Statistical Machine Learning (ENGG*6600*08)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Fall 2023

Dataset

- Consider a dataset $\{x_1, x_2, \ldots, x_n\}$ where $x_i \in \mathbb{R}^d$.
- We have some labels too: $\{y_1, y_2, \dots, y_n\}$ where $y_i \in \mathbb{R}^p$. 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.

• We can consider a map f(.) which maps data to the labels:

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

where $\boldsymbol{\epsilon} \in \mathbb{R}^p$ 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).

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

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).

 One way to do this estimation is to <u>minimize the least squares error</u> between the labels and the estimated model:

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

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

• We can write it in matrix form. Let:



• The cost function is simplified as:

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{2}^{2} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^{\top} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{y}^{\top} - \boldsymbol{\beta}^{\top} \mathbf{X}^{\top}) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$
$$= \underbrace{\mathbf{y}^{\top} \mathbf{y}^{\top}}_{\mathbf{y} (\mathbf{z} + \mathbf{i})} \underbrace{\mathbf{y}^{\top} \mathbf{X} \boldsymbol{\beta}}_{\mathbf{z} (\mathbf{z} + \mathbf{i}) \times \mathbf{i}} + \underbrace{\boldsymbol{\beta}^{\top} \mathbf{X}^{\top} \mathbf{X} \boldsymbol{\beta}}_{\mathbf{z} (\mathbf{z} + \mathbf{i}) \times \mathbf{i}}$$

Taking derivative of the cost function and setting to zero:

$$\frac{\partial}{\partial \beta} (\|\boldsymbol{y} - \boldsymbol{X}\beta\|_{2}^{2}) = -\boldsymbol{X}^{\top}\boldsymbol{y} - \boldsymbol{X}^{\top}\boldsymbol{y} + 2\boldsymbol{X}^{\top}\boldsymbol{X}\beta \stackrel{\text{set}}{=} \boldsymbol{0}$$

$$\implies \underbrace{2\boldsymbol{X}^{\top}\boldsymbol{X}\beta}_{\beta} = \underbrace{2\boldsymbol{X}^{\top}\boldsymbol{y}}_{\beta} \stackrel{\text{set}}{=} \underbrace{\boldsymbol{0}}_{\beta} \stackrel{\text{set}}{=} \underbrace{\boldsymbol{0}$$

• Test phase: The predicted labels for some data (X) is:

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

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

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

where $\beta_j \in \mathbb{R}^p$ and $y_j \in \mathbb{R}^p$. In this case, the optimization problem becomes: $n \not A p$ minimize $\| y - \dot{x} B \|_{F}^p$ where $\|.\|_F$ denotes the Frobenius norm. $d \not A | J \not P$ $n \chi (d \not A)$ $n \not A p$

(7)



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

$$\frac{\partial}{\partial B} (\|\mathbf{Y} - \mathbf{X}B\|_{F}^{2}) = -\mathbf{X}^{\top} \mathbf{Y} - \mathbf{X}^{\top} \mathbf{Y} + 2\mathbf{X}^{\top} \mathbf{X}B \stackrel{\text{set}}{=} \mathbf{0}$$

$$\implies \mathbf{X} \stackrel{\mathsf{T}}{=} \mathbf{X}B = \mathbf{X}^{\top} \mathbf{Y}$$

$$\implies \mathbf{B} = (\mathbf{X}^{\top} \mathbf{X})^{-1} \mathbf{X}^{\top} \mathbf{Y}$$
(8)

• Test phase: The predicted labels for some data (X) is:

$$\overbrace{\boldsymbol{Y}}=\overbrace{\boldsymbol{X}} = \boldsymbol{X}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{Y}.$$
(9)



 $X \cup U_1 \cup U_2 \cup U_3 \rightarrow U_3 = \left[\left\| \hat{y} - y \right\|_F^2 \right]$ $\begin{array}{c} \overbrace{V} = & \overbrace{V}_{1} & \overbrace{V}_{2} & \overbrace{V}_{3} \\ \overbrace{V} = & \overbrace{V}_{1} & \overbrace{V}_{2} & \overbrace{V}_{3} \\ \end{array} = & X & \overbrace{V}_{1} & \bigvee_{2} & \overbrace{V}_{3} \\ \end{array}$ $X f_1(V_1)(f_2)(V_2) f_3(V_3)$



Ridge Linear Regression

• We add ℓ_2 norm regularization term as a penalty on the size of the learnable parameters.

• Case
$$p = 1$$
:
minimize $\sum_{\{\beta_j\}_{j=0}^d}^n (y_i - \beta_0 - \sum_{j=1}^d \beta_j x_{ij})^2 + (\lambda) \sum_{j=0}^d \beta_j^2,$ (10)

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

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

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

Ridge Linear Regression

• The cost function is simplified as:

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_{2}^{2} + \left(\lambda \|\boldsymbol{\beta}\|_{2}^{2}\right) = \mathbf{y}^{\top}\mathbf{y} - \mathbf{y}^{\top}\mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^{\top}\mathbf{X}^{\top}\mathbf{Y} + \boldsymbol{\beta}^{\top}\mathbf{X}^{\top}\mathbf{X}\boldsymbol{\beta} + \left(\lambda\boldsymbol{\beta}^{\top}\boldsymbol{\beta}\right) \mathbf{y}^{\top}\mathbf{y} - \mathbf{y}^{\top}\mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^{\top}\mathbf{y}^{\top}\mathbf{y} + \boldsymbol{\beta}^{\top}\mathbf{y}^{\top}\mathbf$$

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

$$\frac{\partial}{\partial \beta} (\|\mathbf{y} - \mathbf{X}\beta\|_{2}^{2} + \|\beta\|_{2}^{2}) = -\mathbf{X}^{\top}\mathbf{y} - \mathbf{X}^{\top}\mathbf{y} + 2\mathbf{X}^{\top}\mathbf{X}\beta + 2\lambda\beta \stackrel{\text{set}}{=} \mathbf{0}$$

$$\implies (2\mathbf{x}^{\top}\mathbf{X}\beta + 2)\beta = 2\mathbf{X}^{\top}\mathbf{y} \implies \beta(\mathbf{X}^{\top}\mathbf{X} + \lambda\mathbf{I})\beta = 2\mathbf{X}^{\top}\mathbf{y}$$

$$\implies \beta = (\mathbf{X}^{\top}\mathbf{X}+\lambda\mathbf{I})^{-1}\mathbf{X}^{\top}\mathbf{y}.$$
(12)

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

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

Ridge Linear Regression

• More general case where *p* is not necessarily one:

$$\underset{\boldsymbol{B}}{\text{minimize}} \quad \|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{B}\|_{F}^{2} + \lambda \|\boldsymbol{B}\|_{F}^{2}. \tag{14}$$

The cost function is simplified as:

$$\|\mathbf{Y} - \mathbf{X}\mathbf{B}\|_{F}^{2} + \lambda \|\mathbf{B}\|_{F}^{2} = \operatorname{tr}((\mathbf{Y} - \mathbf{X}\mathbf{B})^{\top}(\mathbf{Y} - \mathbf{X}\mathbf{B})) + \lambda \operatorname{tr}(\mathbf{B}^{\top}\mathbf{B})$$

= $\operatorname{tr}((\mathbf{Y}^{\top} - \mathbf{B}^{\top}\mathbf{X}^{\top})(\mathbf{Y} - \mathbf{X}\mathbf{B})) + \lambda \operatorname{tr}(\mathbf{B}^{\top}\mathbf{B})$
= $\operatorname{tr}(\mathbf{Y}^{\top}\mathbf{Y}) - \operatorname{tr}(\mathbf{Y}^{\top}\mathbf{X}\mathbf{B}) - \operatorname{tr}(\mathbf{B}^{\top}\mathbf{X}^{\top}\mathbf{Y}) + \operatorname{tr}(\mathbf{B}^{\top}\mathbf{X}^{\top}\mathbf{X}\mathbf{B}) + \lambda \operatorname{tr}(\mathbf{B}^{\top}\mathbf{B}).$

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

$$\frac{\partial}{\partial B}(\|\mathbf{Y} - \mathbf{X}B\|_{F}^{2} + \lambda \operatorname{tr}(B^{\top}B)) = \underbrace{-\mathbf{X}^{\top}\mathbf{Y} - \mathbf{X}^{\top}\mathbf{Y} + 2\mathbf{X}^{\top}\mathbf{X}B}_{B} + 2\lambda B \stackrel{\text{set}}{=} 0$$

$$\implies \underbrace{2\mathbf{X}^{\top}\mathbf{X}B}_{B} + 2\lambda B \stackrel{\text{set}}{=} \underbrace{2\mathbf{X}^{\top}\mathbf{Y}}_{B} \Longrightarrow \underbrace{2\mathbf{X}^{\top}\mathbf{X} + \lambda I}_{B} \stackrel{\text{set}}{=} \underbrace{2\mathbf{X}^{\top}\mathbf{X} + \lambda I}_{B} \stackrel{\text{s$$

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

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

Linear Regression

ℓ_1 Norm Regularization

As explained before, sparsity is very useful and effective. If x = [x₁,...,x_d][⊤], for having sparsity, we should use subset selection for the regularization:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \widetilde{J}(\mathbf{x};\theta) := J(\mathbf{x};\theta) + \alpha ||\mathbf{x}||_{0}, \tag{17}$$

where:

$$||\mathbf{x}||_{0} := \sum_{j=1}^{d} \mathbb{I}(x_{j} \neq 0) = \begin{cases} 0 & \text{if } x_{j} = 0, \\ 1 & \text{if } x_{j} \neq 0, \end{cases}$$
(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 "l₀" norm (subset selection) is l₁ norm. Therefore, we write the regularized optimization as:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \widetilde{J}(\mathbf{x}; \theta) := J(\mathbf{x}; \theta) + \alpha ||\mathbf{x}||_{1}. \tag{19}$$

The l₁ regularization is also referred to as lasso (least absolute shrinkage and selection operator) regularization [3].

$\beta_1 \chi_1 + (\beta_2 \chi_2) + \dots + \beta_{j+1} \chi_{j+1}$ Lasso Linear Regression

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

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

where $\lambda > 0$ is the regularization parameter and |.| denotes the absolute value function. • We can write it in matrix form. The above optimization problem becomes:

$$\underbrace{\qquad \qquad \text{minimize} \quad \| \mathbf{\beta} - \mathbf{x} \mathbf{\beta} \|_2^2 + \lambda \| \mathbf{\beta} \|_1. \quad \mathbf{a} \qquad (21)$$

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

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

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

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

B1 x1 + B2 x2 + ... + (B: xj) + ... + B_{J+1} x'_{J+1} $= \left(\left| \frac{\beta_{1}}{\beta_{1}} + \dots + \frac{\beta_{n}}{\beta_{n}} + \frac{\beta_{n}}{\beta_{n}} + \dots + \frac{\beta_{n}}{\beta_{n}} + \frac{\beta_{n}}{\beta_{$ $+\left(\beta_{j}, \beta_{j}\right)$

• We had:

ath = bta

We had:

1

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

• The cost is simplified as:

$$|\widehat{\boldsymbol{z}} - \widehat{\beta_j} \widehat{\boldsymbol{x}}_j'||_2^2 + \widehat{\lambda}|\widehat{\beta_j}| = (\widehat{\boldsymbol{z}} - \widehat{\beta_j} \widehat{\boldsymbol{x}}_j')^{\text{T}} (\widehat{\boldsymbol{z}} - \widehat{\beta_j} \widehat{\boldsymbol{x}}_j') + \widehat{\lambda}|\widehat{\beta_j}|$$

= $(\widehat{\boldsymbol{z}}^{\text{T}} - \widehat{\beta_j} \widehat{\boldsymbol{x}}_j'^{\text{T}}) (\widehat{\boldsymbol{z}} - \widehat{\beta_j} \widehat{\boldsymbol{x}}_j') + \widehat{\lambda}|\widehat{\beta_j}| = \underbrace{\widehat{\boldsymbol{z}}^{\text{T}} \widehat{\boldsymbol{z}} - \widehat{\beta_j} \widehat{\boldsymbol{z}}^{\text{T}} \widehat{\boldsymbol{x}}_j' - \widehat{\beta_j} \widehat{\boldsymbol{x}}_j'^{\text{T}} \widehat{\boldsymbol{z}} + \widehat{\beta_j}^2 \widehat{\boldsymbol{x}}_j'^{\text{T}} \widehat{\boldsymbol{x}}_j' + \widehat{\lambda}|\widehat{\beta_j}|,$

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

• Taking derivative of the cost function with respect to β_i and setting to zero:

$$\begin{array}{l} \partial \\ \partial \beta_{j} \\ (||z - \beta_{j} \mathbf{x}_{j}'||_{2}^{2} + \lambda |\beta_{j}|) = \underbrace{\left[\mathbf{z}^{\top} \mathbf{x}_{j}' \right]}_{\mathbf{x}_{j}' \top \mathbf{z}} \mathbf{x}_{j}' \mathbf{z}_{j}' \mathbf{z}_{j} + 2\beta_{j} \mathbf{x}_{j}'^{\top} \mathbf{x}_{j}' + \lambda \operatorname{sign}(\beta_{j}) \\ = \underbrace{\left[\mathbf{x}_{j}'^{\top} \mathbf{z} \right]}_{\mathbf{x}_{j}' \top \mathbf{z}} \mathbf{x}_{j}' \mathbf{z}_{j} + 2\beta_{j} \mathbf{x}_{j}'^{\top} \mathbf{x}_{j}' + \lambda \operatorname{sign}(\beta_{j}) = \underbrace{\left[-2\mathbf{x}_{j}'^{\top} \mathbf{z} + 2\beta_{j} \mathbf{x}_{j}'^{\top} \mathbf{x}_{j}' + \lambda \operatorname{sign}(\beta_{j}) \right]}_{\mathbf{z}' \top \mathbf{z}' - \mathbf{z}'} \mathbf{z}_{j}' \mathbf{z}_{j}$$

We found:

If the cost is (1/2) ||z - β_jx'_j||²₂ + λ|β_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):

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

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

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

• This is a soft-thresholding function:



Comparison of ℓ_2 and ℓ_1 Regularization

- An intuition for why the ℓ₁ norm regularization is sparse is illustrated below (credit if Tibshirani - 1996 [3]).
- The objective J(x; θ) 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 l₂ norm and a diamond bowl (cone) for l₁ 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, 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, "<u>Statistical learning with sparsity</u>: the lasso and generalizations", Chapman and Hall/CRC, 2015 [5].
- Some slides of this slide deck are inspired by teachings of <u>Prof. Mu Zhu</u> at University of Waterloo, Department of Statistics and Prof. <u>Hoda Mohammadzade</u> 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 l₁-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.