

A Sharper Computational Tool for L_2E Regression

Xiaoqian Liu

UT MD Anderson Cancer Center

Collaborators:

Eric Chi, Rice University

Kenneth Lange, UCLA

Oct. 5, 2023

The classical linear regression:

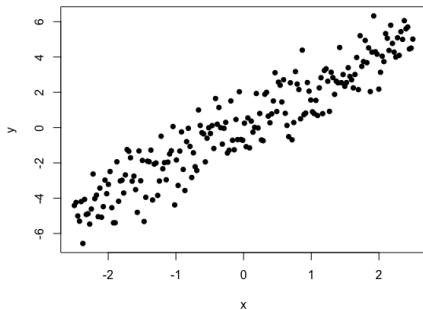
$$y = X\beta + \epsilon$$

- $y \in \mathbb{R}^n, X \in \mathbb{R}^{n \times p}$
- $\epsilon \in \mathbb{R}^n \sim N(0, \tau^{-2}I_n)$
- Least Square estimator / Maximum Likelihood Estimator (MLE):

$$\min_{\beta} \frac{1}{2} \|y - X\beta\|_2^2$$

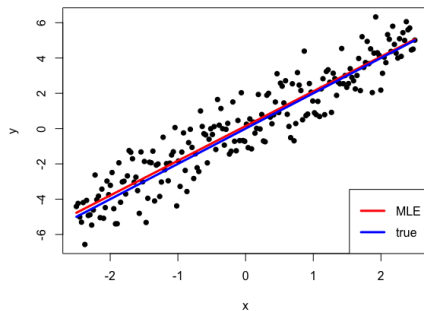
Motivation

An ideal data set:



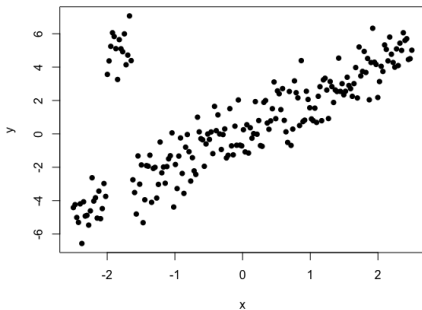
Motivation

An ideal data set:



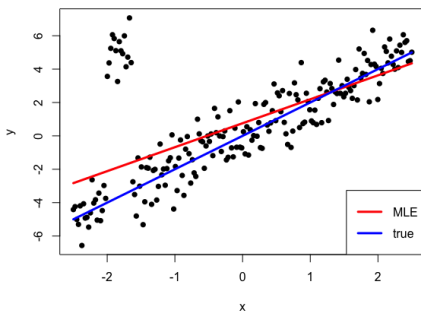
Motivation

In reality, data could be contaminated (**outliers!**).



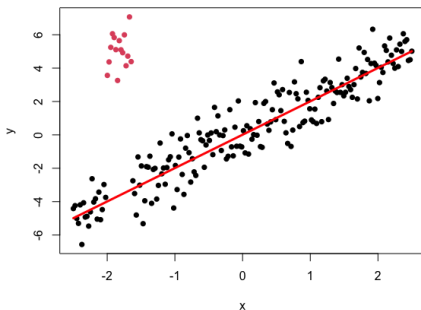
Motivation

In reality, data could be contaminated (outliers!).



Motivation

Our aims: **robust estimation** + **outlier detection** + **structure recovery**



- 1 L2E regression
 - L₂E criterion
 - Structured L₂E model
- 2 Computational framework
 - Updating the vector of coefficients
 - Updating the precision parameter
- 3 Examples
- 4 Discussion

L₂-distance estimation (L₂E) (Scott, 2001)

Seek a parametric model $f(x | \theta)$ under a minimum distance criterion (minimum integrated square error)

$$\min_{\theta} \int [f(x | \theta) - f(x)]^2 dx \quad (1)$$

$$\begin{aligned} & \int [f(x | \theta) - f(x)]^2 dx \\ = & \int f(x | \theta)^2 dx - 2 \int f(x | \theta)f(x) dx + \int f(x)^2 dx \end{aligned}$$

$$\hat{\theta}_{L_2E} = \operatorname{argmin}_{\theta} \int f(x | \theta)^2 dx - \frac{2}{n} \sum_{i=1}^n f(x_i | \theta) \quad (2)$$

Suppose $X \sim N(\mu, 1)$, then

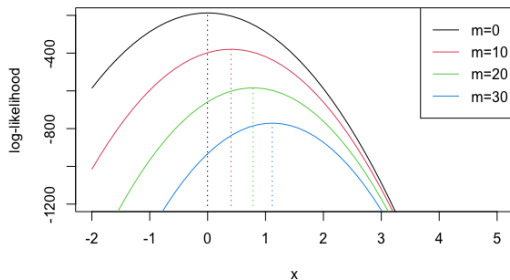
$$\hat{\mu}_{L_2E} = \operatorname{argmin}_{\mu} \frac{1}{2\sqrt{\pi}} - \frac{2}{n} \sum_{i=1}^n f(x_i | \mu)$$

$$\hat{\mu}_{MLE} = \operatorname{argmax}_{\mu} \sum_{i=1}^n \log f(x_i | \mu)$$

- L₂E maximizes the sum of the densities
- MLE maximizes the product of the densities.

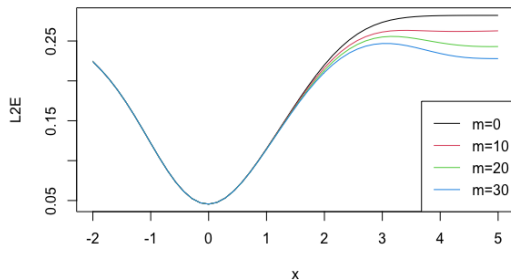
L_2E v.s. MLE

Consider a sample of size 100 from $N(0, 1)$ with m additional data points from a contamination density $N(5, 1)$.



L_2E v.s. MLE

Consider sample of size 100 from $N(0, 1)$ with m additional data points from a contamination density $N(5, 1)$.



Assume a normal model:

- $y_i | X_i = x_i \sim N(x_i^T \beta, \tau^{-2})$
- $\theta = (\beta, \tau)$

$$f(y_i | \beta, \tau) = \frac{\tau}{\sqrt{2\pi}} e^{-\frac{\tau^2 r_i^2}{2}} \quad \text{with } r_i = y_i - \mathbf{x}_i^T \beta$$

L₂E loss:

$$h(\beta, \tau) = \frac{\tau}{2\sqrt{\pi}} - \frac{\tau}{n} \sqrt{\frac{2}{\pi}} \sum_{i=1}^n e^{-\frac{\tau^2 r_i^2}{2}} \quad (3)$$

Structured L_2E regression:

$$\min_{\beta \in \mathbb{R}^p, \tau \in \mathbb{R}_+} h(\beta, \tau), \quad \text{subject to } \beta \in C \quad (4)$$

Examples of C :

- $C = \{\beta \in \mathbb{R}^p : \beta_1 \leq \dots \leq \beta_p\}$ (isotonic regression)
- $C = \{\beta \in \mathbb{R}^p : \|\beta\|_0 \leq k\}$ (sparse regression)

An alternative formulation of (4):

$$\min_{\beta \in \mathbb{R}^p, \tau \in \mathbb{R}_+} h(\beta, \tau) + \psi(\beta), \quad (5)$$

where $\psi(\beta)$ is either the indicator function of C or a non-smooth penalty function such as Lasso.

A computational framework by block descent (Chi and Chi, 2022):

- Update β :

$$\beta^{(k+1)} = \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} h(\beta, \tau^{(k)}) + \psi(\beta)$$

- Update τ :

$$\tau^{(k+1)} = \underset{\tau \in \mathbb{R}^+}{\operatorname{argmin}} h(\beta^{(k+1)}, \tau)$$

Our contributions:

	Chi and Chi (2022)	Our work (Liu et al., 2023)
update β	proximal gradient	(sharp) MM
update τ	proximal gradient	reparameterization & Newton
penalization	convex	distance penalization

Majorization-Minimization (Lange et al., 2000; Lange, 2016)

Goal: Minimize a function $f(x)$

A surrogate function $g(x | \tilde{x})$ majorizes a function $f(x)$ if

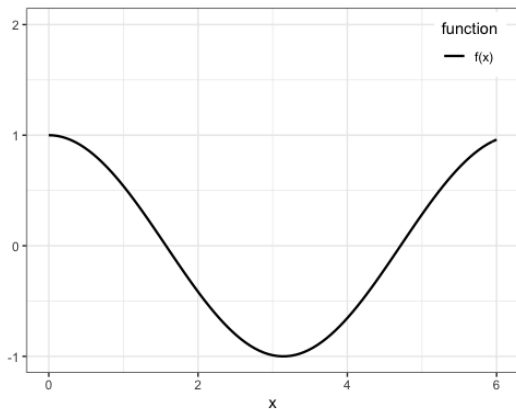
- tangency: $f(\tilde{x}) = g(\tilde{x} | \tilde{x})$
- domination: $f(x) \leq g(x | \tilde{x})$ for all x

The MM iterate:

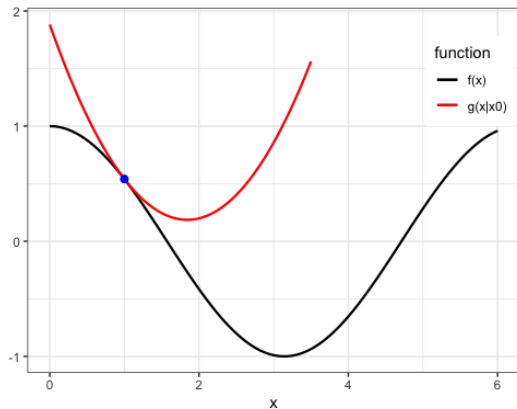
$$x^+ = \underset{x}{\operatorname{argmin}} g(x | \tilde{x})$$

- monotonicity: $f(x^+) \leq g(x^+ | \tilde{x}) \leq g(\tilde{x} | \tilde{x}) = f(\tilde{x})$

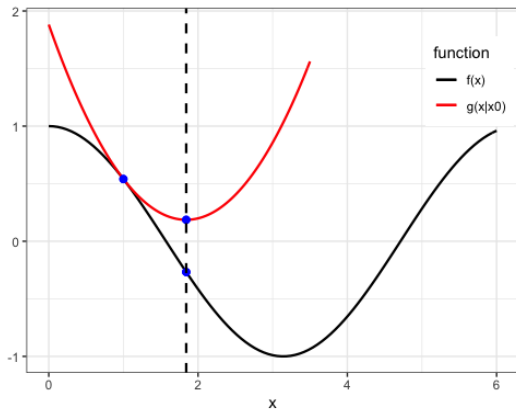
Structured L_2E — Update β



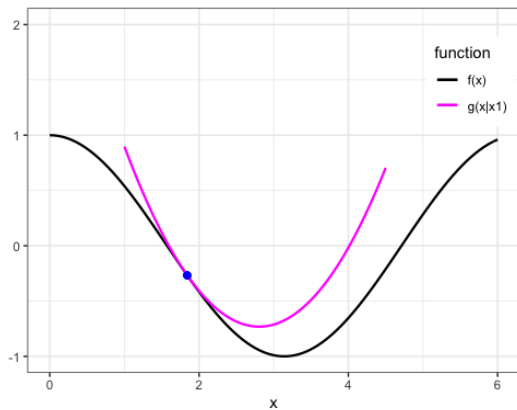
Structured L_2E — Update β



Structured L_2E — Update β

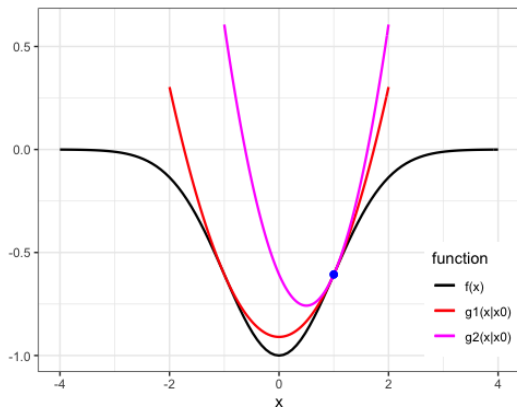


Structured L_2E — Update β



Structured L_2E — Update β

What makes a majorization better than another?



- $g_1(x|x_0) \rightarrow$ a sharp majorization

Structured L_2E — Update β

L_2E loss:

$$h(\beta, \tau) = \frac{\tau}{2\sqrt{\pi}} - \frac{\tau}{n} \sqrt{\frac{2}{\pi}} \sum_{i=1}^n e^{-\frac{\tau^2 r_i^2}{2}}$$

- $f(u) = -\exp(-u)$ is concave

A sharp quadratic univariate majorization w.r.t. r^2 :

$$-\exp\left(-\frac{\tau^2 r^2}{2}\right) \leq -\exp\left(-\frac{\tau^2 \tilde{r}^2}{2}\right) + \frac{\tau^2}{2} \exp\left(-\frac{\tau^2 \tilde{r}^2}{2}\right)(r^2 - \tilde{r}^2)$$

Structured L_2E — Update β

Majorization for the L_2E loss:

$$g(\beta|\tilde{\beta}) = \frac{\tau}{2\sqrt{\pi}} + \frac{\tau^3}{\sqrt{2\pi}n} \sum_{i=1}^n w_i (y_i - x_i^T \beta)^2$$

- $w_i = \exp\left(-\frac{\tau^2(y_i - x_i^T \tilde{\beta})^2}{2}\right)$
- Weights w_i based on residuals from last iterate
- **Effect:** Downweight points as outliers when $\tilde{r}_i = y_i - x_i^T \tilde{\beta}$ is large

Structured L_2E — Update β

Recall for updating β :

$$\operatorname{argmin}_{\beta \in \mathbb{R}^p} h(\beta, \tau) + \psi(\beta)$$

MM iterates for updating β :

$$\beta^+ = \operatorname{argmin}_{\beta} \frac{1}{2} \|\tilde{y} - \tilde{X}\beta\|_2^2 + \psi(\beta)$$

- $\tilde{y} = \sqrt{W}y$, $\tilde{X} = \sqrt{W}X$, W weight matrix
- A penalized least squares problem
- General, simple, and flexible (“plug-and-play”)

Structured L_2E — Update τ

Updating τ :

$$\operatorname{argmin}_{\tau \in \mathbb{R}^+} \frac{\tau}{2\sqrt{\pi}} - \frac{\tau}{n} \sqrt{\frac{2}{\pi}} \sum_{i=1}^n e^{-\frac{\tau^2 r_i^2}{2}}$$

- Reparameterize $\tau = e^\eta \implies$ no constraint on η
- An approximate Newton method

$$\eta_{k+1} = \eta_k - t_k d_k^{-1} \frac{\partial}{\partial \eta} h(\boldsymbol{\beta}, e^{\eta_k}),$$

where $t_k > 0$ is a stepsize parameter chosen via backtracking, and d_k is an approximation of the second derivative $\frac{\partial^2}{\partial \eta^2} h(\boldsymbol{\beta}, e^\eta)$.

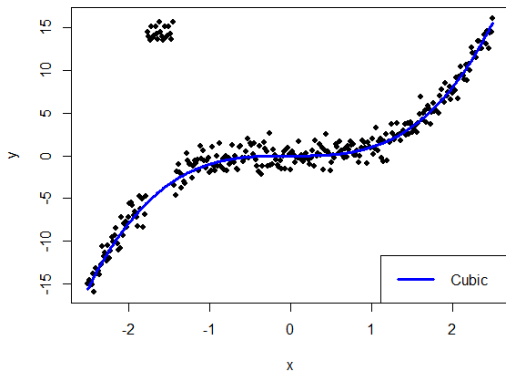
Structured L_2E — computational framework

Algorithm 1 Block descent with MM and approximate Newton

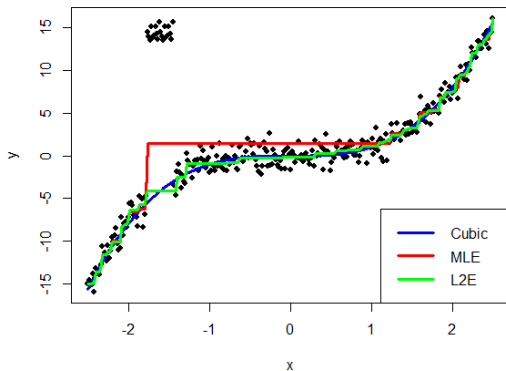
Initialize: $\beta_0 \in \mathbb{R}^p$, $\tau_0 \in \mathbb{R}_+$, N_β , and N_η .

- 1: **for** $k = 1, 2, \dots$ **do**
 - 2: $\beta^+ \leftarrow \beta_{k-1}$
 - 3: **for** $i = 1, \dots, N_\beta$ **do**
 - 4: $\tilde{\mathbf{y}} = \sqrt{\mathbf{W}_+} \mathbf{y}$
 - 5: $\tilde{\mathbf{X}} = \sqrt{\mathbf{W}_+} \mathbf{X}$
 - 6: $\beta^+ = \operatorname{argmin}_{\beta \in \mathbb{R}^p} \frac{1}{2} \|\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\beta\|_2^2 + \lambda\psi(\beta)$ Penalized LS
 - 7: **end for**
 - 8: $\beta_k \leftarrow \beta^+$
 - 9: $\eta^+ \leftarrow \log(\tau_{k-1})$
 - 10: **for** $i = 1, \dots, N_\eta$ **do**
 - 11: $\eta^+ = \eta^+ - t_i d_i^{-1} \frac{\partial}{\partial \eta} h(\beta_k, e^{\eta^+})$ Modified Newton
 - 12: **end for**
 - 13: $\tau_k \leftarrow e^{\eta^+}$
 - 14: **end for**
-

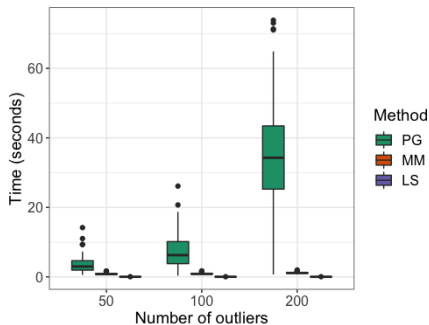
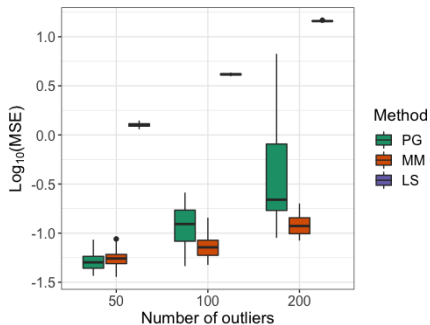
Isotonic regression



Isotonic regression



Isotonic regression



Distance penalization

For a constrained optimization problem

$$\min_{\beta} \ell(\beta) \quad \text{subject to } \beta \in C$$

Distance penalization (Chi et al., 2014; Xu et al., 2017)

$$\psi(\beta) = \frac{1}{2} \text{dist}(\beta, C)^2 = \min_{u \in C} \frac{1}{2} \|\beta - u\|_2^2 \quad (6)$$

The resulting optimization problem:

$$\min_{\beta} \ell(\beta) + \frac{\rho}{2} \text{dist}(\beta, C)^2$$

- If $\rho \rightarrow \infty$, then $\beta \in C$ (recover the constrained solution)
- ρ is assigned a large value in practice

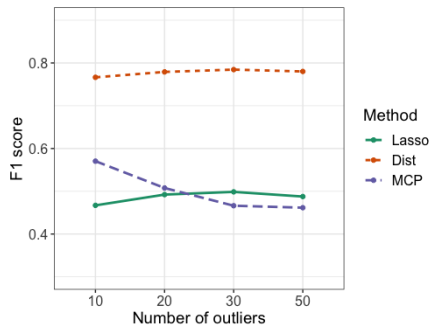
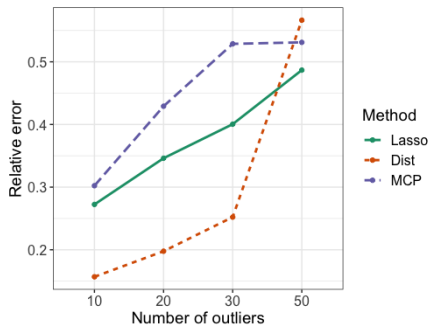
Advantages of distance penalization:

- A general definition
 - diverse structures: sparsity, order constraint, shape constraint
 - multiple constraints: $\frac{1}{2} \sum_{i=1}^l w_i \text{dist}(\beta, C_i)^2$
 - fusion constraint: $L\beta \in C$ (Landeros et al., 2020)
- Only projection onto the constraint set is necessary
 - no requirement that ℓ or C is convex
 - no requirement that ℓ is differentiable
- An efficient proximal distance algorithm (Keys et al., 2019)
 - $\text{dist}(\beta, C)^2 \leq \|\beta - \mathcal{P}_C(\tilde{\beta})\|_2^2$

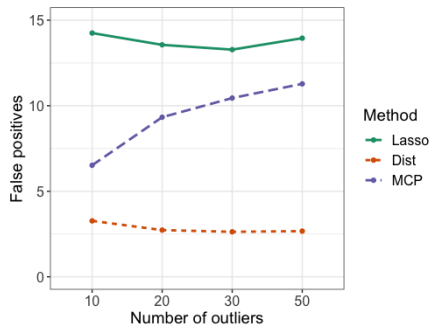
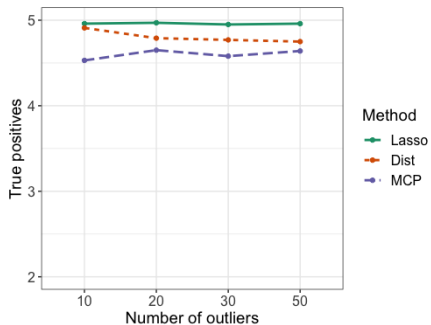
$$y = X\beta + \epsilon$$

- $\beta = (1, 1, 1, 1, 1, 0, \dots, 0)^T \in \mathbb{R}^{50}$
- $X \in \mathbb{R}^{200 \times 50}$ from standard normal distribution
- ϵ standard normal noise
- Shift the first m entries of y and the first m rows of X by 5 to produce outliers

Sparse regression



Sparse regression



Outlier detection

Multivariate regression:

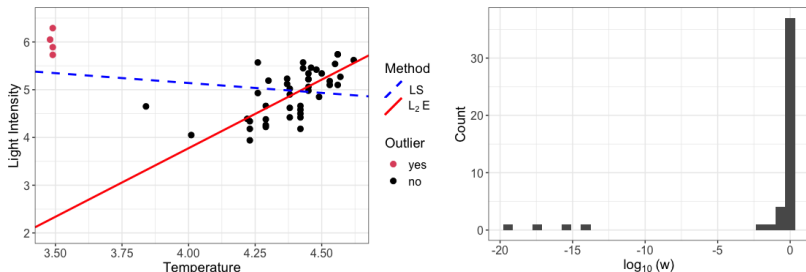


Figure: Fitted regression models from L_2E and LS for the Hertzsprung-Russell Diagram Data (left panel). The four known outliers are successfully detected by the L_2E according to the histogram of the resulting weights (right panel).

Take-home message:

- Structured L_2E regression for
robust estimation + outlier detection + structure recovery
- A sharper computational framework
 - general: various constraints/penalties
 - simple and flexible: “plug-and-play”

Once you have a procedure for solving some structured regression problem, you can use our framework to robustify it!

Thank You!

Reference I

- Chi, E. C., Zhou, H., and Lange, K. (2014). Distance majorization and its applications. *Mathematical Programming*, 146(1):409–436.
- Chi, J. T. and Chi, E. C. (2022). A user-friendly computational framework for robust structured regression with the l_2 criterion. *Journal of Computational and Graphical Statistics*, pages 1–12.
- Keys, K. L., Zhou, H., and Lange, K. (2019). Proximal distance algorithms: Theory and practice. *The Journal of Machine Learning Research*, 20(1):2384–2421.
- Landeros, A., Padilla, O. H. M., Zhou, H., and Lange, K. (2020). Extensions to the proximal distance method of constrained optimization. *arXiv preprint arXiv:2009.00801*.
- Lange, K. (2016). *MM Optimization Algorithms*. SIAM.
- Lange, K., Hunter, D. R., and Yang, I. (2000). Optimization transfer using surrogate objective functions. *Journal of Computational and Graphical Statistics*, 9(1):1–20.

- Liu, X., Chi, E. C., and Lange, K. (2023). A sharper computational tool for regression. *Technometrics*, 65(1):117–126.
- Scott, D. W. (2001). Parametric statistical modeling by minimum integrated square error. *Technometrics*, 43(3):274–285.
- Xu, J., Lange, K., and Chi, E. (2017). Generalized linear model regression under distance-to-set penalties. In *Advances in Neural Information Processing Systems*, pages 1386–1396.