Discussion 8: Linear Regression and Modeling

Data 100, Spring 2020

Suraj Rampure

Friday, March 13th, 2020

Agenda

- Multiple Regression
- Solution to OLS
- Features

As per usual, everything will be posted at

surajrampure.com/teaching/ds100.html

Multiple Regression

Simple Linear Regression y = a + bx

Our model is

intercept
$$\hat{y} = \hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + ... + \theta_p x_p$$
 | feature |

- ullet $x_1, x_2, ..., x_p$ are called "features", "covariates", "explanatory variables", columns"
- $\theta_0, \theta_1, ..., \theta_p$ are called "weights"
- ullet y is called the "response variable"

Our goal is to model the relationship between explanatory variables and our response variable. More concretely, our goal is to find values of $\theta_0, \theta_1, ..., \theta_p$ that minimize our average L_2 loss:

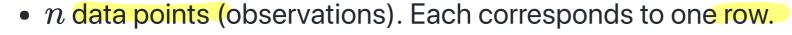
$$L(ec{ heta}) = rac{1}{n} \sum_{i=1}^n (y_i - \hat{y_i})^2$$

We denote the optimal values by $\hat{\theta_0}$, $\hat{\theta_1}$, ..., $\hat{\theta_p}$, or more generally, by $\hat{\theta}$.

Matrix Formulation

Our model is $\hat{y} = X\theta$:

$$egin{bmatrix} \hat{y_1} \ \hat{y_2} \ \hat{y_3} \ \vdots \ \hat{y_n} \end{bmatrix} = egin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & ... & x_{1p} \ 1 & x_{21} & x_{22} & x_{23} & ... & x_{2p} \ 1 & x_{31} & x_{32} & x_{33} & ... & x_{3p} \ \vdots & \vdots & \vdots & \vdots \ 1 & x_{n1} & x_{n2} & x_{n3} & ... & x_{np} \end{bmatrix} egin{bmatrix} heta_0 \ heta_1 \ heta_2 \ heta_3 \ \vdots \ heta_p \end{bmatrix}$$





• Almost always,
$$n > p$$
.



Under this formulation...

$$L(heta) = rac{1}{n} \sum_{i=1}^n (y_i - \hat{y_i})^2 = rac{1}{n} ||y - X heta||_2^2$$

Question: How do we find the optimal $\hat{\theta}$?

From the previous slide, our objective function is the average L_2 loss over our entire dataset (i.e. the mean squared error, MSE).

$$L(heta) = rac{1}{n}||y-X heta||_2^2$$

There are two ways to solve for $\hat{\theta}$:

- 1. Through calculus: take the gradient, set it equal to 0.
 - We are not doing this in our semester, but just know that it can be done.
- 2. Through geometry.
 - This was done in lecture, and we will recap the argument in the worksheet.

Either way, we end up at the same solution:

$$\hat{ heta} = (X^T X)^{-1} X^T y$$

Using our Model
$$p^{nd}$$
: $\alpha = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$, $b = \begin{bmatrix} b_2 \\ b_2 \\ \vdots \\ b_n \end{bmatrix}$, $a \cdot b = a_1 b_1 + a_2 b_2$

Given our initial data X and y, we now know how to solve for the vector θ that minimizes our model's loss. Now, given this vector, how do we use it to make predictions for **unseen data**?

$$\hat{O} = \begin{bmatrix} \hat{O}_{0} \\ \hat{O}_{1} \\ \hat{O}_{2} \end{bmatrix} \qquad \chi = \begin{bmatrix} 1 \\ \text{height} \\ \text{weight} \end{bmatrix}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{weight}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{height}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{height}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{height}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{height}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{height}$$

$$\hat{O} \cdot \chi = \hat{O}_{0} + \hat{O}_{1} \cdot \text{height} + \hat{O}_{2} \cdot \text{height}$$

Features

Note: In order for us to perform linear regression, our model only needs to be linear in terms of the weights.

For example, let's suppose we want to predict mpg given horsepower (x_1) , model year (x_2) and acceleration (x_3) . The following are both linear:

$$\Rightarrow y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$\Rightarrow y = \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 \sin x_3 + \theta_3 e^{x_1 x_2 x_3}$$

$$\Rightarrow \chi_{+} \theta_0 + \theta_1 x_1 + \theta_2 x_1^2 \sin x_3 + \theta_3 e^{x_1 x_2 x_3}$$

$$\Rightarrow \chi_{+} \theta_0 + \xi_0 + \xi_1 \chi_1$$

$$+ \xi_0 \chi_1 + \xi_0 \chi_2 \chi_2 + \xi_0 \chi_3 \chi_4$$

On the other hand, the following model is **not linear**: why?

$$y= heta_0+ heta_1x_1^{ heta_2}+ heta_2^2x_1x_2x_3$$
 could not model with $\hat{y}=\chi \phi$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 \sin x_3 + \theta_3 e^{x_1 x_2 x_3}$$