

Robust High-dimensional Statistics II: Minimum Divergence Approach

Dr. Abhik Ghosh

Indian Statistical Institute, Kolkata, India.

2022



- 1 **Introduction**
- 2 **MD-LASSO: The Initial Attempt**
- 3 **Minimum Penalized Density Power Divergence Estimator (MPDPDE)**
- 4 **Robustness of MPDPDE: Influence functions**
- 5 **Theoretical Properties of MPDPDE**
- 6 **Extending MPDPDE for GLM**
- 7 **Other Minimum Distance Methods**

- 1 **Introduction**
- 2 MD-LASSO: The Initial Attempt
- 3 Minimum Penalized Density Power Divergence Estimator (MPDPDE)
- 4 Robustness of MPDPDE: Influence functions
- 5 Theoretical Properties of MPDPDE
- 6 Extending MPDPDE for GLM
- 7 Other Minimum Distance Methods

Background: Minimum Distance Approach

- **Statistical Divergence Measure** $d(\cdot, \cdot)$ is a non-negative functional of two density (or distribution) functions that equals zero if and only if two arguments are identically equal.
- Examples: Kullback-Leibler (KL) Divergence, (Squared) Hellinger distance, L_2 -Divergence

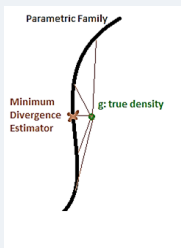
Minimum Divergence estimator

- Parametric inference: Model density family $\mathcal{F} = \{f_\theta : \theta \in \Theta \subseteq \mathbb{R}^p\}$.
- Sample is observed from a population having true density g .
- **Minimum divergence estimator:**

$$\hat{\theta} = \arg \min_{\theta \in \Theta} d(\hat{g}, f_\theta),$$

where \hat{g} is a non-parametric sample-based estimate of g .

- Motivation: MLE is the minimum KL divergence estimator!



Density based divergences has become extremely popular in robust statistical inference!

(Pardo, 2006; Basu et al., 2011)

$$\hat{\theta} = \arg \min_{\theta} \left\{ d(\hat{g}, f_{\theta}) + \sum_{j=1}^p \rho_{\lambda}(|\beta_j|) \right\}$$

where $\rho_{\lambda}(\cdot)$ is a penalty function (lasso, SCAD, MCP, ...).

- If \hat{g} involves kernel, it is extremely difficult to obtain and study the resulting estimator in high-dimensional settings!
- We must use an appropriate distance for which g can be estimated without kernel (e.g., use the fact $gdx = dG(x)$)
- For useful distances, we can write $d(g, f_{\theta}) = \int \rho(\mathbf{z}, \theta)g(\mathbf{z})d\mathbf{z}$ + a term independent of θ . Then the corresponding empirical loss-function has the form:

$$\mathbf{d}(\hat{g}, \mathbf{f}_{\theta}) = \frac{1}{n} \sum_{i=1}^n \rho(\mathbf{z}_i, \theta) + \mathbf{K}(\text{can be ignored}).$$

- 1 Introduction
- 2 MD-LASSO: The Initial Attempt**
- 3 Minimum Penalized Density Power Divergence Estimator (MPDPDE)
- 4 Robustness of MPDPDE: Influence functions
- 5 Theoretical Properties of MPDPDE
- 6 Extending MPDPDE for GLM
- 7 Other Minimum Distance Methods

MD-LASSO for linear regression

Standard **linear regression model** (LRM): $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$,
where $\mathbf{y} = (y_1, \dots, y_n)^T$ are responses, $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_n)^T$ is the design matrix, and
 $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ are the random error components. $f_{\boldsymbol{\beta}}(Y|\mathbf{x}) \equiv N(\mathbf{x}^T \boldsymbol{\beta}, \sigma^2)$

Minimum L_2 divergence Estimator

$$\begin{aligned} \text{Minimize } d_1(g(Y|\mathbf{x}), f_{\boldsymbol{\beta}}(Y|\mathbf{x})) &= \int [g(Y|\mathbf{x}) - f_{\boldsymbol{\beta}}(Y|\mathbf{x})]^2 dY \\ &= \int f_{\boldsymbol{\beta}}(Y|\mathbf{x})^2 dY - 2 \int f_{\boldsymbol{\beta}}(Y|\mathbf{x}) g(Y|\mathbf{x}) dY + \int g(Y|\mathbf{x}) dY \\ &= \frac{1}{\sqrt{2\pi}\sigma} - \frac{2}{\sqrt{2\pi}\sigma} \int e^{-\frac{1}{2\sigma^2}(Y-\mathbf{x}^T \boldsymbol{\beta})^2} dG(Y|\mathbf{x}) + \text{constant}. \end{aligned}$$

Minimum Distance (MD) Loss Function (Lozano et al., 2016)

$$L_n(\boldsymbol{\beta}) = -c \log \left[\sum_{i=1}^n e^{-\frac{1}{2c}(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2} \right], \quad c \text{ is a scaling parameter.}$$

The MD-LASSO (Lozano et al., 2016)

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ -c \log \left[\sum_{i=1}^n e^{-\frac{1}{2c}(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2} \right] + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

Standard **linear regression model** (LRM): $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where $\mathbf{y} = (y_1, \dots, y_n)^T$ are responses, $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_n)^T$ is the design matrix, and $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ are the random error components. $f_{\boldsymbol{\beta}}(Y|\mathbf{x}) \equiv N(\mathbf{x}^T \boldsymbol{\beta}, \sigma^2)$

Definition

Consider any sample of n points (Y, \mathbf{x}) and all possible corrupted samples (Y', \mathbf{x}') that are obtained by replacing m of the original points by arbitrary values. The **breakdown point** of a regression estimator $\hat{\boldsymbol{\beta}}$ at the sample (Y, \mathbf{x}) is defined as

$$\epsilon_n^*(\hat{\boldsymbol{\beta}}; Y, \mathbf{x}) = \min \left\{ \frac{m}{n} : \sup_{(Y', \mathbf{x}')} \|\hat{\boldsymbol{\beta}}(Y', \mathbf{x}')\|_2 = \infty \right\}.$$

Asymptotic breakdown point: $\epsilon^* = \lim_{n \rightarrow \infty} \epsilon_n^*(\hat{\boldsymbol{\beta}}; Y, \mathbf{x})$.

Result (Lozano et al., 2016)

For any finite choice of c , let $Q_c(\boldsymbol{\beta})$ denote the MD-Lasso objective function and consider the non-empty set $B_c = \{\boldsymbol{\beta} : \boldsymbol{\beta} \text{ is a local optimum for the MD-Lasso problem and } Q_c(\boldsymbol{\beta}) \leq Q_c(\mathbf{0})\}$. For every $\alpha \in (0, 1)$, the breakdown point of any solution in B_c is at least α .

The MD-Lasso can tolerate at least αn arbitrarily corrupted observations and still produce estimates having bounded ℓ_2 -norm!

Lozano et al. (2016) derived the consistency of the MD-LASSO under the high-dimensional set-up assuming **Restricted Strong Convexity in a local neighborhood** of true parameter value, when

- MD-LASSO belongs to the same local convexity neighborhood!
- MD-LASSO is computed under additional restriction $\|\beta\|_2 \leq r$ (feasibility)!

- 1 Introduction
- 2 MD-LASSO: The Initial Attempt
- 3 Minimum Penalized Density Power Divergence Estimator (MPDPDE)**
- 4 Robustness of MPDPDE: Influence functions
- 5 Theoretical Properties of MPDPDE
- 6 Extending MPDPDE for GLM
- 7 Other Minimum Distance Methods

The Density Power Divergence for IID Data

Data: $X_1, \dots, X_n \sim g$ (IID), Model: $f_\theta, \theta \in \Theta \subseteq \mathbb{R}^p$;

Density Power Divergence is a generalization of the KL-divergence (Basu et al., 1998)

$$d_\alpha(g, f_\theta) = \int f_\theta^{1+\alpha} - \frac{1+\alpha}{\alpha} \int f_\theta^\alpha g + \frac{1}{\alpha} \int g^{1+\alpha}, \quad \text{if } \alpha > 0,$$
$$d_0(g, f_\theta) = \lim_{\alpha \rightarrow 0} d_\alpha(g, f_\theta) = \int g \log(g/f_\theta) = \text{KL-divergence.}$$

The Minimum Density Power Divergence Estimation (MDPDE)

The MDPDE of θ is the minimizer of $d_\alpha(\hat{g}, f_\theta)$ with respect to θ , or equivalently the minimizer of

$$H_n(\theta) = \int f_\theta^{1+\alpha} - \frac{1+\alpha}{\alpha} \frac{1}{n} \sum_{i=1}^n f_\theta^\alpha(\mathbf{X}_i).$$

Properties of the MDPDE under IID data (Basu et al., 2011)

- MDPDE at $\alpha = 0$ is the MLE, most efficient but highly non-robust (unbounded IF).
- Consistent and Asymptotically normal for all $\alpha \geq 0$.
- Large $\alpha =$ **more robust**, **less efficient**. Small $\alpha =$ **less robust**, **more efficient**.
- The loss in efficiency is not quite significant at small $\alpha > 0$ in most cases.

Independent Non-homogeneous (INH) Set-up

- Data: $X_i \sim g_i$ independently for $i = 1, \dots, n$.
- Model: $X_i \sim f_{i,\theta}$, $\theta \in \Theta \subseteq \mathbb{R}^p$ for $i = 1, \dots, n$.
- Example: Fixed-design regressions.

The MDPDE under INH Set-up (Ghosh and Basu, 2013)

- The MDPDE of θ is defined as the minimizer of average DPD measure $\frac{1}{n} \sum_{i=1}^n d_\alpha(\hat{g}_i, f_{i,\theta})$.
- This leads to the objective function

$$H_n(\theta) = \frac{1}{n} \sum_{i=1}^n \left[\int \mathbf{f}_{i,\theta}^{1+\alpha} - \frac{1+\alpha}{\alpha} \mathbf{f}_{i,\theta}^\alpha(\mathbf{X}_i) \right].$$

Properties are similar to the IID case (under slightly different assumptions).

The DPD loss function for Linear Regression

Standard **linear regression model** (LRM): $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where $\mathbf{y} = (y_1, \dots, y_n)^T$ are responses, $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_n)^T$ is the design matrix, and $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ are the random error components.

For each i , $y_i \sim f_{i,\boldsymbol{\theta}}$, where $f_{i,\boldsymbol{\theta}} \equiv N(\mathbf{x}_i^T \boldsymbol{\beta}, \sigma^2)$ and $\boldsymbol{\theta} = (\boldsymbol{\beta}, \sigma^2)$. So, it belongs to the **INH Set-up!!**

The DPD-based loss function (Ghosh and Basu, 2013)

- According to the INH set-up, the MDPDE should be the minimizer of

$$H_n(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \left[\int f_{i,\boldsymbol{\theta}}^{1+\alpha} - \frac{1+\alpha}{\alpha} f_{i,\boldsymbol{\theta}}^\alpha(X_i) \right].$$

- For Normal error density, it simplifies to

$$L_n^\alpha(\boldsymbol{\beta}, \sigma) = \frac{1}{(2\pi)^{\alpha/2} \sigma^\alpha \sqrt{1+\alpha}} \left[1 - \frac{(1+\alpha)^{3/2}}{\alpha} \frac{1}{n} \sum_{i=1}^n e^{-\alpha \frac{(y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2\sigma^2}} \right] + \frac{1}{\alpha}$$

Durio and Isaia (2011) considered it from random-design perspective: same objective function!

As $\alpha \downarrow 0$, $L_n^\alpha(\boldsymbol{\beta}, \sigma)$ coincides (in a limiting sense) with the negative log-likelihood (plus one).

(why? think L-Hospital's rule.)

Penalized DPD

Standard **linear regression model** (LRM): $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$,
where $\mathbf{y} = (y_1, \dots, y_n)^T$ are responses, $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_n)^T$ is the design matrix, and
 $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T \sim \mathcal{N}_n(\mathbf{0}, \sigma^2 \mathbf{I}_n)$ are the random error components.

Penalized DPD

$$Q_n^\alpha(\boldsymbol{\beta}, \sigma) = L_n^\alpha(\boldsymbol{\beta}, \sigma) + \sum_{j=1}^p \rho_\lambda(|\beta_j|)$$

where $\rho_\lambda(\cdot)$ is a penalty function (lasso, SCAD, MCP, ...).

Minimum Penalized DPD estimator (MPDPD) of $(\boldsymbol{\beta}, \sigma)$ at any given $\alpha \in [0, 1]$

Minimize the penalized DPD objective function $Q_n^\alpha(\boldsymbol{\beta}, \sigma)$ with respect to both $\boldsymbol{\beta}$ and σ !

- As $\alpha \downarrow 0$, this becomes non-concave penalized negative log-likelihood (non-robust).
- At $\alpha = 1$, it has a link to MD-LASSO if the penalty is ℓ_1 (but not exactly coincide - why?)
- Zang et al. (2017): Empirically studied the MPDPDE with grouped lasso penalty (without theory).
- Ghosh and Majumdar (2020): Complete theory, including IF, for general non-concave penalties!

Non-convex minimization problem:

Minimize the penalized DPD objective function $Q_n^\alpha(\beta, \sigma)$ with respect to both β and σ !

Main idea

Starting from $\hat{\beta}, \hat{\sigma}$, **Iteratively** minimize the following:

$$R_\lambda^\alpha(\beta) = L_n^\alpha(\beta, \hat{\sigma}) + \sum_{j=1}^p \rho_\lambda(|\beta_j|),$$

$$S^\alpha(\sigma) = L_n^\alpha(\hat{\beta}, \sigma).$$

Solving Individual Minimization Problems

- Update β using a Concave-Convex Procedure (Yuille and Rangarajan, 2003):

$$\rho_\lambda(|\beta_j|) = \tilde{J}_\lambda(|\beta_j|) + \lambda|\beta_j| \simeq \nabla \tilde{J}_\lambda(|\beta_j^c|)\beta_j + \lambda|\beta_j|$$

where $\tilde{J}_\lambda(\cdot)$ is differentiable and concave, β^c is a current solution.

- Update σ using gradient descent.

$$\hat{\beta}^{(k+1)} = \operatorname{argmin}_{\beta} \left\{ L_n^{\alpha}(\beta, \hat{\sigma}^{(k)}) + \sum_{j=1}^p \left[\nabla \tilde{J}_{\lambda}(|\hat{\beta}_j^{(k)}|) \beta_j + \lambda |\beta_j| \right] \right\};$$

$$\hat{\sigma}^{2(k+1)} = \frac{\left[\sum_{i=1}^n w_i^{(k)} - \frac{\alpha}{(1+\alpha)^{3/2}} \right]}{\left[\sum_{i=1}^n w_i^{(k)} (y_i - \mathbf{x}_i^T \beta^{(k+1)})^2 \right]},$$

$$w_i^{(k)} := \exp \left\{ -\alpha \frac{(y_i - \mathbf{x}_i^T \beta^{(k)})^2}{\sigma^{2(k)}} \right\}.$$

To choose λ , we use a robust High-dimensional BIC:

$$\text{RobHBIC}(\lambda) = \log(\hat{\sigma}^2) + \frac{\log \log(\mathbf{n}) \log \mathbf{p}}{\mathbf{n}} \|\hat{\beta}\|_0, \quad (1)$$

where $\hat{\sigma}^2$ is estimated each time as the MPDPDE of σ^2 .

Select the optimal λ^* that minimizes the HBIC over a pre-determined set of values Λ_n :

$$\lambda^* = \underset{\lambda \in \Lambda_n}{\operatorname{argmin}} \text{RobHBIC}(\lambda).$$

A General formulation with Location-Scale Error

Standard **linear regression model** (LRM): $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$,
where $\mathbf{y} = (y_1, \dots, y_n)^T$ are responses, $\mathbf{X} = (\mathbf{x}_1 \cdots \mathbf{x}_n)^T$ is the design matrix, and
 $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^T$ are the random error components.

We assume the error component $\epsilon_j \sim \frac{1}{\sigma} f\left(\frac{\epsilon}{\sigma}\right)$,

where f is any univariate density with mean 0 and variance 1, such that $M_f^{(\alpha)} = \int f(\epsilon)^{1+\alpha} d\epsilon < \infty$.

Corresponding DPD-based loss function (Ghosh and Majumdar, 2020)

$$L_n^{(\alpha)}(\boldsymbol{\beta}, \sigma) = \frac{1}{\sigma^\alpha} \mathbf{M}_f^{(\alpha)} - \frac{1+\alpha}{\alpha} \frac{1}{n\sigma^\alpha} \sum_{i=1}^n f^\alpha\left(\frac{\mathbf{y}_i - \mathbf{x}_i^t \boldsymbol{\beta}}{\sigma}\right) + \frac{1}{\alpha}. \quad (2)$$

The General MPDPDE

$$\left(\hat{\boldsymbol{\beta}}, \hat{\sigma}\right) = \arg \min_{(\boldsymbol{\beta}, \sigma)} \left\{ L_n^{(\alpha)}(\boldsymbol{\beta}, \sigma) + \sum_{j=1}^p p_\lambda(|\beta_j|) \right\}$$

- 1 Introduction
- 2 MD-LASSO: The Initial Attempt
- 3 Minimum Penalized Density Power Divergence Estimator (MPDPDE)
- 4 Robustness of MPDPDE: Influence functions**
- 5 Theoretical Properties of MPDPDE
- 6 Extending MPDPDE for GLM
- 7 Other Minimum Distance Methods

MPDPE Functional

$$\mathbf{T}_\alpha(\mathcal{G}) = \left(\mathbf{T}_\alpha^\beta, T_\alpha^\sigma \right) = \arg \min_{(\beta, \sigma)} \left\{ \int L_\alpha^*((y, \mathbf{x}); \theta) d\mathcal{G}(y, \mathbf{x}) + \sum_{j=1}^p \rho_\lambda(|\beta_j|) \right\}$$

where

$$L_\alpha^*((y, \mathbf{x}); \theta) = \frac{1}{\sigma^\alpha} M_f^{(\alpha)} - \frac{1 + \alpha}{\alpha} \frac{1}{\sigma^\alpha} f^\alpha \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) + \frac{1}{\alpha}.$$

MPDPDE Functional

$$\mathbf{T}_\alpha(\mathcal{G}) = \left(\mathbf{T}_\alpha^\beta, T_\alpha^\sigma \right) = \arg \min_{(\beta, \sigma)} \left\{ \int L_\alpha^*((y, \mathbf{x}); \theta) d\mathcal{G}(y, \mathbf{x}) + \sum_{j=1}^p \rho_\lambda(|\beta_j|) \right\}$$

where

$$L_\alpha^*((y, \mathbf{x}); \theta) = \frac{1}{\sigma^\alpha} M_f^{(\alpha)} - \frac{1+\alpha}{\alpha} \frac{1}{\sigma^\alpha} f^\alpha \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) + \frac{1}{\alpha}.$$

$$\psi_\alpha((y, \mathbf{x}); \theta) = \nabla L_\alpha^*((y, \mathbf{x}); \theta) = \frac{(1+\alpha)}{\sigma^{\alpha+1}} \begin{bmatrix} \psi_{1,\alpha} \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) \mathbf{x} \\ \psi_{2,\alpha} \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) \end{bmatrix}, \quad (3)$$

where $M_f^{(\alpha)} = \int f^{1+\alpha}$ and

$$\begin{aligned} \psi_{1,\alpha}(s) &= u(s) f^\alpha(s), \\ \psi_{2,\alpha}(s) &= \{s u(s) + 1\} f^\alpha(s) - \frac{\alpha}{\alpha+1} M_f^{(\alpha)}. \end{aligned} \quad (4)$$

Assumptions

Suppose that the **penalty** $\tilde{p}_\lambda(s) = p_\lambda(|s|)$ is **twice differentiable in s** and assume the following:

- $M_f^{(\alpha)}$, $\int \psi_\alpha((y, \mathbf{x}); \theta) dG(y, \mathbf{x})$, $\mathbf{J}_\alpha(G; \theta) = \int \nabla \psi_\alpha((y, \mathbf{x}); \theta) dG(y, \mathbf{x})$ are all finite.
- $\mathbf{J}_\alpha^*(G; \theta) := \left[\mathbf{J}_\alpha(G; \theta) + \text{diag} \left\{ \tilde{\mathbf{P}}_\lambda^{**}(\beta), 0 \right\} \right]$ is invertible at $\theta = \theta_g = (\beta_g, \sigma_g)^T = \mathbf{T}_\alpha(G)$.

Define: $\tilde{\mathbf{P}}_\lambda^{**}(\beta) = \text{diag} \left\{ \tilde{p}_\lambda''(\beta_1), \dots, \tilde{p}_\lambda''(\beta_p) \right\}$.

Main Result (Ghosh and Majumdar, 2020)

Whenever it exists, the influence function of the MPDPDE functional \mathbf{T}_α at G is given by

$$\mathcal{IF}((y_t, \mathbf{x}_t), \mathbf{T}_\alpha, \mathbf{G}) = -\mathbf{J}_\alpha^*(\mathbf{G}; \theta^g)^{-1} \begin{bmatrix} \frac{(1+\alpha)}{\sigma_g^{\alpha+1}} \psi_{1,\alpha} \left(\frac{y_t - \mathbf{x}_t^T \beta_g}{\sigma_g} \right) \mathbf{x}_t + \tilde{\mathbf{P}}_\lambda^*(\beta_g) \\ \frac{(1+\alpha)}{\sigma_g^{\alpha+1}} \psi_{2,\alpha} \left(\frac{y_t - \mathbf{x}_t^T \beta_g}{\sigma_g} \right) \end{bmatrix}. \quad (5)$$

Further, if Θ is compact, then the above IF exists for all (y_t, \mathbf{x}_t) .

Assumptions

Consider the general penalty $p_\lambda(|s|)$ where $p_\lambda(s)$ is twice differentiable in s and assume that true value $\beta_g = (\beta_{1g}^T, \mathbf{0}_{p-s}^T)^T$, where β_{1g} contains all and only s -non-zero elements of β_g .

Define: $\mathbf{P}_\lambda^*(\mathbf{v}) = (p'_\lambda(|v_1|)\text{sign}(v_1), \dots, p'_\lambda(|v_q|)\text{sign}(v_q))$, and $\mathbf{P}_\lambda^{**}(\mathbf{v}) = \text{diag}\{p''_\lambda(|v_1|), \dots, p''_\lambda(|v_q|)\}$, for any $\mathbf{v} = (v_1, \dots, v_q)^T$ with all non-zero elements.

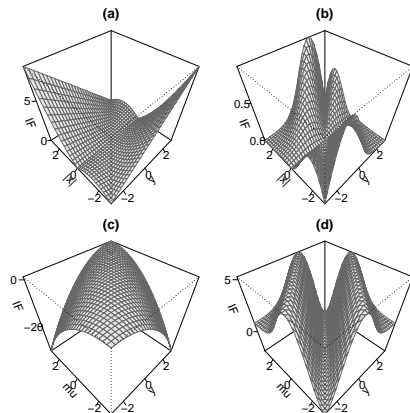
Main Result (Ghosh and Majumdar, 2020)

Let us denote $\theta_g = (\beta_{1g}^T, \mathbf{0}_{p-s}^T, \sigma_g)^T = \mathbf{T}_\alpha(G) = (\mathbf{T}_{1,\alpha}^\beta(G)^T, \mathbf{T}_{2,\alpha}^\beta(G)^T, T_\alpha^\sigma(G))^T$. Then, whenever the associated quantities exists, the influence function of $\mathbf{T}_{2,\alpha}^\beta$ is identically zero at G and that of $(\mathbf{T}_{1,\alpha}^\beta, T_\alpha^\sigma)$ at G is given by

$$\text{IF}((\mathbf{y}_t, \mathbf{x}_t), (\mathbf{T}_{1,\alpha}^\beta, \mathbf{T}_\alpha^\sigma), \mathbf{G}) = -\mathbf{S}_\alpha(\mathbf{G}; \theta_g)^{-1} \begin{bmatrix} \frac{(1+\alpha)}{(\sigma_g)^{\alpha+1}} \psi_{1,\alpha} \left(\frac{y_t - \mathbf{x}_t^T \beta_g}{\sigma_g} \right) \mathbf{x}_{1,t} + \mathbf{P}_\lambda^*(\beta_{1g}) \\ \frac{(1+\alpha)}{(\sigma_g)^{\alpha+1}} \psi_{2,\alpha} \left(\frac{y_t - \mathbf{x}_t^T \beta_g}{\sigma_g} \right) \end{bmatrix},$$

$$\mathbf{S}_\alpha(\mathbf{G}; \theta) = -\frac{(1+\alpha)}{\sigma^{\alpha+2}} E_G \begin{bmatrix} J_{11,\alpha} \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) \mathbf{x}_1 \mathbf{x}_1^T & J_{12,\alpha} \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) \mathbf{x}_1 \\ J_{12,\alpha} \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) \mathbf{x}_1^T & J_{22,\alpha} \left(\frac{y - \mathbf{x}^T \beta}{\sigma} \right) \end{bmatrix} + \begin{bmatrix} \mathbf{P}_\lambda^{**}(\beta_1) & \mathbf{0}_s \\ \mathbf{0}_s^T & 0 \end{bmatrix}.$$

Further, if Θ is compact, then the IF exists in both cases for all (y_t, \mathbf{x}_t) .



Influence function plots for β (panels a and b, $(y_t, \|\mathbf{x}_{1t}\|_1)$ on the (x, y) axes, and ℓ_2 norms of IFs are plotted) and σ (panels c and d, $(y_t, \mathbf{x}_t^T \beta)$ on the axes). We assume \mathbf{x}_{1t} is drawn from $\mathcal{N}_5(\mathbf{0}, \mathbf{I})$, and $\beta_1 = (1, 1, 1, 1, 1)^T$, $\sigma = 1$. Panels a and c are for $\alpha = 0$, while b and d are for $\alpha = 0.5$

- 1 Introduction
- 2 MD-LASSO: The Initial Attempt
- 3 Minimum Penalized Density Power Divergence Estimator (MPDPDE)
- 4 Robustness of MPDPDE: Influence functions
- 5 Theoretical Properties of MPDPDE**
- 6 Extending MPDPDE for GLM
- 7 Other Minimum Distance Methods

Theorem (Ghosh and Majumdar, 2020)

Consider the general penalized DPD loss function $Q_{n,\lambda}^\alpha(\boldsymbol{\theta})$ for a fixed $\alpha \geq 0$ and a general error density f and a penalty function p_λ satisfying Assumption (P). Then, $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\beta}}, \hat{\sigma})$ is a strict local minimizer of $Q_{n,\lambda}^\alpha(\boldsymbol{\theta})$ if and only if

$$\frac{1 + \alpha}{n\hat{\sigma}^{\alpha+1}} \sum_{i=1}^n \psi_{1,\alpha}(r_i(\hat{\boldsymbol{\theta}})) \mathbf{x}_{1i} + \mathbf{P}_\lambda^*(\hat{\boldsymbol{\beta}}_1) = \mathbf{0}, \quad (6)$$

$$\left\| \frac{1 + \alpha}{n\lambda\hat{\sigma}^{\alpha+1}} \sum_{i=1}^n \psi_{1,\alpha}(r_i(\hat{\boldsymbol{\theta}})) \mathbf{x}_{2i} \right\|_\infty < \rho(p_\lambda) = p'_\lambda(0+)/\lambda, \quad (7)$$

$$\frac{1 + \alpha}{n\hat{\sigma}^{\alpha+1}} \sum_{i=1}^n \psi_{2,\alpha}(r_i(\hat{\boldsymbol{\theta}})) = 0, \quad (8)$$

$$\Lambda_{\min} \left(-\frac{1 + \alpha}{n\hat{\sigma}^{\alpha+2}} \sum_{i=1}^n \begin{bmatrix} J_{11,\alpha}(r_i(\hat{\boldsymbol{\theta}})) \mathbf{x}_{1i} \mathbf{x}_{1i}^T & J_{12,\alpha}(r_i(\hat{\boldsymbol{\theta}})) \mathbf{x}_{1i} \\ J_{12,\alpha}(r_i(\hat{\boldsymbol{\theta}})) \mathbf{x}_{1i}^T & J_{22,\alpha}(r_i(\hat{\boldsymbol{\theta}})) \end{bmatrix} \right) > \zeta(p_\lambda, \hat{\boldsymbol{\beta}}_1) \quad (9)$$

$$= \lim_{\epsilon \downarrow 0} \max_{1 \leq j \leq s} \sup_{t_1 < t_2 \in (|\beta_j| - \epsilon, |\beta_j| + \epsilon)} -\frac{p'_\lambda(t_2) - p'_\lambda(t_1)}{t_2 - t_1}.$$

where $\hat{\boldsymbol{\beta}}_1$ contains non-zero components of $\hat{\boldsymbol{\beta}}$, and $\mathbf{x}_i = (\mathbf{x}_{1i}^T, \mathbf{x}_{2i}^T)^T$ is the corresponding partition of \mathbf{x}_i for each i with \mathbf{x}_{1i} having the same dimension as $\hat{\boldsymbol{\beta}}_1$.

Denote the non-zero index set of the true coefficient vector β^* by S .

- **Restricted eigenvalue condition**

$$\frac{\|\mathbf{X}\delta\|^2}{n\|\delta\|^2} \geq \kappa$$

for some $\kappa > 0$ and $\delta \in \mathbb{R}^p$ s.t. $\|\delta_{S^c}\|_1 \leq 3\|\delta_S\|_1$.

- **Our condition** [Ghosh and Majumdar, 2020]

$$\min_{(\delta, \sigma) \in \mathcal{N}_0} \Lambda_{\min} \left[\frac{1}{n} \mathbf{X}_S^T \nabla^2 L_n^\alpha(\delta, \sigma) \mathbf{X}_S \right] \geq c$$

for $c > 0$, and

$$\mathcal{N}_0 = \left\{ (\delta, \sigma) : \delta_{S^c} = \mathbf{0}, \|(\delta_S, \sigma) - (\beta_S^*, \sigma^*)\|_\infty < d_n = \frac{\min_{j \in S} |\beta_j^*|}{2} \right\}$$

Results: Weak Oracle Properties

(A1) $\|\mathbf{x}^{(j)}\| = O(\sqrt{n})$ for $j = 1, \dots, p$, where $\mathbf{x}^{(j)}$ is the j^{th} column of \mathbf{X} .

(A2) For some $C \in (0, 1)$, $\tau_1 \in [0, 0.5]$, we have

$$\left\| \left(\mathbf{X}_S^{*T} \nabla^2 L_n^\alpha(\theta_0) \mathbf{X}_S^* \right)^{-1} \right\|_\infty = O\left(\frac{b_s}{n}\right), \quad \left\| \left(\mathbf{X}_N^{*T} \nabla^2 L_n^\alpha(\theta_0) \mathbf{X}_S^* \right) \left(\mathbf{X}_S^{*T} \nabla^2 L_n^\alpha(\theta_0) \mathbf{X}_S^* \right)^{-1} \right\|_\infty < \min \left\{ \frac{C p'_\lambda(0+)}{p'_\lambda(d_n)}, O(n^{\tau_1}) \right\},$$

$$\max_{(\delta, \sigma) \in \mathcal{N}_0} \max_{1 \leq j \leq p+1} \left\{ \Lambda_{\max} \left(\mathbf{X}_S^* \left[\sum_{i=1}^n \nabla^2_{(\delta, \sigma)} \psi_{\alpha, j} \left((y_i, \mathbf{x}_i), (\delta^T, \mathbf{0}_{p-s}^T, \sigma)^T \right) \right] \mathbf{X}_S^{*T} \right) \right\} = O(n).$$

(A3) Let $s = O(n^{\tau_0})$. For some $\tau \in (0, 0.5]$, define $\tau^* = \min\{0.5, 2\tau - \tau_0\} - \tau_1$. Then,

$$d_n \geq \log n / n^\tau, \quad b_s = o(\min\{n^{1/2-\tau} \sqrt{\log n}, n^\tau / s \log n\}), \quad p'_\lambda(d_n) = o\left(\frac{\log n}{b_s n^\tau}\right), \quad \lambda \geq \frac{(\log n)^2}{n^{\tau^*}}.$$

Also, $\max_{1 \leq j \leq p} \|\mathbf{x}^{(j)}\|_\infty = o(n^{\tau^*} / \sqrt{\log n})$ and $\max_{(\delta, \sigma) \in \mathcal{N}_0} \zeta(p_\lambda; \delta) = o\left(\max_{(\delta, \sigma) \in \mathcal{N}_0} \Lambda_{\min} \left[\frac{1}{n} \mathbf{X}_S^{*T} \nabla^2 L_n^\alpha(\delta, \sigma) \mathbf{X}_S^* \right]\right)$.

(A4) For any $\mathbf{a} \in \mathbb{R}^n$ and $0 < \epsilon < \|\mathbf{a}\| / \|\mathbf{a}\|_\infty$, there is $c_1 > 0$ such that

$$P \left(\left| \frac{1 + \alpha}{\sigma_0^{\alpha+1}} \sum_{i=1}^n a_i \psi_{1, \alpha}(r_i(\theta_0)) \right| > \|\mathbf{a}\| \epsilon \right) \leq 2e^{-c_1 \epsilon^2}.$$

Theorem (Ghosh and Majumdar, 2020)

Let $s = o(n)$, $\log p = O(n^{1-2\tau^*})$ and Assumptions (P), (A1)–(A4) hold for the given α . Then, there exist MNPDPDEs $\hat{\beta} = (\hat{\beta}_S^T, \hat{\beta}_N^T)^T$ of β , with $\hat{\beta}_S \in \mathbb{R}^s$, and $\hat{\sigma}$ of σ such that $(\hat{\beta}, \hat{\sigma})$ is a (strict) local minimizer of $Q_{n, \lambda}^\alpha(\theta)$, with

$$\hat{\beta}_N = \mathbf{0}_{p-s}, \quad \left\| \hat{\beta}_S - \beta_{S0} \right\|_\infty = O\left(\frac{\log n}{n^\tau}\right), \quad |\hat{\sigma} - \sigma_0| = O\left(\frac{\log n}{n^\tau}\right),$$

holding with probability $\geq 1 - (2/n)(1 + s + (p - s) \exp[-n^{1-2\tau^*}])$.

Stronger Assumptions

(A2*) For some $c > 0$, the design matrix \mathbf{X} satisfies

$$\begin{aligned} \min_{(\delta, \sigma) \in \mathcal{N}_0} \Lambda_{\min} \left[\mathbf{X}_S^{*T} \nabla^2 L_n^\alpha(\delta, \sigma) \mathbf{X}_S^* \right] &\geq cn, \quad \left\| \left(\mathbf{X}_N^{*T} \nabla^2 L_n^\alpha(\theta_0) \mathbf{X}_S^* \right) \right\|_{2, \infty} = O(n), \\ \max_{(\delta, \sigma) \in \mathcal{N}_0} \max_{1 \leq j \leq p+1} \Lambda_{\max} \left(\mathbf{X}_S^* \left[\sum_{i=1}^n \nabla_{(\delta, \sigma)}^2 \psi_{\alpha, j} \left((y_i, \mathbf{x}_i), (\delta^T, \mathbf{0}_{p-s}^T, \sigma)^T \right) \right] \mathbf{X}_S^{*T} \right) &= O(n), \\ E \left\| \frac{1 + \alpha}{\sigma_0^{\alpha+1}} \sum_{i=1}^n \psi_{1, \alpha}(r_i(\theta_0)) \mathbf{x}_{Si} \right\|_2^2 &= O(sn), \quad E \left| \frac{1 + \alpha}{\sigma_0^{\alpha+1}} \sum_{i=1}^n \psi_{2, \alpha}(r_i(\theta_0)) \right|^2 = O(n), \end{aligned}$$

(A3*) For some $\tau \in (0, 0.5]$, we have

$$\begin{aligned} \rho'_\lambda(d_n) &= O(n^{-1/2}), \quad d_n \gg \lambda \gg \min \left\{ s^{1/2} n^{-1/2}, n^{\frac{\tau-1}{2}} \sqrt{\log n} \right\}, \\ \max_{1 \leq j \leq p} \|\mathbf{x}^{(j)}\|_\infty &= o(n^{(1-\tau)/2} / \sqrt{\log n}), \quad \max_{(\delta, \sigma) \in \mathcal{N}_0} \zeta(p_\lambda; \delta) = o(1). \end{aligned}$$

Theorem (Ghosh and Majumdar, 2020)

Suppose $s \ll n$ and $\log p = O(n^{\tau^*})$ for some $\tau^* \in (0, 0.5)$ and Assumptions (P), (A1), (A2*), (A3*) and (A4) hold at a fixed $\alpha \geq 0$. Then, there exists a strict minimizer (strict MPDPDE) $(\hat{\beta}^T, \hat{\sigma})^T = \left((\hat{\beta}_S^T, \hat{\beta}_N^T)^T, \hat{\sigma} \right)$, with $\hat{\beta}_S \in \mathbb{R}^s$, that satisfies the following with probability tending to 1 as $n \rightarrow \infty$.

$$\hat{\beta}_N = \mathbf{0}, \quad \left\| \hat{\beta} - \beta_0 \right\| = \mathbf{O}(\sqrt{\mathbf{s}/n}), \quad |\hat{\sigma} - \sigma_0| = \mathbf{O}(n^{-1/2}).$$

Additional Assumption

(A5) The penalty, loss function and design matrix satisfy the following conditions:

$$p'_\lambda(d_n) = O((sn)^{-\frac{1}{2}}), \quad \max_{1 \leq k \leq n} E |\psi_{k,\alpha}(r_i(\theta_0))|^3 = O(1), \quad k = 1, 2,$$

$$\min_{(\delta, \sigma) \in \mathcal{N}_0} \Lambda_{\min} [\mathbf{X}_S^{*T} \Sigma_{\alpha^*} ((\delta^T, \mathbf{0}_{p-s}, \sigma)^T) \mathbf{X}_S^*] \geq cn, \quad \sum_{i=1}^n [\mathbf{x}_{Si}^{*T} (\mathbf{X}_S^{*T} \Sigma^*(\theta_0) \mathbf{X}_S^*)^{-1} \mathbf{x}_{Si}^*]^{3/2} = o(1).$$

where

$$\Sigma_{\alpha^*}(\theta) = \begin{pmatrix} \mathbf{K}_{11}^{(\alpha)}(\theta) & \mathbf{K}_{12}^{(\alpha)}(\theta) \\ \mathbf{K}_{12}^{(\alpha)}(\theta) & \mathbf{K}_{22}^{(\alpha)}(\theta) \end{pmatrix}, \quad \mathbf{K}_{ij,\alpha}(\theta) = \frac{(1+\alpha)^2}{\sigma^{2\alpha+2}} \text{diag} \{ \psi_{i,\alpha}(r_1(\theta))\psi_{j,\alpha}(r_1(\theta)), \dots, \psi_{i,\alpha}(r_n(\theta))\psi_{j,\alpha}(r_n(\theta)) \}.$$

Theorem (Ghosh and Majumdar, 2020)

In addition to the assumptions of the previous theorem, suppose that $s = o(n^{1/3})$ and (A5) holds. Then, the strict MNPDPE $(\hat{\beta}, \hat{\sigma})$ satisfies the following results with probability tending to 1 as $n \rightarrow \infty$.

- 1 $\hat{\beta}_N = \mathbf{0}$, where $\hat{\beta} = (\hat{\beta}_S^T, \hat{\beta}_N^T)^T$ with $\hat{\beta}_S \in \mathbb{R}^s$.
- 2 Let $\mathbf{A}_n \in \mathbb{R}^{q \times (s+1)}$ such that $\mathbf{A}_n \mathbf{A}_n^T \rightarrow \mathbf{G}$ as $n \rightarrow \infty$, where \mathbf{G} is symmetric and positive definite. Then,

$$\mathbf{A}_n (\mathbf{X}_S^{*T} \Sigma_{\alpha^*}(\theta_0) \mathbf{X}_S^*)^{-1/2} (\mathbf{X}_S^{*T} \Sigma_{\alpha}(\theta_0) \mathbf{X}_S^*) \left((\hat{\beta}_S, \hat{\sigma}) - (\beta_{S0}, \sigma_0) \right) \xrightarrow{D} \mathbf{N}_q(\mathbf{0}_q, \mathbf{G}). \quad (10)$$

- 1 Introduction
- 2 MD-LASSO: The Initial Attempt
- 3 Minimum Penalized Density Power Divergence Estimator (MPDPDE)
- 4 Robustness of MPDPDE: Influence functions
- 5 Theoretical Properties of MPDPDE
- 6 Extending MPDPDE for GLM**
- 7 Other Minimum Distance Methods

Generalized linear model (GLM): Given a covariate value $\mathbf{X} = \mathbf{x}$, the response variable Y has density

$$f(y; \mathbf{x}^T \boldsymbol{\beta}) = \exp \{y\theta - b(\theta) + c(y)\}, \quad \text{with } E[Y|\mathbf{x}] = b'(\theta) = g^{-1}(\mathbf{x}^T \boldsymbol{\beta}), \quad (11)$$

where $b(\cdot)$ and $c(\cdot)$ are known functions, g is a known monotone differentiable link function, and the canonical parameter θ is defined via the linear predictor $\eta = \mathbf{x}^t \boldsymbol{\beta}$.

DPD Loss Function (Ghosh and Basu, 2016)

$$L_n^{(\alpha)}(\boldsymbol{\beta}) = \frac{1}{n(1+\alpha)} \sum_{i=1}^n \left[\int \mathbf{f}(\mathbf{y}, \mathbf{x}_i^t \boldsymbol{\beta})^{\alpha+1} d\mathbf{y} - \frac{1+\alpha}{\alpha} \mathbf{f}(\mathbf{y}_i, \mathbf{x}_i^t \boldsymbol{\beta})^\alpha + \frac{1}{\alpha} \right].$$

Minimum Penalized DPD Estimator (MPDPDE)

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ L_n^{(\alpha)}(\boldsymbol{\beta}, \sigma) + \sum_{j=1}^p p_\lambda(|\beta_j|) \right\}.$$

- **(A1)** The l_2 -norm of each column of \mathbf{X} is $O(\sqrt{n})$.

(A2) $\|\boldsymbol{\Psi}_{n,\alpha}(\beta_0)^{-1}\|_\infty = O\left(\frac{b_s}{n}\right)$, $\|\tilde{\boldsymbol{\Psi}}_{n,\alpha}(\beta_0)\boldsymbol{\Psi}_{n,\alpha}(\beta_0)^{-1}\|_\infty < \min\left\{\frac{C\rho'_\lambda(0+)}{\rho'_\lambda(d_n)}, O(n^{\tau_1})\right\}$ and

$\max_{\delta \in \mathcal{N}_0} \max_{1 \leq j \leq p+1} \{\Lambda_{\max}(\mathbf{X}_S [\nabla_\delta^2 \Gamma_{j,\alpha}(\delta)] \mathbf{X}_S^t)\} = O(n)$, for some constants $C \in (0, 1)$,

$\tau_1 \in [0, 0.5]$, where $\mathcal{N}_0 = \{\delta \in \mathbb{R}^s : \|\delta - \beta_{S0}\|_\infty \leq d_n\}$,

$\boldsymbol{\Gamma}_\alpha(\delta) = (\Gamma_{1,\alpha}(\delta), \dots, \Gamma_{p,\alpha}(\delta))^t = \sum_{i=1}^n \psi_\alpha(y_i; \mathbf{x}_{iS}^t \delta)$. and ∇_δ^2 denotes the second-order derivative with respect to δ .

(A3) For some $\tau \in (0, 0.5]$, we have $d_n \geq \log n/n^\tau$ and

$b_s = o(\min\{n^{1/2-\tau} \sqrt{\log n}, n^\tau/s \log n\})$. Further, with $s = O(n^{\tau_0})$, we define

$\tau^* = \min\{0.5, 2\tau - \tau_0\} - \tau_1$. Then, we have $\rho'_\lambda(d_n) = o\left(\frac{\log n}{b_s n^\tau}\right)$, $\lambda \geq (\log n)^2 n^{-\tau^*}$, and

$\max_{\delta \in \mathcal{N}_0} \zeta(p_\lambda; \delta) = o\left(\max_{\delta \in \mathcal{N}_0} \Lambda_{\min}\left[n^{-1} \boldsymbol{\Psi}_{n,\alpha}((\delta^t, \mathbf{0}_{p-s}^t))\right]\right)$, where ζ denote the local concavity of the penalty p_λ as defined in ?, Definition 1. Also, the maximum (in absolute) element of the design matrix \mathbf{X} is of order $o\left(n^{\tau^*}/\sqrt{\log n}\right)$.

(A4) For any $\mathbf{a} \in \mathbb{R}^n$ and $0 < \epsilon < \|\mathbf{a}\|/\|\mathbf{a}\|_\infty$, there exists a $c_1 > 0$ such that

$$P\left(\left|\sum_{i=1}^n a_i \psi_\alpha(y_i, \mathbf{x}_i^t \beta_0)\right| > \|\mathbf{a}\| \epsilon\right) \leq 2e^{-c_1 \epsilon^2}.$$

(A5) $\rho'_\lambda(d_n) = O((sn)^{-1/2})$, $\max_{1 \leq i \leq n} E|\psi_\alpha(y_i, \mathbf{x}_i^t \beta_0)|^3 = O(1)$,

Theorem (Ghosh, 2021+)

Under the set-up of ultra high-dimensional GLMs (11) with $s = o(n)$ and $\log p = O(n^{1-2\tau^*})$, let us assume that Assumptions (A1)–(A5) and (P) hold for some fixed $\alpha \geq 0$.

Then, there exist MNPDPDEs $\hat{\beta} = (\hat{\beta}_S, \hat{\beta}_N)^t$ of β , with $\hat{\beta}_S \in \mathbb{R}^s$, such that $\hat{\beta}$ is a (strict) local minimizer (MPDPDE) and satisfies the following:

a) $\hat{\beta}_N = \mathbf{0}_{p-s}$, and $\|\hat{\beta}_S - \beta_{0S}\|_\infty = O(n^{-\tau} \log n)$,

with probability at least $P_n = 1 - (2/n)(1 + s + (p - s) \exp[-n^{1-2\tau^*}])$.

b) Let $\mathbf{A}_n \in R^{q \times (s+1)}$ such that $\mathbf{A}_n \mathbf{A}_n^t \rightarrow \mathbf{G}$ as $n \rightarrow \infty$, where \mathbf{G} is symmetric and positive definite. Then, with probability tending to one,

$$\mathbf{A}_n \Omega_{n,\alpha}(\beta_0)^{-1/2} \Psi_{n,\alpha}(\beta_0) (\hat{\beta}_S - \beta_{0S}) \xrightarrow{\mathcal{D}} \mathbf{N}_q(\mathbf{0}_q, \mathbf{G}).$$

$$\Psi_{n,\alpha}(\beta) = \left(\mathbf{X}_S^t \Sigma_{n,\alpha}(\beta_0) \mathbf{X}_S \right), \quad \Omega_{n,\alpha}(\beta) = \left(\mathbf{X}_S^t \Sigma_{n,\alpha}^*(\beta_0) \mathbf{X}_S \right),$$

where $\Sigma_{n,\alpha}(\beta) = \text{diag} \left\{ -\frac{\partial}{\partial \eta} \psi_\alpha(y_i, \eta) \Big|_{\eta = \mathbf{x}_i^t \beta} : i = 1, \dots, n \right\}$

and $\Sigma_{n,\alpha}^*(\beta) = \text{diag} \left\{ \psi_\alpha(y_i, \mathbf{x}_i^t \beta)^2 : i = 1, \dots, n \right\}$.

Theorem (Ghosh, 2021+)

Under the ultra high-dimensional GLMs (11) with $s \ll n$, $\log p = O(n^{\tau^*})$ for some $\tau^* \in (0, 0.5)$, let us assume that Assumptions (A1), (A2*), (A3*), (A4), (A5) and (P) hold for some fixed $\alpha \geq 0$. Then, there exist MNPDPDEs $\widehat{\beta} = (\widehat{\beta}_S^t, \widehat{\beta}_N^t)^t$ of β , with $\widehat{\beta}_S \in \mathbb{R}^s$, such that $\widehat{\beta}$ is a (strict) local minimizer (MPDPDE) and satisfies the following results with probability tending to 1 as $n \rightarrow \infty$.

a) $\widehat{\beta}_N = \mathbf{0}_{p-s}$, and $\|\widehat{\beta} - \beta_0\|_2 = O(\sqrt{s/n})$.

b) Let $\mathbf{A}_n \in \mathbb{R}^{q \times (s+1)}$ such that $\mathbf{A}_n \mathbf{A}_n^t \rightarrow \mathbf{G}$ as $n \rightarrow \infty$, where \mathbf{G} is symmetric and positive definite. Then,

$$\mathbf{A}_n \Omega_{n,\alpha}(\beta_0)^{-1/2} \Psi_{n,\alpha}(\beta_0) \left(\widehat{\beta}_S - \beta_{0S} \right) \xrightarrow{\mathcal{D}} \mathbf{N}_q(\mathbf{0}_q, \mathbf{G}).$$

- 1 Introduction
- 2 MD-LASSO: The Initial Attempt
- 3 Minimum Penalized Density Power Divergence Estimator (MPDPDE)
- 4 Robustness of MPDPDE: Influence functions
- 5 Theoretical Properties of MPDPDE
- 6 Extending MPDPDE for GLM
- 7 Other Minimum Distance Methods**

Penalized log-DPD or γ -divergence

The log-DPD or γ -divergence (Jones et al., 2001; Fujisawa and Eguchi, 2008)

$$d_\gamma(g, f_\theta) = \log \int f_\theta^{1+\gamma} - \frac{1+\gamma}{\gamma} \log \int f_\theta^\gamma g + \frac{1}{\gamma} \log \int g^{1+\gamma}, \quad \text{if } \gamma > 0,$$

$$d_0(g, f_\theta) = \lim_{\gamma \rightarrow 0} d_\gamma(g, f_\theta) = \int g \log(g/f_\theta) = \text{KL-divergence}.$$

The log-DPD or γ -divergence loss function for LRM (Fujisawa and Eguchi, 2008)

Standard **linear regression model** (LRM): $\mathbf{y} = \mathbf{X}\beta + \epsilon$.

$$L_n^{(\gamma)}(\theta) = -\frac{1}{\gamma} \log \left(\frac{1}{n} \sum_{i=1}^n f_\theta(\mathbf{y}_i | \mathbf{x}_i) \right) + \frac{1}{1+\gamma} \log \left(\frac{1}{n} \sum_{i=1}^n \int f_\theta(\mathbf{y} | \mathbf{x}_i) d\mathbf{y} \right), \quad f_\theta(\cdot | \mathbf{x}_i) \equiv \mathcal{N}(\mathbf{x}_i^\top \beta, \sigma^2).$$

$$\text{Minimum Penalized log-DPD Estimator } (\hat{\beta}, \hat{\sigma}) = \arg \min_{(\beta, \sigma)} \left\{ L_n^{(\gamma)}(\beta, \sigma) + \sum_{j=1}^p p_\lambda(|\beta_j|) \right\}$$

- Kawashima and Fujisawa (2017): Proposed and studied empirically with ℓ_1 -penalty only! (no theory)
- Castilla et al. (2020): Derived general theory and IF for general non-concave penalties (similarly as DPD)!

Penalized Bregman divergence

The Bregman divergence (Bregman, 1967))

$$Q(\nu, \mu) = -q(\nu) + q(\mu) + (\nu - \mu)q'(\mu), \quad q \text{ concave.} \quad (12)$$

The log-DPD or γ -divergence loss function for GLM (Fujisawa and Eguchi, 2008)

Given a covariate value $\mathbf{X} = \mathbf{x}$, the response variable Y has density

$$f(y; \mathbf{x}^T \boldsymbol{\beta}) = \exp \{y\theta - b(\theta) + c(y)\}, \quad \text{with } E[Y|\mathbf{x}] = b'(\theta) = g^{-1}(\mathbf{x}^T \boldsymbol{\beta}),$$

$$L_n^{(q)}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n Q(\mathbf{y}_i, g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta})).$$

Minimum Penalized Bregman Divergence Estimator (MPBDE)

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \left\{ L_n^{(q)}(\boldsymbol{\beta}) + \sum_{j=1}^p p_{\lambda}(|\beta_j|) \right\}$$

- Zhang et al. (2010): Studied the properties of MPBDE for the cases $p < n$ and $pn!$
- Zhang et al. (2011): Developed an appropriate algorithm for computation of MPBDE with ℓ_1 penalty!
- Jiang and Zhang (2013): Studied the properties of MPBDE for the high-dimensional cases with $p > n!$

8 **Illustrations**

9 **Summary and Conclusion**

Simulation Setup and Performance Metrics

- Obtain rows of \mathbf{X} as $n = 100$ random draws from $\mathcal{N}(0, \Sigma_X)$, where Σ_X is a positive definite with $(i, j)^{\text{th}}$ element given by $0.5^{|i-j|}$.
- Given p , we consider two settings for β :
 - Setting A (strong signal): For $j \in \{1, 2, 4, 7, 11\}$, $\beta_j = j$, otherwise 0;
 - Setting B (weak signal): Set $\beta_1 = \beta_7 = 1.5$, $\beta_2 = 0.5$, $\beta_4 = \beta_{11} = 1$, and 0 otherwise.
- Generate the random errors as $\epsilon \sim N(0, 0.5^2)$, and set $\mathbf{y} = \mathbf{X}\beta + \epsilon$.
- Three outlier settings:
 - Y-outliers: We add 20 to the response variables of a random 10% of samples,
 - X-outliers: We add 20 to each of the elements in the first 10 rows of \mathbf{X} for a random 10% of samples,
 - No outliers.
- Methods compared- RLARS, sLTS, LAD-Lasso, DPD-lasso, log DPD-lasso, DPD-ncv(with SCAD), along with non-robust Lasso, SCAD and MCP.
- RLARS solution is used as our starting point.

Performance Metrics

$$\text{MSEE}(\hat{\beta}) = (1/p)\|\hat{\beta} - \beta_0\|^2, \quad \text{RMSPE}(\hat{\beta}) = \sqrt{\|\mathbf{y}_{\text{test}} - \mathbf{X}_{\text{test}}\hat{\beta}\|^2},$$

$$\text{EE}(\hat{\sigma}) = |\hat{\sigma} - \sigma_0|, \quad \text{MS}(\hat{\beta}) = |\text{supp}(\hat{\beta})|.$$

$$\text{TP}(\hat{\beta}) = \frac{|\text{supp}(\hat{\beta}) \cap \text{supp}(\beta_0)|}{|\text{supp}(\beta_0)|}, \quad \text{TN}(\hat{\beta}) = \frac{|\text{supp}(\hat{\beta}) \cap \text{supp}(\beta_0)|}{|\text{supp}(\beta_0)|},$$

Table of outputs for $\rho = 500$, $n = 100$ and Y-outliers

Setting B: $\beta_1 = \beta_7 = 1.5, \beta_2 = 0.5, \beta_4 = \beta_{11} = 1$, and 0 otherwise.						
Method	MSEE($\hat{\beta}$) ($\times 10^{-4}$)	RMSPE($\hat{\beta}$) ($\times 10^{-2}$)	EE($\hat{\sigma}$)	TP($\hat{\beta}$)	TN($\hat{\beta}$)	MS($\hat{\beta}$)
RLARS	1.1	4.58	0.09	1.00	1.00	6.00
sLTS	6.2	6.06	0.23	1.00	0.93	40.07
LAD-Lasso	68.6	15.65	2.77	0.65	0.99	6.28
DPD-ncv, $\alpha = 0.2$	0.8	4.28	0.06	1.00	1.00	5.00
DPD-ncv, $\alpha = 0.4$	0.8	4.30	0.06	1.00	1.00	5.00
DPD-ncv, $\alpha = 0.6$	0.8	4.50	0.06	1.00	1.00	5.00
DPD-ncv, $\alpha = 0.8$	0.7	4.59	0.06	1.00	1.00	5.00
DPD-ncv, $\alpha = 1$	0.8	4.61	0.06	1.00	1.00	5.00
DPD-Lasso, $\alpha = 0.2$	61.3	15.10	0.05	1.00	0.00	499.08
DPD-Lasso, $\alpha = 0.4$	58.9	14.41	0.17	1.00	0.05	477.15
DPD-Lasso, $\alpha = 0.6$	56.5	14.85	0.14	1.00	0.10	450.22
DPD-Lasso, $\alpha = 0.8$	55.1	14.29	0.02	1.00	0.13	435.72
DPD-Lasso, $\alpha = 1$	54.2	14.16	0.01	1.00	0.13	433.65
LDPD-Lasso, $\alpha = 0.2$	2.1	5.09	0.07	1.00	0.99	10.19
LDPD-Lasso, $\alpha = 0.4$	2.2	5.12	0.09	1.00	0.99	7.97
LDPD-Lasso, $\alpha = 0.6$	2.3	5.14	0.11	1.00	0.99	7.62
LDPD-Lasso, $\alpha = 0.8$	2.3	5.14	0.13	1.00	1.00	7.38
LDPD-Lasso, $\alpha = 1$	2.3	5.15	0.14	1.00	1.00	7.38
Lasso	134.1	22.41	4.54	0.02	1.00	0.24
SCAD	128.6	20.97	3.60	0.32	0.99	8.72
MCP	141.6	21.09	3.69	0.24	0.99	4.52

Table of outputs for $p = 500$, $n = 100$ and X-outliers

Setting B: $\beta_1 = \beta_7 = 1.5, \beta_2 = 0.5, \beta_4 = \beta_{11} = 1$, and 0 otherwise.						
Method	MSEE($\hat{\beta}$) ($\times 10^{-4}$)	RMSPE($\hat{\beta}$) ($\times 10^{-2}$)	EE($\hat{\sigma}$)	TP($\hat{\beta}$)	TN($\hat{\beta}$)	MS($\hat{\beta}$)
RLARS	2.0	4.2	0.14	1.00	0.99	12.00
sLTS	8.7	5.3	0.24	1.00	0.92	42.50
LAD-Lasso	108.0	20.4	2.87	0.38	0.99	7.71
DPD-ncv, $\alpha = 0.2$	1.2	4.1	0.08	1.00	1.00	7.00
DPD-ncv, $\alpha = 0.4$	1.1	4.0	0.10	1.00	1.00	7.00
DPD-ncv, $\alpha = 0.6$	1.1	4.2	0.12	1.00	1.00	7.00
DPD-ncv, $\alpha = 0.8$	1.4	4.2	0.14	1.00	1.00	7.00
DPD-ncv, $\alpha = 1$	1.5	4.2	0.15	1.00	1.00	7.00
DPD-Lasso, $\alpha = 0.2$	59.5	13.8	0.05	1.00	0.01	495.26
DPD-Lasso, $\alpha = 0.4$	48.6	10.8	0.20	1.00	0.16	420.56
DPD-Lasso, $\alpha = 0.6$	35.3	9.2	0.28	1.00	0.35	329.12
DPD-Lasso, $\alpha = 0.8$	27.6	8.6	0.13	1.00	0.45	278.17
DPD-Lasso, $\alpha = 1$	25.7	9.2	0.01	1.00	0.47	267.29
LDPD-Lasso, $\alpha = 0.2$	1.9	5.0	0.06	1.00	0.98	15.14
LDPD-Lasso, $\alpha = 0.4$	1.8	5.0	0.07	1.00	0.98	14.04
LDPD-Lasso, $\alpha = 0.6$	1.8	5.1	0.07	1.00	0.98	14.03
LDPD-Lasso, $\alpha = 0.8$	1.8	5.0	0.07	1.00	0.98	14.47
LDPD-Lasso, $\alpha = 1$	1.8	5.0	0.07	1.00	0.98	13.90
LASSO	22.6	10.3	0.13	0.99	0.87	70.32
SCAD	45.8	13.8	0.55	0.81	0.98	16.25
MCP	45.2	12.8	0.49	0.81	0.97	16.45

Table of outputs for $p = 500$, $n = 100$ and no outliers

Setting B: $\beta_1 = \beta_7 = 1.5, \beta_2 = 0.5, \beta_4 = \beta_{11} = 1$, and 0 otherwise.						
Method	MSEE($\hat{\beta}$) ($\times 10^{-4}$)	RMSPE($\hat{\beta}$) ($\times 10^{-2}$)	EE($\hat{\sigma}$)	TP($\hat{\beta}$)	TN($\hat{\beta}$)	MS($\hat{\beta}$)
RLARS	1.4	4.73	0.12	1.00	0.99	10.00
sLTS	7.9	5.65	0.24	1.00	0.93	42.00
LAD-Lasso	4.7	3.90	0.42	1.00	1.00	7.30
DPD-ncv, $\alpha = 0.2$	1.4	4.73	0.12	1.00	0.99	10.00
DPD-ncv, $\alpha = 0.4$	1.4	4.73	0.12	1.00	0.99	10.00
DPD-ncv, $\alpha = 0.6$	1.4	4.73	0.12	1.00	0.99	10.00
DPD-ncv, $\alpha = 0.8$	1.4	4.73	0.12	1.00	0.99	10.00
DPD-ncv, $\alpha = 1$	1.4	4.73	0.12	1.00	0.99	10.00
DPD-Lasso, $\alpha = 0.2$	79.1	14.56	0.10	1.00	0.00	499.00
DPD-Lasso, $\alpha = 0.4$	58.2	12.98	0.25	1.00	0.14	429.70
DPD-Lasso, $\alpha = 0.6$	44.9	10.18	0.25	1.00	0.31	348.80
DPD-Lasso, $\alpha = 0.8$	19.9	8.86	0.05	1.00	0.57	215.60
DPD-Lasso, $\alpha = 1$	21.1	11.46	0.00	1.00	0.59	208.70
LDPD-Lasso, $\alpha = 0.2$	1.9	3.94	0.06	1.00	0.97	17.50
LDPD-Lasso, $\alpha = 0.4$	2.0	4.05	0.09	1.00	0.98	15.50
LDPD-Lasso, $\alpha = 0.6$	2.0	4.23	0.09	1.00	0.98	16.90
LDPD-Lasso, $\alpha = 0.8$	2.0	4.16	0.08	1.00	0.98	16.00
LDPD-Lasso, $\alpha = 1$	2.0	4.10	0.08	1.00	0.98	16.40
Lasso	2.1	3.59	0.33	1.00	0.98	12.90
SCAD	0.3	3.71	0.21	1.00	0.99	9.70
MCP	0.3	3.69	0.20	1.00	1.00	6.80

8 Illustrations

9 **Summary and Conclusion**

- We discussed the need minimum distance approach for robust estimation and variable selection in high-dimensional data analyses
- We discussed different minimum distance methods, namely the procedures based on L_2 -divergence, DPD, log-DPD and Bregman divergence
- The DPD based procedures are discussed in great detail, including their IF, oracle consistency and asymptotic normality results, under high-dimensional linear regression model.
- Extension of the minimum penalized DPD estimator is also discussed for ultra-high dimensional GLM.
- Numerical illustrations are provided for comparison of different minimum distance methods, along with an interesting real data example.

- Basu, A., Harris, I. R., Hjort, N. L., and Jones, M. C. (1998). Robust and efficient estimation by minimising a density power divergence. *Biometrika*, **85**, 549–559.
- Basu, A., Shioya, H. and Park, C. (2011). *Statistical Inference: The Minimum Distance Approach*. Chapman & Hall/CRC, Boca de Raton.
- Castilla, E., Ghosh, A., Jaenada, M., and Pardo, L. (2020). On regularization methods based on Rényi's pseudodistances for sparse high-dimensional linear regression models. *arXiv preprint*, arXiv:2007.15929.
- Durio, A. and Isaia, E. D. (2011). The minimum density power divergence approach in building robust regression models. *Informatica*, **22**(1), 43–56.
- Fujisawa, H. and Eguchi, S. (2008). Robust Parameter Estimation with a Small Bias Against Heavy Contamination. *J. Mult. Anal.*, **99**, 2053–2081.
- Ghosh, A. (2021+). Robustness Concerns in High-dimensional Data Analysis and Potential Solutions. To appear in *Big Data Analytics in Chemoinformatics and Bioinformatics (with applications to computer-aided drug design, cancer biology, emerging pathogens and computational toxicology)*, Basak, S.C. and Vracko, M. Eds., Elsevier.
- Ghosh, A., and Basu, A. (2013). Robust estimation for independent non-homogeneous observations using density power divergence with applications to linear regression. *Electron. J. Stat.*, **7**, 2420–2456.
- Ghosh, A., and Basu, A. (2016). Robust Estimation in Generalized Linear Models : The Density Power Divergence Approach. *Test*, **25**(2), 269–290.
- Ghosh, A. and Majumdar, S. (2020). Ultrahigh-dimensional Robust and Efficient Sparse Regression using Non-Concave Penalized Density Power Divergence. *IEEE Trans. Info. Theory*, **66**(12), 7812–7827.

- Jiang, Y., and Zhang, C. (2013). High-dimensional regression and classification under a class of convex loss functions. *Stat. Its Interface*, **6(2)**, 285-299.
- Jones, M. C., Hjort, N. L., Harris, I. R., and Basu, A. (2001). A comparison of related density-based minimum divergence estimators. *Biometrika*, **88(3)**, 865–873.
- Kawashima, T. and Fujisawa, H. (2017). Robust and Sparse Regression via γ -Divergence. *Entropy*, **19(11)**, 608.
- Lozano, A. C. ; Meinshausen, N- and Yang, E. (2016). Minimum Distance Lasso for robust high-dimensional regression. *Electron. J. Stat.*, **10**, 1296–1340.
- Pardo, L. (2006). *Statistical inference based on divergence measures*. CRC press.
- Yuille, A. L., and Rangarajan, A. (2003). The concave-convex procedure. *Neural Comput.*, **15(4)**, 915–936.
- Zang, Y., Zhao, Q., Zhang, Q., et al. (2017). Inferring gene regulatory relationships with a high-dimensional robust approach. *Genet. Epidemiol.*, **41(5)**, 437–454.
- Zhang, C., Jiang, Y., and Chai, Y. (2010). Penalized Bregman divergence for large-dimensional regression and classification. *Biometrika*, **97(3)**, 551–566.
- Zhang, C., Zhang, Z., and Chai, Y. (2011). Penalized Bregman divergence estimation via coordinate descent. *JIRSS*, **10(2)**, 125–140

THANK YOU!

Contact me at:

abhik.ghosh@isical.ac.in