

1 Introduction

Given an example $x = \{x^1, x^2, \dots, x^d\}$, where x^d is the d-th feature, the linear model is defined as

$$f(x) = w^1 x^1 + w^2 x^2 + \dots + w^d x^d + b. \quad (1)$$

Its' vectorized formulation can be written:

$$f(x) = \mathbf{w}^T X + b, \quad (2)$$

where $\mathbf{w} = \{w^1; w^2; \dots; w^d\}$.

2 Linear Regression

Problem Def. Given the dataset $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, where $x_n = \{x_n^1, x_n^2, \dots, x_n^d\}$ and $y \in \mathbb{R}$, the linear regression is to learn a model to predict the value y .

For convenience, we consider $d=1$ such that $x_n \in \mathbb{R}$. With Eq. 1, the linear regression can be constructed as:

$$f(x_n) = wx_n + b \text{ s.t. } f(x_n) = y_n. \quad (3)$$

How to calculate w and b ? The problem lies in the measure between the predict and the true value. The mean squared error (MSE) is commonly used in the regression task:

$$\begin{aligned} Error &= \sum_i^n (f(x_n) - y_n)^2, \\ Error &= \sum_i^n (wx_n + b - y_n)^2. \end{aligned} \quad (4)$$

Least squares method. ¹

$$\begin{aligned} \frac{\partial E}{\partial w} &= 2(w \sum_i^n x_i^2 - \sum_i^n (y_i - b)x_i), \\ \frac{\partial E}{\partial b} &= 2(nb - \sum_i^n (y_i - wx_i)). \end{aligned} \quad (5)$$

The closed-form:

$$\begin{aligned} w &= \frac{\sum_i^n y_i (x_i - \bar{x})}{\sum_i^n x_i^2 - 1/m(\sum_i^n x_i)^2}, \\ b &= 1/m \sum_i^n (y_i - wx_i). \end{aligned} \quad (6)$$

Generalizing to the multivariate linear regression, in which $x_n = \{x_n^d, \dots, x_n^d, 1\}$, the error can be written as:

$$\begin{aligned} Error &= |\mathbf{w}X - Y|^2 \\ &= (\mathbf{w}^T X - Y)^T (\mathbf{w}X - Y) \\ &= (X^T \mathbf{w}^T - Y^T) (\mathbf{w}X - Y) \\ &= X^T \mathbf{w}^T \mathbf{w}X - X^T \mathbf{w}^T Y - Y^T \mathbf{w}X + Y^T Y. \end{aligned} \quad (7)$$

The derivative:

$$\frac{\partial E}{\partial w} = 2X^T (X\mathbf{w} - Y). \quad (8)$$

The closed-form:

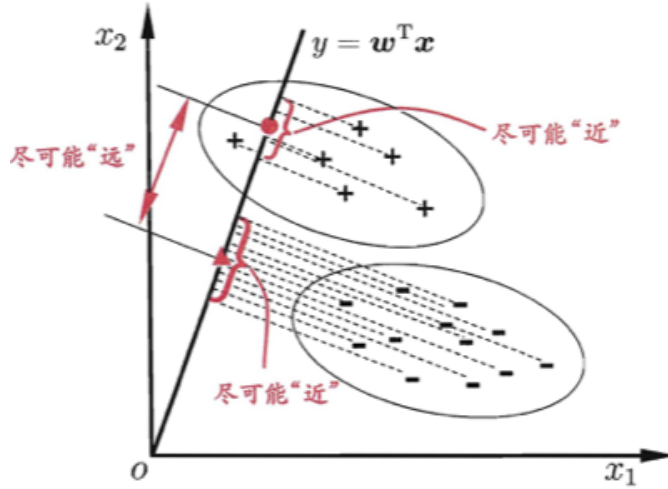
$$\mathbf{w} = (X^T X)^{-1} X^T Y, \quad (9)$$

where $(X^T X)$ has to be a full-rank matrix.

¹https://colab.research.google.com/drive/1h_1gtTb86X4KKrmBbgw13mMOT-3j62jf#scrollTo=jYkD1fYw85w7

3 Linear Model for Classification

Problem Def. Given the dataset $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, where $x_n = \{x_n^1, x_n^2, \dots, x_n^d\}$ and $y = \{0, \dots, C\}$, the classification model is to predict the class $y_n = C_i$.



We transform the linear regression problem to the classification as the figure shown. This is also called Linear Discriminant Analysis (LDA), which assumes that each class come from an individual Gaussian distribution. The objective can be formulated as minimizing the cross-distance among intra-class projection while maximizing the inter-classes projection:

$$Error = -\frac{\|\mu_0 - \mu_1\|^2}{\Sigma_0 + \Sigma_1}, \quad (10)$$

where μ denotes the mean and Σ is the variance.

Formally, the variance is defined as within-class scatter matrix:

$$S_w = \Sigma_0 + \Sigma_1 = \sum_0 (x - \mu_0)(x - \mu_0)^T + \sum_1 (x - \mu_1)(x - \mu_1)^T \quad (11)$$

The difference between class is between-class scatter matrix:

$$S_b = (\mu_0 - \mu_1)(\mu_0 - \mu_1)^T. \quad (12)$$

After mapping to the line, the objective used for optimization is:

$$J = \frac{\mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}}. \quad (13)$$

The derivative:

$$\frac{\partial J}{\partial \mathbf{w}} = \frac{2S_b \mathbf{w} \times \mathbf{w}^T S_w \mathbf{w} - 2S_w \mathbf{w} \times \mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}} = 0 \quad (14)$$

The closed-form:

$$\begin{aligned} S_b \mathbf{w} \times \mathbf{w}^T S_w \mathbf{w} &= S_w \mathbf{w} \times \mathbf{w}^T S_b \mathbf{w} \\ S_b \mathbf{w} &= S_w \mathbf{w} \left(\frac{\mathbf{w}^T S_b \mathbf{w}}{\mathbf{w}^T S_w \mathbf{w}} \right) \\ (S_w^{-1} S_b) \mathbf{w} &= \lambda \mathbf{w} \end{aligned} \quad (15)$$

If $(S_w^{-1} S_b)$ is non-singular matrix, \mathbf{w} is the Eigenvector of $(S_w^{-1} S_b)$.

4 Recommended Reading

- [1] https://en.m.wikipedia.org/wiki/Linear_discriminant_analysis.
- [2] https://en.m.wikipedia.org/wiki/Logistic_regression
- [3] <https://stats.stackexchange.com/questions/95247/logistic-regression-vs-lda-as-two-class-classifiers>
- [4] https://en.m.wikipedia.org/wiki/Generalized_linear_model