
CS189/289A SPRING 2023: LINEAR REGRESSION

LECTURE NOTE

Chawin Sitawarin *

University of California, Berkeley

Demo: <https://colab.research.google.com/drive/1Pv1BUScUd6fHm2s2pQ9oLPHwg60QSYi8>

1 Optimization View of Linear Regression

One very nice thing about linear regression is that it can be tackled or studied from various perspectives and domains. One can simply look at linear regression as a standalone optimization problem. Specifically, given some matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and vector $\mathbf{y} \in \mathbb{R}^n$, we wish to find the parameters $\mathbf{w} \in \mathbb{R}^d$ that minimize the loss function $loss(\mathbf{w}; \mathbf{X}, \mathbf{y})$ which is chosen to be the mean (or sum) of squared errors (MSE):

$$\mathbf{w}^* = \arg \min_{\mathbf{w} \in \mathbb{R}^d} loss(\mathbf{w}; \mathbf{X}, \mathbf{y}) \quad (1)$$

$$\text{where } L(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \mathbf{w})^2 = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2. \quad (2)$$

A simple interpretation is that we want to best predict or “reconstruct” any variable y_i given only its feature vector \mathbf{x}_i by making an assumption that their relationship is linear.

2 Probabilistic View of Linear Regression

A different way to motivate linear regression is as a probabilistic model for the data. Here, we assume that the data is generated by a linear model with additive Gaussian noise. Or more specifically, we assume that the data is generated by the following model:

$$p(y \mid \mathbf{x}, \mathbf{w}) = \mathcal{N}(y \mid \mathbf{x}^\top \mathbf{w}, \sigma^2), \quad (3)$$

or in other words,

$$y = \mathbf{x}^\top \mathbf{w} + \epsilon \quad \text{where } \epsilon \sim \mathcal{N}(0, \sigma^2) \quad (4)$$

- This view shows another example of MLE.
- This view will be useful when we study extensions of linear regression such as ridge regression and LASSO, particularly as a solution of the Maximum A Posteriori (MAP) problem.

In this setup, we can view linear regression as a way to estimate the parameters \mathbf{w} and σ^2 from the data via *Maximum Likelihood Estimation* (MLE).

*Corresponding email: chawins@berkeley.edu

We write the Likelihood function as

$$L(\mathbf{w}, \sigma \mid \mathbf{X}, \mathbf{y}) = \prod_{i=1}^n p(y_i \mid \mathbf{x}_i, \mathbf{w}, \sigma^2) \quad (5)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^\top \mathbf{w})^2\right) \quad (6)$$

$$= (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2\right) \quad (7)$$

We can maximize the likelihood directly, but more often than not, it is more convenient to minimize the *Negative Log-Likelihood* (NLL) instead. This term is frequently used in machine learning as well as most of the deep learning models.

$$\arg \max_{\mathbf{w}, \sigma} L(\mathbf{w}, \sigma \mid \mathbf{X}, \mathbf{y}) = \arg \max_{\mathbf{w}, \sigma} \log L(\mathbf{w}, \sigma \mid \mathbf{X}, \mathbf{y}) \quad (\text{Log-Likelihood}) \quad (8)$$

$$= -\arg \min_{\mathbf{w}, \sigma} \log L(\mathbf{w}, \sigma \mid \mathbf{X}, \mathbf{y}) \quad (\text{Negative Log-Likelihood}) \quad (9)$$

$$= \arg \min_{\mathbf{w}, \sigma} \frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \frac{n}{2} \log(2\pi\sigma^2) \quad (10)$$

Most of the times, we focus more on finding \mathbf{w}^* and ignore σ^* . To find \mathbf{w}^* we set the gradient to zero, i.e.,

$$\nabla_{\mathbf{w}} L(\mathbf{w}^*, \sigma \mid \mathbf{X}, \mathbf{y}) = 0 \quad (11)$$

$$\frac{\partial}{\partial \mathbf{w}} \left(\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{X}\mathbf{w}^*\|_2^2 + \frac{n}{2} \log(2\pi\sigma^2) \right) = 0 \quad (12)$$

$$\frac{1}{\sigma^2} \mathbf{X}^\top (\mathbf{X}\mathbf{w}^* - \mathbf{y}) = 0 \quad (13)$$

$$\mathbf{X}^\top \mathbf{y} - \mathbf{X}^\top \mathbf{X}\mathbf{w}^* = 0 \quad (14)$$

$$\mathbf{X}^\top \mathbf{X}\mathbf{w}^* = \mathbf{X}^\top \mathbf{y} \quad (15)$$

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (16)$$

We call the final MSE produced by \mathbf{w} *Residual Sum of Squares* (RSS).

$$RSS(\mathbf{w}) = \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 = \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \quad (17)$$

Exercise 1

Show that the MLE solution for σ (i.e., σ^*) is given by $\sigma^* = \frac{1}{n} \|\mathbf{y} - \mathbf{X}\mathbf{w}^*\|_2^2$.

2.1 Generative Model View of Linear Regression

Instead of viewing linear regression as a way to estimate $p(y \mid \mathbf{x})$ using the parameters \mathbf{w} and σ^2 , we can also view it as a “generative model” for the joint distribution $p(y, \mathbf{x})$ instead. In this case, we make an assumption that $p(y, \mathbf{x})$ is *jointly Gaussian*, or precisely, we assume that they are generated from the following distribution:

$$p(y, \mathbf{x}) = \mathcal{N}\left(\begin{bmatrix} y \\ \mathbf{x} \end{bmatrix} \mid \begin{bmatrix} \mu_y \\ \boldsymbol{\mu}_x \end{bmatrix}, \begin{bmatrix} \Sigma_{yy} & \Sigma_{xy} \\ \Sigma_{xy} & \Sigma_{xx} \end{bmatrix}\right) \quad (18)$$

Then we will estimate all the parameters $\mu_y, \boldsymbol{\mu}_x, \Sigma_{yy}, \Sigma_{xy}, \Sigma_{xx}$ from the data with MLE. Given that the joint distribution is Gaussian, we know that the marginal distribution for each variable is also Gaussian. As a result, we know

that the following empirical means and covariances are the MLE solutions:

$$\hat{\mu}_y = \frac{1}{n} \sum_{i=1}^n y_i \quad (19)$$

$$\hat{\mu}_x = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad (20)$$

$$\hat{\Sigma}_{yy} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{\mu}_y)^2 \quad (21)$$

$$\hat{\Sigma}_{xy} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_x) (y_i - \hat{\mu}_y)^\top \quad (22)$$

$$\hat{\Sigma}_{xx} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_x) (\mathbf{x}_i - \hat{\mu}_x)^\top \quad (23)$$

Now we will show that $\mathbb{E}[y \mid \mathbf{x}]$ takes the same form as the predicted value \hat{y} in the previous section.

$$\mathbb{E}[y \mid \mathbf{x}] = \mu_{y|x} \quad (24)$$

$$= \mu_y + \Sigma_{xy}^\top \Sigma_{xx}^{-1} (\mathbf{x} - \mu_x) \quad (25)$$

Then plugging in the MLE solutions, we get:

$$\hat{y} := f(\mathbf{x}) = \hat{\mu}_y + \hat{\Sigma}_{xy}^\top \hat{\Sigma}_{xx}^{-1} (\mathbf{x} - \hat{\mu}_x) \quad (26)$$

$$= \hat{\Sigma}_{xy}^\top \hat{\Sigma}_{xx}^{-1} \mathbf{x} + \left(\hat{\mu}_y - \hat{\Sigma}_{xy}^\top \hat{\Sigma}_{xx}^{-1} \hat{\mu}_x \right) \quad (27)$$

$$= \mathbf{w}_1^\top \mathbf{x} + w_0 \quad (28)$$

$$= \begin{bmatrix} \hat{\Sigma}_{xy}^\top \hat{\Sigma}_{xx}^{-1} \\ \hat{\mu}_y - \hat{\Sigma}_{xy}^\top \hat{\Sigma}_{xx}^{-1} \hat{\mu}_x \end{bmatrix}^\top \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \quad (29)$$

It takes some algebra to convince that $[\mathbf{w}_1; w_0] = \mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ in the previous section if the rows of \mathbf{X} are $[\mathbf{x}_i \ 1]^\top$'s (appended by 1).

Acknowledgement

This note is, for the most parts, based on [Murphy \[2022\]](#).

References

K. P. Murphy. *Probabilistic Machine Learning: An introduction*. MIT Press, 2022. URL [probml.ai](#). 3