}$ |
| EN | Eqn (6) | $(\lambda n_e) \sum_{j=1}^p \sum_{j \neq k} |\theta_j| + (\sigma^2 n_e) \sum_{j=1}^p \sum_{j \neq k} \theta_j^2$ |
| adaptive | Eqn (7) | $(\lambda n_e) \sum_{j=1}^p \sum_{j \neq k} |\theta_j||\hat{\theta}_j|^{-\gamma}$, |
| SCAD | Eqn (8) | $\sum_{j=1}^p \left\{ (n_e \lambda |\theta_j| - (a+1)\theta^2/(2a^2 n_e))) \cdot 1_{|\theta_j| \leq \lambda n_e} + \right.$ $\left. (a\lambda n_e|\theta_j| - (\lambda n_e)^2/2 - \theta_j^2(1 - 1/(2a^2))) (a-1)^{-1} \cdot 1_{|\theta_j| \in (\lambda n_e, a\lambda n_e]} \right\}$ |
| group lasso | Eqn (9) | $(\lambda n_e) \sum_{l=1}^g \sqrt{q_l} ||\boldsymbol{\theta}_l||_2$ |

When $Y$ is non-Gaussian, the achieved regularization effects $P(\boldsymbol{\theta})$ in Eqns (5) to (11) are second-order approximate. For example, Table 2 lists the analytical form of $P(\boldsymbol{\theta})$ for the lasso-type noise ($\gamma = 1$ in Eqn (5)). For all the regression types, in addition to the $l_1$ penalty, there is an additional big-O term on $\boldsymbol{\theta}$, which is arbitrarily small under some regularity conditions (more details are provided in Sec 2.3).

Table 2: Expected penalty term in PANDA with lasso-type noise $e_j \sim N(0, \lambda|\theta_j|^{-\gamma})$

| $Y$ | $P(\boldsymbol{\theta})$ |
|---|---|
| Bernoulli | $\frac{\lambda n_e}{2} \frac{\exp(\theta_0)}{(1+\exp(\theta_0))^2} \sum_j |\theta_j| + O(\lambda^2 n_e ||\boldsymbol{\theta}||_2^2) + C$ |
| Exponential | $\frac{\lambda n_e}{2} \exp(\boldsymbol{\theta}_0) \sum_j |\theta_j| + O(\lambda^2 n_e ||\boldsymbol{\theta}||_2^2) + C$ |
| Poisson | $\frac{\lambda n_e}{2} \exp(\theta_0) \sum_j |\theta_j| + O(\lambda^2 n_e ||\boldsymbol{\theta}||_2^2) + C$ |
| Negative Binomial | $\frac{\lambda n_e}{2} \frac{r \exp(\boldsymbol{\theta}_0)}{(r+\exp(\boldsymbol{\theta}_0))} \sum_j |\theta_j| + O(\lambda^2 n_e ||\boldsymbol{\theta}||_2^2) + C$ ($r$ is the # of failures) |

For relatively small $n_e$, especially when $n_e < p$, PANDA promotes sparsity on $\boldsymbol{\theta}$ by imposing $n_e$ linear constraints on $\boldsymbol{\theta}$. Applying the second-order approximation at $\mathbf{e}_i^T \boldsymbol{\theta} = 0$, we have

$$l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}) = -\left\{ \sum_{i=1}^{n+n_e} \left( h(\tilde{y}_i) + \left( \theta_0 + \sum_j \theta_j \tilde{x}_{ij} \right) \tilde{y}_i \right) - B_j \left( \theta_0 + \sum_j \theta_j \tilde{x}_{ij} \right) \right\}$$

$$\approx l(\boldsymbol{\theta}|\mathbf{x}) + C_1 \sum_{i=1}^{n_e} \left( \sum_j \theta_j e_{ij} \right)^2 + C, \tag{12}$$

where $C_1$ and $C$ the same as in Eqn (4). The regularization effect obtained in Eqn (12) with fixed $n_e$ is different from the regularization presented in Proposition 1 in the sense that it takes effect by promoting the orthogonality between $\boldsymbol{\theta}$ and $\mathbf{e}_{x,i}, i = 1 \ldots, n_e$ rather than penalizing the individual parameters. The formal results are given in Proposition 2. The proof is provided in Sec S.2 of the supplementary materials.

**Proposition 2 (orthogonal regularization effect of PANDA for GLM with fixed $n_e$).** With fixed $n_e$ and the approximate loss function in Eqn (12), PANDA estimates $\boldsymbol{\theta}$ by solving

$$\hat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}) \approx \arg\min_{\boldsymbol{\theta}} (l(\boldsymbol{\theta}|\mathbf{x}) + C_1 \sum_{i=1}^{n_e} \left( \mathbf{e}_i^T \boldsymbol{\theta} \right)^2 \tag{13}$$

in each iteration, which is conceptually equivalent to the constrained optimization problem

$\min l(\boldsymbol{\theta}|\mathbf{x})$ subject to

$$\sum_{i=1}^{n_e}(\mathbf{e}_i^T\boldsymbol{\theta})^2 \leq \sum_{i=1}^{n_e}\left(\mathbf{e}_i^T\hat{\boldsymbol{\theta}}\right)^2 \text{ or equivalently,} \tag{14}$$

$$\exists\, 0 < d_i < (\sum_{i=1}^{n_e}(\mathbf{e}_i^T\hat{\boldsymbol{\theta}})^2)^{1/2} \text{ such that } |\mathbf{e}_i^T\boldsymbol{\theta}| \leq d_i, \text{ for } i = 1, \ldots, n_e. \tag{15}$$

$\hat{\boldsymbol{\theta}}$ in Eqns (14) and (15) is the solution from Eqn (13). Proposition 2 suggests that the (unconstrained) optimization problem PANDA solves in each iteration is equivalent to a constrained optimization problem with $n_e$ linear constraints on $\boldsymbol{\theta}$. When $n_e < p$, the $n_e$ constraints in Eqn (15) only affect a subset of the $p$ parameters. For the $l_0$ penalty ($\gamma = 2$ in Eqn 5)) and when $\lambda$ is large, the constraints take effect on exactly $n_e$ parameters. In other words, the following two optimization problems are equivalent.

$$\text{Problem 1:} \quad \bar{\boldsymbol{\theta}} = \mathrm{E}_{\mathbf{e}}(\hat{\boldsymbol{\theta}}) = \mathrm{E}_{\mathbf{e}}\left\{\arg\min_{\boldsymbol{\theta}}(l(\boldsymbol{\theta}|\mathbf{x}) + C_1\sum_{i=1}^{n_e}\left(\sum_j \theta_j e_{ij}\right)^2\right\} \tag{16}$$

$$\text{Problem 2:} \quad \hat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} l(\boldsymbol{\theta}|\mathbf{x}), \text{subject to } \sum_{j=1}^{p}\mathbf{1}(\theta_j \neq 0) = p - n_e \tag{17}$$

Figure 2 plots the heat maps of the constrained region on $\boldsymbol{\theta} = (\theta_1, \theta_2)$ as suggested by Eqn (14) for $n_e = 1, 2$ and 10, respectively when $\hat{\boldsymbol{\theta}} = (0.01, 1)$ (the upper panel) and $\hat{\boldsymbol{\theta}} = (0.01, 0.01)$ (the bottom panel), with the $l_0$ penalty. Specifically, each heat map is made of $50,000$ "dots" uniformly distributed in the $[2, 2] \times [2, 2]$ solution region (for plotting purposes, we focus on the region of $[2, 2]^2$; in theory, the region can be as large as $(-\infty, \infty)^2$). The relative density of a particular constraint on $\boldsymbol{\theta}$ out the $5,000$ repeats is proportional to the grayness of the dot. In the upper panel, with $n_e = 1$, the chance of constraining $\theta_1$ at 0 is much higher than at any non-zero values. As $n_e$ increases from 1 to 2 to 10, the constrained region for $\boldsymbol{\theta}$ shrinks (to 0 for $\theta_1$ and to within $[-1, 1]$ for $\theta_2$). In the bottom panel, setting $n_e = 1$ still lead a substantial chance of getting non-zero $\boldsymbol{\theta}$. As $n_e$ increases to 2 and 10, the chance of $\boldsymbol{\theta} = 0$ drastically increases and is almost certain at $n_e = 10$.

As $n_e$ further increases, the regularization effect moves away from promoting the orthogonality between $\mathbf{e}_x$ and $\theta$ to focusing more on the individual parameter regularization, and eventually converges to those in Proposition 1. In practice, $n_e$ can be pre-specified if users have prior knowledge on the sparsity of $\boldsymbol{\theta}$, otherwise be regarded as a tuning parameter, chosen by the CV procedure or an information criterion (AIC or BIC) for model selection.

## 2.2   Algorithmic Steps of PANDA

The practical implementation of PANDA starts with some initial values for $\boldsymbol{\theta}$. The estimates of $\boldsymbol{\theta}$ and the variance terms of the pre-specified NGD are updated iteratively until convergence. The detailed steps are listed in Algorithm 1, along with some remarks on specifying the algorithmic parameters and convergence criterion (Remarks 1 to 5).

**Remark 1 (convergence criterion).** We provide several choices to evaluate the convergence of the PANDA algorithm. First, we may eyeball the trace plots of $\bar{l}^{(t)}$, which is often sufficient. Second, we can apply a cutoff value, say $\tau$ on the absolute percentage change on $\bar{l}^{(t)}$ from two consecutive iterations: if $|\bar{l}^{(t+1)} - \bar{l}^{(t)}|/\bar{l}^{(t)} < \tau$, then we may declare convergence. $\tau$ is supposed to be close-to-0 upon convergence, but being arbitrarily close to 0 would be difficult to achieve given the fluctuation around $\bar{l}^{(t)}$ with finite $m$ or $n_e$ due to the randomness of the augmented noises from iteration to iteration. Finally, we develop a formal statistical test for convergence
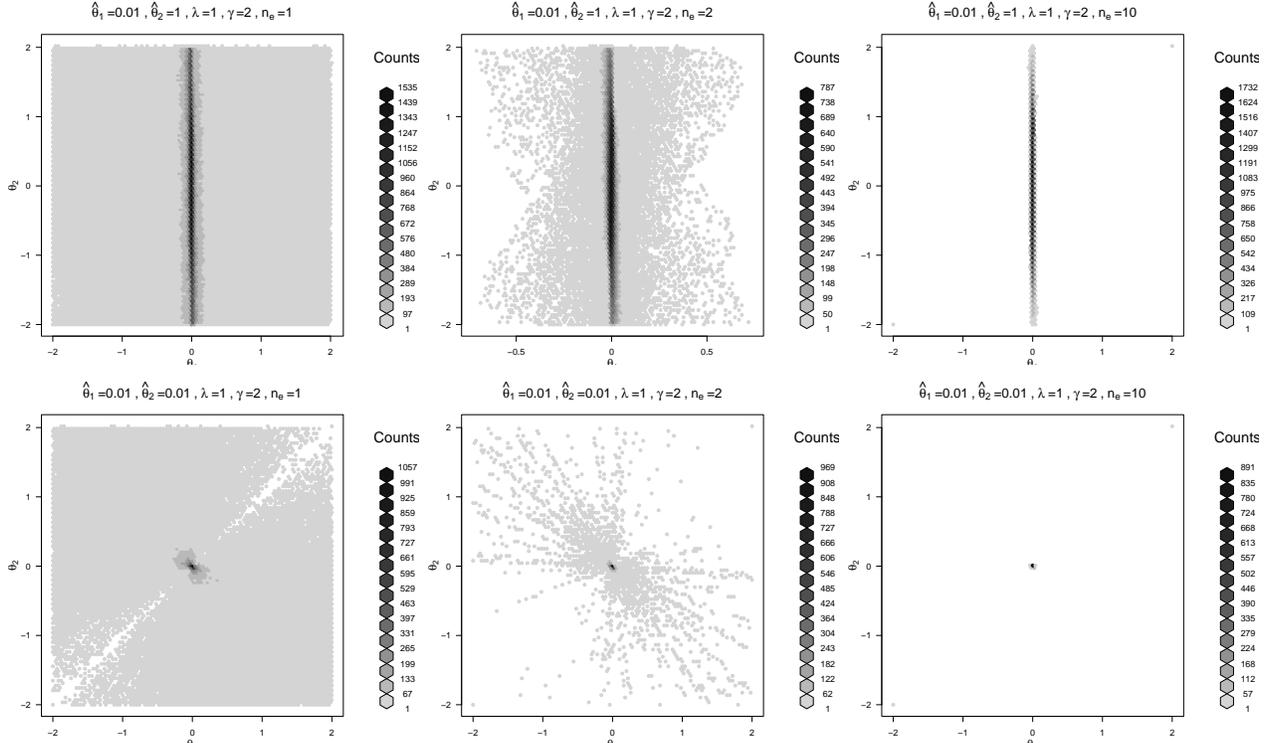
Figure 2: Heat maps of the $l_0$ constraint through the orthogonal regularization in PANDA with fixed $n_e$ for 2-dimensional $\boldsymbol{\theta} = (\theta_1, \theta_2)$

based on $\bar{l}^{(t)}$; but the test should be used with caution as it tends to claim non-convergence. The details of the test are provided in Sec S.7 of the supplementary materials.

**Remark 2** (**maximum iteration** $T$)**.** $T$ should be set at a number large enough so to allow the algorithm to reach convergence within a reasonable time period. When $n_e$ is large, we expect the algorithm to converge with a relatively small $T$ (in the examples in Sec 4, convergence is achieved for $T \leq 20$ with a large $n_e$). If $n_e$ is small, especially when PANDA is used to realize the $l_0$ regularization, $T$ should be set a large number for convergence.

**Remark 3** (**$m$ and $r$**)**.** In practical implement, $n_e$, no matter how large, is still finite. In addition, one might not want to set $n_e$ at a very large value as it will slow down the per-iteration computation. With a finite $n_e$, there is random fluctuation around the loss function and parameter estimates since each iteration is based on a different set of finite samples, even when the PANDA algorithm converges. To mitigate the random fluctuation, we can take the moving averages of the estimated parameters over multiple ($m$) iterations. The same rationale applies to the banking of $r$ estimates after convergence. In addition, taking the averages of the estimates obtained from the multiple augmented data sets also leads to a small generalization error due to the ensemble-learning type of effect PANDA brings (see Sec 3.4 for more details). In our empirical studies, $r = O(10)$ seems to be sufficient.

**Remark 4** (**Bounding at** $\tau_0$)**.** The bounding at $\tau_0$ is necessary. Despite the fact that estimates of zero-valued $\theta$ can get arbitrarily close to 0 (see Sec 3.1 for the almost sure convergence of the minimizers in PANDA), being exactly 0 cannot be achieved computationally in practice due to the numerical nature of PANDA. In addition, after the convergence of the

---

**Algorithm 1** PANDA for GLM

---

1: **Input**: initial parameter estimates $\bar{\boldsymbol{\theta}}^{(0)} = (\bar{\boldsymbol{\theta}}_1^{(0)}, \ldots, \bar{\boldsymbol{\theta}}_p^{(0)})$; NGD; maximum iteration $T$; noisy data size $n_e$ and moving average (MA) window width $m$ and number of banked parameter estimates $r$ after convergence; threshold $\tau_0$ .

2: **Output**: regularized parameter estimates $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\theta}}_1, \ldots, \hat{\boldsymbol{\theta}}_p)$

3: Centerize the observed independent variables $\mathbf{x}$.

4: $t \leftarrow 0$; convergence $\leftarrow 0$

5: **While** $t < T$ and convergence $= 0$

6:     $t \leftarrow t + 1$

7:     Generate $\mathbf{e}_x^{(t)}$ from NGD $\mathrm{N}(0, \mathrm{V}(\bar{\boldsymbol{\theta}}^{(t-1)}))$ and set $\mathbf{e}_y \equiv \bar{y}$ ($\sim$Bern$(\hat{p})$ for Bernoulli $Y$).

8:     Combine $(\mathbf{y}, \mathbf{x})$ with $(\mathbf{e}_y^{(t)}, \mathbf{e}_x^{(t)})$ to obtain the augmented data $(\tilde{\mathbf{y}}^{(t)}, \tilde{\mathbf{x}}^{(t)})$

9:     Run GLM on $(\tilde{\mathbf{y}}^{(t)}, \tilde{\mathbf{x}}^{(t)})$ and obtain MLE $\hat{\boldsymbol{\theta}}^{(t)}$. For linear regression with Gaussian $Y$, ordinary least-squares estimates are obtained.

10:     If $t > m$, calculate MA $\bar{\boldsymbol{\theta}}^{(t)} = m^{-1} \sum_{l=t-m+1}^{t} \hat{\boldsymbol{\theta}}^{(l)}$; otherwise $\bar{\boldsymbol{\theta}}^{(t)} = \hat{\boldsymbol{\theta}}^{(t)}$. Calculate the loss function $l^{(t)}$ with $\bar{\boldsymbol{\theta}}^{(t)}$ plugged in Eqn (3).

11:     Calculate the averaged loss function $\bar{l}^{(t)} = m^{-1} \sum_{l=t-m+1}^{t} l^{(l)}$.

12:     Let convergence $\leftarrow 1$ if $\bar{l}^{(t)}$ satisfies one of the convergence criteria (Remark 1).

13: **End While**

14: Run Lines 6 and 9 above for additional $m + r$ iterations, and record $\bar{\boldsymbol{\theta}}^{(l)}$ for $l = t + m + 1, \ldots, t + m + r$. Let $\bar{\boldsymbol{\theta}} = (\bar{\boldsymbol{\theta}}_1, \ldots, \bar{\boldsymbol{\theta}}_p)$, where $\bar{\boldsymbol{\theta}}_j = (\bar{\boldsymbol{\theta}}_j^{(t+m+1)}, \ldots, \bar{\boldsymbol{\theta}}_j^{(t+m+r)})$ for $k = 1, \ldots, p$.

15: Set $\hat{\theta}_j = 0$ if $\max\{|\bar{\boldsymbol{\theta}}_j|\} < \tau_0$; and $\hat{\theta}_j = r^{-1} \sum_{l=t+m+1}^{t+m+r} \bar{\boldsymbol{\theta}}_j^{(l)}$ o.w.

---

PANDA algorithm, there is still mild fluctuation around the parameter estimates due to the randomness of the augmented noise, especially when $n_e$ or $m$ is not large. We suggest bounding the absolute maximum of the estimates over a sequence of iterations as given in Algorithm 1, which seems to be a robust criterion in the empirical studies in Sec 4.

**Remark 5 (non-convex regularizers).** PANDA optimizes a convex objective function in each iteration of the GLM on the augmented data even when the targeted regularizer itself is non-convex, such as the SCAD or $l_0$. As such, PANDA does not run into the same type of computational difficulties that gradient-based techniques often experience for non-convex optimization (e.g., getting stuck in a local optimum). As a matter of fact, due to the stochastic nature of PANDA from iteration to iteration, it can escape from a local optimum especially if it is unstable, and lands at a more stable local optimum or even the global optimum. That being said, the initial values used in PANDA would also affect the final solutions when the targeted regularizer is non-convex.

## 2.3   $n_e$ **vs.** $m$

Upon convergence, the expected regularization in Proposition 1 can be realized either by letting $m \to \infty$ suggested by $\lim_{m\to\infty} m^{-1} \sum_{t=1}^{m} \sum_{i=1}^{n_e} (e_i^{(t)} - \sum_j e_{ik}^{(t)} \theta_j)^2$ or by letting $n_e \to \infty$ suggested by $n_e \lim_{n_e\to\infty} n_e^{-1} \sum_{i=1}^{n_e} (e_i^{(t)} - \sum_j e_{ik}^{(t)} \theta_j)^2$ under the constraint $n_e \mathrm{V}(e_j) = O(1) \ \forall \ \theta_j$. The constraint $n_e \mathrm{V}(e_j) = O(1)$ guarantees that injected noise $\mathbf{e}$ does not over-regularize

or overwhelm the information about $\boldsymbol{\theta}$ contained in the observed data $\mathbf{x}$ even when $n_e$ is large. For example, $\mathrm{V}(e_j) = \lambda|\theta_j|^{-1}$ in the case of the lasso-type noise, and $n_e\lambda$ can be treated as one tuning parameter. The targeted regularization implied by the lower-order term $n_e\big(C_1\sum_j \theta_j^2 \mathrm{V}(e_j)\big)$ in Eqn (4) can be approximated arbitrarily well as $n_e \to \infty$ with $n_e V(e_j) = O(1)$. When $m \to \infty$ and $n_e$ is fixed, there exists, more or less, other type of regularization on $\boldsymbol{\theta}$ on top of the targeted regularization given that the higher-order term $O\big(\sum_j \big(\theta_j^4 n_e \mathrm{V}^2(e_j)\big)\big)$ in Eqn (4) does not disappear. If we also require $n_e V(e_j) = O(1)$ in the large $m$ and small $n_e$ case, then $O\big(\sum_j\big(\theta_j^4 n_e \mathrm{V}^2(e_j)\big)\big) = O\big(\sum_j\big(\theta_j^4 \mathrm{V}(e_j)\big)\big)$, then the high-order term would also be ignorable if $\theta_j^4 V(e_j)$ is small.

Figure 3 illustrates the differences between the realized regularization effect $P(\boldsymbol{\theta})$, when the targeted regularization is lasso ($P(\boldsymbol{\theta}) = |\boldsymbol{\theta}|$), by letting $n_e \to \infty$ ($m$ is small) vs $m \to \infty$ ($n_e$ is small) and its relationships with $\boldsymbol{\theta}$ for several types of GLM (the regularization effects when $Y$ follows an exponential distribution are similar to when $Y$ is Poisson and the results from the former are not provided). For $n_e \to \infty$ ($\lambda n_e = 1$ fixed at 1 and $m = 50$), the realized penalty is identical to the targeted lasso in all four regression types, and is very close lasso at $n_e = 100$ except for some very mild random fluctuation. The realized regularization on $\boldsymbol{\theta}$ at $m \to \infty$ and small $n_e$ varies by regression type. When $|\boldsymbol{\theta}|$ is small, the target regularization is realized as the higher-order term that involves $|\boldsymbol{\theta}|$ in Eqn (4) is ignorable. As $|\theta|$ increases, the the higher-order term becomes less ignorable and regularization deviates from lasso, except for linear regression where the higher-order term is analytically 0. Specifically, the realized regularization is sub-linear for logistic and NB regression, and super-linear for Poisson regression.

In summary, to achieve the expected regularization effect in Proposition 1, one can set either $m$ or $n_e$ at a large number. Computationally, a large $n_e$ often requires less iterations even when $m$ is as small as 1. On the other hand, a very large $n_e$ slows down the computation per iteration. Taken together, the actual time taken to reach convergence might not differ that much between the two cases. In some sense, the choices on $n_e$ and $m$ more or less depends on each other. If a large $n_e$ still results in noticeable fluctuation around $\hat{\boldsymbol{\theta}}$, then a large $m$ can be used to speed up the convergence. For a small $n_e$, a relatively large $m$ should be used to yield stable penalty. Fig 4 shows the parameter estimation trajectories of zero-valued regression coefficients across with $\lambda$ in linear regression and Poisson regression on simulated data when the lasso-type noise is used in PANDA There are 30 predictors ($p = 30$) and $n = 100$ in each case. In the linear regression, the predictors were simulated from $\mathrm{N}(0, 1)$; in the Poisson regression, the predictors were simulated from $\mathrm{Unif}(-0.3, 0.5)$. Out of the 30 coefficients, 9 were set at 0, and the other 21 non-zero coefficients ranged from 0.5 to 1. The estimation trajectories for the 9 zero-valued parameters look very similar between large $n_e$ vs. large $m$ in both regression settings.

If the targeted penalty is $l_0$ (Proposition 2) and $n > p$, $n_e$ can be tuned within $[1, p]$. There are other considerations regarding the choices of $m$ and $n_e$ when using PANDA to obtain inference on $\boldsymbol{\theta}$. More details are provided in Sec 3.3.
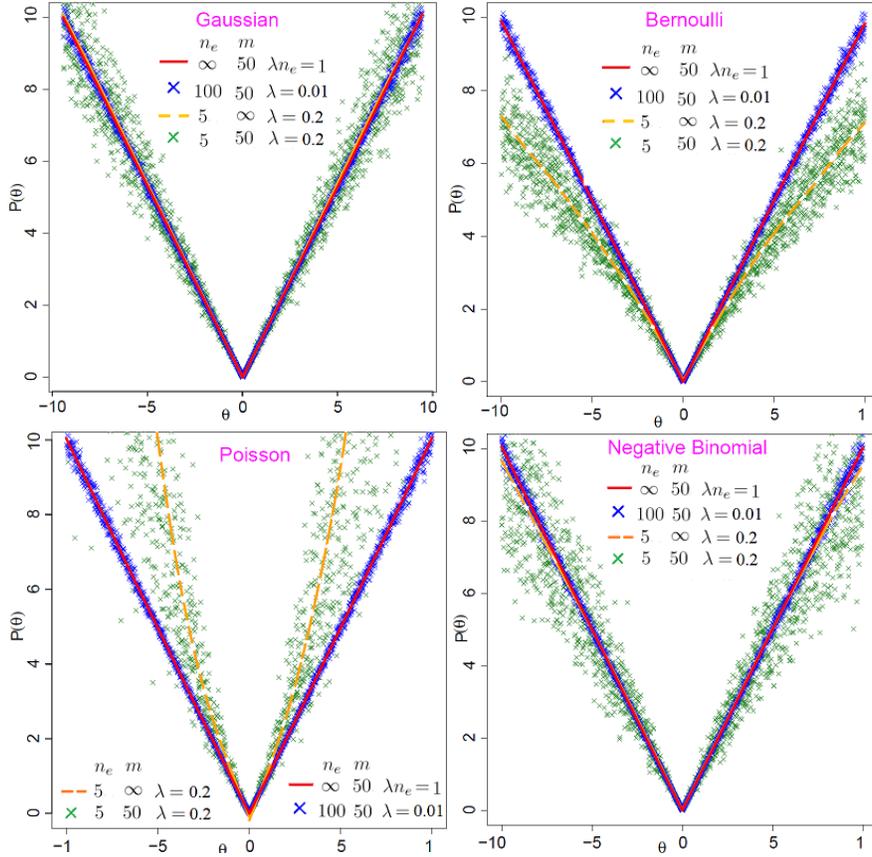
Figure 3: Realized regularization by PANDA in different GLMs for the targeted penalty $P(\boldsymbol{\theta}) = |\boldsymbol{\theta}|$.
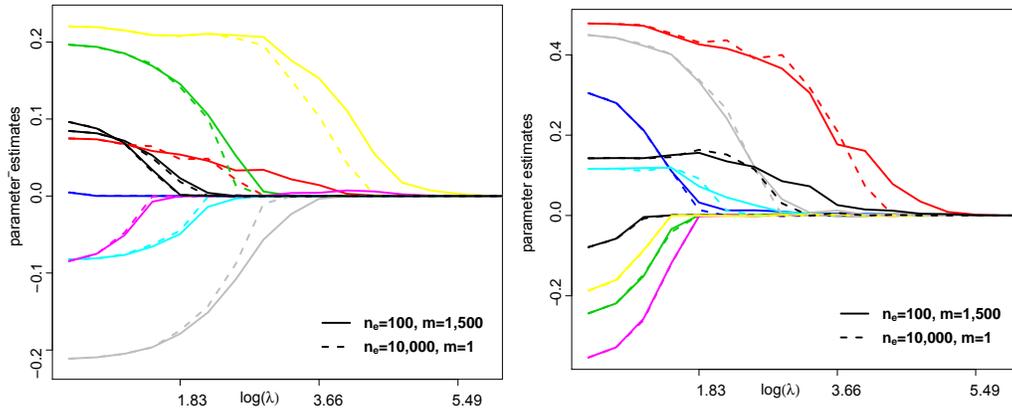


Figure 4: Estimation trajectories of zero-valued regression coefficients across $\lambda$ in linear regression (left) and Poisson regression (right) with lasso-type noise in PANDA

# 3 Theoretical Properties and Statistical Inferences

In this section, we establish the almost sure (a. s.) convergence of the data augmented loss function to its expectation and the a. s. convergence of the minimizer of the former to the minimizer of the expected loss function as $n_e \to \infty$ or $m \to \infty$ (Sec 3.1). We also examine the Fisher information of the parameters in noise-augmented data (Sec 3.2) and statistical inferences of the parameters via PANDA (Sec 3.3), and claim that PANDA exhibits ensemble learning behavior (Sec 3.4).

## 3.1 Almost sure convergence of noise augmented loss function and its minimizer

Let $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ denote the average loss function over $m$ iterations of the PANDA algorithm upon convergence. Theorem 1 presents the asymptotic properties of $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ under two scenarios: 1) $n_e \to \infty$ while $n_e \mathrm{V}(e_j) = O(1)$ for a given $\theta_j$ and $m\ (\geq 1)$ is fixed at a constant; 2) $m \to \infty$ when $n_e (> p)$ takes a finite constant.

**Theorem 1. (asymptotic properties of the noise-augmented loss function and its minimizer for PANDA)** Assume $\boldsymbol{\theta}$ belongs to a compact set. Let $l_p(\boldsymbol{\theta}|\mathbf{x}) = \mathrm{E}_{\mathbf{e}}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}))$.
1) If $n_e \to \infty$ while $n_e \mathrm{V}(e_j) = O(1)$ for any given $\theta_j$ and $m \geq 1$ is held at a constant, then

$$n_e^{1/2} C_1^{-1} \left( \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x}) \right) \xrightarrow{d} N(0,1) \tag{18}$$

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{a.s.} l_p(\boldsymbol{\theta}|\mathbf{x}) \overset{n_e \to \infty}{\longrightarrow} l(\boldsymbol{\theta}|\mathbf{x}) + P(\boldsymbol{\theta}) + C \tag{19}$$

$$\arg\inf_{\boldsymbol{\theta}} \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{a.s.} \arg\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x}), \tag{20}$$

2) If $m \to \infty$ while $n_e (> p)$ is fixed, then

$$m^{1/2} C_2^{-1} \left( \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x}) \right) \xrightarrow{d} N(0,1) \tag{21}$$

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{a.s.} l_p(\boldsymbol{\theta}|\mathbf{x}) \overset{m \to \infty}{\longrightarrow} l(\boldsymbol{\theta}|\mathbf{x}) + P(\boldsymbol{\theta}) + C \tag{22}$$

$$\arg\inf_{\boldsymbol{\theta}} \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{a.s.} \arg\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x}). \tag{23}$$

$P(\boldsymbol{\theta})$ in Eqns (19) and (22) is the same as defined in Proposition 1. $C_1$ and $C_2$ are functions of $\boldsymbol{\theta}$ and take different forms for different types of $Y$.

The proof of Theorem 1 is provided in Sec S.3 of the supplementary materials. There are two important takeaways. First, Theorem 1 states that $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ follows a Gaussian distribution at the rate of $\sqrt{n_e}$ and $\sqrt{m}$ under the two scenarios, respectively, implying that the augmented loss function in PANDA is trainable for practical implementation. The fluctuation of $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ around its expected value due to noise augmentation is controlled and the tail of the distribution of $d = \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x})$ decays to 0 exponentially fast in $n_e$ and $m$ as $\Pr(d > t) \leq \exp(-n_e t^2/2C^2)$ and $\Pr(d > t) \leq \exp(-mt^2/2C^2)$ for any $t > 0$. Second, $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ converges a.s. to its expectation (the penalized loss function given the observed data $(\mathbf{x}, \mathbf{y})$ with the targeted penalty term), guaranteeing that PANDA does what it is designed to do.

When there exists multicollinearity among $\mathbf{X}$, the loss function minimized in PANDA has an optimum region rather than a single optimum point. To examine the asymptotic properties in this case, we define the *optimum parameter set* (Definition 1) and show that the parameters learned by PANDA fall in the optimum parameter set asymptotically (Proposition 3).

**Definition 1. (optimum parameter set)** Let the expected loss function $l_p(\boldsymbol{\theta}|\mathbf{x})$ be a continuous function in $\boldsymbol{\theta}$. The optimum set is defined as $\boldsymbol{\Theta}^0 = \{ \boldsymbol{\theta}^0 \in \boldsymbol{\Theta} \mid l_p(\boldsymbol{\theta}^0|\mathbf{x}) \leq l_p(\boldsymbol{\theta}|\mathbf{x}), \forall \boldsymbol{\theta} \in \boldsymbol{\Theta} \}$, where $\boldsymbol{\Theta}$ is the set containing all possible parameter values. The distance from $\boldsymbol{\theta} \in \boldsymbol{\Theta}$ to $\boldsymbol{\Theta}^0$ is defined as $d(\boldsymbol{\theta}, \boldsymbol{\Theta}^0) = \min_{\boldsymbol{\theta}^0 \in \boldsymbol{\Theta}^0} ||\boldsymbol{\theta} - \boldsymbol{\theta}^0||_2$.

**Proposition 3. (consistency of parameter estimate in presence of multicollinearity)** Let $\hat{\boldsymbol{\theta}}_p^0 = \arg\min_{\boldsymbol{\theta}} \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$. Given

$$\sup_{\boldsymbol{\theta}} \left| \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - \bar{l}_p(\boldsymbol{\theta}|\mathbf{x}) \right| \to 0 \text{ as } n_e \to \infty \bigcap n_e \mathrm{V}(e_j) = O(1) \ \forall j = 1, ..., p \text{ or } m \to \infty; \tag{24}$$

and assume $\boldsymbol{\theta}$ is compact, then $\Pr\left(\limsup_{m\to\infty \text{ or } n_e\to\infty} d(\hat{\boldsymbol{\theta}}_p^0, \boldsymbol{\Theta}^0) \leq \delta\right) = 1 \ \forall \ \delta > 0$.

The proof is given in Sec S.4 of the supplementary materials. Multicollinearity does not affect the convergence of the loss functions in PANDA; therefore, Eqn (24) holds per the proof of Theorem 1.

## 3.2 Fisher Information in Noise Augmented Data

The augmented noise in PANDA brings endogenous information to observed data $\mathbf{x}$ to regularize the estimation of $\boldsymbol{\theta}$. The expected regularization can be achieved by letting $(n_e \to \infty) \cap (n_e \mathrm{V}(e_j) = O(1))$. At the first sight, it seems that a large amount of augmented noisy data could potentially overshadow the information about the parameters in the observed data, leading to over-regularization. We claim that this is not the case because of the constraint $n_e \mathrm{V}(e_j) = O(1) \ \forall \ j$. In other words, $n_e$ combined with the tuning parameters from the NGD variance term is treated as a single tuning parameter. For example, with the lasso-type noise, $n_e \lambda$ is treated one tuning parameter: if $n_e$ is large, then $\lambda$ takes a small value so to keep $n_e \lambda = O(1)$. Proposition 4 provides the theoretical justification that, as long as $n_e \mathrm{V}(e_j) = O(1)$ for any given $\theta_j$, the amount of regularization brought by the augmented data to $\theta_j$ remains as constant even for $n_e \to \infty$. Proposition 4 is established in the context of the bridge-type noise; the same conclusion can be obtained for other noise types in a similar fashion. The proof is provided in Sec S.5 of the supplementary materials.

**Proposition 4** (**Fisher information in noise augmented data**). The regularization on the coefficients $\boldsymbol{\theta}$ in GLM introduced through the augmented bridge-type noise is proportional to $n_e \lambda |\boldsymbol{\theta}|^{-\gamma}$. Specifically, $I_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(\boldsymbol{\theta})$, the Fisher information on $\boldsymbol{\theta}$ contained in the augmented data $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ is the summation of $I_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta})$, the Fisher information on $\boldsymbol{\theta}$ contained in the observed data, and $I_{\mathbf{e}}(\boldsymbol{\theta})$, the amount of regularization on $\boldsymbol{\theta}$.

$$I_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(\boldsymbol{\theta}) = I_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta}) + (\lambda n_e) B''(\theta_0 + 0) \mathrm{Diag}\{|\theta_1|^{-\gamma}, \ldots, |\theta_p|^{-\gamma}\} + O(\lambda n_e^{1/2}) J, \tag{25}$$

where $J$ is a $p \times p$ matrix with all elements at 1. The higher-order term $O(\lambda n_e^{1/2})$ becomes $O(\lambda^{1/2})$ if $\lambda n_e = O(1)$ and is ignorable if $\lambda$ is small. Eqn (25) suggests that the information about $\boldsymbol{\theta}$ does not increase with $n_e$ as along as $\lambda n_e |\boldsymbol{\theta}|^{-\gamma}$ is kept at a constant. In addition, the closer $|\boldsymbol{\theta}|$ is to 0, the stronger the regularization the augmented information brings to $\boldsymbol{\theta}$.

## 3.3 Asymptotic Distribution of Regularized Parameters via PANDA

Proposition 5 presents the asymptotic distribution of the estimated $\hat{\boldsymbol{\theta}}$ via PANDA, based on which we can obtain statistical inferences for $\boldsymbol{\theta}$. The proof is given in Sec S.6 of the supplementary materials.

**Proposition 5** (**asymptotic distribution of parameter estimates via PANDA**). Let $\hat{\boldsymbol{\theta}}^{(t)}$ denote the estimate of $\boldsymbol{\theta}$ in iteration $t$ of the PANDA algorithm. The final estimate for $\boldsymbol{\theta}$ is denoted by $\bar{\boldsymbol{\theta}} = r^{-1} \sum_{t=1}^{r} \hat{\boldsymbol{\theta}}^{(t)}$ from $r \geq 1$ iterations after convergence. Assume $n_e \mathrm{V}(e) = o(\sqrt{n}) \ \forall \boldsymbol{\theta}$.

$$\sqrt{n}(\hat{\boldsymbol{\theta}}^{(t)} - \boldsymbol{\theta}) \xrightarrow{d} N(\mathbf{0}, \Sigma^{(t)}) \text{ as } n \to \infty, \tag{26}$$

$$\sqrt{n}(\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}) \xrightarrow{d} N(\mathbf{0}, \bar{\Sigma} + \Lambda) \text{ as } n \to \infty; r \to \infty, \tag{27}$$

where $\Sigma^{(t)} = I_{\tilde{\mathbf{x}}^{(t)}, \tilde{\mathbf{y}}}(\boldsymbol{\theta})^{-1} I_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta}) I_{\tilde{\mathbf{x}}^{(t)}, \tilde{\mathbf{x}}}(\boldsymbol{\theta})^{-1}$ in iteration $t$, $\bar{\Sigma} = r^{-1} \sum_{t=1}^{r} \Sigma^{(t)}$, and $\Lambda = \mathrm{V}(\hat{\boldsymbol{\theta}}^{(t)})$ is the between iteration variability of $\hat{\boldsymbol{\theta}}^{(t)}$.

The regularity condition $n_e \mathrm{V}(e) = o(\sqrt{n})$ takes different forms for different NGDs (e.g., for the bridge-type noise, it is $\lambda n_e = o(\sqrt{n})$). The asymptotic variance of $\hat{\boldsymbol{\theta}}^{(t)}$ involves the inverse of $I_{\tilde{\mathbf{x}}^{(t)}, \tilde{\mathbf{y}}}(\boldsymbol{\theta})$, which always exists given the augmented data. Eqn (27) suggests the overall variance on $\bar{\boldsymbol{\theta}}$ is the summation of two variance components, $\bar{\Sigma}$, the per-iteration variance of $\hat{\boldsymbol{\theta}}^{(t)}$, and $\Lambda$, the between-iteration variance of $\hat{\boldsymbol{\theta}}^{(t)}$. $\bar{\Sigma}$ contains the unknown $\boldsymbol{\theta}$ and can be estimated by plugging in $\hat{\boldsymbol{\theta}}^{(t)}$, with the caveat that the uncertainty around $\hat{\boldsymbol{\theta}}^{(t)}$ is not accounted for. $\Lambda$ can be estimated by the sample variance of $\hat{\boldsymbol{\theta}}^{(t)}$ over $r$ iterations; that is, $(r-1)^{-1} \sum_{t=1}^{r} (\hat{\boldsymbol{\theta}}^{(t)} - \bar{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}}^{(t)} - \bar{\boldsymbol{\theta}})^T$.

In the case of linear regression, the asymptotic distribution of $\hat{\boldsymbol{\theta}}^{(t)}$ in Eqn (26) becomes

$$\sqrt{n}(\hat{\boldsymbol{\theta}}^{(t)} - \boldsymbol{\theta}) \xrightarrow{d} N\left(\mathbf{0}, \sigma^2 (\boldsymbol{M}^{(t)})^{-1} (\mathbf{x}^T \mathbf{x}) (\boldsymbol{M}^{(t)})^{-1}\right), \tag{28}$$

where $\mathbf{M}^{(t)} = (\mathbf{x}^T \mathbf{x} + n_e \mathrm{diag}(\mathrm{V}(\mathbf{e})))$. The asymptotic variance in Eqn (28) contains the unknown $\sigma^2$, which can be estimated by $\hat{\sigma}^2 = \mathrm{SSE}/(n - \nu) = (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\theta}}^{(t)})^T (\mathbf{y} - \mathbf{x}\hat{\boldsymbol{\theta}}^{(t)})/(n - \nu)$ with the degree of freedom $\nu = \mathrm{tr}(\mathbf{x}(\boldsymbol{M}^{(t)})^{-1} \mathbf{x}^T)$. $\hat{\sigma}^2$ converges to $\sigma^2 \chi_{n-\nu}^2$ in distribution.

When applying PANDA to obtain inference in GLMs, we should set $n_e$ at a small number and $m$ at a large number to achieve valid inference and targeted regularization effect simultaneously. We recommend $n_e = o(n)$ as long as $n_e + n > p$ (e.g., one order of magnitude smaller than $n$), especially when $n$ is small. This is different from when the main goal is variable selection (except for $l_0$), regularized estimation, or prediction without uncertainty quantification, where a large $n_e$ can be used to achieve the targeted regularization effect with fewer iterations per Proposition 1. The reason is that a large $n_e$ (relative to $n$) tends to underestimate $\bar{\Sigma} + \Lambda$, the asymptotic variance of $\bar{\boldsymbol{\theta}}$, resulting in a lower-than-nominal coverage rate and an inflated type I error rate. As mentioned above, $\bar{\Sigma} = r^{-1} \sum_{t=1}^{r} \Sigma^{(t)}$ is estimated by plugging in $\hat{\boldsymbol{\theta}}^{(t)}$ for $t = 1, \ldots, r$ upon convergence, pretending that it is the true parameter value and thus ignoring the uncertainty around it. Though this issue exists regardless of whether a large or a small $n_e$ is used, using a small $n_e$ helps to re-capture this lost variability with the between-iteration variability $\Lambda$. Specifically, $\hat{\boldsymbol{\theta}}^{(t)}$ is a regularized estimate from minimizing a loss function summed over the data component $\mathbf{x}$ and a penalty term, or equivalently, a summation of the loss functions constructed from the data component $\mathbf{x}$ and the augmented data component $\mathbf{e}$ in the context of PANDA. Instead of focusing on how $\hat{\boldsymbol{\theta}}^{(t)}$ changes with sample data $\mathbf{x}$, which is fixed throughout iterations, we quantify how it changes with $\mathbf{e}$. If a large $n_e$ is used, the ignored sampling variability around $\hat{\boldsymbol{\theta}}^{(t)}$ can hardly be recovered through $\Lambda$ as it is close to 0, which is easy to understand as the realized regularization effect with a large $n_e$ is close to its expectation and it is almost like solving the same analytical constrained optimization at every iteration, leading to very similar $\hat{\boldsymbol{\theta}}^{(t)}$ across iterations upon convergence.

## 3.4   Ensemble Learning Behavior of PANDA with Fixed $n_e$

Ensemble learning methods combine multiple learners to achieve better predictive performance than that from an individual learner. Let $Y$ be the observed outcome and $\bar{Y}$ be its prediction

from an ensemble method. Brown et al. (2005) suggest that the generalization error of the ensemble method made of $M$ learners, $E(\bar{Y} - Y)^2$, can be decomposed as

$$M^{-2}\Big[\big(\textstyle\sum_i(\mathrm{E}(\hat{Y}_i) - Y)\big)^2 + \sum_i \mathrm{E}\big(\hat{Y}_i - E(\hat{Y}_i)\big)^2 + \sum_i \sum_{j \neq i} E\big((\hat{Y}_i - \mathrm{E}(\hat{Y}_i))(\hat{Y}_j - E(\hat{Y}_j))\big)\Big], \quad (29)$$

where $\hat{Y}_i$ refers to the prediction from the $i$-th learner in the ensemble for $i = 1, \ldots, M$. The success of ensemble methods, in part, can be attributed to the diversity term among the $M$ learners that is captured by the third term (covariance) in Eqn (29): as the diversity increase, the covariance decrease, and the overall generalization error decreases. The diversity can be achieved by perturbing the training data such as taking a subset of observation, or a subset of attributes to train the learners.

We show that PANDA, in addition to achieving the targeted regularization effects, also exhibits some ensemble learning behavior with a fixed $n_e$, which may propel it to edge out the existing constrained regularization approaches with smaller generalization error in prediction. Intuitively, upon convergence, the final estimates of $\boldsymbol{\theta}$ are averages over the estimates trained from different sets of noise augmented data from $r$ iterations, generating the diversity among the learners needed for the ensemble learning.

**Claim 1** (**Ensemble learning behavior of PANDA with fixed $n_e$**). Upon convergence, the average estimates over the sets of parameter estimates from multiple iterations of PANDA with a fixed $n_e$ can be regarded an ensemble learner.

If the diversity brought by PANDA with a fixed (small) $n_e$ and a large $m$ surpasses the increase in MSE (the sum of the first two terms in Eqn (29)), PANDA would lead to a smaller generalization error compared to the existing constrained optimization approaches for penalized GLM regression that don't promote diversity.

# 4 Numerical Examples

## 4.1 $l_0$ Regularization via PANDA

We demonstrate the PANDA-$l_0$ regularization in linear regression using the prostate cancer dataset and in logistic regression using the kyphosis dataset (Tibshirani, 1996). The prostate cancer dataset consists of 8 $X$'s and 97 observations. We standardized the $X$'s and centralized $Y$ prior to the application of the PANDA algorithm. The kyphosis dataset consists of 81 observations ($64 : 17$ for $Y = 0 : 1$). We included both the linear and the quadratic terms of the three standardized $X$'s ($X_4, X_5, X_6$ are the quadratic terms of $X_1, X_2, X_3$, respectively), in the logistic regression following Tibshirani (1996). We examine the regression coefficient estimation trajectories as $n_e$ increases from 1 to $p$ and as $\lambda$ increases while holding $n_e$ constant. For comparison, we also run the lasso regression in each case via the R package `glmnet`.

The results are presented in Fig 5. Column A shows that the PANDA-$l_0$ regularization shrinks only $n_e$ coefficients towards 0, leaving the other coefficients unregularized, but lasso shrinks all coefficients simultaneously. The observations are consistent with Proposition 2 and Eqn (17), which state the number of selected variables through PANDA-$l_0$ is $p - n_e$ for $n_e < p$. In the logistic regression case, due to the high correlations (0.957, 0.969 and 0.974) between the linear and the quadratic terms, the shrinkage occurs roughly around the same $\lambda$ for a fixed $n_e$ for each linear+quadratic pair in the trajectory. The plots in column B examine the effect of $\lambda$ on the estimation trajectory fixing $n_e$ at $p - p_0$ in PANDA-$l_0$, where $p_0$ is number

of variables selected by lasso ($p_0 = 3$ in linear regression and $p_0 = 4$ in logistic regression). As $\lambda$ increases, $n_e$ coefficients shrinks to 0. Further increasing $\lambda$ has no regularization effect on the remaining non-zero coefficients, despite some minor fluctuation around the non-zero parameter estimates. The plots in column C are similar to column B but $n_e$ is fixed at $p$. For small $\lambda$, the estimation trajectories are similar to using $n_e < p$ as in column B; as $\lambda$ continues to increase, the non-zero coefficients eventually get shrunk to 0, but in a different manner than lasso in the sense that its shrinkage process is not gradual but rather abruptly. For $n_e > p$, the estimation trajectories would be somewhere between column C and the lasso trajectories (e.g. Fig 4), and eventually become the lasso trajectories as $n_e$ becomes very large.
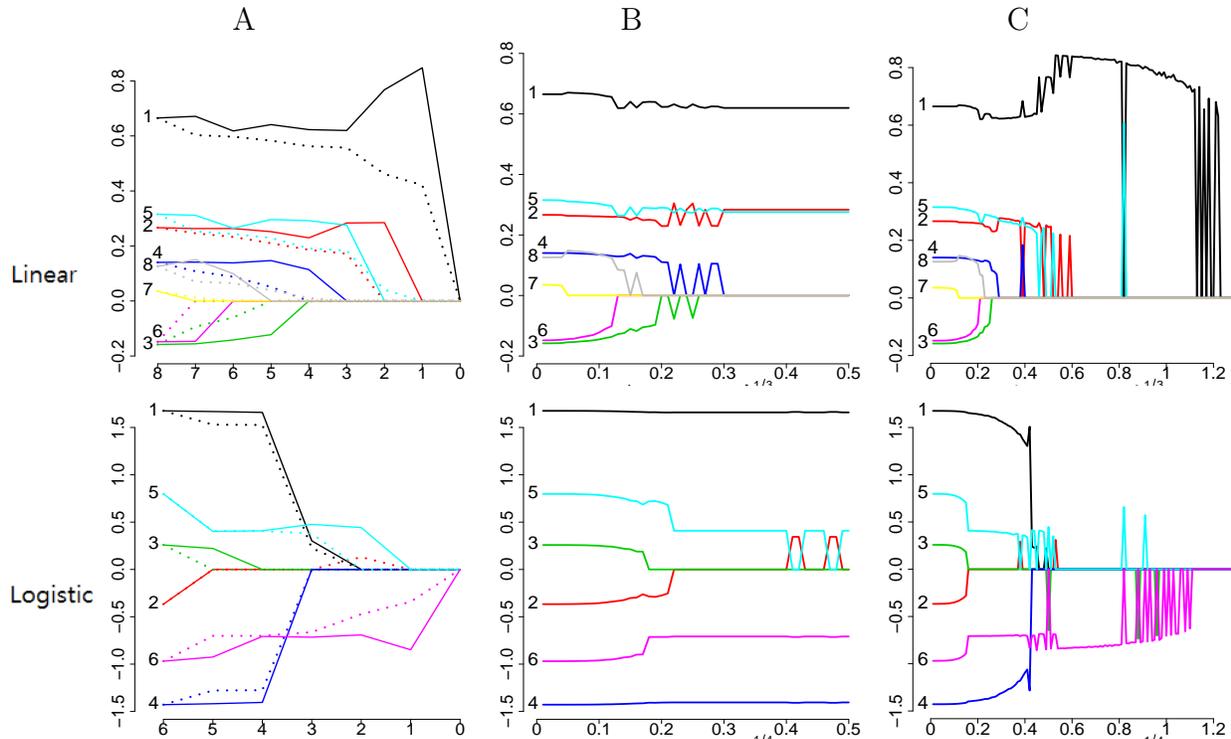


Figure 5: Estimation trajectory in linear and logistic regression as $n_e$ changes (column A), $\lambda$ changes at fixed $n_e < p$ (column B), and $\lambda$ changes at $n_e = p$ (column C). The solid lines in column A are from the PANDA-$l_0$ regularization by varying $n_e$ from 1 to $p$ and the dash lines represent the lasso regression via R package `glmnet` with the smallest $\lambda$ that yields $p - n_e$ non-zero estimates.

## 4.2 Inference for GLM parameters via PANDA

We investigate the inferential validity for GLM coefficients based on the asymptotic distributions in Proposition 5 via simulation studies. We examine Gaussian ($\sigma^2 = 1$), Poisson, Bernoulli, exponential (exp), and Negative Binomial (NB) (number of failure fixed at $r = 5$) outcomes with $p = 30$ in each case. For the Gaussian and NB outcomes, the predictors were simulated from $N(0, 1)$; for the Bernoulli, exp, and Poisson outcomes, the predictors were simulated from $\text{Unif}(-3, 3)$, $\text{Unif}(-1, 2)$, and $\text{Unif}(-0.3, 0.5)$, respectively. We examine three sample size scenarios $n = 50, 70, 100$, with 500 repetitions in each simulation case. The bridge-type noise is employed with $\gamma = 1, n_e = n, \lambda n_e \in (1.5, 7)$ in logistic regression and $\gamma = 2, n_e = 9, \lambda = n/10$ in the other GLMs. The achieved regularization effect is lasso in the logistic regression and $l_0$ in the other GLMs as $n_e = 9$ is set at the number of zero coefficients. In each repetition, we calculate the 95% CIs for the 30 regression coefficients (21 are non-

zero and 9 are zero) and examine the coverage probabilities (CP) and the CI widths. Tables 3 presents the results, benchmarked against the post-lasso inferential procedure (Lee et al., 2016; Taylor and Tibshirani, 2017) implemented via the R package `selectiveInference`.

Table 3: Empirical CP and CI width via PANDA and post-selection procedures with lasso penalty

| | zero coefficients (9) | | | | | | non-zero coefficients (21) | | | | | |
| | PANDA | | | post selection | | | PANDA | | | post selection | | |
| sample size | 50 | 70 | 100 | 50 | 70 | 100 | 50 | 70 | 100 | 50 | 70 | 100 |
| | mean CP (%) among the 9 coefficients | | | | | | mean CP (%) among the 21 coefficients | | | | | |
| Gaussian | 98.2 | 99.5 | 99.9 | NA | NA | NA | 91.4 | 96.5 | 97.7 | 92.3 | 93.2 | 94.2 |
| Bernoulli | 100 | 99.5 | 96.6 | NA | NA | NA | 97.3 | 88.4 | 92.6 | 65.9 | 75.2 | 82.9 |
| Exp | 95.1 | 95.5 | 96.5 | - | - | - | 87.1 | 99.5 | 94.4 | - | - | - |
| Poisson | 92.2 | 95.8 | 98.5 | - | - | - | 87.0 | 87.5 | 94.1 | - | - | - |
| NB | 95.8 | 99.4 | 100 | - | - | - | 83.1 | 95.4 | 99.9 | - | - | - |
| | mean CI width among the 9 coefficients | | | | | | mean CI width among the 21 coefficients | | | | | |
| Gaussian | 0.28 | 0.15 | 0.08 | NA | NA | NA | 0.91 | 0.74 | 0.57 | 29.6 | 2.01 | 1.26 |
| Bernoulli | 14.6 | 1.32 | 0.93 | NA | NA | NA | 24.8 | 2.15 | 1.46 | 22.0 | 10.6 | 4.64 |
| Exp | 0.39 | 0.23 | 0.14 | - | - | - | 1.07 | 0.95 | 0.77 | - | - | - |
| Poisson | 0.76 | 0.44 | 0.25 | - | - | - | 1.28 | 1.08 | 0.91 | - | - | - |
| NB | 0.54 | 0.28 | 0.15 | - | - | - | 1.19 | 1.11 | 1.05 | - | - | - |

NA: Not Available. R Package `selectiveInference` does not provide inference for coefficients whose estimates are 0. In addition, it only produces CIs for linear and logistic regression with the lasso regularization. CIs obtained by `selectiveInference` that have infinite lower/upper bounds are excluded from the summary ($4 \sim 18\%$).

For true zero-valued coefficients, PANDA maintains the nominal 95% coverage for all the examined outcome types and sample sizes. The R `selectiveInference` package does not provide inference for coefficients whose estimates are 0 (that is, not selected by lasso in the first place). Among these 9 zero-valued coefficients, lasso only selected some of them a few times out of the 500 repetitions. When the true coefficients are not 0, the CIs from PANDA have better coverage with much narrower CIs than the post-selection procedure (except for logistic regression at $n=50$). The post-selection procedure experiences severe under-coverage in the logistic regression for all $n$. We also examined the case of a larger $n_e$ ($n_e = 2n$ in the logistic regression and $n_e=n$ for the other GLMs). There was some under-coverage (CP $\geq\sim 90\%$ for zero coefficients; $\geq\sim 80\%$ for non-zero coefficients), which improved as $n$ increased.

## 4.3 Comparison with Existing Regularization in Linear and Logistic Regression

To examine the regularization effects by PANDA in linear and logistic regression, we use the same simulation setting as Examples 4.1 and 4.3 in Fan and Li (2001). In the linear regression, $Y = \mathbf{x}\boldsymbol{\beta} + \epsilon$, where $\boldsymbol{\beta}^T = (3, 1.5, 0, 0, 2, 0, 0, 0)$ ($p = 8$), $x_j \sim N(0, 1)$ for $j = 1, \ldots, p$ with $\text{corr}(x_j, x_{j'}) = 0.5^{|j-j'|}$, and $\epsilon \sim N(0, \sigma^2)$. Three sets of $(n, \sigma)$ were examined: (40, 3), (40, 1), and (60, 1). For the logistic regression, $n$ was set at 200; $Y \sim \text{Ber}\left(e^{\mathbf{X}^T\boldsymbol{\beta}}/(1 + e^{\mathbf{X}^T\boldsymbol{\beta}})\right)$, where the first six components of $\mathbf{X}$ were the same as those in linear regression and the last two components $x$ were drawn from Bernoulli(0.5) independently.

The medians of relative model error (MRME) and the number of correctly and incorrectly identified zero coefficients (out of 5) over 100 repetitions were obtained in each regression case. The estimates from the ridge, lasso, adaptive lasso and EN regressions via the existing approaches were obtained from R package `glmnet` and those from SCAD were from R package

`ncvreg`. We examine two scenarios of PANDA: large $n_e$/small $m$ and large $n_e$/small $m$. The specific values of $n_e$ and $m$, along with other PANDA algorithmic parameters are summarized in Table **??** in the supplementary materials. The results are presented in Table 4 for the linear regression and in Table 5 for the logistic regression. In summary, PANDA is either consistent with or performs better (due to its additional ensemble behaviors) than existing approaches for the same type of regularizer, per the MRME and the true 0/false 0 counts. The superiority of PANDA is specially obvious in the logistic regression. In general, PANDA-SCAD and PANDA-$l_0$ have the best performance.

Table 4: PANDA vs. Existing Approaches in Penalized Linear Regression with Various Regularizers

| | ridge | lasso | adaptive lasso | EN | SCAD | $l_0$ | ridge | lasso | adaptive lasso | EN | SCAD | $l_0$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | \multicolumn{6}{MRME} | | | # of correctly/incorrectly identified zero coefficients | | | | |
| | \multicolumn{12}{$n = 40, \sigma = 3$} | | | | | | | | | | |
| Existing | 80.09 | 67.86 | 68.47 | 68.24 | 72.79 | | 0/0 | 2.78/0.04 | 2.80/0 | 2.66/0.03 | 3.59/0.09 | |
| PANDA | 80.06 | 67.70 | 67.18 | 68.31 | 72.50 | 78.99 | 0.01/0 | 2.37/0.01 | 2.69/0.01 | 2.50/0.01 | 4.01/0.17 | 3.83/0.13 |
| | \multicolumn{12}{$n = 40, \sigma = 1$} | | | | | | | | | | |
| Existing | 94.60 | 68.03 | 68.32 | 69.45 | 44.42 | | 0/0 | 2.87/0 | 2.83/0 | 2.56/0.03 | 4.72/0 | |
| PANDA | 95.24 | 67.38 | 63.58 | 68.40 | 44.87 | 45.11 | 0.13/0 | 2.69/0 | 3.07/0 | 2.62/0 | 4.91/0 | 4.86/0 |
| | \multicolumn{12}{$n = 60, \sigma = 1$} | | | | | | | | | | |
| Existing | 97.40 | 66.40 | 68.34 | 67.92 | 44.91 | | 0/0 | 2.61/0 | 2.66/0 | 2.55/0.03 | 4.96/0 | |
| PANDA | 97.62 | 66.22 | 61.48 | 67.02 | 44.82 | 44.77 | 0.19/0 | 2.55/0 | 3.06/0 | 2.43/0 | 5.00/0 | 5.00/0 |

Table 5: PANDA vs. Existing Approaches in Penalized Logistic Regression with Various Regularizers

| | ridge | lasso | adaptive lasso | EN | SCAD | $l_0$ | ridge | lasso | adaptive lasso | EN | SCAD | $l_0$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | \multicolumn{6}{MRME} | | | # of correctly/incorrectly identified zero coefficients | | | | |
| Existing | 85.16 | 68.67 | 67.96 | 69.71 | 48.14 | | 0.07/0 | 2.10/0 | 2.05/0 | 2.09/0 | 4.31/0 | |
| PANDA | 76.50 | 61.15 | 58.60 | 62.14 | 34.87 | 37.66 | 0.17/0 | 2.43/0 | 2.83/0 | 2.44/0 | 4.97/0.02 | 4.89/0 |

## 4.4  Sports Article Objectivity Data

We implemented PANDA in a real-life dataset that contains 1000 sports articles that are labeled "objective" or "subjective". The data set is available for download from the UCI Machine Learning Repository (Rizk and Awad, 2018). There are 59 variables in the original data. The independent variables $X$'s are the extracted features from the articles such as the frequencies of different types of words, (e.g., the objective and subjective SENTIWORDNET scores, foreign words, subordinating preposition or conjunction) and frequencies of different types of punctuation (e.g., questions marks, exclamation marks), and text complexity score, among others. After removing the redundant features (perfectly linear dependent variables) and the highly imbalanced features (e.g, $>99\%$ in one category), and adjusting for the total word counts, we kept 48 $X$'s plus $Y$ (365 "subjective" and 635 "objective"). We split the 1000 cases into 800 training samples and 200 testing samples (100 subjective vs. 100 objective).

We learned the logistic regression parameters based on the 800 training samples and make predictions for the 200 testing samples via the trained model. We run the logistic regression with lasso, ridge, EN, and adaptive lasso penalties via the R package `glmnet`, and with the SCAD penalty via the `ncvreg` package, and obtained the regularized regression with the same types of penalty listed the above using PANDA. For the existing approaches, the 10-fold CV was used for hyper-parameter tuning. For PANDA, we run 100 iterations with $n_e = 1000$ and

$m = 10$. The algorithm converged after $10 \sim 15$ iterations, and the final parameter estimates were averaged over the last $r = 20$ iterations with $\tau_0 = 0.01$.

Table 6 presents the results on the MSE, classification accuracy rate, and computational time. Compared to the MLE from the non-regularized logistic regression, the prediction MSE and the accuracy rate on the testing samples via PANDA are similar or slightly better with the regularizers realized with the R packages `glmnet` and `ncvreg`. Specifically, the prediction MSE decreases by $\approx 10\%$; the accuracy increases by 1.5% to 2% for the same regularizer types. The number of zero coefficients ranges about 10 to 20 (out of a total of 48), depending on which regularizer is used per `glmnet` and `ncvreg`. PANDA took about 1.5 to 2 seconds to run 50 iterations. However, 25 iterations (costing 0.7 to 1 seconds) would also be sufficient for this application. Suppose $t$ values of tuning parameters are used in a $K$-fold cross-validation. The total time including the hyper-parameter tuning would be around $0.7tK$ to $tK$ seconds. Say $t = 10$ and $K = 10$, then it will take about 1 to 1.5 mins for PANDA. This is significantly longer than the existing method, which is expected since PANDA involves random sampling of data points and running GLM for every iteration.

Table 6: PANDA vs Existing Approaches in the Sports Article Objectivity Data

| penalty | ridge | lasso | EN | adaptive lasso | SCAD | $l_0$ |
|---|---|---|---|---|---|---|
| Prediction MSE (0.1573 with MLE) | | | | | | |
| Existing | 0.1539 | 0.1561 | 0.1544 | 0.1561 | 0.1629 | |
| PANDA | 0.1312 | 0.1260 | 0.1277 | 0.1280 | 0.1340 | 0.1448 |
| Accuracy Rate/Sensitivity/Specificity (%): 78.5/94/63 with MLE | | | | | | |
| Existing | 78.5/95/62 | 77.5/95/60 | 77.5/95/60 | 78/95/61 | 77/94/60 | |
| PANDA | 83.5/91/76 | 82.5/87/78 | 84.5/90/79 | 82/86/78 | 81.5/85/78 | 79.5/92/67 |
| # of zero-valued coefficients (0 with MLE) | | | | | | |
| Existing | 0 | 8 | 4 | 10 | 25 | |
| PANDA‡ | 1 | 8 | 4 | 10 | 25 | 19 |
| Computational Time (sec)$^{\|}$ | | | | | | |
| Existing | 0.7 $\sim$ 0.8 | | | | | |
| PANDA | 0.3 $\sim$ 0.4 per 10 iterations | | | | | $\sim$2.5 |

‡ hyperparameters were tuned to match the # of zero-valued coefficients in existing methods.
$^{\|}$ V1.1.463 on PC (Intel Core i7-7660U CPU @ 2.50 GHz)

We also run PANDA using the same tuning parameters selected by the R packages for the existing approaches for the same type of regularizer. PANDA performs better than the existing approaches with smaller RMSE, slightly better accuracy rates, and doubled zero coefficients in most cases.

# 5   Discussion

PANDA is a regularization technique through noise augmentation. PANDA effectively regularizes parameter estimation and allows valid inferences for GLMs, and displays ensemble learning behavior in certain cases. We establish the Gaussian tail of the noise-augmented loss function and the almost sure convergence to its expectation – a penalized loss function with the targeted regularizer, providing the theoretical justification for PANDA as a regularization technique and that the noise-augmented loss function is trainable. For a pre-fixed $n_e < p$, we show that PANDA is equivalent to imposing $n_e$ linear constraints on parameters and can lead to the $l_0$ regularization. PANDA is straightforward to implement. There is no need for sophis-

ticated optimization techniques as PANDA can leverage existing functions or procedures for running GLMs in any statistical software. In terms of the computational time, large $n_e$ usually leads to convergence with a small number of iterations, but the per-iteration computational cost can be high. If PANDA is applied to yield the $l_0$ regularization or to obtain inference in GLMs on top of variable selection, a small $n_e$ with a relatively large $m$ should be used.

The PANDA algorithm calculates $\bar{\boldsymbol{\theta}}$, the average of $m$ minimizers of $l(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ from the latest $m$ iterations, so to leverage the existing software for running GLM and to maintain its computational advantage over the existing approaches that employ sophisticated optimization techniques. Proposition 1 suggests the average of $m$ noise-augmented loss function $l(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ yields a single minimizer $\hat{\boldsymbol{\theta}}$, the Monte Carlo version of $\mathrm{E}_{\mathbf{e}}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ as $m \to \infty$. We establish in Corollary S.1 in the supplementary materials that $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$ are first-order equivalent for large $m$ and $n_e$ for PANDA in linear regression, We also present simulation results in the linear regression and Poisson regression settings to illustrate the similarity between $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$.

For linear regression, the OLS estimator obtained from the noise-augmented data in each iteration of the PANDA algorithm is a *weighted ridge estimator* on the observed data. Compared to a regular ridge estimator, where the same constant is added to all the diagonal elements of $\mathbf{x}^T\mathbf{x}$, different constants are used for different diagonal elements in weighted ridge regression. The formal results and the proof are provided in Sec S.9 of the supplementary materials.

PANDA and the noise augmentation technique, in general, can be extended to regularize other types of learning problems such as undirected graphical models, where some of the existing techniques are GLM-based. The realized $l_0$ penalty by noise augmentation can be used to regularize learning problems where such penalty is desired but hard to realize due to computational constraints. Regarding the ensemble learning behavior of PANDA, it is worthwhile to study further the underlying theory and run more empirical studies to quantify the benefits of the diversity term enabled by PANDA in generalization error reduction.

# References

Berk, R., Brown, L., Buja, A., Zhang, K., Zhao, L., et al. (2013). Valid post-selection inference. *The Annals of Statistics*, 41(2):802–837.

Brown, G., Wyatt, J. L., and Tino, P. (2005). Managing diversity in regression ensembles. *Journal of Machine Learning Research*, 6:1621–2650.

Dicker, L., Huang, B., and Lin, X. (2013). Variable selection and estimation with the seamless-$l_0$ penalty. *Statistica Sinica*, 23:929–962.

Efron, B. (2014). Estimation and accuracy after model selection. *Journal of the American Statistical Association*, 109(507):991–1007.

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

Frank, I. E. and Friedman, J. H. (1993). A statistical view of some chemometrics regression tools. *Technometrics*, 35:109–148.

Javanmard, A. and Montanari, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research*, 15:2869–2909.

Lee, J., Sun, D., Sun, Y., and Taylor, J. (2016). Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, 44(3):907–927.

Leeb, H. and Pötscher, B. M. (2005). Model selection and inference: Facts and fiction. *Econometric Theory*, 21(1):21–59.

Leeb, H., Pötscher, B. M., et al. (2006). Can one estimate the conditional distribution of post-model-selection estimators? *The Annals of Statistics*, 34(5):2554–2591.

Liu, Z. and Li, G. (2016). Efficient regularized regression with $l_0$ penalty for variable selection and network construction. *Computational and Mathematical Methods in Medicine*, 3456153.

Lockhart, R., Taylor, J., Tibshirani, R. J., and Tibshirani, R. (2014). A significance test for the lasso. *Annals of statistics*, 42(2):413.

Reid, S., Taylor, J., and Tibshirani, R. (2017). Post-selection point and interval estimation of signal sizes in gaussian samples. *The Canadian Journal of Statistics*, 45(2):128–148.

Rizk, Y. and Awad, M. (2018). Sports articles for objectivity analysis data set. https://archive.ics.uci.edu/ml/datasets/Sports+articles+for+objectivity+analysis.

Simon, N., Friedman, J., Hastie, T., , and Tibshirani, R. (2013). *SGL: Fit a GLM (or cox model) with a combination of lasso and group lasso regularization*. R package version 1.1.

Taylor, J. and Tibshirani, R. (2017). Post-selection inference for $l_1$-penalized likelihood models. *The Canadian Journal of Statistics*, 46(1):41–61.

Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B*, 58(1):267–288.

Tibshirani, R., Saunders, M., Rosset, S., Zhu, J., and Knight, K. (2005). Sparsity and smoothness via the fused lasso. *Journal of the Royal Statistical Society, Series B*, 67(1):91–108.

Tibshirani, R. J., Taylor, J., Lockhart, R., and Tibshirani, R. (2016). Exact post-selection inference for sequential regression procedures. *Journal of the American Statistical Association*, 111(154):600–620.

Van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202.

Yuan, M. and Lin, Y. (2014). Model selection and estimation in regression with grouped variables. *Journal of the Royal Statistical Society, Series B*, 1:685–693.

Zhang, C. and Zhang, S. S. (2013). Confidence intervals for low-dimensional parameters in high-dimensional linear models. *Journal of the Royal Statistical Society Statistical Methodology Series B*, 76(1):217–242.

Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101(1):1418–1429.

Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society, Series B*, 62(2):301–320.

# Supplementary Materials to

## *Adaptive Noisy Data Augmentation for Regularized Estimation and Inference of Generalized Linear Models*

Yinan Li, Fang Liu

Department of Applied and Computational Mathematics & Statistics

University of Notre Dame, Notre Dame, IN 46556, U.S.A.

## S.1 Proof of Proposition 1

We take the Taylor expansion of $l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$, which is the negative log-likelihood, around $e_i = 0$ for $i = 1, \ldots, n_e$, and evaluate its expectation over the distribution of $e_i$.

$$l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})) = l(\boldsymbol{\theta}|\mathbf{x}) + l_p(\boldsymbol{\theta}|\mathbf{e}) = l(\boldsymbol{\theta}|\mathbf{x})) + \sum_{i=1}^{n_e} l_i(\boldsymbol{\theta}|e_i)$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) - \sum_{i=1}^{n_e} \left( h(e_i) + e_i \left( \theta_0 + \sum_j \theta_j e_{ij} \right) - B \left( \theta_0 + \sum_j \theta_j e_{ij} \right) \right)$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + \sum_{i=1}^{n_e} l_i(\boldsymbol{\theta}|e_i)|_{\mathbf{e}_i=0} - \sum_{i=1}^{n_e} \left\{ e_i \sum_j (\theta_j e_{ij}) - \sum_j (\theta_j e_{ij}) B' \left( \theta_0 + \sum_j \theta_j e_{ij}|_{e_{ij}=0} \right) \right.$$

$$\left. - \sum_{d=2}^{\infty} (d!)^{-1} \sum_j (\theta_j e_{ij})^d B^{(d)} \left( \theta_0 + \sum_j \theta_j e_{ij}|_{e_{ij}=0} \right) \right\}$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + C + \sum_{i=1}^{n_e} \left[ (B'(\theta_0) - e_i) \sum_j (\theta_j e_{ij}) + \sum_{d=2}^{\infty} (d!)^{-1} B^{(d)}(\theta_0) \sum_j (\theta_j e_{ij})^d \right],$$

where $C = B(\theta_0) - \sum_{i=1}^{n_e} (h(e_i) + e_i \theta_0)$, a constant independent of $\boldsymbol{\theta}$. The expectation of $l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ over the distribution of $e_{ij} \sim N(0, \mathrm{V}(e_j))$ is

$$\mathrm{E}_\mathbf{e}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})) = l(\boldsymbol{\theta}|\mathbf{x}) + C + n_e \left( \tfrac{1}{2} B''(\theta_0) \sum_j \theta_j^2 \mathrm{V}(e_j) \right) + O \left( n_e \sum_j \left( \theta_j^4 \mathrm{E}(e_j^4) \right) \right)$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + C + n_e \left( \tfrac{1}{2} B''(\theta_0) \sum_j \theta_j^2 \mathrm{V}(e_j) \right) + O \left( n_e \sum_j \left( \theta_j^4 V^2(e_j) \right) \right)$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + n_e \left( C_1 \sum_j \theta_j^2 \mathrm{V}(e_j) \right) + C + O \left( n_e \sum_j \left( \theta_j^4 V^2(e_j) \right) \right), \text{ where } C_1 = 2^{-1} B''(\theta_0). \quad (1)$$

There are two ways to realize the expectation in Eqn (1) empirically. One way is to approximate $l_p(\boldsymbol{\theta}|\mathbf{e})$ by $\lim_{m \to \infty} m^{-1} \sum_{t=1}^m \sum_{i=1}^{n_e} l_i(\boldsymbol{\theta}|\mathbf{e}_i^{(t)})$. The other way, under the constraint $n_e \mathrm{V}(e_j) = O(1)$, is to let $n_e \to \infty$, in which case the second term in Eqn (1) becomes $n_e C_1 \sum_j \theta_j^2 \mathrm{V}(e_{ij}) = n_e C_1 \sum_j \left( \theta_j^2 \lim_{n_e \to \infty} n_e^{-1} \sum_{i=1}^{n_e} \mathbf{e}_{ij}^2 n_e \mathrm{V}(e_j) \right)$. Between the two approaches, letting $n_e \to \infty \cap [n_e \mathrm{V}(e_j) = O(1)]$ also leads to the big-$O$ term $O \left( \sum_j \left( \theta_j^4 n_e V^2(e_j) \right) \right) \to 0$ in Eqn (1); in other words, the second order Taylor approximation of $\mathrm{E}_\mathbf{e}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}))$ is arbitrarily close to $\mathrm{E}_\mathbf{e}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}))$.

If $\theta_0 = 0$, then $C_{1j} = B''(0)$ and Eqn (1) can be simplified to

$$l(\boldsymbol{\theta}|\mathbf{x}) + C_2 \sum_j \theta_j^2 (n_e \mathrm{V}(e_j)) + C + O \left( \sum_j \left( \theta_j^4 n_e V^2(e_j) \right) \right). \quad (2)$$

For linear regression, the expectation of $l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \sum_{i=1}^{n+n_e} \left( \tilde{y}_i - \sum_j \tilde{x}_{ij} \theta_j \right)^2$ over the distribution of noise $\mathbf{e}$ is

$$\mathrm{E}_\mathbf{e}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})) = \sum_{i=1}^n \left( x_{ij} - \sum_j x_{ij} \theta_j \right)^2 + \mathrm{E}_\mathbf{e} \left( \sum_{i=1}^{n_e} \left( e_i - \sum_j e_{ij} \theta_j \right)^2 \right) \quad (3)$$

$$= \sum_{i=1}^n \left( x_{ij} - \sum_j x_{ij} \theta_j \right)^2 = \sum_{i=1}^{n_e} \mathrm{E}_\mathbf{e} \left( \sum_j e_{ij} \theta_j \right)^2 = l(\boldsymbol{\theta}|\mathbf{x}) + n_e \sum_j \theta_j^2 \mathrm{V}(e_{ij}). \quad (4)$$

1

## S.2 Proof of Proposition 2

Define Optimization Problem 1 $\hat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \arg\min_{\boldsymbol{\theta}} l(\boldsymbol{\theta}|\mathbf{x}) + C_1 \sum_{i=1}^{n_e} \left(\mathbf{e}_i^T \boldsymbol{\theta}\right)^2$. Due to its convexity in $\boldsymbol{\theta}$, $\hat{\boldsymbol{\theta}}$ can be solved directly from $\nabla_{\boldsymbol{\theta}} f(\hat{\boldsymbol{\theta}}) = 0$.

Define a constrained optimization Problem 2

$$\hat{\boldsymbol{\theta}} = \arg\min_{\boldsymbol{\theta}} l(\boldsymbol{\theta}|\mathbf{x}),$$
$$\text{s.t. } \sum_{i=1}^{n_e}(\mathbf{e}_i^T \boldsymbol{\theta})^2 \le d, \tag{5}$$

the corresponding Lagrangian for which is $L(\boldsymbol{\theta}) = l(\boldsymbol{\theta}|\mathbf{x}) + \lambda_L \left( \sum_{i=1}^{n_e} \left(\mathbf{e}_i^T \boldsymbol{\theta}\right)^2\right) - d\right)$, and the KKT conditions are

$$\nabla L(\boldsymbol{\theta}^*) = 0$$
$$\lambda_L \ge 0$$
$$\lambda_L \left( \sum_{i=1}^{n_e} \left(\mathbf{e}_i^T \boldsymbol{\theta}^*\right)^2 - d \right) = 0. \tag{6}$$

For Problem 2 to have the same solution as Problem 1, that is, $\boldsymbol{\theta}^* = \hat{\boldsymbol{\theta}}$, we set $\lambda_L = C_1$ and $d = \sum_{i=1}^{n_e} \left(\mathbf{e}_i^T \hat{\boldsymbol{\theta}}\right)^2$. The constraint in Eqn (5) now becomes

$$\sum_{i=1}^{n_e}(\mathbf{e}_i^T \boldsymbol{\theta})^2 \le \sum_{i=1}^{n_e} \left(\mathbf{e}_i^T \hat{\boldsymbol{\theta}}\right)^2 \tag{7}$$

Given that $n_e$ noise data points are independent per the PANDA procedure, Eqn (7) can also be regarded as $n_e$ linear constraints on $\boldsymbol{\theta}$

$$\exists\, 0 < d_i < \left( \sum_{i=1}^{n_e}(\mathbf{e}_i^T \hat{\boldsymbol{\theta}})^2 \right)^{1/2} : |\mathbf{e}_i^T \boldsymbol{\theta}| \le d_i, i = 1, \dots, n_e \tag{8}$$

## S.3 Proof of Theorem 1

We prove Theorem 1 for linear regression ($Y_i$ is Gaussian), Poisson regression, exponential regression, negative binomial regression, and logistic regression (when $Y_i$ is Bernoulli), respectively. WLOG, we use the bridge-type noise $e_{ij} \sim N(0, \lambda|\boldsymbol{\theta}|^{-\gamma})$ to demonstrate the proofs, which can be easily extended to other types of noises. Prior to the proof of Theorem 1, we state a theoretical result in Claim 2, on which the subsequent proofs rely on.

**Claim 2.** If $l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ and $l_p(\boldsymbol{\theta}|\mathbf{x})$ are convex functions w.r.t. $\boldsymbol{\theta}$ and share the same parameter space $\boldsymbol{\theta}$, then

$$\left| \inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - \inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x}) \right| \le \sup_{\boldsymbol{\theta}} |l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x})|$$

**Proof**: Since both $\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ and $\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x})$ are convex optimization problems, each has a global optimum, denoted by $\hat{\boldsymbol{\theta}}$ and $\tilde{\boldsymbol{\theta}}$, respectively. Therefore, $\left| \inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - \inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x}) \right|$ $= \left| l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) \right|$. Consider the following two scenarios,

i). if $l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \ge l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x})$, then $l_p(\tilde{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \ge l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \ge l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x})$ and

$$\left| l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) \right| = l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) \le l_p(\tilde{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) = \left| l_p(\tilde{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) \right|$$

2

ii). if $l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) < l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x})$, then $l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) < l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) < l_p(\hat{\boldsymbol{\theta}}|\mathbf{x})$ and

$$\left| l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) \right| = l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) - l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \leq l_p(\hat{\boldsymbol{\theta}}|\mathbf{x}) - l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \left| l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\hat{\boldsymbol{\theta}}|\mathbf{x}) \right|.$$

All taken together, $\left| l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) \right| \leq \max\left( \left| l_p(\tilde{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})\mathbf{1} - l_p(\tilde{\boldsymbol{\theta}}|\mathbf{x}) \right|, \left| l_p(\hat{\boldsymbol{\theta}}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\hat{\boldsymbol{\theta}}|\mathbf{x}) \right| \right)$
$\leq \sup_{\boldsymbol{\theta}} |l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x})|.$

## S.3.1 Linear regression

In this case the regularization effects with $n_e \to \infty$ and $m \to \infty$ are the same. The loss function upon convergence is

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = \sum_{i=1}^{n}\left(y_i - \sum_j x_{ij}\theta_j\right)^2 + m^{-1}\sum_{t=1}^{m}\sum_{i=1}^{n_e}\left(e_i - \sum_j e_{ij}^{(t)}\theta_j\right)^2.$$

Since $e_{ij}^{(t)} = \sqrt{\lambda|\boldsymbol{\theta}|^{-\gamma}}z_{ij}^{(t)}$, where $z_{ij}^{(t)} \sim N(0,1)$. Therefore,

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = l(\boldsymbol{\theta}|\mathbf{x}) + m^{-1}\sum_{t=1}^{m}\sum_{i=1}^{n_e}\left(\sum_j \frac{\lambda\theta_j^2}{|\theta_j|^\gamma}z_{ij}^{(t)2} + 2\sum_{j<k}\frac{\lambda\theta_j\theta_k}{|\theta_j\theta_k|^{\frac{\gamma}{2}}}z_{ik}^{(t)}z_{ij}^{(t)}\right)$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + m^{-1}\sum_{t=1}^{m}\sum_j\left(\frac{\lambda\theta_j^2}{|\theta_j|^\gamma}\sum_{i=1}^{n_e}z_{ij}^{(t)2}\right) + 2m^{-1}\sum_{t=1}^{m}\sum_{j<k}\left(\frac{\lambda\theta_j\theta_k}{|\theta_j\theta_k|^{\frac{\gamma}{2}}}\sum_{i=1}^{n_e}z_{ij}^{(t)}z_{ik}^{(t)}\right).$$

Since $\sum_{i=1}^{n_e}z_{ij}^{(t)2} \sim \Gamma\left(\frac{n_e}{2}, 2\right)$ and $\Gamma\left(\frac{n_e}{2}, 2\right) \approx N(n_e, 2n_e) = n_e + (2n_e)^{1/2}N(0,1) = n_e + (2n_e)^{1/2}z_1$ as $n_e \to \infty$; $\sum_{i=1}^{n_e}z_{il}^{(t)}z_{iv}^{(t)} \sim \Gamma\left(\frac{n_e}{2}, 2\right) - \Gamma\left(\frac{n_e}{2}, 2\right) \approx N(0, 4n_e) = 2n_e^{1/2}N(0,1) = 2n_e^{1/2}z_2$ as $n_e \to \infty$, where $z_1 \sim N(0,1)$ and $z_2 \sim N(0,1)$. Therefore, the distribution of $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ can be approximated by

$$l(\boldsymbol{\theta}|\mathbf{x}) + \sum_j n_e\lambda|\theta_j|^{2-\gamma} \tag{9}$$

$$+ \sum_j\left(n_e\lambda|\theta_j|^{2-\gamma}2^{1/2}n_e^{-1/2}\left(\frac{1}{m}\sum_{t=1}^{m}z_1^{(t)}\right)\right) + \sum_{j<k}\left(\frac{\lambda n_e\theta_j\theta_k}{|\theta_j\theta_k|^{\frac{\gamma}{2}}}(2n_e^{-1/2})\left(\frac{1}{m}\sum_{t=1}^{m}z_2^{(t)}\right)\right)$$

$$= l_p(\boldsymbol{\theta}|\mathbf{x}) + (mn_e)^{-1/2}C_1N(0,1) \text{ where } C_1 = n_e\lambda\left(2\left|\left|\left(\frac{\boldsymbol{\theta}}{|\boldsymbol{\theta}|^{\frac{\gamma}{2}}}\right)\left(\frac{\boldsymbol{\theta}}{|\boldsymbol{\theta}|^{\frac{\gamma}{2}}}\right)^T\right|\right|_2^2\right)^{1/2}, \tag{10}$$

where $l_p(\boldsymbol{\theta}|\mathbf{x}) = l(\boldsymbol{\theta}|\mathbf{x}) + \sum_j n_e\lambda|\theta_j|^{2-\gamma} = \mathrm{E}_{\mathbf{e}}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}))$. Exactly the same Eqn (10) can be obtained by letting $m \to \infty$ rather than $n_e \to \infty$.

Per the strong law of large numbers (LLN), Eqn (10) suggests $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ converges almost surely to its mean for all $\boldsymbol{\theta} \in \boldsymbol{\theta}$ as $m \to \infty$ or $n_e \to \infty$ (with $n_e\lambda = O(1)$), assuming $|\theta_j|$ belongs to a compact parameter space and is bounded by $B$. Consequently, $\sup_{\boldsymbol{\theta}}\left|\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x})\right| \xrightarrow{\text{a.s.}}$
$0$ as $m \to \infty$ or $n_e \to \infty$. Per Claim 2, $\inf_{\boldsymbol{\theta}}\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{\text{a.s.}} \inf_{\boldsymbol{\theta}}l_p(\boldsymbol{\theta}|\mathbf{x})$, and $\arg\inf_{\boldsymbol{\theta}}\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{\text{a.s.}}$
$\arg\inf_{\boldsymbol{\theta}}l_p(\boldsymbol{\theta}|\mathbf{x})$ due to the convexity of the loss function.

## S.3.2 Poisson Regression

The averaged noise-augmented loss function over $m$ iterations upon convergence is

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = l(\boldsymbol{\theta}|\mathbf{x}) - \frac{1}{m}\sum_{t=1}^{m}\sum_{i=1}^{n_e}\left(e_i\left(\theta_0 + \sum_j e_{ij}^{(t)}\theta_j\right) - \log(e_i!) - \exp\left(\theta_0 + \sum_j e_{ij}^{(t)}\theta_j\right)\right) \tag{11}$$

3

$$=l(\boldsymbol{\theta}|\mathbf{x})-\frac{1}{m}\sum_{t=1}^{m}e_i\sum_{j}\sum_{i=1}^{n_e}\theta_j e_{ij}^{(t)}+\frac{1}{m}\sum_{t=1}^{m}\sum_{i=1}^{n_e}\exp\left(\theta_0+\sum_{j}\theta_j e_{ij}^{(t)}\right)+C$$

$$=l(\boldsymbol{\theta}|\mathbf{x})\boxed{-\frac{1}{m}\sum_{t=1}^{m}e_i\sum_{j}\left(\frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}}\sum_{i=1}^{n_e}z_{ij}^{(t)}\right)+\frac{1}{m}\sum_{t=1}^{m}\sum_{i=1}^{n_e}\exp\left(\theta_0+\sum_{j}\frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}}z_{ij}^{(t)}\right)}+C, \quad (12)$$

$$=l(\boldsymbol{\theta}|\mathbf{x})+P(\boldsymbol{\theta})+C,$$

where $P(\boldsymbol{\theta})$ refers to the boxed expression in Eqn (12), $z_{ij}^{(t)}\sim N(0,1)$, $e_i\equiv n^{-1}\sum_{i=1}^{n}y_i$ that is a constant, and $C$ is a constant not related to $\boldsymbol{\theta}$. The regularizer $P(\boldsymbol{\theta})$ is different for $n_e\to\infty$ vs $m\to\infty$. We thus consider each case separately.

*Case 1: $n_e\to\infty$, $n_e\lambda=O(1)$ and fixed $m$*
Assume $m=1$ WLOG, then $z_{ij}^{(t)}$ can be abbreviated as $z_{ij}$. $n_e\to\infty$ and $\lambda n_e=O(1)$ implies that $\lambda\to 0$, therefore, $\sum_j\frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}}z_{ij}\to 0$ in Eqn (12). Apply the second order Taylor expansion around $\sum_j\theta_j z_{ij}=0$ to Eqn (12), as $n_e\to\infty$,

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})\to l(\boldsymbol{\theta}|\mathbf{x})-e_i\sum_j\left(\frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}}\sum_{i=1}^{n_e}z_{ij}\right)+\exp(\theta_0)\sum_{i=1}^{n_e}\sum_j\frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}}z_{ij}$$

$$+\tfrac{1}{2}\exp(\theta_0)\sum_{i=1}^{n_e}\left(\sum_j\frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}}z_{ij}\right)^2+O\left(n_e^{-1}\right)C_1(\boldsymbol{\theta})N(1,1)+C \quad (13)$$

$$\approx l(\boldsymbol{\theta}|\mathbf{x})+\tfrac{1}{2}\exp(\theta_0)\sum_j\left(\frac{\lambda\theta_j^2}{|\theta_j|^{\gamma}}\sum_{i=1}^{n_e}z_{ij}^2\right)+\exp(\theta_0)\sum_{j<k}\left(\frac{\lambda\theta_j\theta_k}{|\theta_j\theta_k|^{\frac{\gamma}{2}}}\sum_{i=1}^{n_e}z_{ij}z_{ik}\right)$$

$$+O\left(n_e^{-1}\right)C_1(\boldsymbol{\theta})N(1,1)+C \quad (14)$$

$$\to l(\boldsymbol{\theta}|\mathbf{x})+\frac{\lambda n_e}{2}\exp(\theta_0)\sum_k|\theta_j|^{2-\gamma}+O\left(n_e^{-0.5}\right)C_2(\boldsymbol{\theta})N(0,1)+O\left(n_e^{-1}\right)C_1(\boldsymbol{\theta})N(1,1)+C. \quad (15)$$

For Poisson regression, $e_i\equiv n^{-1}\sum_{i=1}^{n}y_i$, the average of the observations in the outcome node (the log of which estimates $\theta_0$) with the canonical log link function. In other words, when $n_e\to\infty$ $e_i=\exp(\theta_0)$; therefore, the second and third terms in Eqn (13) cancel out. $C_1(\boldsymbol{\theta})$ and $C_2(\boldsymbol{\theta})$ are functions of $\boldsymbol{\theta}$ and the standard deviations associated with the two asymptotic normality terms in Eqn (15) which result from the summation over $n_e$ noise terms per the CLT, and the $C_2(\boldsymbol{\theta})$ term is the rate-limiting term and

$$C_2(\boldsymbol{\theta})=\frac{\lambda n_e}{2}\left(2\exp(2\theta_0)\left\|\left(|\boldsymbol{\theta}|^{1-\frac{\gamma}{2}}\right)\left(|\boldsymbol{\theta}|^{1-\frac{\gamma}{2}}\right)^T\right\|_2^2\right)^{1/2} \text{ where } \lambda n_e=O(1). \quad (16)$$

Note that $l(\boldsymbol{\theta}|\mathbf{x})+\frac{\lambda n_e}{2}\exp(\theta_0)\sum_k|\theta_j|^{2-\gamma}$ in Eqn (15) is $l_p(\boldsymbol{\theta}|\mathbf{x})=\mathrm{E}_{\mathbf{e}}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})$ per Proposition 1 and Appendix S.1. As $n_e\to\infty$ and $\lambda n_e=O(1)$, per the strong LLN and Eqn (15), $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})$ converges almost surely to $l_p(\boldsymbol{\theta}|\mathbf{x})$. Given the convexity of the loss function and per Claim 2, $\arg\inf_{\boldsymbol{\theta}}\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})\xrightarrow{a.s.}\arg\inf_{\boldsymbol{\theta}}l_p(\boldsymbol{\theta}|\mathbf{x})$.

*Case 2: $m\to\infty$ and fixed $n_e$*
The 2nd term in Eqn (12) is the summation of Gaussian variables, and the 3rd term follows a log-normal distribution. Therefore, we can rewrite Eqn (12) as

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})=l(\boldsymbol{\theta}|\mathbf{x})-e_i\sum_j\frac{\sqrt{\lambda}n_e\theta_j}{\sqrt{m}|\theta_j|^{\frac{\gamma}{2}}}N(0,1)+m^{-1}\sum_{t=1}^{m}\sum_{i=1}^{n_e}\mathrm{LogN}\left(\theta_0,\sum_j\frac{\lambda\theta_j^2}{|\theta_j|^{\gamma}}\right)+C. \quad (17)$$

Applying the CLT to Eqn (17) as $m \to \infty$,

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \to l(\boldsymbol{\theta}|\mathbf{x}) - e_i \sum_j \frac{\sqrt{\lambda} n_e \theta_j}{\sqrt{m}|\theta_j|^{\frac{\gamma}{2}}} N(0,1) \tag{18}$$

$$+ \left\{ \frac{n_e}{m} \left( \exp\left( \sum_j \frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} \right)^2 - 1 \right) \exp\left( 2\theta_0 + \left( \sum_j \frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} \right)^2 \right) \right\}^{1/2} N(0,1) + C,$$

suggesting that $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ follows a Gaussian distribution asymptotically. Per the strong LLN as $m \to \infty$, Eqn (18) converges almost surely to

$$\mathrm{E}_{\mathbf{e}}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})) = l_p(\boldsymbol{\theta}|\mathbf{x}) = l(\boldsymbol{\theta}|\mathbf{x}) + P(\boldsymbol{\theta}) + C$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + n_e \exp(\theta_0) \exp\left( \frac{1}{2}\lambda \left( \sum_k |\theta_j|^{1-\frac{\gamma}{2}} \right)^2 \right) + C \tag{19}$$

for all $\boldsymbol{\theta} \in \boldsymbol{\Theta}$ assuming $\boldsymbol{\Theta}$ to be compact. Per claim 2, $\sup_{\boldsymbol{\theta}} \left| \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x}) \right| \xrightarrow{\text{a.s.}} 0$ as $m \to \infty$, which leads to $\inf_{\boldsymbol{\theta}} \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{\text{a.s.}} \inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x}) \Rightarrow \arg\inf_{\boldsymbol{\theta}} \bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{\text{a.s.}} \arg\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x})$ given the convexity of the loss function.

### S.3.3   Exponential Regression

The averaged noise-augmented loss function over $m$ iterations upon convergence is

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = l(\boldsymbol{\theta}|\mathbf{x}) - m^{-1} \sum_{t=1}^{m} \sum_{i=1}^{n_e} \left( \theta_0 + \sum_j e_{ij}^{(t)} \theta_j - e_i \exp\left( \theta_0 + \sum_j e_{ij}^{(t)} \theta_j \right) \right),$$

where $e_i = n^{-1} \sum_{i=1}^{n} y_i$. The above loss function is equivalent to the loss function in Eqn (11) in the PGM case except for the constant term that does not involve $\boldsymbol{\theta}$. Therefore, the proof for PGM also applies in the case of EGM.

### S.3.4   Negative Binomial Regression

The averaged noise-augmented loss function over $m$ iterations upon convergence is

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = l(\boldsymbol{\theta}|\mathbf{x}) - \frac{1}{m} \sum_{t=1}^{m} \sum_{i=1}^{n_e} \left( \log\left( \frac{\Gamma(e_i+r)r^r}{\Gamma(e_i+1)\Gamma(r)} \right) + e_i \sum_j e_{ij}^{(t)} \theta_j \right.$$

$$\left. - (r+e_i) \log\left( r + \exp\left( \theta_0 + \sum_j e_{ij}^{(t)} \theta_j \right) \right) \right) \tag{20}$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + C - \frac{1}{m} \sum_{t=1}^{m} \sum_{i=1}^{n_e} e_i \sum_j e_{ij}^{(t)} \theta_j + \frac{1}{m} \sum_{t=1}^{m} \sum_{i=1}^{n_e} (r+e_i) \log\left( r + \exp\left( \theta_0 + \sum_j e_{ij}^{(t)} \theta_j \right) \right) \tag{21}$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + C \boxed{- \frac{1}{m} \sum_{t=1}^{m} e_i \sum_j \left( \frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} \sum_{i=1}^{n_e} z_{ij}^{(t)} \right) + \frac{1}{m} \sum_{t=1}^{m} \sum_{i=1}^{n_e} (r+1) \log\left( r \exp\left( \theta_0 + \sum_j \frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} z_{ij}^{(t)} \right) \right)} \tag{22}$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + P(\boldsymbol{\theta}) + C = l_p(\boldsymbol{\theta}|\mathbf{x}) + C,$$

where $P(\boldsymbol{\theta})$ refers to the boxed expression in Eqn (22), $z_{ij}^{(t)} \sim N(0,1)$, $e_i \equiv n^{-1} \sum_{i=1}^{n} y_i$ is a constant, and $C$ is a constant not related to $\boldsymbol{\theta}$. The regularizer $P(\boldsymbol{\theta})$ is different for $n_e \to \infty$ vs $m \to \infty$. We thus consider each case separately.

*Case 1: $n_e \to \infty$ and $n_e\lambda = O(1)$ and fixed $m$*

Let $m = 1$ WLOG, thus $z_{ij}^{(t)}$ can be abbreviated as $z_{ij}$. Since $n_e \to \infty$ and $\lambda n_e = O(1)$, implying $\lambda \to 0$ and thus $\exp\left(\sum_j \frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} z_{ij}\right) \to 1$. Applying the second order Taylor expansion around $\sum_j \theta_j z_{ij} = 0$ to Eqn (12), we have

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = l(\boldsymbol{\theta}|\mathbf{x}) - e_i \sum_j \left(\frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} \sum_{i=1}^{n_e} z_{ij}\right) + \frac{(r+e_i)\exp(\theta_0)}{r+\exp(\theta_0)} \sum_{i=1}^{n_e} \sum_j \frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} z_{ij}$$

$$+ \frac{1}{2}\sum_{i=1}^{n_e} \frac{(r+e_i)r\exp(\theta_0)}{(r+\exp(\theta_0))^2}\left(\sum_j \frac{\sqrt{\lambda}\theta_j}{|\theta_j|^{\frac{\gamma}{2}}} z_{ij}\right)^2 + O\left(n_e^{-1}\right) N(1, C_1(\boldsymbol{\theta})) + C \qquad (23)$$

$$\to l(\boldsymbol{\theta}|\mathbf{x}) + \frac{1}{2}\sum_j \frac{r\exp(\theta_0)}{r+\exp(\theta_0)}\left(\frac{\lambda\theta_j^2}{|\theta_j|^\gamma} \sum_{i=1}^{n_e} z_{ij}^2\right) + \sum_{k<l} \frac{r\exp(\theta_0)}{r+\exp(\theta_0)}\left(\frac{\lambda\theta_j\theta_k}{|\theta_j\theta_k|^{\frac{\gamma}{2}}} \sum_{i=1}^{n_e} z_{ij}e_{0il}\right)$$

$$+ O\left(n_e^{-1}\right) N(1, C_1(\boldsymbol{\theta})) + C \qquad (24)$$

$$\to l(\boldsymbol{\theta}|\mathbf{x}) + \frac{1}{2}\sum_j \frac{r\exp(\theta_0)}{r+\exp(\theta_0)}\left(\frac{\lambda\theta_j^2}{|\theta_j|^\gamma} \sum_{i=1}^{n_e} z_{ij}^2\right) + O(n_e^{-1}) N(1, C_1(\boldsymbol{\theta})) + O(n_e^{-0.5}) C_2(\boldsymbol{\theta}) N(0,1) + C \quad (25)$$

In NB regression, the logarithm of the average of the observations in the outcome node estimates $\theta_0$ with the canonical log link function. In other words, when $n_e \to \infty$ $e_i = \exp(\theta_0)$, and $r+\exp(\theta_0) = r+e_i$; therefore, the second and third terms in Eqn (23) cancel out and the forth term can be simplied as shown above. $C_1(\boldsymbol{\theta})$ and $C_2(\boldsymbol{\theta})$ are functions of $\boldsymbol{\theta}$ and the standard deviations associated with the two asymptotic normality terms in Eqn (25) that result from the summation over $n_e$ noise terms per the CLT, and the $C_2(\boldsymbol{\theta})$ term is the rate-limiting term and

$$C_2(\boldsymbol{\theta}) = \frac{\lambda n_e}{2}\left(2\left(\frac{r\exp(\theta_0)}{r+\exp(\theta_0)}\right)^2 \left|\left|\left(|\boldsymbol{\theta}|^{1-\frac{\gamma}{2}}\right)\left(|\boldsymbol{\theta}|^{1-\frac{\gamma}{2}}\right)^T\right|\right|_2^2\right)^{1/2}. \qquad (26)$$

Note that $l(\boldsymbol{\theta}|\mathbf{x}) + \frac{1}{2}\sum_j \frac{r\exp(\theta_0)}{r+\exp(\theta_0)}\left(\frac{\lambda\theta_j^2}{|\theta_j|^\gamma} \sum_{i=1}^{n_e} z_{ij}^2\right)$ in Eqn (25) is $l_p(\boldsymbol{\theta}|\mathbf{x}) = \mathrm{E}_{\mathbf{e}}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ per Proposition 1 and Appendix S.1. As $n_e \to \infty$ and $\lambda n_e = O(1)$, per the strong LLN and Eqn (25), $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ converges almost surely to $l_p(\boldsymbol{\theta}|\mathbf{x})$. Given the convexity of the loss function and Claim 2,

$$\arg\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{a.s.} \arg\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x}).$$

*Case 2: $m \to \infty$ and fixed $n_e$*
The second term in Eqn (21) is the summation over Gaussian variables, therefore, the equation can be written as

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = l(\boldsymbol{\theta}|\mathbf{x}) - e_i \sum_j \frac{\sqrt{\lambda}n_e\theta_j}{\sqrt{m}|\theta_j^{(t-1)}|^{\frac{\gamma}{2}}} N(0,1) + \frac{1}{m}\sum_{t=1}^m \sum_{i=1}^{n_e} U_i^{(t)} + C,$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) - e_i \sum_j \frac{\sqrt{\lambda}n_e\theta_j}{\sqrt{m}|\theta_j^{(t-1)}|^{\frac{\gamma}{2}}} N(0,1) + \frac{n_e}{m}\sum_{t=1}^m U^{(t)} + C, \qquad (27)$$

where $U_i^{(t)} = (r+e_i)\log\left(r+\exp\left(\sum_j e_{ij}^{(t)}\theta_j\right)\right)$. The second equation holds because $U_i^{(t)}$ is the same for all $i = 1, \ldots, n_e$. Applying the CLT to the $U$-term in Eqn (27) as $m \to \infty$,

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \to l(\boldsymbol{\theta}|\mathbf{x}) - e_i \sum_j \frac{\sqrt{\lambda}n_e\theta_j}{\sqrt{m}|\theta_j|^{\frac{\gamma}{2}}} N(0,1) + n_e\mathrm{E}\left(U^{(t)}\right) + \frac{n_e}{\sqrt{m}} N(0, \sigma_U)$$

$$= l(\boldsymbol{\theta}|\mathbf{x}) + n_e E(U^{(t)}) - e_i \sum_j \frac{\sqrt{\lambda}n_e\theta_j}{\sqrt{m}|\theta_j|^{\frac{\gamma}{2}}} N(0,1) + \frac{n_e}{\sigma_U\sqrt{m}} N(0,1), \qquad (28)$$

where $\sigma_U$ is the standard deviation of $U^{(t)}$. Since $\log(r+\exp(*))\to\max\{\log(r), *\}$, as $*\to\pm\infty$, $\sigma_U$ is a finite. Eqn (28) suggests that $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ follows a Gaussian distribution as $m\to\infty$.

Additionally, applying the strong LLN to Eqn (21), $\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$ converges almost surely to its mean $l_p(\boldsymbol{\theta}|\mathbf{x}) = \mathrm{E}(l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}))$ for all $\boldsymbol{\theta}\in\boldsymbol{\theta}$ as $m\to\infty$, assuming $\boldsymbol{\theta}$ to be compact; that is,

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \to l_p(\boldsymbol{\theta}|\mathbf{x}) + C = l(\boldsymbol{\theta}|\mathbf{x}) + n_e\mathrm{E}(U_i^{(t)}) + C. \tag{29}$$

It follows that $\sup_{\boldsymbol{\theta}}\left|\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p(\boldsymbol{\theta}|\mathbf{x})\right| \xrightarrow{\text{a.s.}} 0$ as $m\to\infty \Rightarrow \inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{\text{a.s.}} \inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x}) \Rightarrow$ $\arg\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) \xrightarrow{\text{a.s.}} \arg\inf_{\boldsymbol{\theta}} l_p(\boldsymbol{\theta}|\mathbf{x})$ given the convexity of the loss function.

### S.3.5 Binomial Regression

The averaged noise-augmented loss function over $m$ iterations upon convergence is

$$\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) = l(\boldsymbol{\theta}|\mathbf{x}) - \frac{1}{m}\sum_{t=1}^{m}\sum_{i=1}^{n_e}\left(e_i\sum_j e_{ij}^{(t)}\theta_j - \log\left(1+\exp\left(\theta_0 + \sum_j e_{ij}^{(t)}\theta_j\right)\right)\right),$$

which is a special case of Eqn (20) when $r=1$, and the proof for NBGM also applies to BGM.

## S.4 Proof of Proposition 3

In the case of multicollinearity, PANDA with sparsity regularization might experience difficulty in learning minimizer $\hat{\boldsymbol{\theta}}_p^{(n_e)}$ (or $\hat{\boldsymbol{\theta}}_p^{(m)}$) when $n_e$( or $m$) $\to\infty$. In such a case, we prove that there exists $\epsilon > 0$ and a sub-sequence $[n_e]_i$ (or $[m]_i$), such that letting $\boldsymbol{\theta}_p^i \triangleq \hat{\boldsymbol{\theta}}_p^{[n_e]_i}$ (or $\hat{\boldsymbol{\theta}}_p^{[m]_i}$), then $d(\boldsymbol{\theta}_p^i, \boldsymbol{\Theta}^0) > \epsilon$, where $\boldsymbol{\Theta}^0$ is the optimum parameter set. Denote $l_p^i = l_p(\boldsymbol{\theta}_p^i|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})$, then by Eqn (24), there exists a sub-sequence $[i]$, such that,

$$\Pr\left(\sup_{\boldsymbol{\theta}}\left|\bar{l}_p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - \bar{l}_p(\boldsymbol{\theta}|\mathbf{x})\right| > \delta\right) < k^{-1}, k\in N. \tag{30}$$

Since $\boldsymbol{\theta}$ is compact, the sub-sequence $[i]$ converges to a point $\hat{\boldsymbol{\theta}}^* \in \boldsymbol{\Theta}$, $d\left(\hat{\boldsymbol{\theta}}^*, \boldsymbol{\Theta}^0\right) \geq \epsilon$, $\hat{\boldsymbol{\theta}}^* \notin \boldsymbol{\Theta}^0$. On the other hand, for any $\boldsymbol{\theta}\in\boldsymbol{\Theta}$, we have

$$\begin{aligned}\bar{l}_p(\hat{\boldsymbol{\theta}}^*|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - \bar{l}_p(\boldsymbol{\theta}|\mathbf{x}) =& (\bar{l}_p(\hat{\boldsymbol{\theta}}^*|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - \bar{l}_p(\boldsymbol{\theta}_p^{[i]}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}})) + (\bar{l}_p(\boldsymbol{\theta}_p^{[i]}|\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) - l_p^{[i]}(\boldsymbol{\theta}_p^{[i]})) \\ & + (l_p^{[i]}(\boldsymbol{\theta}_p^{[i]}) - l_p^{[i]}(\boldsymbol{\theta})) + (l_p^{[i]}(\boldsymbol{\theta}) - \bar{l}_p(\boldsymbol{\theta}|\mathbf{x})).\end{aligned}$$

By the continuity of the loss function and $\lim_{i\to\infty}\boldsymbol{\theta}_p^{[i]} = \hat{\boldsymbol{\theta}}^*$, the first term in the above equation is arbitrarily small with $i\to\infty$; by equation (30), the second and forth terms are arbitrarily small with $i\to\infty$, and the third term is non-positive. Since $\boldsymbol{\theta}\in\boldsymbol{\Theta}$ is arbitrary , we must have $\hat{\boldsymbol{\theta}}^* \in \boldsymbol{\Theta}^0$, which is a contradiction. The Proposition is proved.

## S.5 Proof of Proposition 4

WLOG, we derive the Fisher information with the bridge-type noise. The proofs for other types of noise are similar. The Fisher information matrix $I_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(\boldsymbol{\theta})$ on the augmented data is obtained by taking the expectation of the negative second derivative of the noise-augmented loss function in Eqn (3) over the distribution of data $\mathbf{x}$ and augmented noise $\mathbf{e}$.

$$I_{\tilde{\mathbf{x}}, \tilde{\mathbf{y}}}(\boldsymbol{\theta}) = \mathrm{E}_{\mathbf{x}}\left(\mathbf{x}^T\boldsymbol{B}''(\mathbf{x})\mathbf{x}\right) + \mathrm{E}_{\mathbf{e}}\left(\mathbf{e}_{\mathbf{x}}^T\boldsymbol{B}''(\mathbf{e}_{\mathbf{x}})\mathbf{e}_{\mathbf{x}}\right) = I_{\mathbf{x}, \mathbf{y}}(\boldsymbol{\theta}) + \mathrm{E}_{\mathbf{e}}\left(\sum_{i=1}^{n_e}\mathbf{e}_{\mathbf{x},i}^T B''(\mathbf{e}_i\boldsymbol{\theta})\mathbf{e}_{\mathbf{x},i}\right),$$

where $\boldsymbol{B}''(\mathbf{x}) = \mathrm{diag}\{B''(\mathbf{x}_1\boldsymbol{\theta}),\ldots,B''(\mathbf{x}_n\boldsymbol{\theta})\}$ and $\boldsymbol{B}(\mathbf{e_x}) = \mathrm{diag}\{B''(\mathbf{e}_{\mathbf{x},1}\boldsymbol{\theta}),\ldots,B''(\mathbf{e}_{\mathbf{x},n_e}\boldsymbol{\theta})\}$. Let $\lambda n_e = O(1)$ and $\mathrm{V}(\mathbf{e}_{\mathbf{x},i})$ denote the covariance matrix of $\mathbf{e}_{\mathbf{x},i}$; take the second-order Taylor expansion around $\mathbf{e}_{\mathbf{x},i}\boldsymbol{\theta} = 0$, we have

$$I_{\tilde{\mathbf{x}},\tilde{\mathbf{y}}}(\boldsymbol{\theta}) = I_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta}) + n_e B''(0)\mathrm{V}(\mathbf{e}_{\mathbf{x},i}) + O(\lambda n_e^{1/2})J_p$$

$$= I_{\mathbf{x},\mathbf{y}}(\boldsymbol{\theta}) + (\lambda n_e)B''(0)\mathrm{diag}\{|\boldsymbol{\theta}_{j1}|^{-\gamma},\ldots,|\boldsymbol{\theta}_{jp}|^{-\gamma}\} + O(\lambda n_e^{1/2})J_p,$$

where $J_p$ is a $p \times p$ matrix with all elements equal to 1.

# S.6    Proof of Proposition 5

Given $n^{-1/2}l'(\boldsymbol{\theta}|\mathbf{x}) \xrightarrow{d} N(0, I_1^{-1}(\boldsymbol{\theta}))$, where $l'(\boldsymbol{\theta}|\mathbf{x})$ is the first derivative of the negative log-likelihood function given the observed data $\mathbf{x}$ and $I_1(\boldsymbol{\theta})$ is the information matrix over one observation. It follows that

$$n^{-1/2}(l'(\boldsymbol{\theta}|\mathbf{x}) + l'(\boldsymbol{\theta}|\mathbf{e})) = n^{-1/2}l'(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}}) \xrightarrow{d} N(n^{-1/2}l'(\boldsymbol{\theta}|\mathbf{e}), I_1(\boldsymbol{\theta})) \tag{31}$$

where $\mathbf{e}$ is the augmented noise and $l'(\boldsymbol{\theta}|\mathbf{e}) = \sum_{i=1}^{n_e} l'(\boldsymbol{\theta}|\mathbf{e}_i)$. Let $\phi(\mathbf{e}) = n^{-1/2}l'(\boldsymbol{\theta}|\mathbf{e})$ and it expectation over the distribution of $\mathbf{e}$ can be worked out for different types of noise. For example, with the bridge-type noise, $\phi(\mathbf{e}) = n^{-1/2}l'(\boldsymbol{\theta}|\mathbf{e})$ and $\mathrm{E}_{\mathbf{e}}(\phi) = \frac{\lambda n_e}{\sqrt{n}}\sigma^2\mathrm{sgn}(\theta_0)$ for Gaussian outcome nodes, $\frac{\lambda n_e}{8\sqrt{n}}\mathrm{sgn}(\theta_0) + \frac{\lambda^2 n_e}{\sqrt{n}}O(|\theta_0|)$ for Bernoulli outcome nodes, $\frac{\lambda n_e}{2\sqrt{n}}\mathrm{sgn}(\theta_0) + \frac{\lambda^2 n_e}{n}O(|\theta_0|)$ for exponential and Poisson outcome nodes and $\frac{\lambda n_e r}{2(r+1)n}\mathrm{sgn}(\theta_0) + \frac{\lambda^2 n_e}{\sqrt{n}}O(|\theta_0|)$ for NB outcome nodes. If $\lambda n_e = o(\sqrt{n})$, then $\mathrm{E}_{\mathbf{e}}(\phi) \to 0$ as $n \to \infty$.

Upon the convergence of the PANDA algorithm, the MLE of $\boldsymbol{\theta}$ based on $(\tilde{\mathbf{x}},\tilde{\mathbf{y}})$ is the minimizer $\hat{\boldsymbol{\theta}}_{j,\mathbf{e}}$ from solving $l'(\hat{\boldsymbol{\theta}}_{j,\mathbf{e}}) = 0$, its first-order Taylor expansion around $\boldsymbol{\theta}$ is $l'(\hat{\boldsymbol{\theta}}_{\mathbf{e}}) \approx l'(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}}) + l''(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})(\hat{\boldsymbol{\theta}}_{\mathbf{e}} - \boldsymbol{\theta}) = 0$. Therefore, $\hat{\boldsymbol{\theta}}_{\mathbf{e}} - \boldsymbol{\theta} = -(l''(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}}))^{-1}l'(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})$ and $\sqrt{n}\left(\hat{\boldsymbol{\theta}}_{\mathbf{e}} - \boldsymbol{\theta}\right) = -(n^{-1}l''(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}}))^{-1}\left(n^{-1/2}l'(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})\right)$, where $l''(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})$ is the Hessian matrix and $l''(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}}) \to I_p(\boldsymbol{\theta})$ as $n \to \infty$. Taken together with Eqn (31), assume $\lambda n_e = o(\sqrt{n})$, by the Slutsky's theorem, as $n \to \infty$

$$\sqrt{n}\left(\hat{\boldsymbol{\theta}}_{\mathbf{e}} - \boldsymbol{\theta}\right) = (n^{-1}l''(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}}))^{-1}\left(n^{-1/2}l'(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})\right) \xrightarrow{d} N\left(\mathbf{0}, I_p(\boldsymbol{\theta})^{-1}I(\boldsymbol{\theta})I_p(\boldsymbol{\theta})^{-1}\right) \triangleq N(\mathbf{0}, \Sigma_{\mathbf{e}}).$$

When the mean of $m > 1$ estimates over consecutive iteration are taken as the final estimate for $\boldsymbol{\theta}$, that is $\bar{\boldsymbol{\theta}} = m^{-1}\sum_{t=1}^m \hat{\boldsymbol{\theta}}_{\mathbf{e}}^{(t)}$, the variability among the $m$ consecutive estimates will need to be accounted for and be reflected in the variance of the final estimate. It is easy to establish this in the Bayesian framework. Specifically,

$$\mathrm{E}(\boldsymbol{\theta}|\mathbf{x}) = \mathrm{E}_{\mathbf{e}}(\mathrm{E}(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})) = \mathrm{E}_{\mathbf{e}}(\hat{\boldsymbol{\theta}}_{\mathbf{e}}) = m^{-1}\sum_{t=1}^m \hat{\boldsymbol{\theta}}_{\mathbf{e}}^{(t)} \triangleq \bar{\boldsymbol{\theta}} \text{ as } m \to \infty$$

$$V(\boldsymbol{\theta}|\mathbf{x}) = \mathrm{E}_{\mathbf{e}}(V(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})) + V_{\mathbf{e}}(\mathrm{E}(\boldsymbol{\theta}|\tilde{\mathbf{x}},\tilde{\mathbf{y}})) = \mathrm{E}_{\mathbf{e}}(\Sigma_{\mathbf{e}}) + V_{\mathbf{e}}(\hat{\boldsymbol{\theta}}_{\mathbf{e}}) \triangleq \bar{\Sigma} + \Lambda$$

$$= m^{-1}\sum_{t=1}^m \Sigma_{\mathbf{e}}^{(t)} + (m-1)^{-1}\sum_{t=1}^m \left(\hat{\boldsymbol{\theta}}^{(t)} - \bar{\boldsymbol{\theta}}\right)\left(\hat{\boldsymbol{\theta}}_{\mathbf{e}}^{(t)} - \bar{\boldsymbol{\theta}}\right)' \text{ as } m \to \infty.$$

Per the large-sample Bayesian theory, the posterior mean and variance of $\boldsymbol{\theta}$ given $\mathbf{x}$ are asymptotically equivalent ($n \to \infty$) to the MLE for $\boldsymbol{\theta}$ and the inverse information matrix of $\boldsymbol{\theta}$ contained in $\mathbf{x}$. In other words,

$$\sqrt{n}(\bar{\boldsymbol{\theta}} - \boldsymbol{\theta}) \to N\left(\mathbf{0}, \bar{\Sigma} + \Lambda\right).$$

In the case of a finite $m$ (as in practical application), $V(\boldsymbol{\theta}|\mathbf{x})$ is estimated by $\bar{\Sigma} + (1 + m^{-1})\Lambda$

with the correction for the finite $m$. Applying Proposition 5 with lasso-type noise, we have
$$\sqrt{n}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \to N\left(n^{-1/2}\lambda n_e \text{sgn}(\boldsymbol{\theta})M^{-1}, \sigma^2 M^{-1}(\mathbf{y}'\mathbf{x})M^{-1}\right),$$
where $M = (\mathbf{y}'\mathbf{x} + \text{diag}(\lambda n_e|\boldsymbol{\theta}|^{-1}))$ and $\sigma^2$ is the variance of the error term in the linear regression, and is estimated by
$$\begin{aligned}\hat{\sigma}^2 &= \text{SSE}(n - \nu)^{-1} = (n - \nu)^{-1}(\mathbf{x}\boldsymbol{\theta} + \epsilon)'(I - H)(\mathbf{x}\boldsymbol{\theta} + \epsilon)\\ &= (n - \nu)^{-1}\epsilon'(I - H)\epsilon + (n - \nu)^{-1}\left(\boldsymbol{\theta}'\mathbf{x}'(I - H_{\mathbf{x}}\boldsymbol{\theta} + 2\boldsymbol{\theta}'\mathbf{x}'(I - H)\epsilon\right)\end{aligned}$$
where $H = \mathbf{x}(\mathbf{y}'\mathbf{x} + \text{diag}(\lambda n_e|\boldsymbol{\theta}|^{-1}))^{-1}\mathbf{y}'$ and $\nu = \text{trace}(H)$.

## S.7 A Formal Test on the Convergence of the PANDA Algorithm

When presenting the PANDA algorithm in Sec 2.2, we state that a formal statistical test can be used to test convergence. This test is based on the assumption of $n_e \to \infty$ or $m \to \infty$ and should work well when either $n_e$ or $m$ is large in practice. WLOG, we establish the test below for $n_e \to \infty$; the procedure is similar for $m \to \infty$ by replacing $n_e$ with $m$.

Theorem 1 shows that as $n_e \to \infty$, the distribution of the loss function in iteration $t$ converges to a Gaussian distribution (Eqn (18)). The asymptotic Gaussian distribution involves $C_1(\boldsymbol{\theta})$, which is unknown and can be estimated by plugging the $\hat{\boldsymbol{\theta}}^{(t)}$ from the current iteration $t$. Specifically,
$$C_1^{(t)} = \tfrac{\lambda n_e}{2}\left(\kappa\left|\left|\left(\hat{\boldsymbol{\theta}}^{(t)}|\hat{\boldsymbol{\theta}}^{(t)}|^{-\gamma/2}\right)\left(\hat{\boldsymbol{\theta}}^{(t)}|\hat{\boldsymbol{\theta}}^{(t)}|^{-\gamma/2}\right)^T\right|\right|_2^2\right)^{1/2},$$
where $\kappa$ is a constant that depends on the type of $Y$ ($\kappa = 8$ for Gaussian, $2\exp(2\theta_0)/(1 + \exp(2\theta_0))^4$ for Bernoulli, $2\exp(2\theta_0)$ for Poisson, $2$ for Exponential, and $2r^2\exp(2\theta_0)/(r + \exp(\theta_0))^2$ for NB; see Eqns (10), (16) and (26)). Let $d^{(t)} = \bar{l}_p(\tilde{\mathbf{x}}^{(t+1)}, \tilde{\mathbf{y}}) - \bar{l}_p(\tilde{\mathbf{x}}^{(t)}, \tilde{\mathbf{y}})$ denote the difference in the loss function from two consecutive iterations of the PANDA algorithm, which is $n_e^{-1/2}(C_1^{(t+1)}z^{(t+1)} - C_1^{(t)}z^{(t)})$ per Eqn (18). If the PANDA algorithm converges, the estimates $\hat{\boldsymbol{\theta}}^{(t)}$ stabilizes, so does $C_1^{(t)}$; in other words, $C_1^{(t+1)} \approx C_1^{(t)}$ and a nonzero $d^{(t)}$ is mostly due to the randomness of the injected noise with an expected mean of 0; that is,
$$z^{(t)} = d^{(t)}/\sqrt{n_e^{-1}\left[C_1^{(t)2} + C_1^{(t+1)2}\right]}. \tag{32}$$

Since $z^{(t)}$ is independent from $z^{(t+1)}$ (augmented noises are drawn independently across iteration). If $|z^{(t)}| > z_{1-\alpha/2}$, then we may claim the PANDA algorithm has not converged at iteration $t$ at the significance level of $\alpha$.

The denominator in Eqn (32) assumes $C_1^{(t)}$ and $C_1^{(t+1)}$ are independent whey are likely to positively correlated as both use the original data $(\mathbf{x}, \mathbf{y})$. With the under-estimated variance, $z^{(t)}$ would be over-estimated, and convergence is likely to rejected more often than necessary.

## S.8 Minimizer of Averaged noise-augmented Loss Function vs Averaged minimizer of Noise-augmented Loss Functions

Per Proposition 1, one would take the average over $m$ noise-augmented loss function $l(\Theta|\mathbf{x}, \mathbf{e})$ to yield a single minimizer $\hat{\boldsymbol{\theta}}$, which is the Monte Carlo version of $E_{\mathbf{e}}(l_p(\boldsymbol{\theta}|\mathbf{x}, \mathbf{e})$ as $m \to \infty$. However, PANDA would lose its computational edge. To maintain the computational

advantage for PANDA, we instead calculate $\bar{\boldsymbol{\theta}}$, the average of $m$ minimizers of $l(\Theta|\mathbf{x}, \mathbf{e})$ from the latest $m$ iterations, which is the approach that the PANDA algorithm uses. We establish in Corollary S.1 that $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$ are equivalent under some regularity conditions. We also present some numerical examples below to illustrate the similarity between $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$.

**Corollary S.1** (**First-order equivalence between minimizer of averaged noise-augmented loss functions vs averaged minimizers of single noise-augmented loss functions**). The average $\bar{\boldsymbol{\theta}}$ of $m$ minimizers of the $m$ perturbed loss functions upon convergence is first-order equivalent to the minimizer $\hat{\boldsymbol{\theta}}$ of the averaged $m$ noise-augmented loss functions as $m \to \infty$ or as $n_e \to \infty$ while $V(\theta_j n_e) = O(1)$. In addition, The higher-order difference between $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$ also approaches 0 as $n_e \to \infty$ while $V(\theta_j n_e) = O(1)$.

Proof: WLOG, we work with the bridge-type noise. in this proof. The average of the minimizers of the $m$ loss functions is

$$\bar{\boldsymbol{\theta}} = m^{-1}\sum_{t=1}^{m} \left(\mathbf{x}'\mathbf{x} + \sum_{i=1}^{n_e} \mathbf{e}_{i,\mathbf{x}}^{(t)'}\mathbf{e}_{i,\mathbf{x}}^{(t)}\right)^{-1}\mathbf{x}'\mathbf{y}, \tag{33}$$

where $e_{ij} \sim N(0, \lambda|\theta_j|^{-1})$. Let $\sum_{i=1}^{n_e}\mathbf{e}_{i,\mathbf{x}}^{(t)'}\mathbf{e}_{i,\mathbf{x}}^{(t)} = \mathrm{E}\left(\sum_{i=1}^{n_e}\mathbf{e}_{i,\mathbf{x}}^{(t)'}\mathbf{e}_{i,\mathbf{x}}^{(t)}\right) + A^{(t)} = \mathrm{diag}(\lambda n_e|\boldsymbol{\theta}|^{-\gamma}) + \bar{A}^{(t)}$. $A^{(t)}$ can be regarded as the sample deviation of $\sum_{i=1}^{n_e}\mathbf{e}_{i,\mathbf{x}}^{(t)'}\mathbf{e}_{i,\mathbf{x}}^{(t)}$ from its mean. Let $\bar{A} = m^{-1}\sum_{t=1}^{m} \bar{A}^{(t)}$, the elements of which are

$$\begin{cases} \bar{A}[j,j] = m^{-1}\sum_{t=1}^{m}\sum_{i=1}^{n_e} e_{ij}^{(t)2} - \lambda n_e|\theta_j|^{-1} & \sim \lambda|m\theta_j|^{-1}(\chi_{n_e m}^2 - n_e m) \\ \bar{A}[j,k] = m^{-1}\sum_{t=1}^{m}\sum_{i=1}^{n_e} e_{ij}^{(t)}e_{ij}^{(t)} & \sim \lambda|\theta_j\theta_j|^{-\frac{1}{2}}m^{-1}\sum_{t=1}^{m}\sum_{i=1}^{n_e} z_{ti}z_{ti}' \end{cases}, \tag{34}$$

where $z_{ti} \sim N(0,1)$ and $z_{ti}' \sim N(0,1)$ independently. Let $S = (\mathbf{x}'\mathbf{x} + \mathrm{diag}(\lambda n_e|\boldsymbol{\theta}|^{-1}))^{-1}$. The Taylor expansion of the inverse of the sum of two matrices, assuming $A^{(t)}$ to be a small increment, is $(S^{-1} + A^{(t)})^{-1} = S - SA^{(t)}S + SA^{(t)}SA^{(t)}S + \ldots$ Therefore, Eqn (33) becomes

$$\bar{\boldsymbol{\theta}} = S\mathbf{x}'\mathbf{y} - S\left(\bar{A} + O(\lambda^2 n_e)\right)S\mathbf{x}'\mathbf{y}. \tag{35}$$

On the other hand, the minimizer of the average of $m$ loss functions is

$$\hat{\boldsymbol{\theta}} = \left(\mathbf{x}'\mathbf{x} + \sum_{i=1}^{n_e m}\hat{\mathbf{e}}_{ij}'\hat{\mathbf{e}}_{ij}\right)^{-1}\mathbf{x}'\mathbf{y} = \left(\mathbf{x}'\mathbf{x} + \mathrm{diag}(\lambda n_e|\boldsymbol{\theta}|^{-1}) + \hat{A}\right)^{-1}\mathbf{x}'\mathbf{y},$$

$$= S\mathbf{x}'\mathbf{y} - S\left(\hat{A} + O(\lambda^2 n_e)\right)S\mathbf{x}'\mathbf{y}, \tag{36}$$

where $e_{ij} \sim N(0, \lambda|m\theta_j|^{-1})$ for the sake of yielding the same regularization effect as imposed on $\bar{\boldsymbol{\theta}}$; and $\hat{A}$ is defined in a similar manner as $\bar{A}$, the elements of which are
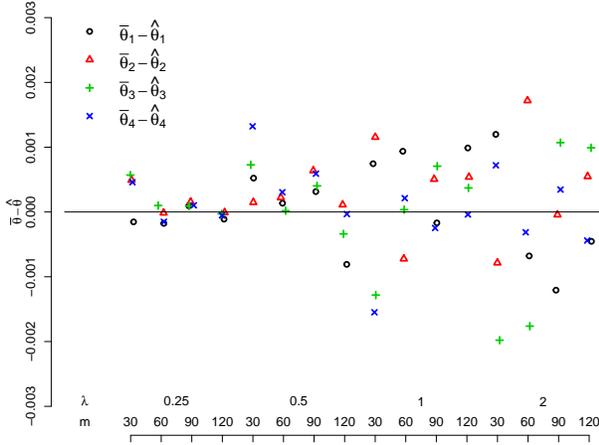
$$\begin{cases} \hat{A}[j,j] = \sum_{i=1}^{n_e m} e_{ij}^2 - \lambda n_e|\theta_j|^{-1} & \sim \lambda|m\theta_j|^{-1}(\chi_{n_e m}^2 - n_e m) \\ \hat{A}[j,k] = \sum_{i=1}^{n_e m} e_{ij}e_{ij} & \sim \frac{\lambda}{m}|\theta_j\theta_k|^{-\frac{1}{2}}\sum_{i=1}^{n_e m} z_i z_i' \end{cases}, \tag{37}$$

where $z_i \sim N(0,1)$ and $z_i' \sim N(0,1)$ independently. $\bar{A}$ and $\hat{A}$ in Eqn (34) and (37) follow the same distribution. The expected values of $\bar{A}[j,j]$, $\bar{A}[j,k]$, $\hat{A}[j,j]$, and $\hat{A}[j,k]$ are all equal to zero; the variance of $\bar{A}[j,j]$ and $\hat{A}[j,j]$ is $\lambda^2|m\theta_{jk}|^{-2}2n_e m = 2\lambda(\lambda n_e)|\theta_{jk}|^{-2}2/m$, and that of $\bar{A}[j,k]$ and $\hat{A}[j,k]$ is $\lambda^2 m^{-2}|\theta_{jk}\theta_{jl}|^{-1}n_e m = \lambda(\lambda n_e)|\theta_{jk}|^{-2}2/m$. As $m$ increases, both variance terms shrink to 0. As $n_e$ increases while $O(n_e\lambda) = 1$, then both variance terms shrinks to 0 as well. In other words, we expect $\bar{A}$ and $\hat{A}$ to be very similar. As such, $\bar{\boldsymbol{\theta}}$ in Eqn (35) and $\hat{\boldsymbol{\theta}}$ in Eqn (36) are also very similar. In addition, as $n_e$ increases and $\lambda n_e = O(1)$, the higher-order terms also goes to 0.

To first illustrate the similarity between $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$, we simulated data ($n = 30$) from linear

regression and a Poisson regression models, where the linear predictor is $\mathbf{X}^T\boldsymbol{\theta} = X_1 + 0.75X_2 + 0.5X_3 + 0X_4$. $\mathbf{X}$ and the error in the linear regression was simulated from $\mathrm{N}(0,1)$ independently. The PANDA augmented noises $\mathbf{e}$ in both cases were drawn from $\mathrm{N}(0,\lambda^2)$ with $n_e = 200$. We examined $m = 30, 60, 90, 120$ and $\lambda^2 = 0.25, 0.5, 1, 2$, calculated $\hat{\boldsymbol{\theta}}$ and $\bar{\boldsymbol{\theta}}$, and plotted their differences in Figure S.1. The results show minimal difference between $\hat{\boldsymbol{\theta}}$ and $\bar{\boldsymbol{\theta}}$.
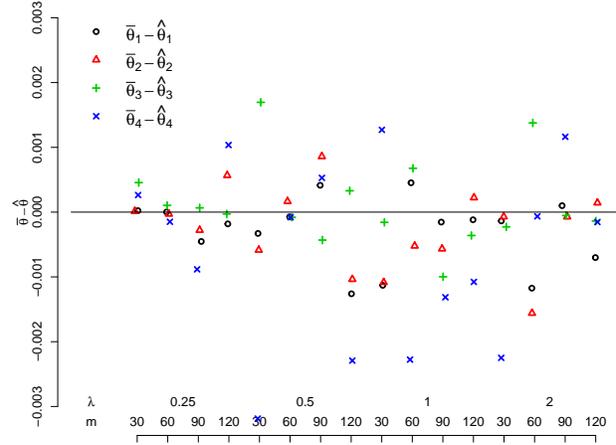


Figure S.1: Differences between $\bar{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}$ in linear regression (top) and in Poisson regression (bottom)

## S.9 PANDA in Each Iteration Realizes Weighted Ridge in Linear Regression

**Corollary S.2 (PANDA and weighted ridge regression).** The OLS estimator in each iteration of PANDA on the noise augmented data is equivalent to the weighted ridge estimator
$\hat{\boldsymbol{\theta}} = \left(\mathbf{x}^T\mathbf{x} + \mathbf{e}^T\mathbf{e}\right)^{-1}\mathbf{x}^T\mathbf{y}$.

The proof is straightforward. Let $\tilde{\mathbf{x}} = (\mathbf{x}, \mathbf{e}_x)^T$. In each iteration of PANDA, the OLS estimator $\hat{\boldsymbol{\theta}} = (\tilde{\mathbf{x}}^T\tilde{\mathbf{x}})^{-1}\tilde{\mathbf{x}}^T(\mathbf{y},\mathbf{0}) = (\mathbf{x}^T\mathbf{x} + \mathbf{e}_x^T\mathbf{e}_x)^{-1}\mathbf{xy}$, leading to Corollary S.2. If $n_e \to \infty$, then $\mathbf{e}_x^T\mathbf{e}_x \to n_e\mathrm{V}(\mathbf{e}_x)$. For example, if the NGD is $\mathrm{N}(0,\lambda|\theta|_j^{-\gamma})$, then $n_e\mathrm{V}(\mathbf{e}_x) = \mathrm{diag}((n_e\lambda)|\theta|_j^{-\gamma})$; and $(\lambda n_e)$ can be tuned as one single tuning parameter.