# IEMS 304 Lecture 2: Simple Linear Regression

Yiping Lu

yiping.lu@northwestern.edu
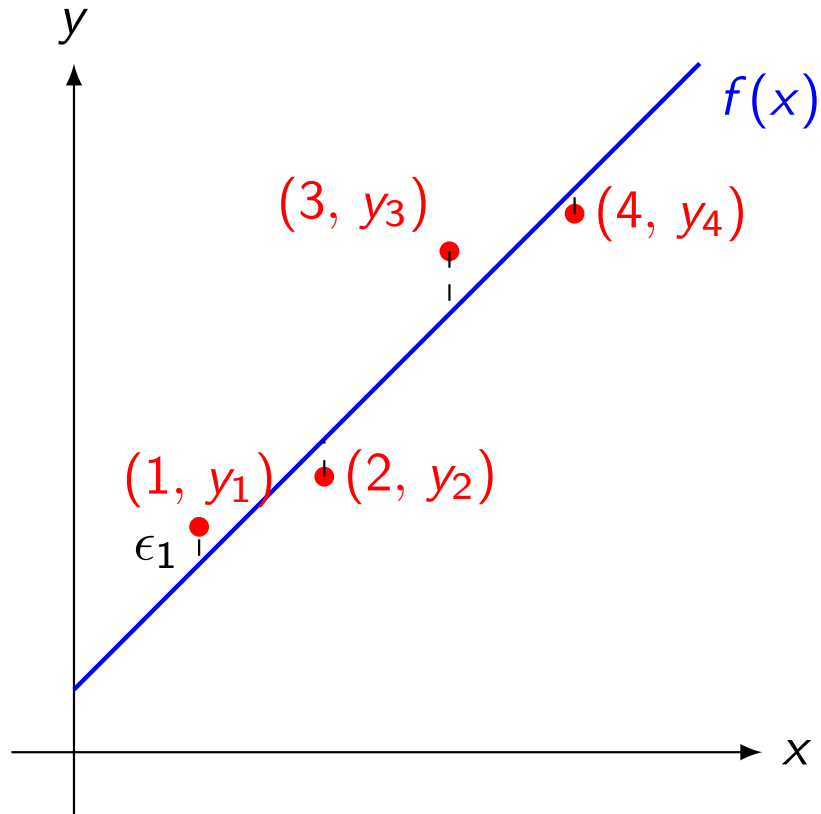
*Industrial Engineering & Management Sciences*

*Northwestern University*

NORTHWESTERN
UNIVERSITY

# Simple Linear Regression

Dataset $(X_1, Y_1), (X_2, Y_2), \ldots$

real number    real number

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

☐ $X$ has an arbitrary distribution, possibly deterministic.

☐ If $X = x$, then $Y = \beta_0 + \beta_1 x + \varepsilon$, with $\beta_0, \beta_1$ being the *coefficients*, and $\varepsilon$ being the *noise* variable.

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

One option to estimate the unknown quantities is to find the optimal fit to $L_2$ loss be precise here, minimize the mean squared error (MSE):

*because I'm usy*

*minimize*

*$L_2$ loss for a single prediction, will return the mean*

prediction

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

*Variable to optimize*

objective function (population)

$$\mathbb{E}[Y | x=x] = \beta_0 + \beta_1 X$$

$$\mathbb{E}[\varepsilon | x=x] = 0$$

☐ How to access $\mathbb{E}$?

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

Only Thing I can compute.

$$(\widehat{\beta}_0, \widehat{\beta}_1) := \arg\min_{(b_0, b_1)} \frac{1}{n} \sum_{i=1}^{n} \left[\left(Y_i - (b_0 + b_1 x_i)\right)^2\right]$$

*Empirical Objective Function*

# Monte Carlo Methods

## How to Estimate $\pi$ ?

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

❐ Hence $\boxed{\pi \approx 4 \times \frac{\text{points in circle}}{\text{total points}}}$.

We minimize in-sample, empirical MSE: (mean square error)

$$(\hat{\beta}_0, \hat{\beta}_1) = \arg \min_{(b_0, b_1)} \underbrace{\frac{1}{n} \sum_{i=1}^{n} (Y_i - (b_0 + b_1 X_i))^2}_{\widehat{\mathrm{MSE}}(b_0, b_1)} \, .$$

**Next.** $\hat{\beta}_0, \hat{\beta}_1$ has closed form solution!

**How ?**

$$f(x) = g_1(g_2(x)) \qquad \frac{\partial f}{\partial x} = \frac{\partial g_1(y)}{\partial y}\Big|_{y=g_2(x)} \cdot \frac{\partial g_2(x)}{\partial x}$$

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

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

$$f(b_0, b_1) = \frac{1}{n}\sum_{i=1}^{n}\left[\underbrace{Y_i - (b_0 + b_1 X_i)}_{g_2(b_0, b_1)}\right]^{\overbrace{2}^{g_1(b_0, b_1)}}$$

$$\nabla_{b_0} f(b_0, b_1) = -\frac{1}{n}\sum_{i=1}^{n}\underbrace{2(Y_i - (b_0 + b_1 X_i))}_{\partial g_1} \cdot \underbrace{1}_{\partial g_2} = 0$$

$$\nabla_{b_1} f(b_0, b_1) = -\frac{1}{n}\sum_{i=1}^{n}\underbrace{2(Y_i - (b_0 + b_1 X_i))}_{\partial g_1} \cdot \underbrace{X_i}_{\partial g_2} = 0$$

linear Eq. n.s.t. $b_0, b_1$

6

$$\nabla_{b_0} f = 0 \implies \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \underbrace{(b_0 + b_1 x_i)}_{\text{residual}} \right) \cdot 1 = 0$$

The error of linear regression on training data

$$\nabla_{b_1} f = 0 \implies \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - (b_0 + b_1 x_i) \right) \cdot x_i = 0$$

① The residual /error on training data is mean zero!

$$Cov(X, Y) = \frac{1}{n} \sum_{i=1}^{n} x_i \cdot Y_i$$

② The residual /error on training data is independent to the data!

$$b_0 = \frac{1}{n} \sum_{i=1}^{n} (Y_i - b_1 x_i) = \bar{Y} - b_1 \bar{x} \qquad (\triangle)$$

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^{n} Y_i, \qquad \bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

Plug $(\triangle)$ into $\nabla_{b_1} f = 0$

$$\frac{1}{n} \sum_{i=1}^{n} \left( Y_i - (\bar{Y} - b_1 \bar{x}) - b_1 x_i \right) x_i = 0$$

$$\implies \frac{1}{n} \sum_{i=1}^{n} \left( (Y_i - \bar{Y}) - b_1 (x_i - \bar{x}) \right) x_i = 0 \qquad (\text{✦})$$

This is using $((x_i - \bar{x}), (Y_i - \bar{Y}))$ as dataset to fit the simple linear regression.

Computing Eq $(\text{✦})$

$$\frac{1}{n} \sum_{i=1}^{n} x_i (Y_i - \bar{Y}) - \frac{1}{n} \sum_{i=1}^{n} x_i (x_i - \bar{x}) b_1 = 0$$

$$-\bar{x} \cdot \left( \frac{1}{n} \sum_{i=1}^{n} (Y_i - \bar{Y}) \right) = 0 \qquad -\bar{x} b_1 \left( \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x}) \right) = 0$$

$$\implies \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(Y_i - \bar{Y}) - \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})(x_i - \bar{x}) b_1 = 0$$

$$\hat{\beta}_1 = \frac{c_{XY}}{s_X^2}, \; = \; \frac{\text{Covariance } (X, Y)}{\text{Covariance } (X, x)}$$

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

$$\underbrace{\frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{x})(Y_i - \bar{y})}_{\text{Covariance } (X, Y)}, \quad \underbrace{\frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{x})^2}_{\text{Var } (x), \; \text{Covariance}(X, x)}$$

$$c_{XY} = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{x})(Y_i - \bar{y}), s_X^2 = \frac{1}{n} \sum_{i=1}^{n} (X_i - \bar{x})^2.$$

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

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

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

❑ 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

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

❒ Bias apply the law of total expectation:

$$\mathbb{E}[\hat{\beta}_1] = \mathbb{E}\left[\mathbb{E}[\hat{\beta}_1 \mid X_1, \ldots, X_n]\right] = \mathbb{E}[\beta_1] = \beta_1.$$

❒ Variance apply the law of total variance:

$$\mathrm{Var}(\hat{\beta}_1) = \mathbb{E}\left[\mathrm{Var}(\hat{\beta}_1 \mid X_1, \ldots, X_n)\right] + \mathrm{Var}\left(\mathbb{E}[\hat{\beta}_1 \mid X_1, \ldots, X_n]\right)$$

$$= \mathbb{E}\left[\frac{\sigma^2}{ns_X^2}\right] + \mathrm{Var}(\beta_1) = \frac{\sigma^2}{n}\mathbb{E}\left[\frac{1}{s_X^2}\right].$$

# Go Beyond Point Estimation

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

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

❒ What happens when the noise variance, $\sigma^2$, increases?

❒ What happens when the number of samples, $n$, increases?

❒ What influences the variance of our predictions?

❒ What happens when we predict at $x$ that is very close to $\overline{x}$? How about very far?

Using the simple linear regression model,
$$\mathbb{E}[(Y - (\beta_0 + \beta_1 X))^2] = \sigma^2. \quad \textit{(convince yourself why.)}$$

Then, a natural estimator for $\sigma^2$ 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 $\sigma^2$. (Later)

# Residual and Error

$$\begin{aligned}
\text{(residual)} \quad e_i &= Y_i - (\hat{\beta}_0 + \hat{\beta}_1 X_i) \\
\text{(noise)} \quad \varepsilon_i &= Y_i - (\beta_0 + \beta_1 X_i)
\end{aligned}$$

# Remark

- The sum of noise variables cannot equal zero all the time, because $\mathrm{Var}(\sum_{i=1}^{n} \varepsilon_i) = n\sigma^2$.

- 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

❒ 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: $\text{SSE} = \sum_{i=1}^{n} e_i^2$;
- total sum of squares: $\text{SST} = \sum_{i=1}^{n} (y_i - \bar{y})^2$;
- regression sum of squares: $\text{SSR} = \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2$.

$$\bar{y} = \arg\min_c \sum_{i=1}^{n} (c - y_i)^2 \text{ is the best constant fit of } \{y_i\}_{i=1}^{n}!$$

❒ We can decompose $\text{SST}$ as

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

# $R^2$ Statistics and Correlation

$R^2$ **(Coefficient of Determination):**
$$R^2 = \frac{\text{SSR}}{\text{SST}}, \quad \text{where} \quad \text{SSR} = \sum(\hat{y}_i - \bar{y})^2, \quad \text{SST} = \sum(y_i - \bar{y})^2.$$

**Theorem**

*Recall* *Pearson correlation coefficient:* $r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}}$, *then we have*

$$R^2 = r^2$$

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,

$$R^2 = \frac{\text{SSR}}{\text{SST}} = \frac{\left(\sum(x_i - \bar{x})(y_i - \bar{y})\right)^2}{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2} = r^2.$$

# Error

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

$$Y = b_0 + b_1 X + \varepsilon$$

$$Y = \max(b_0 + b_1 X + \varepsilon, 0)$$

**Pipeline of Machine Learning**

The model looks similar,

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

with modified assumptions:

☐ $X$ has an arbitrary distribution, possibly deterministic.

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

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

☐ (stronger) $\varepsilon$ is *independent* across observations.

---

**Question.** What is $p(Y_i | X_i; b_0, b_1, s^2)$? $Y_i = b_0 + b_1 X_i + \varepsilon_i,$ $\varepsilon_i \sim N(0, s^2)$

observes a data $(X_i, Y_i)$ $\varepsilon_i = \overbrace{(Y_i - b_0 - b_1 X_i)}^{\text{Residual}}$ $\to$ means $P(\varepsilon_i) = \frac{1}{\sqrt{2\pi s^2}} \exp\left\{\frac{-1}{2s^2} \cdot \varepsilon^2\right\}$

what is the probability that $Y_i$ is the value I observe?. residual

21

# Log-Likelihood

*max likelihood $(\Rightarrow)$ minimize for $(residual)^2$*

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

*is the probability that $Y_i$ is the value I demand*

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

*$\star$ $\exp\left(\begin{smallmatrix}negative\\constant\end{smallmatrix} (residual)^2\right)$*

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

---

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

**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

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

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

**Step 2. Log-Likelihood:**

For $n$ independent observations, the log-likelihood function is

$$\ell(\beta) = \sum_{i=1}^{n} \left[ y_i \log\left( \frac{1}{1 + e^{-\beta^T x_i}} \right) + (1 - y_i) \log\left( 1 - \frac{1}{1 + e^{-\beta^T x_i}} \right) \right].$$

**Step 3. Estimation:**

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

☹ No closed-form solution.

# Gradient Descent

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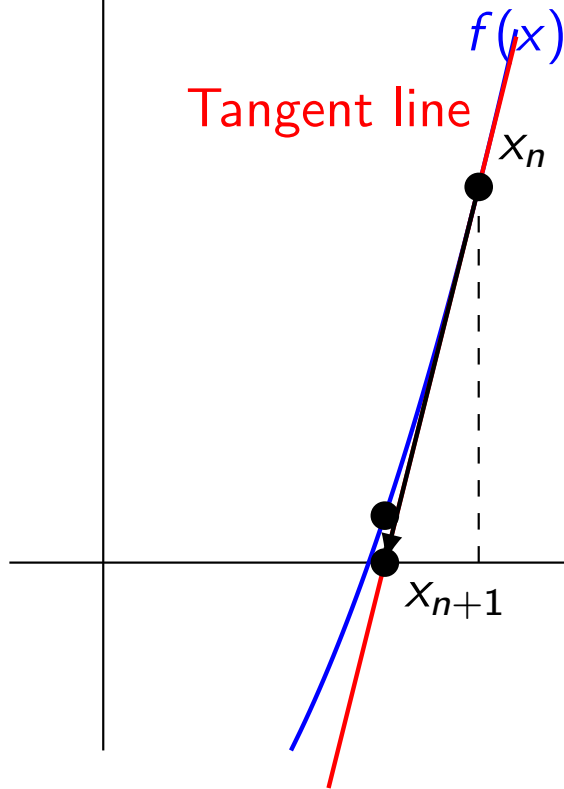
# Ill Conditioned Problems

- 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 $x_1$ direction, leading to oscillatory (zigzag) behavior.

Newton's method is an iterative technique for finding a root of a nonlinear equation $F(x) = 0$ via
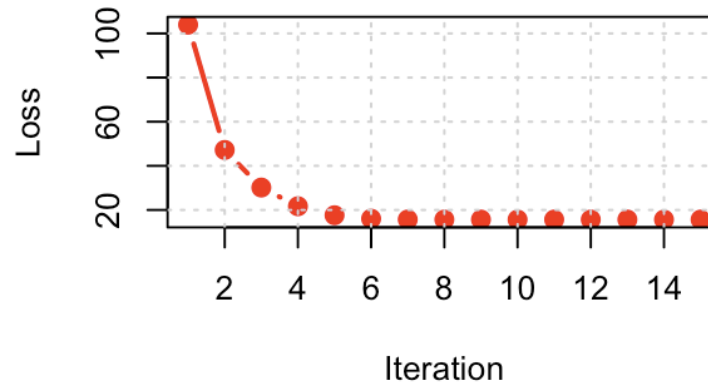
$$x_{n+1} = x_n - F'(x_n)^{-1} F(x_n).$$

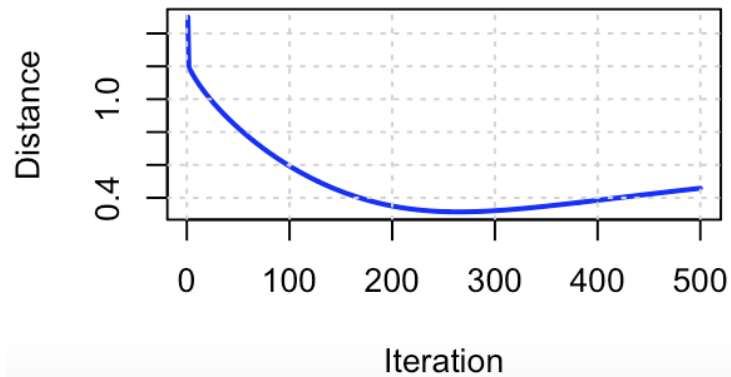What happens if one optimize
$f(x_1, x_2) = 10x_1^2 + x_2^2$?