

# IEMS 304 Lecture 2: Simple Linear Regression

---

Yiping Lu

yiping.lu@northwestern.edu

*Industrial Engineering & Management Sciences  
Northwestern University*

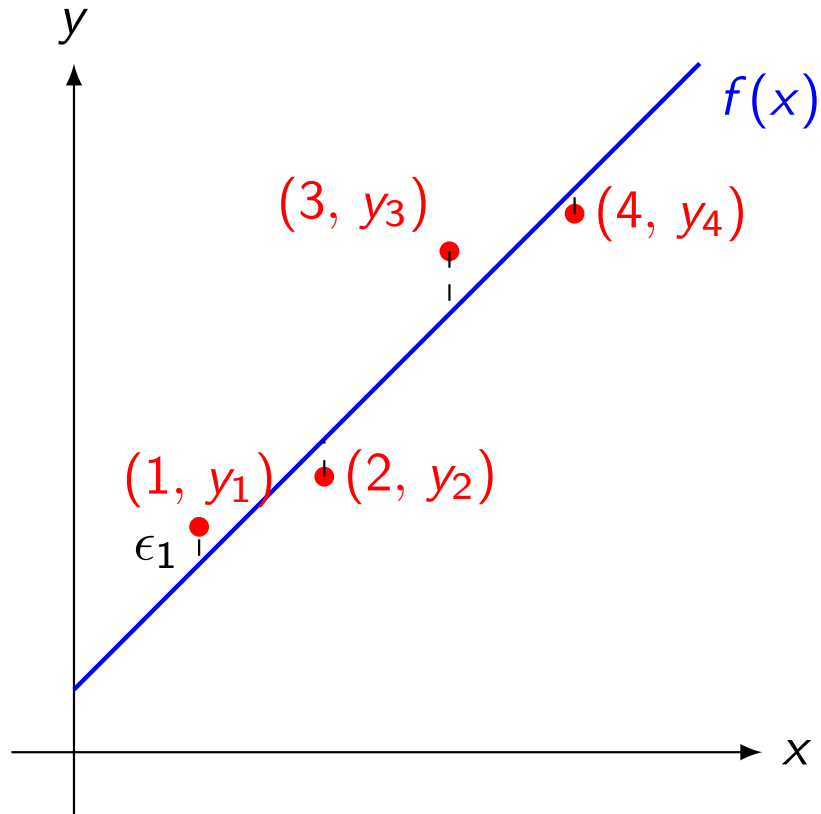


NORTHWESTERN  
UNIVERSITY

# Simple Linear Regression

---

# Linear Regression



Data set  $(x_1, y_1), (x_2, y_2), \dots$   
           $\uparrow$            $\uparrow$   
      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.
- $\mathbb{E}[\varepsilon|X = x] = 0$ ,  $\text{Var}(\varepsilon|X = x) = \sigma^2$ .

# Least Squares Estimators

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):

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

Variable to optimize

objective function (population)

because I'm using  $L_2$  loss, minimize  $L_2$  loss for a single prediction, will return the mean

□ How to access  $\mathbb{E}$ ?

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

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

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

only Thing I can compute.

$$(\hat{\beta}_0, \hat{\beta}_1) := \arg \min_{(b_0, b_1)}$$

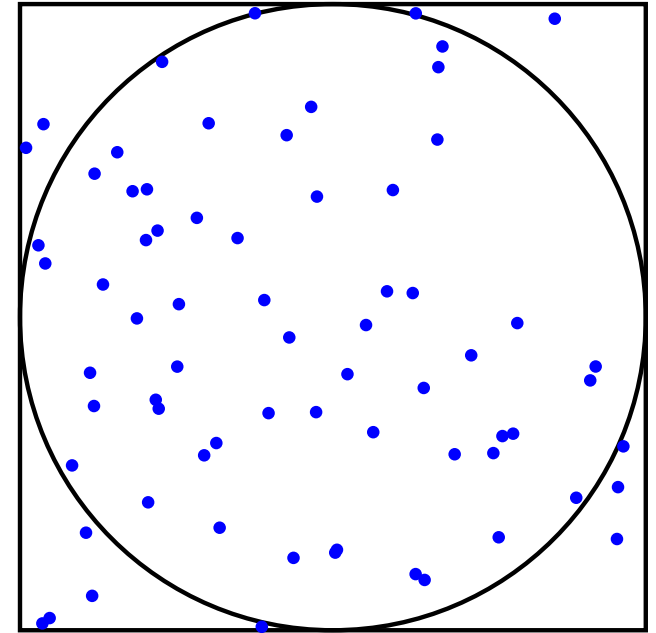
$$\frac{1}{n} \sum_{i=1}^n [(Y_i - (b_0 + b_1 X_i))^2]$$

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



# Find $\beta_0, \beta_1$

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

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

---

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

How ?

# How to find the **Minimizer** of a Function

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

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

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

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

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

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

linear Eq. r.s.f.  $b_0, b_1$

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

The error of linear regression on training data

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

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

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

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

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

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

Plug (Δ) into  $\nabla_{b_1} f = 0$

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

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

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

Computing Eq (★)

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

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

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



## Find $\beta_0, \beta_1$

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

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

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

*Covariance(X, Y)                      Var(X), Covariance(X, X)*

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

$$0 = c_{XY} - \hat{\beta}_1 s_X^2$$

# How accurate is the Model?– Bias

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

---

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

# Model Fitting

□ 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

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

# Unconditioning on $X$

- ❑ **Bias** apply the law of total expectation:

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

- ❑ **Variance** apply the law of total variance:

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

# Go Beyond Point Estimation

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

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

# Exercise

- ❑ 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  $\bar{x}$ ? How about very far?

# How to Estimate $\sigma$ ?

Using the simple linear regression model,

$$\mathbb{E}[(Y - (\beta_0 + \beta_1 X))^2] = \sigma^2. \quad (\text{convince yourself why.})$$

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

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

$$\text{(noise)} \quad \varepsilon_i = Y_i - (\beta_0 + \beta_1 X_i)$$

# Remark

- The sum of noise variables cannot equal zero all the time, because  $\text{Var}(\sum_{i=1}^n \varepsilon_i) = n\sigma^2$ .
- 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.  $\sum_{i=1}^n (X_i - \bar{x})e_i = 0$ .

# Assessing the Fit

---

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

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

□ We can decompose SST as

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

# $R^2$ Statistics and Correlation

$R^2$  (Coefficient of Determination):

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

## 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

$$R^2 = r^2$$

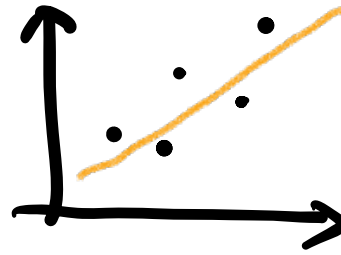
## Prove $R^2 = r^2$

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,

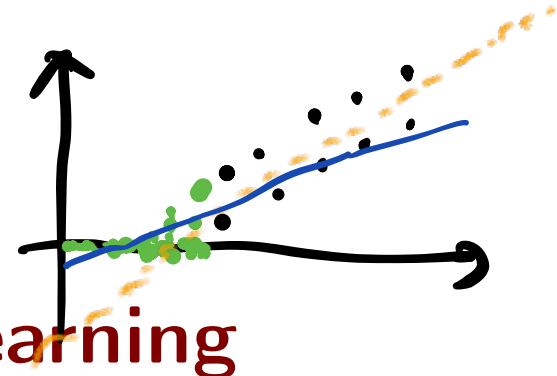
$$R^2 = \frac{SSR}{SST} = \frac{(\sum(x_i - \bar{x})(y_i - \bar{y}))^2}{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2} = r^2.$$

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

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



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



Pipeline of Machine Learning

---



# Log-Likelihood

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.

$$P(\varepsilon_1 = \dots, \varepsilon_2 = xxx) = P(\varepsilon_1 = \dots) \cdot P(\varepsilon_2 = xxx)$$

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

$$Y_i = b_0 + b_1 X_i + \varepsilon_i, \quad (\varepsilon_i \sim N(0, s^2))$$

observes a data  $(X_i, Y_i)$   $\varepsilon_i = \overbrace{(Y_i - b_0 - b_1 X_i)}^{\text{Residual}}$   $\rightarrow$  means  $P(\varepsilon_i) = \frac{1}{\sqrt{2\pi}s^2} \exp\left\{-\frac{1}{2s^2} \underbrace{\varepsilon_i^2}_{\text{Residual}^2}\right\}$   
What is the probability that  $Y_i$  is the value I observe?

# Log-Likelihood

max likelihood  $\Leftrightarrow$  minimize for (residual)<sup>2</sup>.

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

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

is the probability that  $Y_i$  is the value I find

$$\exp(a) \exp(b) = \exp(a+b)$$

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

likelihood of second data

exp (negative constant (residual)<sup>2</sup>)

likelihood of

first data

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

max  $\sum_{\text{all data}} \log$  of the likelihood of each data  $\Leftrightarrow$  min  $\sum_{\text{all data}} -\log$  of likelihood

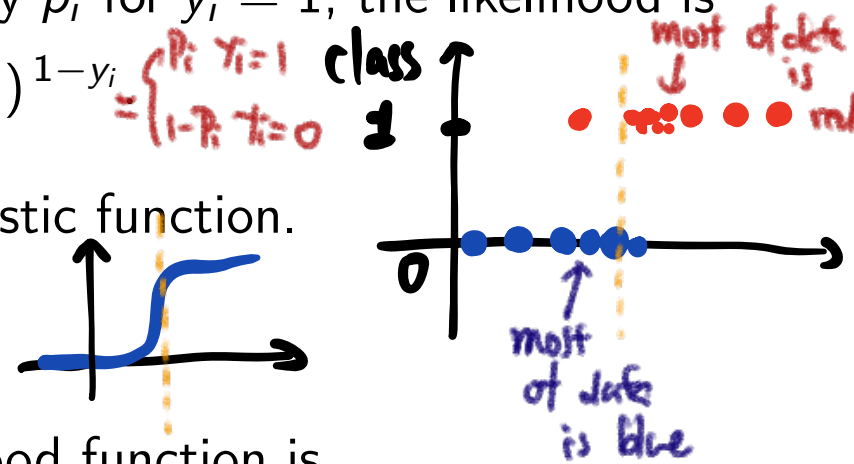
# Logistic regression

## 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

$$L(\mathbf{y} | \mathbf{p}) = p_i^{y_i} (1 - p_i)^{1-y_i}$$

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

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

## 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.

Basic Idea:  $\max \underbrace{P(Y|x)}_{\text{called likelihood}}$

Idea 1:  $\max P(Y|x) \Leftrightarrow \min -\log P(Y|x)$

Idea 2:  $P(Y|x) = \prod_{i=1}^n \underbrace{P(Y_i | x_i)}$

because every data is independent / experience

Fact:  $\log(ab) = \log(a) + \log(b)$ ,  $\log\left(\prod_{i=1}^n P_i\right) = \sum_{i=1}^n \log(P_i)$

Then  $\max P(Y|x)$

$\Leftrightarrow \min -\log P(Y|x) = \min -\log\left(\prod_{i=1}^n P(Y_i | x_i)\right)$

$= \min \sum_{i=1}^n -\log(P(Y_i | x_i))$

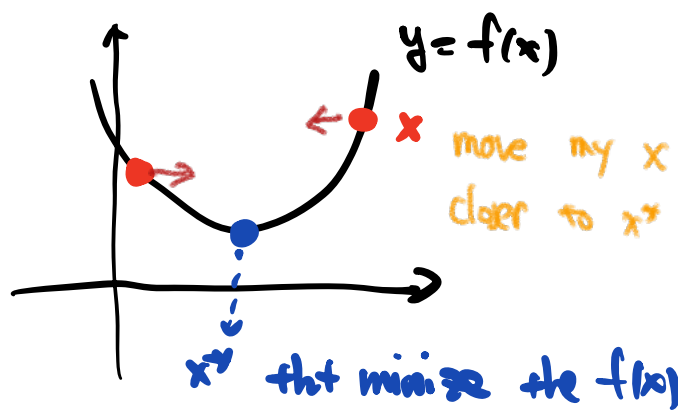
loss function for data  $(x_i, Y_i)$

Example. Assume

$Y_i = f(x_i) + \varepsilon_i$ ,  $\varepsilon_i$  is Gaussian

$\Rightarrow -\log(P(Y_i | x_i)) = (Y_i - f(x_i))^2$

Optimization: Iterative Procedure to minimize a function



Case 1  $\nabla f(x) > 0$

Case 2  $\nabla f(x) < 0$

$t$  is the time of iterative procedure

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

here  $\alpha > 0$  is called learning rate / step size.

$\nabla f$  means for two dimensional function  $f(x)$ ,  $x \in \mathbb{R}^2$

$$\nabla f(x) = \begin{bmatrix} df/dx_1 \\ df/dx_2 \end{bmatrix} \in \mathbb{R}^2$$

It means Taylor expansion holds:

$$f(z) \approx f(y) + \nabla f(y) \cdot (z - y)$$

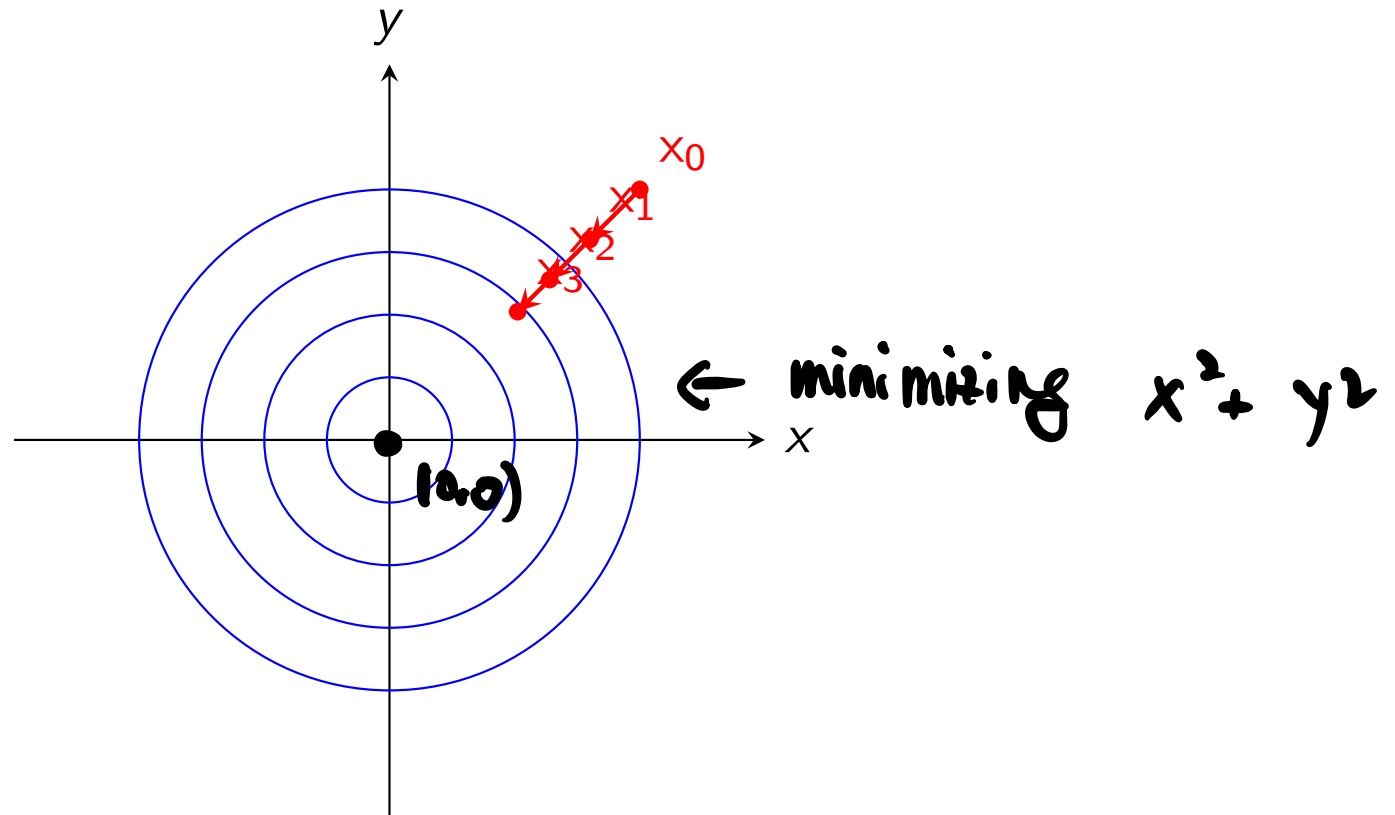
If  $f$  is decaying faster in  $x_1$  direction, then  $z_1 - y_1$  is larger

$$\begin{bmatrix} \frac{df}{dx_1} \\ \frac{df}{dx_2} \end{bmatrix} \begin{bmatrix} z_1 - y_1 \\ z_2 - y_2 \end{bmatrix} = \frac{df}{dx_1} \cdot (z_1 - y_1) + \frac{df}{dx_2} (z_2 - y_2)$$

$= -\alpha \frac{df}{dx_1}$

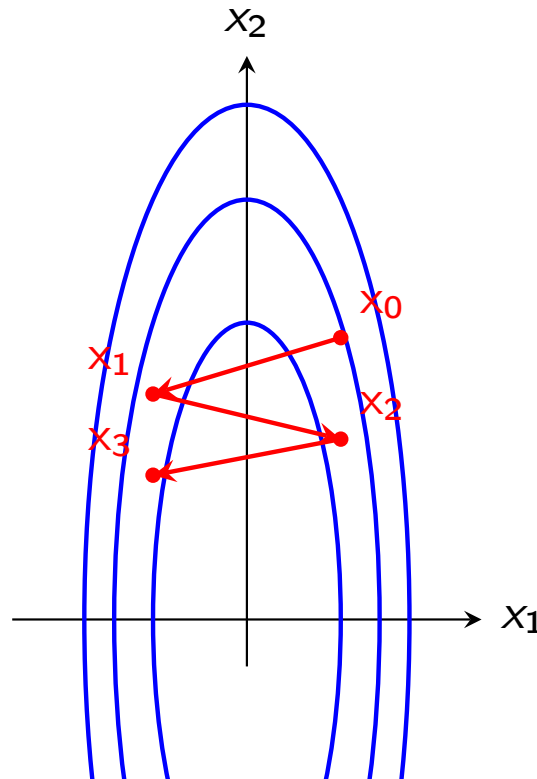
# Gradient Descent

- **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  $\alpha$  is the learning rate.

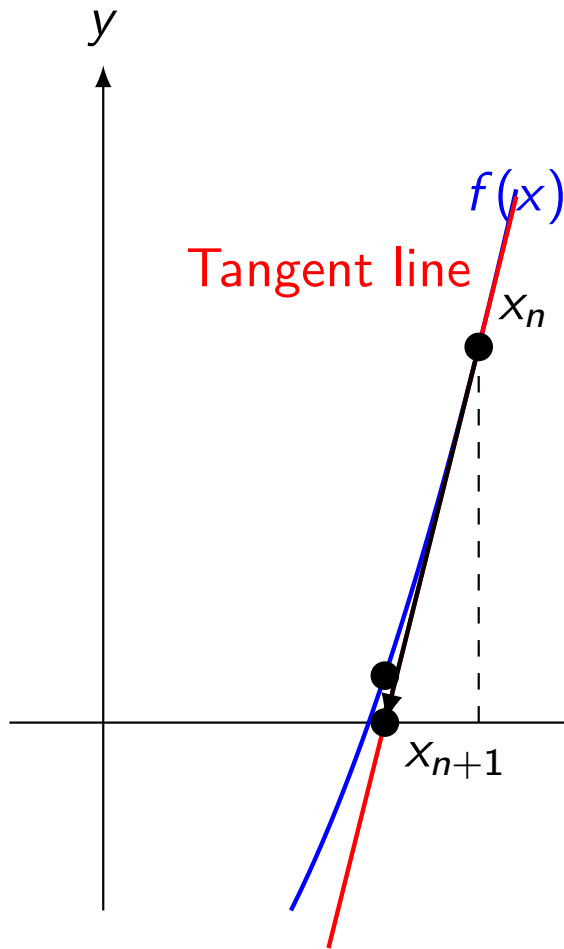


# III Conditioned Problems

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



# Newton Methods



Newton's method is an iterative technique for finding a root of a nonlinear equation  $F(x) = 0$  via

$$x_{n+1} = x_n - F'(x_n)^{-1} F(x_n).$$

$$F(x) = 0$$

$$\Downarrow$$

$$F(x_n) + F'(x_n) \cdot (x - x_n) = 0$$

$\Downarrow$  Solve linear approximation

$$x_{n+1} = x_n - F'(x_n)^{-1} F(x_n)$$

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

In optimization  
 $\min f(x)$

$\Leftrightarrow$  Solve nonlinear Eq

$$\nabla f(x) = 0$$

$\Downarrow$

$$x_{n+1} = x_n - Hf^{-1} \cdot \nabla f$$

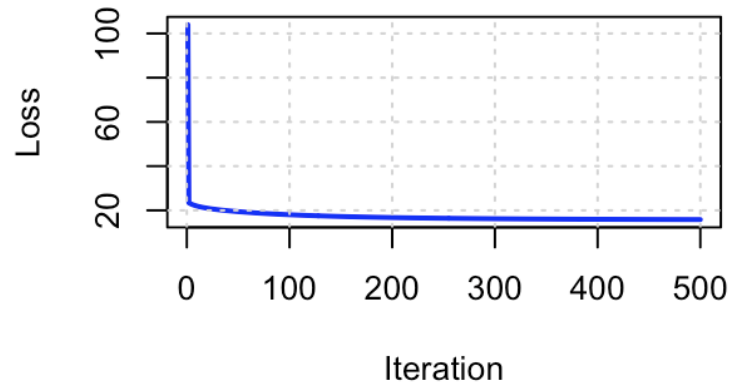
Hessian<sup>-1</sup> · gradient

Newton's method.

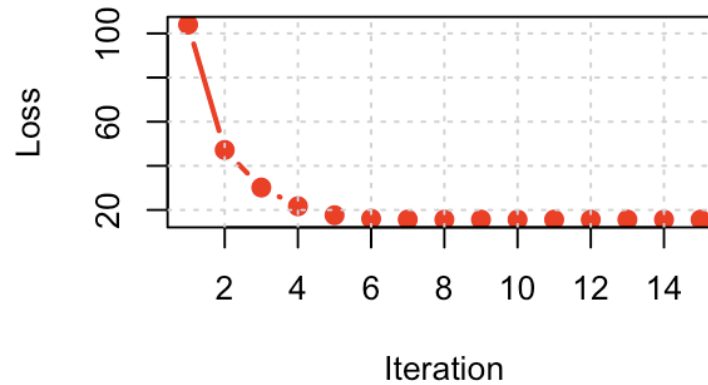


# Homework

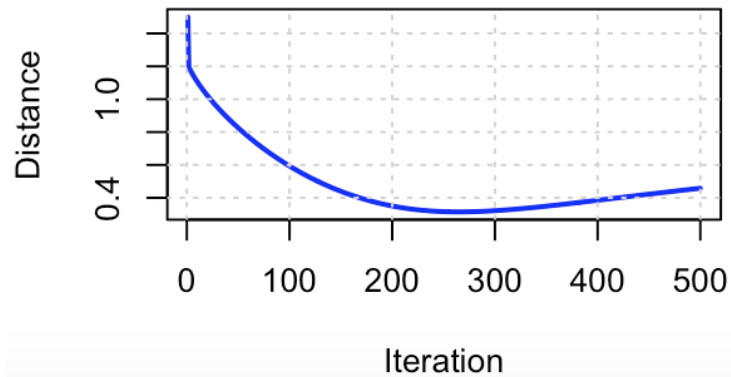
**Gradient Descent: Loss**



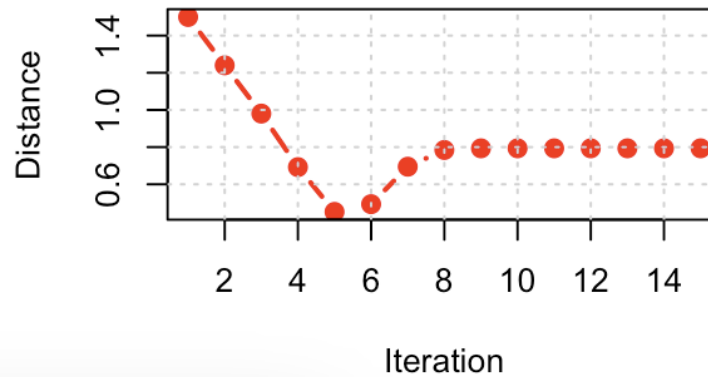
**Newton's Method: Loss**



**Gradient Descent:  $\|\text{beta} - \text{true\_beta}\|$**



**Newton's Method:  $\|\text{beta} - \text{true\_beta}\|$**



# Pipeline of Machine Learning