

# Linear Regression

Statistical Machine Learning (ENGG\*6600\*02)

School of Engineering,  
University of Guelph, ON, Canada

Course Instructor: Benyamin Ghogh  
Summer 2023

## Linear Regression

# Dataset

- Consider a dataset  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$  where  $\mathbf{x}_i \in \mathbb{R}^d$ .
- We have some labels too:  $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$  where  $\mathbf{y}_i \in \mathbb{R}^p$ . Usually  $p = 1$  (but not always).
- The labels are not necessarily discrete but can be continuous.
- Example: data can be the data of weather temperature, longitude, latitude, etc, of the city. The label can be the pollution of the city.

# Linear Regression

- We can consider a map  $f(\cdot)$  which maps data to the labels:

$$\mathbf{y} = f(\mathbf{x}) + \epsilon, \quad (1)$$

where  $\epsilon \in \mathbb{R}^p$  is noise.

- In regression, we want to estimate this map  $f$ .
- In linear regression, we want to estimate this map  $f$  by a line (or affine function).
- First, consider the case  $p = 1$ :

$$\mathbb{R} \ni f(\mathbf{x}) = \beta_0 + \sum_{j=1}^d \beta_j x_j \quad (2)$$

where  $x_j$  is the  $j$ -th element of  $\mathbf{x}$  and  $\{\beta_j \in \mathbb{R}\}_{j=0}^d$  are the learnable parameters and  $\beta_0 \in \mathbb{R}$  is specifically for learning the bias (intercept).

# Linear Regression

- One way to do this estimation is to minimize the least squares error between the labels and the estimated model:

$$\underset{\{\beta_j\}_{j=0}^d}{\text{minimize}} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2 \stackrel{(2)}{=} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right)^2, \quad (3)$$

where  $x_{ij}$  is the  $j$ -th element of  $\mathbf{x}_i$ , i.e.,  $\mathbf{x}_i = [x_{i1}, \dots, x_{id}]^\top$ .

- We can write it in matrix form. Let:

$$\mathbb{R}^{d+1} \ni \boldsymbol{\beta} := [\beta_0, \beta_1, \dots, \beta_d]^\top, \quad \mathbb{R}^{n \times (d+1)} \ni \mathbf{X} := \begin{bmatrix} 1, \mathbf{x}_1^\top \\ 1, \mathbf{x}_2^\top \\ \vdots \\ 1, \mathbf{x}_n^\top \end{bmatrix} = \begin{bmatrix} 1, x_{11}, \dots, x_{1d} \\ 1, x_{21}, \dots, x_{2d} \\ \vdots \\ 1, x_{n1}, \dots, x_{nd} \end{bmatrix},$$

$$\mathbb{R}^n \ni \mathbf{y} := [y_0, y_1, \dots, y_n]^\top.$$

The above optimization problem becomes:

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2. \quad (4)$$

# Linear Regression

- The cost function is simplified as:

$$\begin{aligned}\|\mathbf{y} - \mathbf{X}\beta\|_2^2 &= (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) = (\mathbf{y}^\top - \beta^\top \mathbf{X}^\top)(\mathbf{y} - \mathbf{X}\beta) \\ &= \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\beta - \beta^\top \mathbf{X}^\top \mathbf{y} + \beta^\top \mathbf{X}^\top \mathbf{X}\beta.\end{aligned}$$

- Taking derivative of the cost function and setting to zero:

$$\begin{aligned}\frac{\partial}{\partial \beta}(\|\mathbf{y} - \mathbf{X}\beta\|_2^2) &= -\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\beta \stackrel{\text{set}}{=} \mathbf{0} \\ \implies 2\mathbf{X}^\top \mathbf{X}\beta &= 2\mathbf{X}^\top \mathbf{y} \\ \implies \beta &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}\tag{5}$$

- Test phase: The predicted labels for some data  $\mathbf{X}$  is:

$$\mathbf{y} = \mathbf{X}\beta = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.\tag{6}$$

# Linear Regression

- In a more general form, if  $p$  is not necessarily one, we have:

$$\mathbb{R}^{(d+1) \times p} \ni \mathbf{B} := [\beta_0, \beta_1, \dots, \beta_d]^\top, \quad \mathbb{R}^{n \times p} \ni \mathbf{Y} := \begin{bmatrix} \mathbf{y}_1^\top \\ \mathbf{y}_2^\top \\ \vdots \\ \mathbf{y}_n^\top \end{bmatrix},$$

where  $\beta_j \in \mathbb{R}^p$  and  $\mathbf{y}_j \in \mathbb{R}^p$ .

- In this case, the optimization problem becomes:

$$\underset{\mathbf{B}}{\text{minimize}} \quad \|\mathbf{Y} - \mathbf{XB}\|_F^2, \quad (7)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm.

# Linear Regression

- The cost function is simplified as:

$$\begin{aligned}\|Y - XB\|_F^2 &= \text{tr}((Y - XB)^\top (Y - XB)) = \text{tr}((Y^\top - B^\top X^\top)(Y - XB)) \\ &= \text{tr}(Y^\top Y) - \text{tr}(Y^\top XB) - \text{tr}(B^\top X^\top Y) + \text{tr}(B^\top X^\top XB).\end{aligned}$$

- Taking derivative of the cost function and setting to zero:

$$\begin{aligned}\frac{\partial}{\partial B}(\|Y - XB\|_F^2) &= -X^\top Y - X^\top Y + 2X^\top XB \stackrel{\text{set}}{=} 0 \\ \implies 2X^\top XB &= 2X^\top Y \\ \implies B &= (X^\top X)^{-1}X^\top Y.\end{aligned}\tag{8}$$

- Test phase: The predicted labels for some data  $X$  is:

$$Y = XB = X(X^\top X)^{-1}X^\top Y.\tag{9}$$



## Ridge Linear Regression

# Ridge Linear Regression

- We add  $\ell_2$  norm regularization term as a penalty on the size of the learnable parameters.
- Case  $p = 1$ :

$$\underset{\{\beta_j\}_{j=0}^d}{\text{minimize}} \quad \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right)^2 + \lambda \sum_{j=0}^d \beta_j^2, \quad (10)$$

where  $\lambda > 0$  is the regularization parameter.

- We can write it in matrix form. The above optimization problem becomes:

$$\underset{\boldsymbol{\beta}}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2. \quad (11)$$

# Ridge Linear Regression

- The cost function is simplified as:

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_2^2 = \mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + \lambda\boldsymbol{\beta}^\top \boldsymbol{\beta}.$$

- Taking derivative of the cost function and setting to zero:

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\beta}}(\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda\|\boldsymbol{\beta}\|_2^2) &= -\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + 2\lambda\boldsymbol{\beta} \stackrel{\text{set}}{=} \mathbf{0} \\ \implies 2\mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} + 2\lambda\boldsymbol{\beta} &= 2\mathbf{X}^\top \mathbf{y} \implies 2(\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})\boldsymbol{\beta} = 2\mathbf{X}^\top \mathbf{y} \\ \implies \boldsymbol{\beta} &= (\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}\tag{12}$$

- It is strengthening the main diagonal of  $\mathbf{X}^\top \mathbf{X}$  so it makes  $\mathbf{X}^\top \mathbf{X}$  full rank and non-singular. In other words, it makes  $\mathbf{X}^\top \mathbf{X}$  invertible.
- Test phase: The predicted labels for some data  $\mathbf{X}$  is:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda\mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.\tag{13}$$

# Ridge Linear Regression

- More general case where  $p$  is not necessarily one:

$$\underset{\mathbf{B}}{\text{minimize}} \quad \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2. \quad (14)$$

- The cost function is simplified as:

$$\begin{aligned} \|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \|\mathbf{B}\|_F^2 &= \text{tr}((\mathbf{Y} - \mathbf{XB})^\top (\mathbf{Y} - \mathbf{XB})) + \lambda \text{tr}(\mathbf{B}^\top \mathbf{B}) \\ &= \text{tr}((\mathbf{Y}^\top - \mathbf{B}^\top \mathbf{X}^\top)(\mathbf{Y} - \mathbf{XB})) + \lambda \text{tr}(\mathbf{B}^\top \mathbf{B}) \\ &= \text{tr}(\mathbf{Y}^\top \mathbf{Y}) - \text{tr}(\mathbf{Y}^\top \mathbf{XB}) - \text{tr}(\mathbf{B}^\top \mathbf{X}^\top \mathbf{Y}) + \text{tr}(\mathbf{B}^\top \mathbf{X}^\top \mathbf{XB}) + \lambda \text{tr}(\mathbf{B}^\top \mathbf{B}). \end{aligned}$$

- Taking derivative of the cost function and setting to zero:

$$\begin{aligned} \frac{\partial}{\partial \mathbf{B}} (\|\mathbf{Y} - \mathbf{XB}\|_F^2 + \lambda \text{tr}(\mathbf{B}^\top \mathbf{B})) &= -\mathbf{X}^\top \mathbf{Y} - \mathbf{X}^\top \mathbf{Y} + 2\mathbf{X}^\top \mathbf{XB} + 2\lambda \mathbf{B} \stackrel{\text{set}}{=} \mathbf{0} \\ \implies 2\mathbf{X}^\top \mathbf{XB} + 2\lambda \mathbf{B} &= 2\mathbf{X}^\top \mathbf{Y} \implies 2(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})\mathbf{B} = 2\mathbf{X}^\top \mathbf{Y} \\ \implies \mathbf{B} &= (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y}. \end{aligned} \quad (15)$$

- Test phase: The predicted labels for some data  $\mathbf{X}$  is:

$$\mathbf{Y} = \mathbf{XB} = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y}. \quad (16)$$

## **Lasso Linear Regression**

# $\ell_1$ Norm Regularization

- As explained before, sparsity is very useful and effective. If  $\mathbf{x} = [x_1, \dots, x_d]^\top$ , for having sparsity, we should use **subset selection** for the regularization:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \tilde{J}(\mathbf{x}; \theta) := J(\mathbf{x}; \theta) + \alpha \|\mathbf{x}\|_0, \quad (17)$$

where:

$$\|\mathbf{x}\|_0 := \sum_{j=1}^d \mathbb{I}(x_j \neq 0) = \begin{cases} 0 & \text{if } x_j = 0, \\ 1 & \text{if } x_j \neq 0, \end{cases} \quad (18)$$

is “ $\ell_0$ ” norm, which is not a norm (so we used “.” for it) because it does not satisfy the norm properties [1]. The “ $\ell_0$ ” norm counts the number of non-zero elements so when we penalize it, it means that we want to have sparser solutions with many zero entries.

- According to [2], the convex relaxation of “ $\ell_0$ ” norm (subset selection) is  $\ell_1$  norm. Therefore, we write the regularized optimization as:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \tilde{J}(\mathbf{x}; \theta) := J(\mathbf{x}; \theta) + \alpha \|\mathbf{x}\|_1. \quad (19)$$

- The  $\ell_1$  regularization is also referred to as **lasso (least absolute shrinkage and selection operator)** regularization [3].

# Lasso Linear Regression

- We add  $\ell_1$  norm regularization term as a penalty on the size of the learnable parameters.
- Case  $p = 1$ :

$$\underset{\{\beta_j\}_{j=0}^d}{\text{minimize}} \quad \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right)^2 + \lambda \sum_{j=0}^d |\beta_j|, \quad (20)$$

where  $\lambda > 0$  is the regularization parameter and  $|\cdot|$  denotes the absolute value function.

- We can write it in matrix form. The above optimization problem becomes:

$$\underset{\beta}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1. \quad (21)$$

- Let  $\mathbf{x}'_k \in \mathbb{R}^n$  denote the  $k$ -th column of  $\mathbf{X} \in \mathbb{R}^{n \times (d+1)}$ , i.e.,  $\mathbf{X} = [\mathbf{x}'_1, \dots, \mathbf{x}'_{d+1}]$ . The above cost function can be restated as:

$$\underset{\beta=[\beta_1, \dots, \beta_{d+1}]^\top}{\text{minimize}} \quad \left\| \mathbf{y} - \sum_{k=1}^{d+1} \beta_k \mathbf{x}'_k \right\|_2^2 + \lambda \|\beta\|_1. \quad (22)$$

- This cost function can be further restated as below by taking out the  $j$ -th element from the summation:

$$\underset{\beta=[\beta_1, \dots, \beta_{d+1}]^\top}{\text{minimize}} \quad \left\| \mathbf{y} - \sum_{k=1, k \neq j}^{d+1} \beta_k \mathbf{x}'_k - \beta_j \mathbf{x}'_j \right\|_2^2 + \lambda \|\beta\|_1. \quad (23)$$

# Lasso Linear Regression

- We had:

$$\underset{\boldsymbol{\beta}=[\beta_1,\dots,\beta_{d+1}]^\top}{\text{minimize}} \quad \left\| \mathbf{y} - \sum_{k=1, k \neq j}^{d+1} \beta_k \mathbf{x}'_k - \beta_j \mathbf{x}'_j \right\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

- We define:

$$\mathbb{R}^n \ni \mathbf{z} := \mathbf{y} - \sum_{k=1, k \neq j}^{d+1} \beta_k \mathbf{x}'_k. \quad (24)$$

Therefore the cost function becomes:

$$\underset{\boldsymbol{\beta}=[\beta_1,\dots,\beta_{d+1}]^\top}{\text{minimize}} \quad \|\mathbf{z} - \beta_j \mathbf{x}'_j\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1.$$

- We use coordinate descent for optimization. Considering the  $j$ -th element of  $\boldsymbol{\beta}$ , we have:

$$\underset{\beta_j}{\text{minimize}} \quad \|\mathbf{z} - \beta_j \mathbf{x}'_j\|_2^2 + \lambda |\beta_j|,$$

where  $\mathbf{z}$  is a constant with respect to  $\beta_j$ .



# Lasso Linear Regression

- We had:

$$\underset{\beta_j}{\text{minimize}} \quad \|\mathbf{z} - \beta_j \mathbf{x}'_j\|_2^2 + \lambda |\beta_j|.$$

- The cost is simplified as:

$$\begin{aligned} \|\mathbf{z} - \beta_j \mathbf{x}'_j\|_2^2 + \lambda |\beta_j| &= (\mathbf{z} - \beta_j \mathbf{x}'_j)^\top (\mathbf{z} - \beta_j \mathbf{x}'_j) + \lambda |\beta_j| \\ &= (\mathbf{z}^\top - \beta_j \mathbf{x}'_j{}^\top)(\mathbf{z} - \beta_j \mathbf{x}'_j) + \lambda |\beta_j| = \mathbf{z}^\top \mathbf{z} - \beta_j \mathbf{z}^\top \mathbf{x}'_j - \beta_j \mathbf{x}'_j{}^\top \mathbf{z} + \beta_j^2 \mathbf{x}'_j{}^\top \mathbf{x}'_j + \lambda |\beta_j|, \end{aligned}$$

where it is noticed in calculations that  $\beta_j$  is a scalar.

- Taking derivative of the cost function with respect to  $\beta_j$  and setting to zero:

$$\begin{aligned} \frac{\partial}{\partial \beta_j} (\|\mathbf{z} - \beta_j \mathbf{x}'_j\|_2^2 + \lambda |\beta_j|) &= -\mathbf{z}^\top \mathbf{x}'_j - \mathbf{x}'_j{}^\top \mathbf{z} + 2\beta_j \mathbf{x}'_j{}^\top \mathbf{x}'_j + \lambda \text{sign}(\beta_j) \\ &= -\mathbf{x}'_j{}^\top \mathbf{z} - \mathbf{x}'_j{}^\top \mathbf{z} + 2\beta_j \mathbf{x}'_j{}^\top \mathbf{x}'_j + \lambda \text{sign}(\beta_j) = -2\mathbf{x}'_j{}^\top \mathbf{z} + 2\beta_j \mathbf{x}'_j{}^\top \mathbf{x}'_j + \lambda \text{sign}(\beta_j) \stackrel{\text{set}}{=} 0 \\ \Rightarrow \beta_j &= \begin{cases} \frac{\mathbf{x}'_j{}^\top \mathbf{z}}{\mathbf{x}'_j{}^\top \mathbf{x}'_j} - \frac{\lambda}{2} & \text{if } \beta_j \geq 0 \\ \frac{\mathbf{x}'_j{}^\top \mathbf{z}}{\mathbf{x}'_j{}^\top \mathbf{x}'_j} + \frac{\lambda}{2} & \text{if } \beta_j < 0. \end{cases} \end{aligned}$$

# Lasso Linear Regression

- We found:

$$\beta_j = \begin{cases} \frac{\mathbf{x}_j'^\top \mathbf{z}}{\mathbf{x}_j'^\top \mathbf{x}_j'} - \frac{\lambda}{2} & \text{if } \beta_j \geq 0 \\ \frac{\mathbf{x}_j'^\top \mathbf{z}}{\mathbf{x}_j'^\top \mathbf{x}_j'} + \frac{\lambda}{2} & \text{if } \beta_j < 0. \end{cases} \quad (25)$$

- If the cost is  $(1/2)\|\mathbf{z} - \beta_j \mathbf{x}_j'\|_2^2 + \lambda|\beta_j|$ , then the  $(1/2)$  multipliers of  $\lambda$  will go away (if you put  $(1/2)$  multiplied by the norm and calculate the derivative as was done, you will see why):

$$\beta_j = \begin{cases} \frac{\mathbf{x}_j'^\top \mathbf{z}}{\mathbf{x}_j'^\top \mathbf{x}_j'} - \lambda & \text{if } \beta_j \geq 0 \\ \frac{\mathbf{x}_j'^\top \mathbf{z}}{\mathbf{x}_j'^\top \mathbf{x}_j'} + \lambda & \text{if } \beta_j < 0. \end{cases} \quad (26)$$

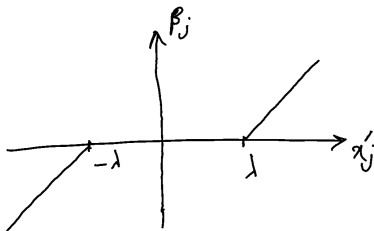
- If the columns of matrix  $\mathbf{X}$  are normalized to have unit length, then  $\mathbf{x}_j'^\top \mathbf{x}_j' = \|\mathbf{x}_j'\|_2^2 = 1$  and it would simplify the solution further to:

$$\beta_j = \begin{cases} \mathbf{x}_j'^\top \mathbf{z} - \lambda & \text{if } \beta_j \geq 0 \\ \mathbf{x}_j'^\top \mathbf{z} + \lambda & \text{if } \beta_j < 0. \end{cases} \quad (27)$$

# Lasso Linear Regression

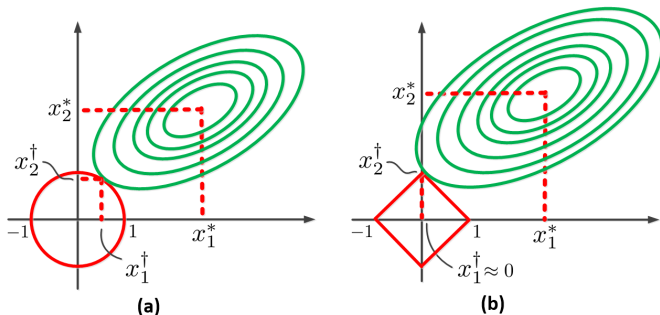
- This is a soft-thresholding function:

$$\beta_j = \begin{cases} \mathbf{x}_j'^\top \mathbf{z} - \lambda & \text{if } \beta_j \geq 0 \\ \mathbf{x}_j'^\top \mathbf{z} + \lambda & \text{if } \beta_j < 0. \end{cases}$$



# Comparison of $\ell_2$ and $\ell_1$ Regularization

- An intuition for why the  $\ell_1$  norm regularization is sparse is illustrated below (credit if Tibshirani - 1996 [3]).
- The objective  $J(\mathbf{x}; \theta)$  has some contour levels like a bowl (if it is convex). The regularization term is also a norm ball, which is a sphere bowl (cone) for  $\ell_2$  norm and a diamond bowl (cone) for  $\ell_1$  norm [1].
- For  $\ell_2$  norm regularization, the objective and the penalty term contact at a point where some of the coordinates might be small; however, for  $\ell_1$  norm, the contact point can be at some point where some variables are exactly zero. This again shows the reason of sparsity in  $\ell_1$  norm regularization.



# Acknowledgment

- For more information on linear regression, see the book: Trevor Hastie, Robert Tibshirani, Jerome H. Friedman, Jerome H. Friedman. "The elements of statistical learning: data mining, inference, and prediction". Vol. 2. New York: springer, 2009 [4].
- Another textbook suitable for sparsity in machine learning is: Robert Tibshirani, Martin Wainwright, Trevor Hastie, "Statistical learning with sparsity: the lasso and generalizations", Chapman and Hall/CRC, 2015 [5].
- Some slides of this slide deck are inspired by teachings of Prof. Mu Zhu at University of Waterloo, Department of Statistics and Prof. Hoda Mohammadzade at Sharif University of Technology, Department of Electrical Engineering.

# References

- [1] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [2] D. L. Donoho, “For most large underdetermined systems of linear equations the minimal  $\ell_1$ -norm solution is also the sparsest solution,” *Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences*, vol. 59, no. 6, pp. 797–829, 2006.
- [3] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [4] T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*, vol. 2. Springer, 2009.
- [5] R. Tibshirani, M. Wainwright, and T. Hastie, *Statistical learning with sparsity: the lasso and generalizations*. Chapman and Hall/CRC, 2015.