IEMS 304 Lecture 2: Simple Linear Regression

Yiping Lu yiping.lu@northwestern.edu

Industrial Engineering & Management Sciences Northwestern University



1

Simple Linear Regression



$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i,$$

- X has an arbitrary distribution, possibly deterministic.
- □ If X = x, then $Y = \beta_0 + \beta_1 x + \varepsilon$, with β_0, β_1 being the *coefficients*, and ε being the *noise* variable.

$$\square \mathbb{E}[\varepsilon|X=x] = 0, \ \operatorname{Var}(\varepsilon|X=x) = \sigma^2.$$

One option to estimate the unknown quantities is to find the optimal fit , to be precise here, minimize the mean squared error (MSE):

$$(\beta_0, \beta_1) = \arg\min_{(b_0, b_1)} \mathbb{E}[(Y - (b_0 + b_1 X)^2)].$$

\square How to access \mathbb{E} ?

• The data we may consider are $\{(X_1, Y_1), \ldots, (X_n, Y_n)\}$.

Monte Carlo Methods

How to Estimate π ?

- □ Draw a square of side length 2 (from −1 to +1) and inscribe a circle of radius 1.
- Randomly sample the points within the square.

□ Count how many points fall inside the circle.

The expectation of fraction of points in the circle is $\frac{\text{the circle's area}}{\text{total points' area}} \approx \frac{\pi r^2}{(2r)^2} = \frac{\pi}{4}$.

 \Box Hence $\pi \approx 4 \times \frac{\text{points in circle}}{\text{total points}}$





<u>Next</u>. $\hat{\beta}_0, \hat{\beta}_1$ has closed form solution!

How ?

How to find the Minimizer of a function $x^* = \arg \min_x f(x)$?

Solve the equation $\nabla f(x^*) = 0$

Find β_0, β_1

$$\hat{\beta}_1 = \frac{c_{XY}}{s_X^2},$$

where c_{XY}, s_X^2 are the sample covariance between X, Y and the sample variance of X respectively. As a reminder,

$$c_{XY}=\frac{1}{n}\sum_{i=1}^{n}(X_{i}-\overline{x})(Y_{i}-\overline{y}),s_{X}^{2}=\frac{1}{n}\sum_{i=1}^{n}(X_{i}-\overline{x})^{2}.$$

$$0 = \overline{xy} - (\overline{y} - \hat{\beta}_1 \overline{x}) \overline{x} - \hat{\beta}_1 \overline{x}^2$$
$$0 = c_{XY} - \hat{\beta}_1 s_X^2$$

How accurate is the Model?- Bias

$$\hat{\beta}_1 = \beta_1 + \frac{1}{ns_X^2} \sum_{i=1}^n (X_i - \overline{x}) \varepsilon_i.$$

<u>Statement</u>: $\hat{\beta}_1$ is unbiased, i.e. $\mathbb{E}[\hat{\beta}_1] = \beta_1$.

\square Find $(\hat{\beta}_0, \hat{\beta}_1)$ that minimize the least square

$$Q = \sum_{i=1}^{n} (y_i - \underbrace{(\hat{\beta}_0 + \hat{\beta}_1 x_i)}_{\hat{y}_i})^2.$$

• Denote
$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$
 as the fitted value;

• Denote $e_i = y_i - \hat{y}_i$ as the **residual**.

Therefore, minimizing the least square can be understood as fitting y_i 's to minimize residuals as good as possible.

How accurate is the Model?– Variance

$$\operatorname{Var}(\hat{\beta}_1) = \operatorname{Var}\left(\beta_1 + \frac{1}{ns_X^2}\sum_{i=1}^n (X_i - \overline{x})\varepsilon_i\right) = \frac{\sigma^2}{ns_X^2}.$$

Unconditioning on X

 $\hfill\square$ Bias apply the law of total expectation:

$$\mathbb{E}[\hat{\beta}_1] = \mathbb{E}\Big[\mathbb{E}[\hat{\beta}_1 \mid X_1, \dots, X_n]\Big] = \mathbb{E}[\beta_1] = \beta_1.$$

□ Variance apply the law of total variance:

$$\begin{aligned} \operatorname{Var}(\hat{\beta}_1) &= \mathbb{E}\Big[\operatorname{Var}(\hat{\beta}_1 \mid X_1, \dots, X_n)\Big] + \operatorname{Var}\Big(\mathbb{E}[\hat{\beta}_1 \mid X_1, \dots, X_n]\Big) \\ &= \mathbb{E}\Big[\frac{\sigma^2}{ns_X^2}\Big] + \operatorname{Var}(\beta_1) = \frac{\sigma^2}{n} \mathbb{E}\Big[\frac{1}{s_X^2}\Big]. \end{aligned}$$

Go Beyond Point Estimation

Fact.
$$\mathbb{E}[\hat{f}(x)] = \beta_0 + \beta_1 x.$$
 and $\operatorname{Var}(\hat{f}(x)) = \frac{\sigma^2}{n} \left(1 + \frac{(x - \overline{x})^2}{s_x^2}\right).$

What is the the standard error of an estimator ? $\operatorname{se}(\hat{\beta}_1) = \frac{\sigma}{\sqrt{ns_{\chi}^2}}$.

 \square What happens when the noise variance, σ^2 , increases?

 \square What happens when the number of samples, *n*, increases?

- □ What influences the variance of our predictions?
- \square What happens when we predict at x that is very close to $\overline{x}?$ How about very far?

Using the simple linear regression model,

$$\mathbb{E}[(Y - (eta_0 + eta_1 X))^2] = \sigma^2.$$
 (convince yourself why.)

Then, a natural estimator for σ^2 would be

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{f}(X_i))^2.$$

Notice that this is a biased estimator. Moreover $s^2 = \frac{n}{n-2}\hat{\sigma}^2$ is an unbiased estimator of σ^2 . (Later)

(residual)
$$e_i = Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i)$$

(noise) $\varepsilon_i = Y_i - (\beta_0 + \beta_1 X_i)$

- The sum of noise variables cannot equal zero all the time, because $Var(\sum_{i=1}^{n} \varepsilon_i) = n\sigma^2$.
- The sum of residuals is *always* zero, i.e. $\sum_{i=1}^{n} e_i = 0$.
- The sample correlation between the residuals and X_i 's is also 0, i.e. $\sum_{i=1}^{n} (X_i \overline{x}) e_i = 0.$

Assessing the Fit

Assessing the Fit

 \square As in simple regression, we calculate

- fitted values: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$;
- residuals: $e_i = y_i \hat{y}_i$;
- error sum of squares: $SSE = \sum_{i=1}^{n} e_i^2$;
- total sum of squares: SST = $\sum_{i=1}^{n} (y_i \bar{y})^2$;
- regression sum of squares: $SSR = \sum_{i=1}^{n} (\hat{y}_i \bar{y})^2$.

$$ar{y} = rgmin_c \sum_{i=1}^n (c-y_i)^2$$
 is the best constant fit of $\{y_i\}_{i=1}^n$

 \square We can decompose SST as

$$\underbrace{\sum_{i=1}^{n} (y_i - \bar{y})^2}_{\text{SST}} = \underbrace{\sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2}_{\text{SSR}} + \underbrace{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}_{\text{SSE}}$$

R² Statistics and Correlation

$$R^{2} \text{ (Coefficient of Determination):}$$

$$R^{2} = \frac{\text{SSR}}{\text{SST}}, \text{ where } \text{SSR} = \sum (\hat{y}_{i} - \bar{y})^{2}, \text{ SST} = \sum (y_{i} - \bar{y})^{2}.$$
Theorem
Recall Pearson correlation coefficient: $r = \frac{\sum (x_{i} - \bar{x})(y_{i} - \bar{y})}{\sqrt{\sum (x_{i} - \bar{x})^{2} \sum (y_{i} - \bar{y})^{2}}}, \text{ then we have}$

$$R^{2} = r^{2}$$

Prove
$$R^2 = r^2$$

Since
$$\hat{\beta}_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2} = r \frac{s_y}{s_x}$$
, we have $SSR = \frac{(\sum(x_i - \bar{x})(y_i - \bar{y}))^2}{\sum(x_i - \bar{x})^2}$. Thus,
 $R^2 = \frac{SSR}{SST} = \frac{(\sum(x_i - \bar{x})(y_i - \bar{y}))^2}{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2} = r^2$.

Prove:
$$s^2 = \frac{n}{n-2}\hat{\sigma}^2$$
 is an *unbiased* estimator of σ^2

Pipeline of Machine Learning

The model looks similar,

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i,$$

with modified assumptions:

 \square X has an arbitrary distribution, possibly deterministic.

□ If X = x, then $Y = \beta_0 + \beta_1 x + \varepsilon$, with β_0, β_1 being the coefficients, and ε being the noise variable.

(stronger)
$$\varepsilon \sim N(0, \sigma^2)$$
, and is independent of X.

 \Box (stronger) ε is *independent* across observations.

Question. What is $p(Y_i|X_i; b_0, b_1, s^2)$?

Given the data, the likelihood under this set of assumption is a function of the unknown parameters, defined as

$$L(b_0, b_1, s^2) = \prod_{i=1}^n p(Y_i | X_i; b_0, b_1, s^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi s^2}} \exp\left\{-\frac{1}{2s^2}(Y_i - (b_0 + b_1 X_i))^2\right\}$$

$$\log(ab) = \log(a) + \log(b)$$

$$\log L(b_0, b_1, s^2) \stackrel{\text{def}}{=} \ell(b_0, b_1, s^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log s^2 - \frac{1}{2s^2} (Y_i - (b_0 + b_1 X_i))^2.$$

Logistic regression

Step 1. Likelihood for a Logistic Binary Outcome:

For each observation $y_i \in \{0, 1\}$ with probability p_i for $y_i = 1$, the likelihood is

$$L(p_i | y_i) = p_i^{y_i} (1 - p_i)^{1 - y_i}.$$

where probability $p_i = \frac{1}{1+e^{-\beta^T x_i}}$ using the logistic function.

Step 2. Log-Likelihood:

For n independent observations, the log-likelihood function is

$$\ell(eta) = \sum_{i=1}^n igg[y_i \logigg(rac{1}{1+e^{-eta^ au_{oldsymbol{x}_i}}igg) + (1-y_i) \logigg(1-rac{1}{1+e^{-eta^ au_{oldsymbol{x}_i}}igg)igg].$$

Step 3. Estimation:

Maximizing $\ell(\beta)$ with respect to β gives the maximum likelihood estimates, leading to the logistic regression model.

Solution No closed-form solution.

Gradient Descent

- **Gradient Descent** is an iterative optimization method to find local minima of a function.
- The update rule is $x_{n+1} = x_n \alpha \nabla f(x_n)$, where α is the learning rate.



III Conditioned Problems

- The function $f(x_1, x_2) = 10x_1^2 + x_2^2$ has very different curvatures along x_1 and x_2 .
- Its level sets are ellipses elongated along the x_2 -axis.
- With a fixed learning rate, gradient descent can overshoot in the steep x₁ direction, leading to oscillatory (zigzag) behavior.





Homework



Gradient Descent: ||beta - true_beta||

Newton's Method: ||beta - true_beta||



Pipeline of Machine Learning