

Linear Regression

Statistical Machine Learning (ENGG*6600*02)

School of Engineering,
University of Guelph, ON, Canada

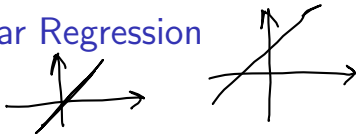
Course Instructor: Benyamin Ghogh
Summer 2023

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



$$x \mapsto y$$

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

$$y = \boxed{f(x)} + \epsilon_i$$

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 $p = 1$:

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

$\beta_0, \beta_1, \dots, \beta_d$

where x_j is the j -th element of 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).

$$\underline{x} = [x_1, \dots, x_d]^T$$

Linear Regression

$$\{(x_i, y_i)\}_{i=1}^n \quad \begin{matrix} \frac{1}{c} \begin{matrix} 1, \dots, 1 \\ x_{11}, x_{12}, \dots, x_{1d} \end{matrix}^T \end{matrix}$$

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

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

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

- We can write it in matrix form. Let:

$$\begin{aligned} \mathbb{R}^{(d+1)} \ni \beta &:= [\beta_0, \beta_1, \dots, \beta_d]^T, & \mathbb{R}^{n \times (d+1)} \ni \mathbf{X} &:= \begin{bmatrix} 1, x_{11}^T \\ 1, x_{21}^T \\ \vdots \\ 1, x_{n1}^T \end{bmatrix} = \begin{bmatrix} 1, x_{11}, \dots, x_{1d} \\ 1, x_{21}, \dots, x_{2d} \\ \vdots \\ 1, x_{n1}, \dots, x_{nd} \end{bmatrix}, \\ \mathbb{R}^n \ni \mathbf{y} &:= [y_1, \dots, y_n]^T. \end{aligned}$$

The above optimization problem becomes:

$$\star \quad \underset{\beta}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{X}\beta\|_2^2 \quad \begin{matrix} \nearrow n \times (d+1) \\ \nearrow (d+1) \times 1 \end{matrix} \quad (4)$$

Linear Regression

$$\|a\|_2^2 = a^T a$$

- The cost function is simplified as:

$$\begin{aligned} \star \quad \|y - X\beta\|_2^2 &= (y - X\beta)^T (y - X\beta) = (y^T - \beta^T X^T)(y - X\beta) \\ &= y^T y - \underbrace{y^T X\beta}_{\text{scalar}} - \underbrace{\beta^T X^T y}_{\text{scalar}} + \beta^T X^T X \beta. \end{aligned}$$

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

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

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

$$y = X\beta = X(X^T X)^{-1} X^T y. \quad (6)$$

Linear Regression

$$\beta \in \begin{bmatrix} \cdot & \cdot & \cdot \\ \vdots & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}_{(d+1) \times p}$$

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

$$\mathbb{R}^{(d+1) \times p} \ni \mathbf{B} := [\beta_0, \beta_1, \dots, \beta_d]^\top,$$

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

$$p=1 \rightarrow \beta_j \in \mathbb{R} \\ \mathbf{y}_j \in \mathbb{R}$$

$$\mathbb{R}^{n \times p} \ni \mathbf{Y} := \begin{bmatrix} \mathbf{y}_1^\top \\ \mathbf{y}_2^\top \\ \vdots \\ \mathbf{y}_n^\top \end{bmatrix},$$

- In this case, the optimization problem becomes:

$$\underset{\mathbf{B}}{\text{minimize}} \quad \|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_F \quad (7)$$

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

Linear Regression

$$\|X\|_F^2 = \text{tr}(X^T X)$$

- The cost function is simplified as:

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

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

$$\begin{aligned} \frac{\partial}{\partial B} (\|Y - XB\|_F^2) &= \underbrace{-X^T Y}_{-2X^T Y} - \underbrace{X^T Y}_{-2X^T Y} + \underbrace{2X^T XB}_{\text{set to } 0} \\ \Rightarrow 2X^T XB &= 2X^T Y \\ \Rightarrow B &= (X^T X)^{-1} X^T Y. \end{aligned} \tag{8}$$

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

$$\star Y = XB = X(X^T X)^{-1} X^T Y; \tag{9}$$

Ridge Linear Regression

Ridge Linear Regression

$$\|\beta\|_2 = \sqrt{\sum_{j=1}^d \beta_j^2}$$

- We add ℓ_2 norm regularization term as a penalty on the size of the learnable parameters.
- Case $p = 1$:

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

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

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

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

$$\begin{bmatrix} \lambda & 0 & 0 & 0 & 0 \\ 0 & \lambda & 0 & 0 & 0 \\ 0 & 0 & \lambda & 0 & 0 \\ 0 & 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & 0 & \lambda \end{bmatrix} \quad \leftarrow (X^T X + \lambda I)^{-1} X^T y \quad \Bigg| \quad (X^T X)^{-1} X^T y$$

Ridge Linear Regression

- The cost function is simplified as:

$$\star \quad \|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 = \underbrace{y^\top y}_{-2X^\top y} - \underbrace{y^\top X\beta}_{-2X^\top y} - \underbrace{\beta^\top X^\top y}_{-2X^\top y} + \underbrace{\beta^\top X^\top X\beta}_{2X^\top X\beta} + \underbrace{\lambda \beta^\top \beta}_{2\lambda\beta}$$

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

$$\begin{aligned} \frac{\partial}{\partial \beta} (\|y - X\beta\|_2^2 + \lambda \|\beta\|_2^2) &= \underbrace{-X^\top y}_{-2X^\top y} - \underbrace{X^\top y}_{-2X^\top y} + \underbrace{2X^\top X\beta}_{2X^\top X\beta} + \underbrace{2\lambda\beta}_{2\lambda\beta} \stackrel{\text{set}}{=} 0 \\ \Rightarrow 2X^\top X\beta + 2\lambda\beta &= 2X^\top y \Rightarrow \cancel{2}(X^\top X + \lambda I)\beta = \cancel{2}X^\top y \\ \Rightarrow \boxed{\beta = (X^\top X + \lambda I)^{-1}X^\top y} \end{aligned} \quad (12)$$

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

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

Ridge Linear Regression

- More general case where p is not necessarily one:

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

- The cost function is simplified as:

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

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

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

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

$$\underbrace{Y = XB = X(X^\top X + \lambda I)^{-1} X^\top Y}_{\text{}} \quad (16)$$

Lasso Linear Regression

ℓ_1 Norm Regularization

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

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

where:

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

is “ ℓ_0 ” norm, which is not a norm (so we used “.” for it) because it does not satisfy the norm properties [1]. The “ ℓ_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 “ ℓ_0 ” norm (subset selection) is ℓ_1 norm. Therefore, we write the regularized optimization as:

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

- The ℓ_1 regularization is also referred to as lasso (least absolute shrinkage and selection operator) regularization [3].

Lasso Linear Regression

- We add ℓ_1 norm regularization term as a penalty on the size of the learnable parameters.
- Case $p = 1$:

$$\star \underset{\{\beta_j\}_{j=0}^d}{\text{minimize}} \underbrace{\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij} \right)^2}_{\text{RSS}} + \underbrace{\lambda \sum_{j=0}^d |\beta_j|}_{\text{penalty}} \rightarrow \|\beta\|_1 \quad (20)$$

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:

$$\star \underset{\beta}{\text{minimize}} \underbrace{\|y - X\beta\|_2^2}_{\text{RSS}} + \underbrace{\lambda \|\beta\|_1}_{\text{penalty}} \quad (21)$$

- 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., $\mathbf{X} = [\mathbf{x}'_1 \cdots \mathbf{x}'_{d+1}]$. The above cost function can be restated as:

$$\underset{\beta = [\beta_0, \dots, \beta_{d+1}]^T}{\text{minimize}} \underbrace{\left\| y - \sum_{k=0}^{d+1} \beta_k \mathbf{x}'_k \right\|_2^2}_{\text{RSS}} + \underbrace{\lambda \|\beta\|_1}_{\text{penalty}} \quad (22)$$

Handwritten notes: $\beta_0 x'_0 + \beta_1 x'_1 + \dots + \beta_{d+1} x'_{d+1}$ and $\beta_j x'_j$ with arrows pointing to the summation terms in the equation.

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

$$\underset{\beta = [\beta_0, \dots, \beta_{d+1}]^T}{\text{minimize}} \underbrace{\left\| y - \sum_{k=0, k \neq j}^{d+1} \beta_k \mathbf{x}'_k - \beta_j \mathbf{x}'_j \right\|_2^2}_{\text{RSS}} + \underbrace{\lambda \|\beta\|_1}_{\text{penalty}} \quad (23)$$

Lasso Linear Regression

- We had:

$$\underset{\beta=[\beta_1, \dots, \beta_{d+1}]^T}{\text{minimize}} \quad \underbrace{\left\| \mathbf{y} - \sum_{k=1, k \neq j}^{d+1} \beta_k \mathbf{x}'_k - \beta_j \mathbf{x}'_j \right\|_2^2}_{\mathcal{E}} + \lambda \|\beta\|_1.$$

- 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:

$$\underset{\beta=[\beta_1, \dots, \beta_{d+1}]^T}{\text{minimize}} \quad \underbrace{\left\| \mathbf{z} - \beta_j \mathbf{x}'_j \right\|_2^2}_{\mathcal{E}} + \underbrace{\lambda \sum_{j=0}^{d+1} |\beta_j|}_{|\beta_0| + |\beta_1| + \dots + |\beta_{d+1}|}.$$

- We use coordinate descent for optimization. Considering the j -th element of β , we have:

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

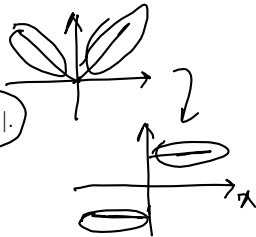
where \mathbf{z} is a constant with respect to β_j .

$$\frac{\partial}{\partial \beta_j} (\beta_j + c)$$

Lasso Linear Regression

- We had:

✖ ✖ $\underset{\beta_j}{\text{minimize}} \quad \underbrace{\|z - \beta_j x'_j\|_2^2}_{\text{RSS}} + \underbrace{\lambda |\beta_j|}_{\text{penalty}}$



- The cost is simplified as:

✖
$$\begin{aligned} \|z - \beta_j x'_j\|_2^2 + \lambda |\beta_j| &= (z - \beta_j x'_j)^T (z - \beta_j x'_j) + \lambda |\beta_j| \\ &= (z^T - \beta_j x_j'^T)(z - \beta_j x'_j) + \lambda |\beta_j| \\ &= \underbrace{z^T z}_{\text{RSS}} - \underbrace{\beta_j z^T x'_j}_{\text{penalty}} - \underbrace{\beta_j x_j'^T z}_{\text{penalty}} + \underbrace{\beta_j^2 x_j'^T x'_j}_{\text{RSS}} + \lambda |\beta_j| \end{aligned}$$

where it is noticed in calculations that β_j is a scalar.

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

✖
$$\begin{aligned} \frac{\partial}{\partial \beta_j} (\|z - \beta_j x'_j\|_2^2 + \lambda |\beta_j|) &= \underbrace{-z^T x'_j - x_j'^T z}_{\text{penalty}} + \underbrace{2\beta_j x_j'^T x'_j}_{\text{RSS}} + \underbrace{\lambda \text{sign}(\beta_j)}_{\text{penalty}} \\ &= \underbrace{-x_j'^T z - x_j'^T z}_{\text{penalty}} + \underbrace{2\beta_j x_j'^T x'_j}_{\text{RSS}} + \lambda \text{sign}(\beta_j) = \underbrace{-2x_j'^T z + 2\beta_j x_j'^T x'_j + \lambda \text{sign}(\beta_j)}_{\text{set to 0}} \\ \Rightarrow \beta_j &= \begin{cases} \frac{x_j'^T z}{x_j'^T x'_j} - \frac{\lambda}{2 x_j'^T x'_j} & \text{if } \beta_j \geq 0 \\ \frac{x_j'^T z}{x_j'^T x'_j} + \frac{\lambda}{2 x_j'^T x'_j} & \text{if } \beta_j < 0 \end{cases} \end{aligned}$$

Handwritten notes and simplifications:

- $-2x_j'^T z + 2\beta_j x_j'^T x'_j + \lambda \text{sign}(\beta_j) = 0$
- $2\beta_j x_j'^T x'_j = 2x_j'^T z - \lambda$
- $\beta_j = \frac{x_j'^T z}{x_j'^T x'_j} - \frac{\lambda}{2 x_j'^T x'_j}$

Lasso Linear Regression

$$\|x_j'\|_2 \leq 1 \rightarrow \|x_j'\|_2^2 = 1 \quad \hookrightarrow \quad x_j'^T x_j' = 1$$

- We found:

$$\beta_j = \begin{cases} \frac{x_j'^T z}{x_j'^T x_j'} - \frac{\lambda}{2x_j'^T x_j'} & \text{if } \beta_j \geq 0 \\ \frac{x_j'^T z}{x_j'^T x_j'} + \frac{\lambda}{2x_j'^T x_j'} & \text{if } \beta_j < 0. \end{cases} \quad (25)$$

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

$$\beta_j = \begin{cases} \frac{x_j'^T z}{x_j'^T x_j'} - \lambda & \text{if } \beta_j \geq 0 \\ \frac{x_j'^T z}{x_j'^T x_j'} + \lambda & \text{if } \beta_j < 0. \end{cases} \quad (26)$$

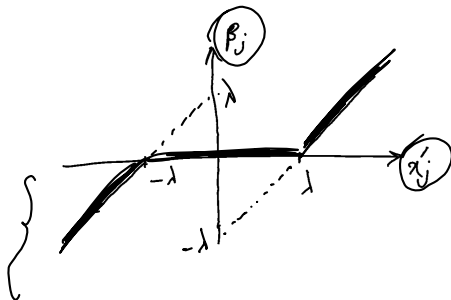
- If the columns of matrix \mathbf{X} are normalized to have unit length, then $x_j'^T x_j' = \|x_j'\|_2^2 = 1$ and it would simplify the solution further to:

$$\star \beta_j = \begin{cases} x_j'^T z - \lambda & \text{if } \beta_j \geq 0 \\ x_j'^T z + \lambda & \text{if } \beta_j < 0. \end{cases} \quad (27)$$

Lasso Linear Regression

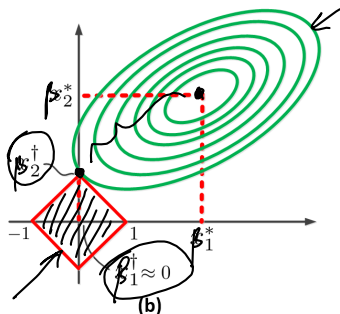
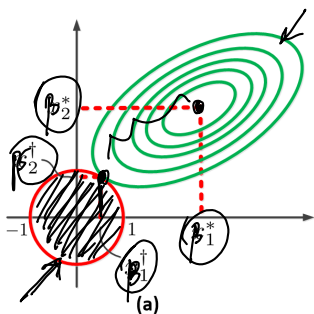
- This is a soft-thresholding function:

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



Comparison of ℓ_2 and ℓ_1 Regularization

- An intuition for why the ℓ_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 ℓ_2 norm and a diamond bowl (cone) for ℓ_1 norm [1].
- For ℓ_2 norm regularization, the objective and the penalty term contact at a point where some of the coordinates might be small; however, for ℓ_1 norm, the contact point can be at some point where some variables are exactly zero. This again shows the reason of sparsity in ℓ_1 norm regularization.



Acknowledgment

- For more information on linear regression, see the book: Trevor Hastie, Robert Tibshirani, Jerome H. Friedman, Jérôme 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 ℓ_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.