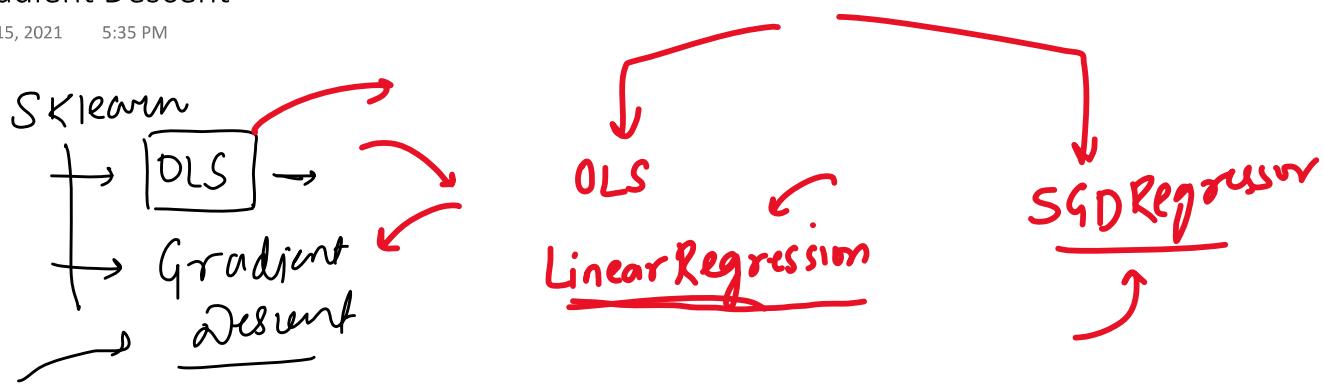


# Why Gradient Descent

Saturday, May 15, 2021 5:35 PM



## Code From Scratch

Monday, May 17, 2021 12:15 PM

$$\beta = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

$(100, 3) \quad (100)$   
 $\downarrow$   
 $100 \quad (100, 1)$

$X \rightarrow \text{matrix}$

cgpa	iq	gender	lpa
1	-	-	-
1	-	-	-

$X \rightarrow X_{\text{train}}$

$Y \rightarrow Y_{\text{train}}$   $\xrightarrow{\text{diabetes}} \text{sklearn} \rightarrow R^2 - \text{same}$

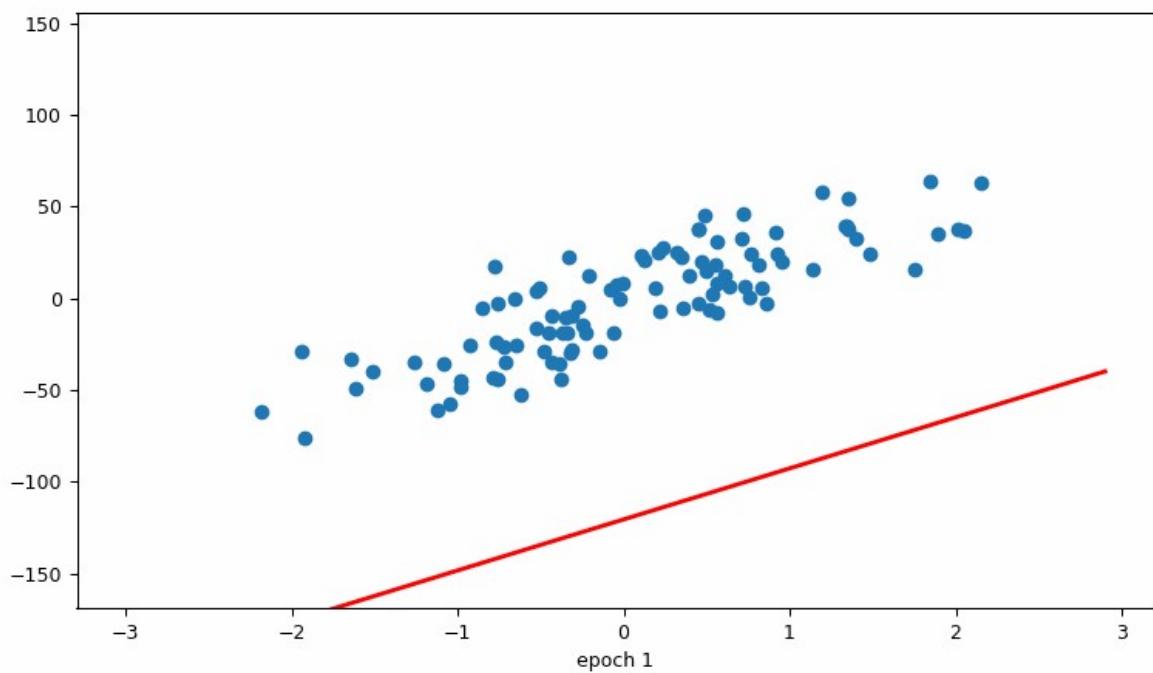
$X_{\text{test}}$   $\beta_0 \quad \beta_1 \rightarrow \beta_3$

$\underbrace{3}_{3}$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3$$

$$= \beta_0 + \underbrace{\beta_1 x_1 + \dots + \beta_3 x_3}_{\text{np dot}} + (\beta, X_{\text{test}})$$

$$\begin{aligned} & X_{\text{test}} \quad (\text{coeff}) \\ & \underbrace{(89, 10)}_{(89, 1) + \beta_0} \quad (10, 1) \\ & 89 \rightarrow ① \quad (89, 1) \\ & \boxed{y_{\text{pred}}} \end{aligned}$$



$$\sum (y_i - mx_i - b)^2$$

$$\sum -2 (y_i - mx_i - b) x_i$$

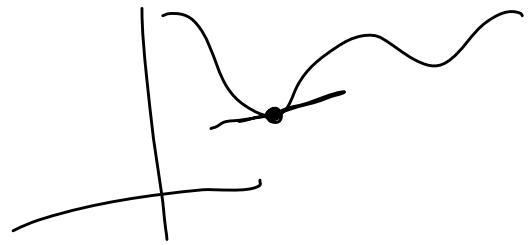
- 2

# What is Gradient Descent?

Thursday, May 20, 2021 1:46 PM

Linear Reg  
Logistic Reg  
Tree

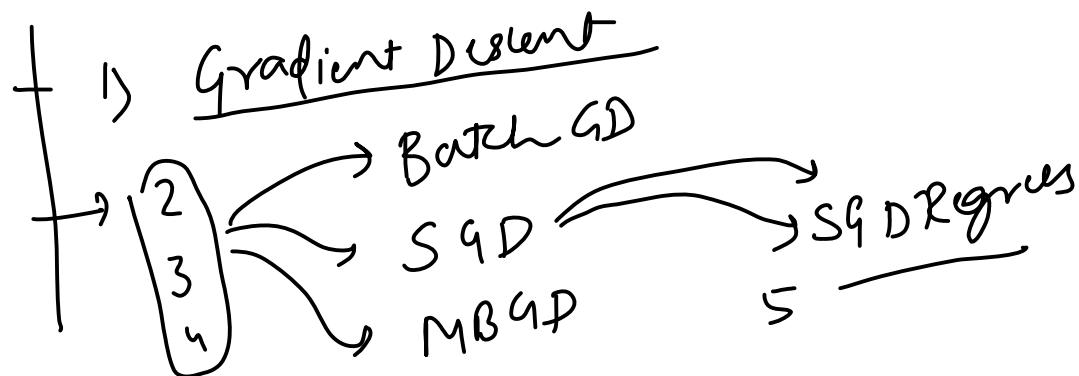
→ Deep Learning



# The Plan

Thursday, May 20, 2021 1:46 PM

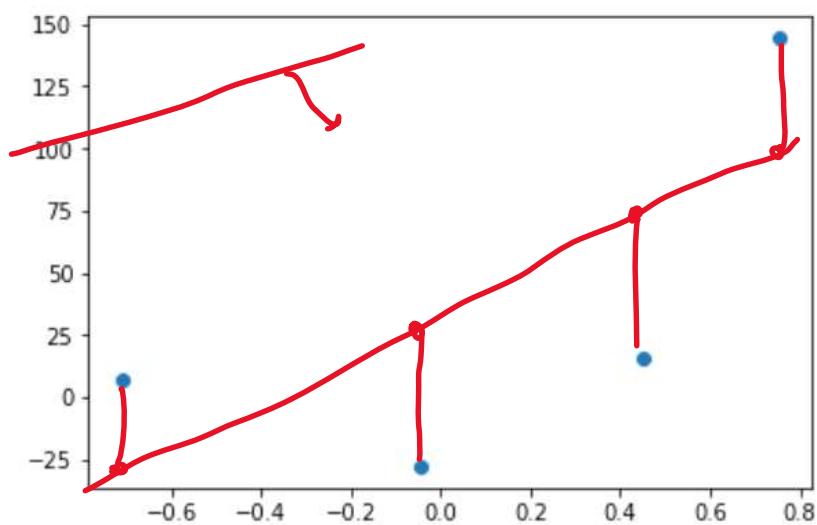
5 videos



# Intuition

Thursday, May 20, 2021 1:46 PM

2 cols  
4 rows



$\text{cgpa} | \text{gpa}$

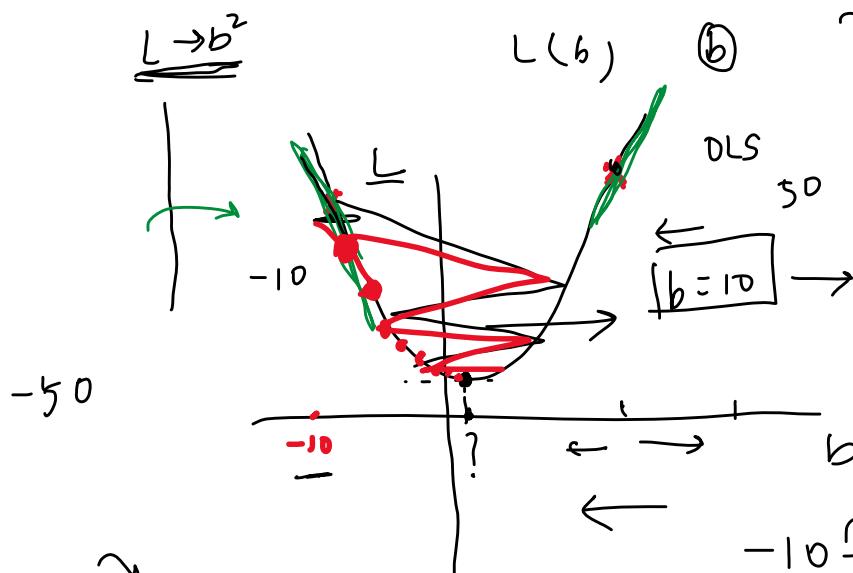
(DLS)  $\rightarrow$

$$\hat{y}_i = mx_i + b$$

$$m = 78.35$$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$L = \sum_{i=1}^n (y_i - mx_i - b)^2$$



$$L = \sum_{i=1}^n (y_i - 78.35 * x_i - b)^2$$

Step 1 - select a random  $b_{\min}$   
 $b_{\min}$   
b increment  
b decrement

$$b = -10 \quad x = 5$$

$$b = -10 \quad \text{Slope} = -ve$$

$$b_{\text{new}} = b_{\text{old}} - \eta \text{ Slope}$$

$$1 \quad b_{\text{new}} = b_{\text{old}} - \eta \text{ Slope}$$

$$b_{\text{new}} = -9.5 - (0.01 * -40) = -9.5 + 0.4 = -9.1$$

$$-10 = (-50)$$

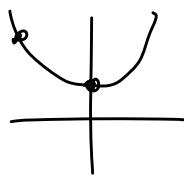
$$40 \\ -40 \\ 10 - (50)$$

$$0.01, 0.0001 \\ -10 \\ x = 50$$

$$b_{\text{new}} = -10 + (0.01 * 50) \\ = -10 + 0.5 = -9.5$$

$$b_{\text{new}} - b_{\text{old}} = \frac{0.0001}{\square}$$

$$\frac{\text{diff } b_{\text{old}} - b_{\text{new}}}{b_{\text{old}}} > 0.0001$$



$$\frac{\text{diff } b_{\text{old}} - b_{\text{new}}}{b_{\text{old}}} > 0.0001$$

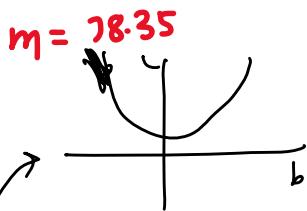
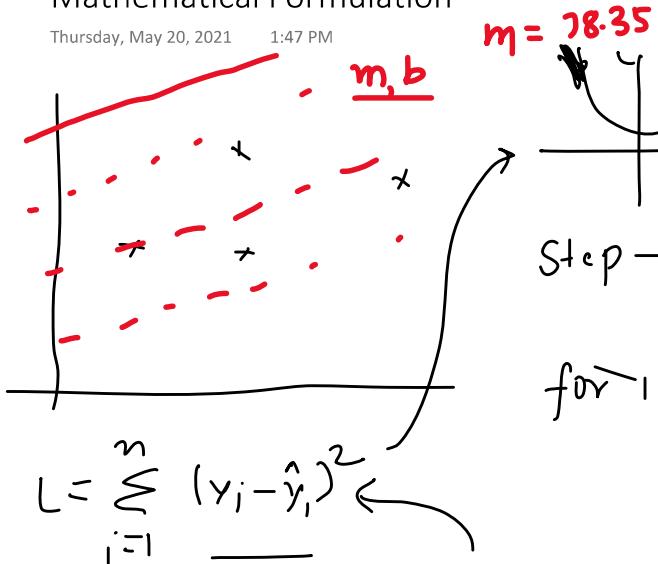
$\Rightarrow 0.000$

2) Iteration  $\rightarrow$  1000, 100,  
epochs

$$\boxed{b_{new} - b_{old} = 0}$$

## Mathematical Formulation

Thursday, May 20, 2021 1:47 PM



Step → Start with a random  
for  $i$  in epochs;

$i = 0$   $\eta, \gamma$

$$\begin{aligned} b &= b \\ b_{\text{new}} &= b_{\text{old}} - \eta \times \text{slope}(b=0) \end{aligned}$$

$$\begin{aligned} \frac{dL}{db} &= \frac{d}{db} \left( \sum_{i=1}^n (y_i - \hat{y}_i)^2 \right) = 2 \sum_{i=1}^n (y_i - mx_i - b) (-1) \\ \frac{d}{db} \sum_{i=1}^n (y_i - mx_i - b)^2 &\xrightarrow{-1} \text{slope} = \boxed{-2 \sum_{i=1}^n (y_i - mx_i - b)} \\ &= -2 \sum_{i=1}^n (y_i - 78.35x_i - 0) \end{aligned}$$

Epoch

$$b_{\text{new}} = b_{\text{old}} - \boxed{\eta \text{slope}(b=0)}$$

$i = 1$

$\text{slope}(b=0)$

$b_{\text{new}}$

Steps SBC

# Example

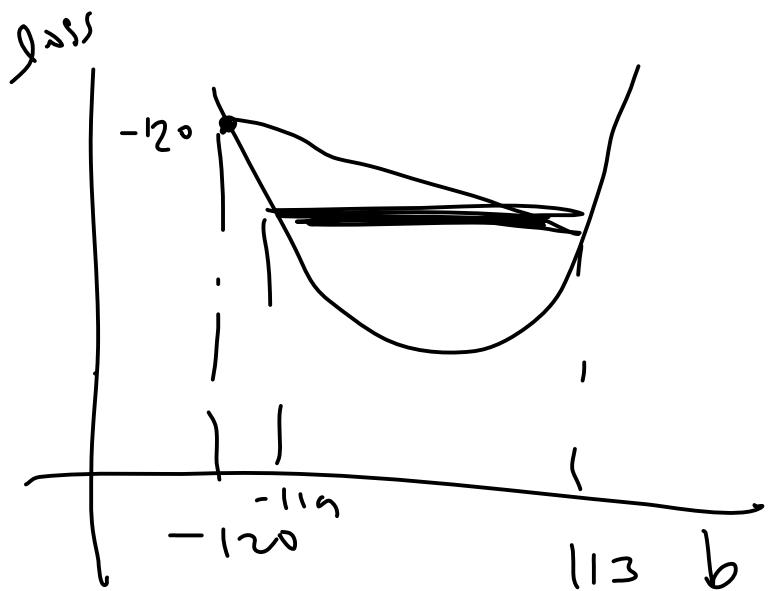
Thursday, May 20, 2021 1:47 PM

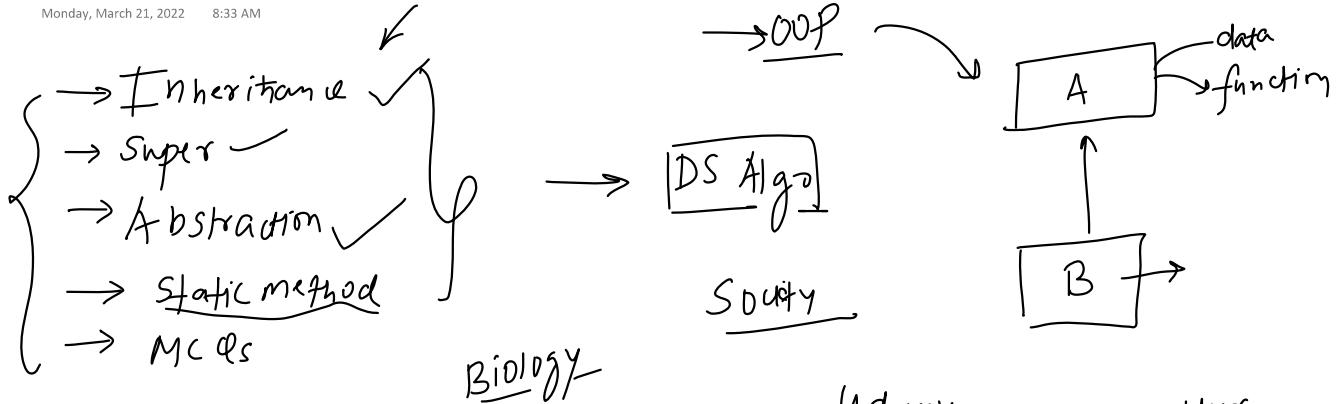
# Code from Scratch

Thursday, May 20, 2021 1:48 PM

# Visualization 1

Thursday, May 20, 2021 1:52 PM



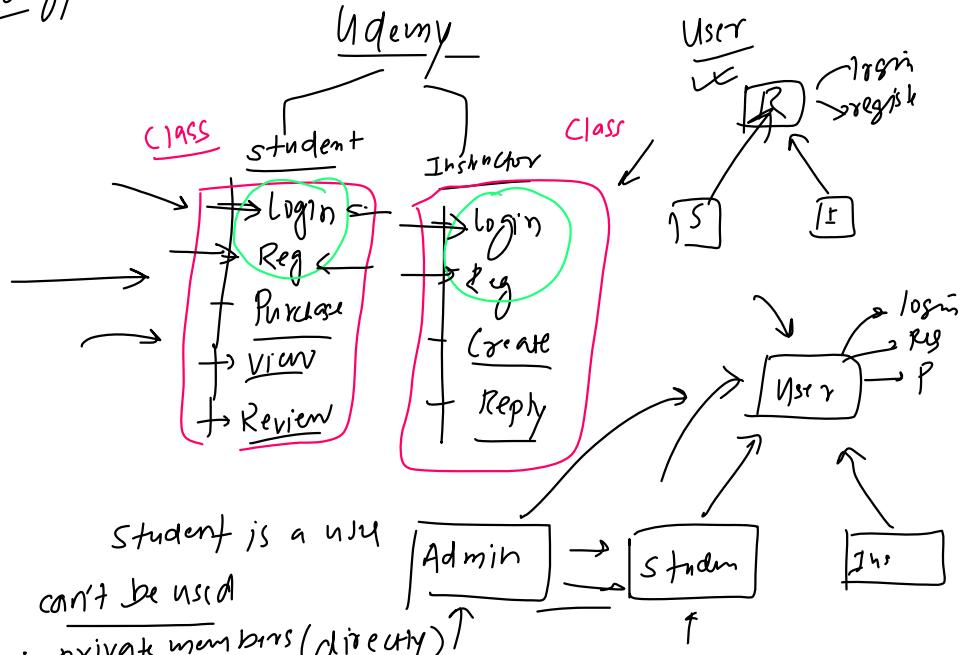


database

→ DRY ✓

Child IS A Parent

mobile is A product

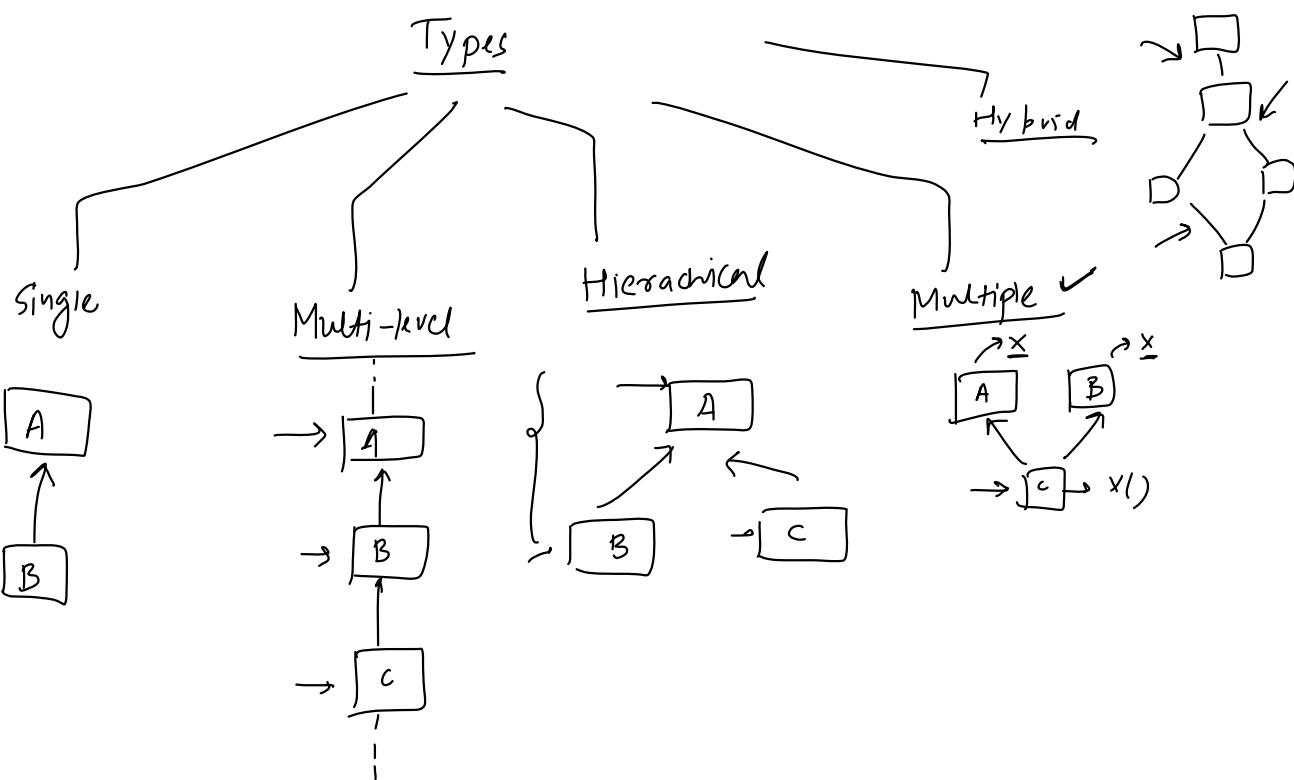


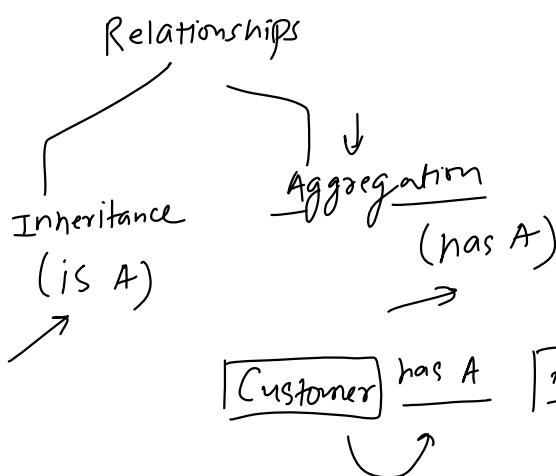
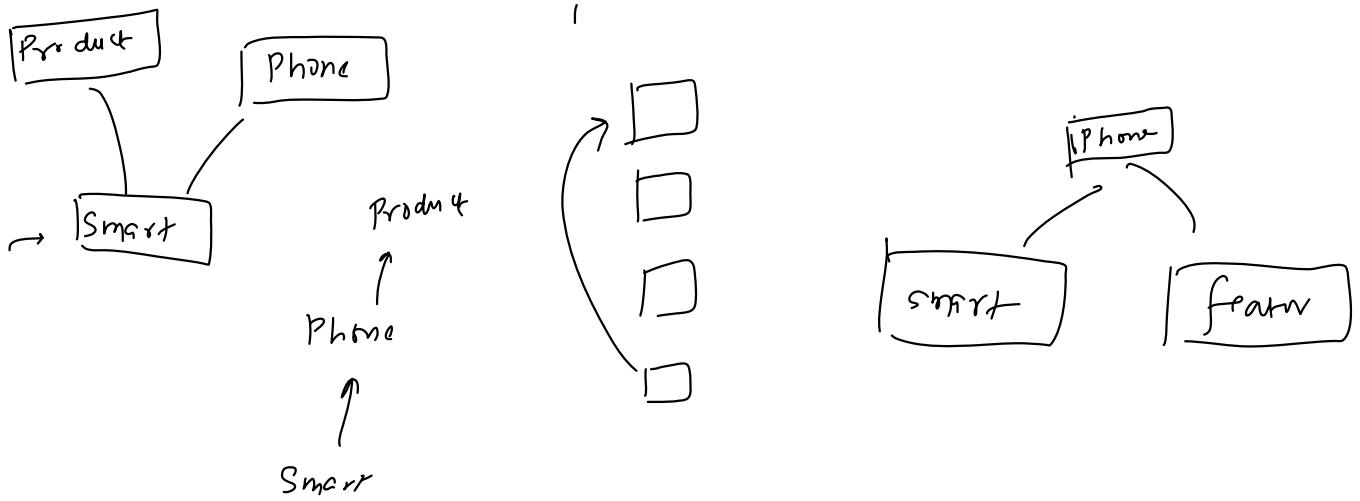
can be used

→ data (property)

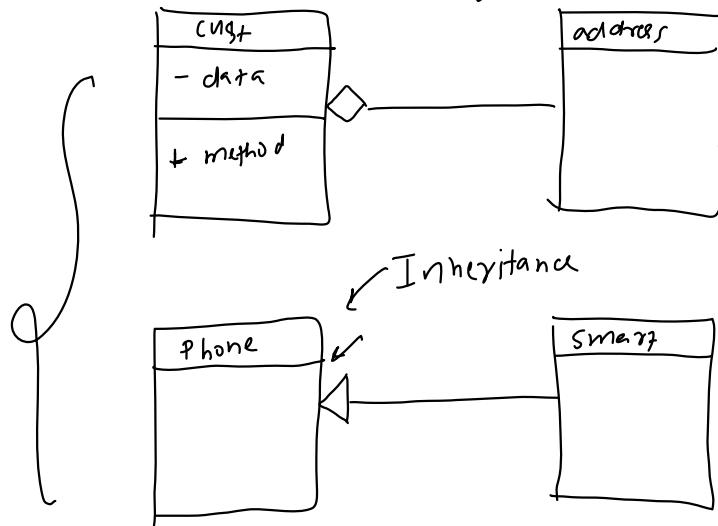
→ function

→ constructor

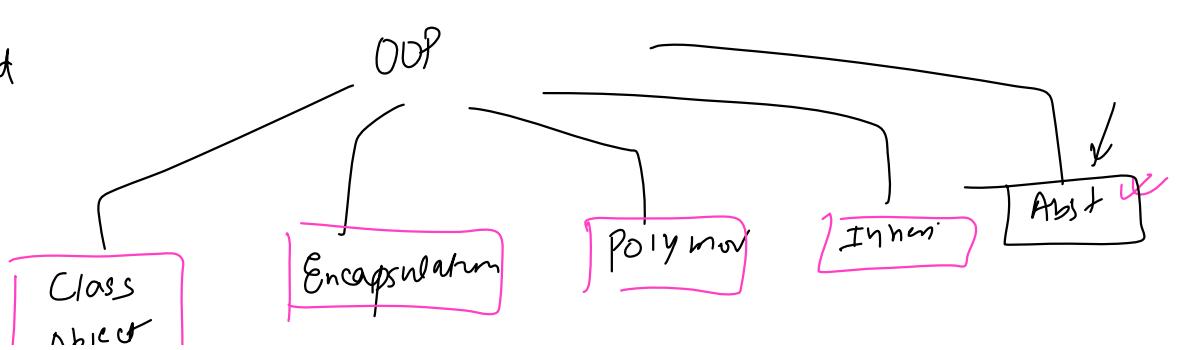




Class diagram



{  
Construct  
static method  
inst vs sm  
super}



Class  
Objec<sup>t</sup>

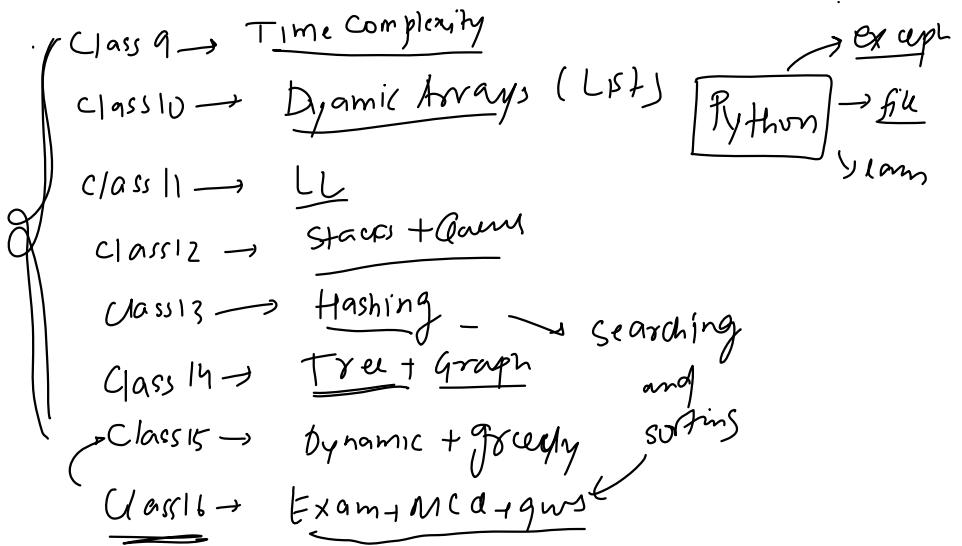
Encapsulation

Principles

Object

(2D) → DS Algo → ⑧ class

{ 17 18 1920 }  
SOL

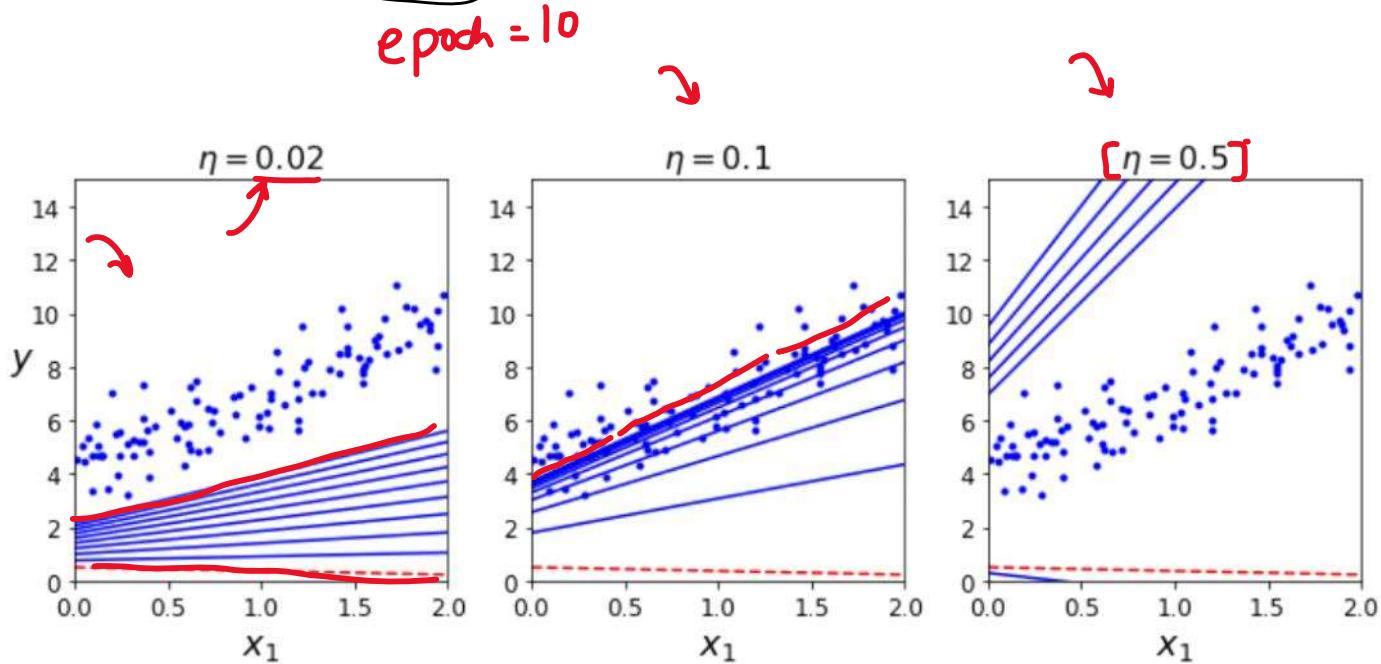


## Few Discussions

Thursday, May 20, 2021 4:47 PM

### 1. Effect of Learning rate

### 2. The universality of Gradient Descent



$$\hat{b} = 0$$

$$b = b_{\text{old}} - \eta \text{ Slope}$$

$$\frac{d L}{d b} = \left[ \sum (y_i - \hat{y}_i)^2 \right] \quad (\text{LR})$$

LDR  $\rightarrow$  [function]

## Adding m into the mix

Thursday, May 20, 2021 1:48 PM

8 steps

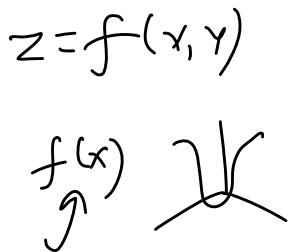
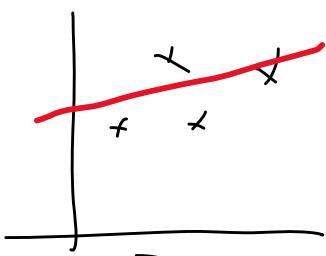
1) init random val for  $m$  and  $b$   
 $m = 1$  and  $b = 0$

2) epochs = 100,  $\eta = 0.01$

for i in epochs:

$$b = b - \eta \boxed{\text{slope}}$$

$$m = m - \eta \boxed{\text{slope}}$$

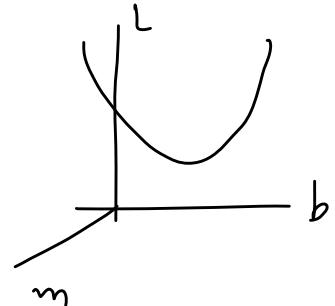


$$b=0$$

$$\frac{L(m, b)}{dL} \quad \boxed{b \text{ slope}} = \frac{\partial L}{\partial b}$$

$$m \text{-slope} = \frac{\partial L}{\partial m}$$

$$\sum (y_i - mx_i - b)^2$$



2D graph of  $L(m, b)$

$$\begin{aligned} \frac{\partial L}{\partial b} &= -2 \sum (y_i - mx_i - b) \\ &= -2 \sum (y_i - mx_i - b) \end{aligned}$$

$\downarrow$

$$= \text{slope}_b \text{ at } b=0$$

$$\begin{aligned} \frac{\partial L}{\partial m} &= 2 \sum (y_i - mx_i - b) \\ &= -2 \sum (y_i - mx_i - b) x_i \\ &\quad \oplus \end{aligned}$$

$\downarrow$

$$\text{slope}_m \text{ at } \boxed{m=1}$$

# Code

Thursday, May 20, 2021 1:48 PM

# Visualization 2

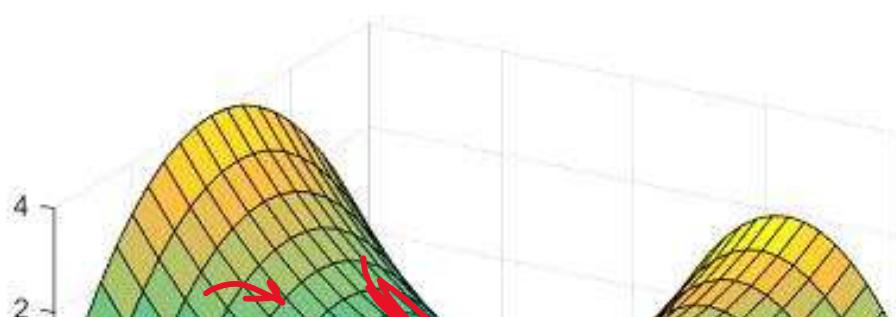
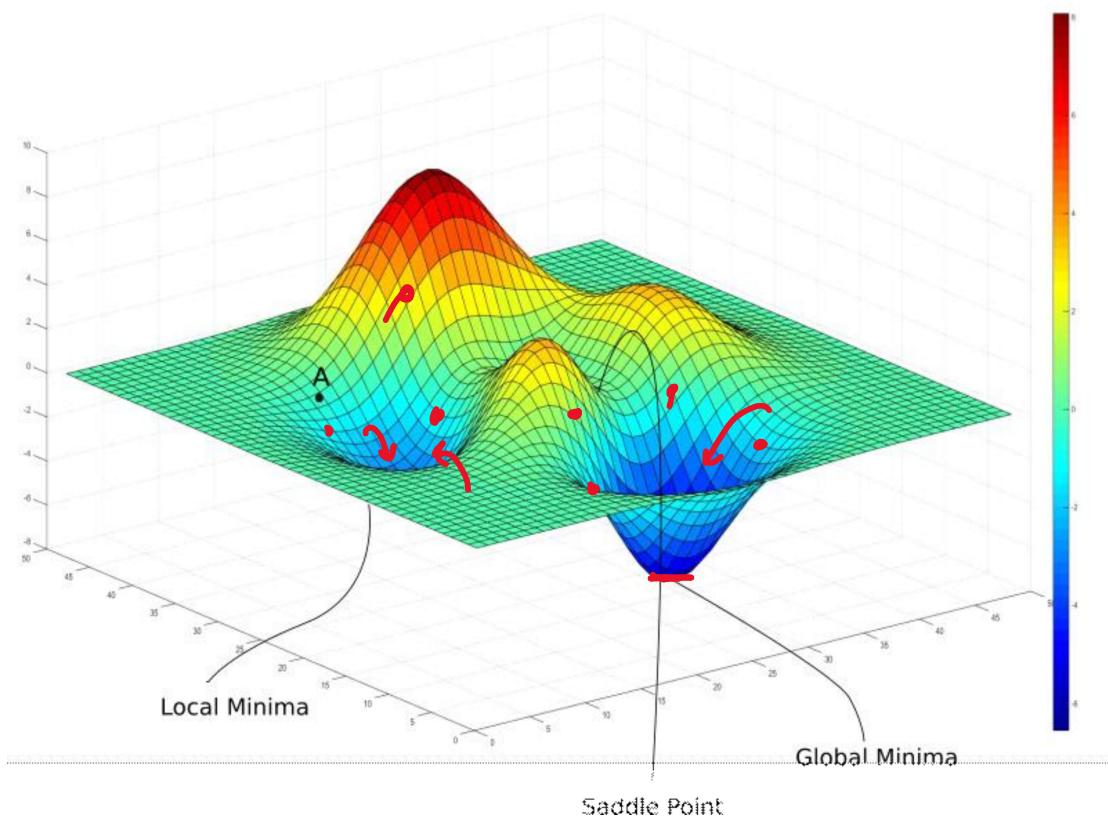
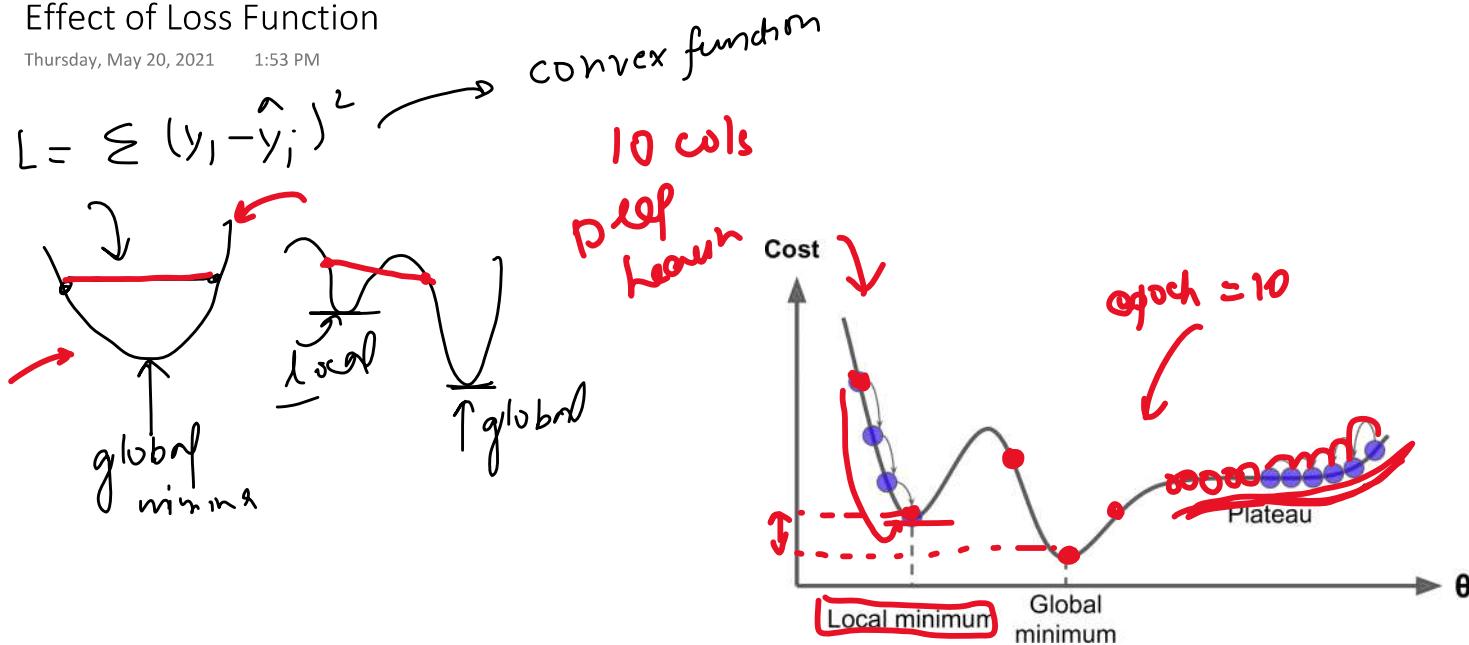
Thursday, May 20, 2021 1:52 PM

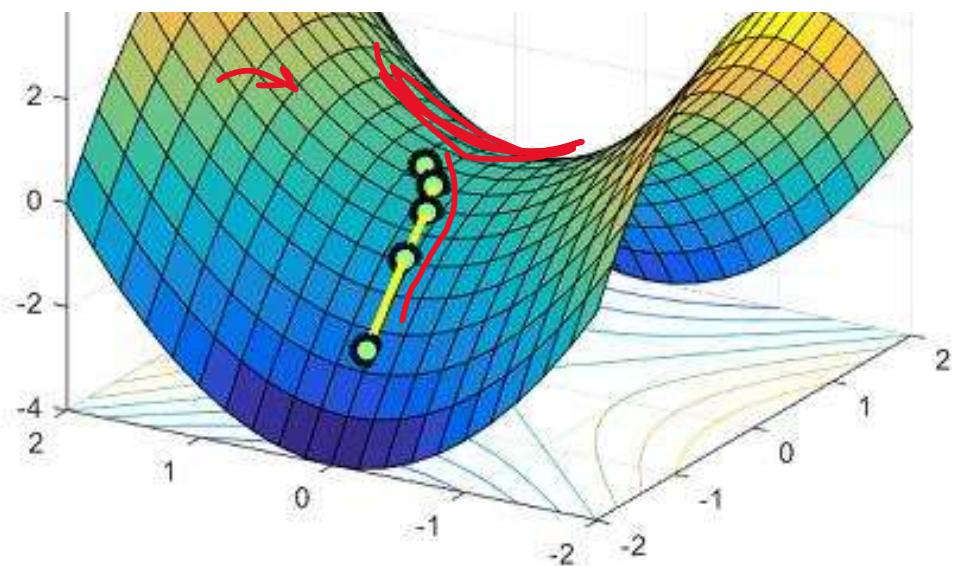
# Effect of Learning Data

Thursday, May 20, 2021 1:53 PM

## Effect of Loss Function

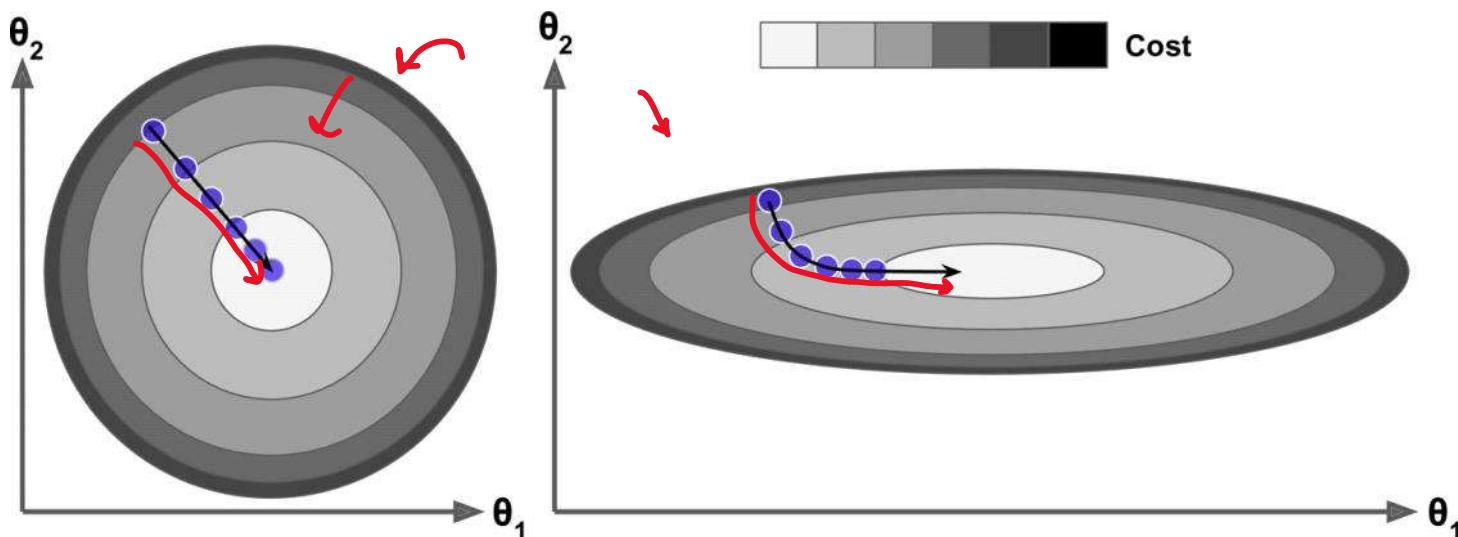
Thursday, May 20, 2021 1:53 PM





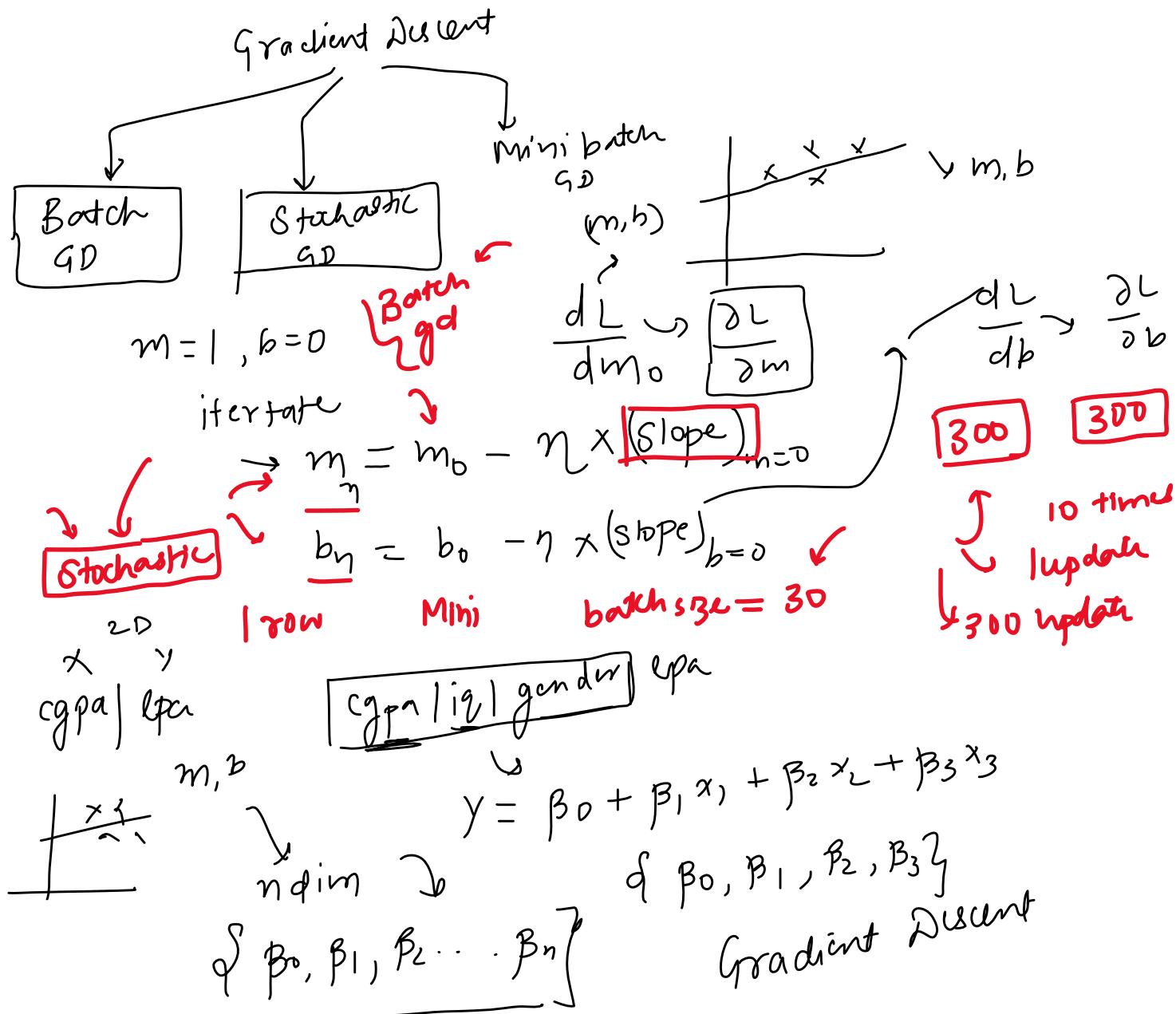
## Effect of Data

Thursday, May 20, 2021 1:53 PM



# Types of Gradient Descent

Saturday, May 22, 2021 1:30 PM



$n$ -dim - dataset 3-cols

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

(lpa) (cgpa) (iq)

$\{\beta_0, \beta_1, \beta_2\}$

$\{m, b\}$

1) Random values

$$\beta_0 = 0, \beta_1, \beta_2 = 1$$

2) epoch = 100, lr = 0.1

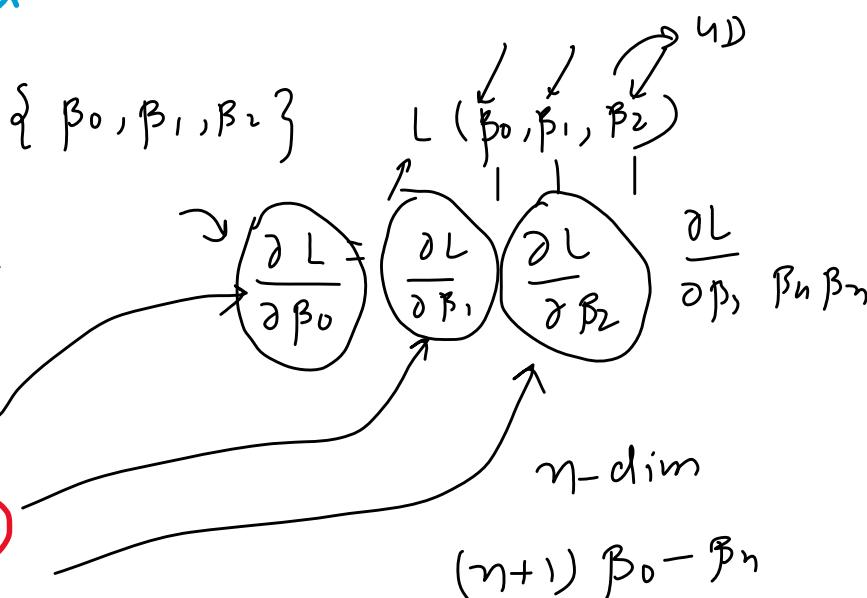
$$\begin{cases} \beta_0 = \beta_0 - \eta \text{ slope} \\ \beta_1 = \beta_1 - \eta \text{ slope} \\ \beta_2 = \beta_2 - \eta \text{ slope} \end{cases}$$

MSE

$$L = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$\{\text{row} = 2, \text{cols} = 2+1\}$

$$\hat{y}_i = \beta_0 +$$



$$L = \frac{1}{2} \left[ (y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2 \right]$$

$x_1$	$x_2$	$y$
8.1	9.3	3.2
7.5	9.5	3.5

$$\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2$$

$$\frac{\partial L}{\partial \beta_0} = \frac{1}{2} \left[ 2(y_1 - \hat{y}_1)(-1) + 2(y_2 - \hat{y}_2)(-1) \right]$$

$$\hat{y}_1 = \beta_0 + \beta_1 x_{11} + \beta_2 x_{12}$$

$$\frac{\partial L}{\partial \beta_0} = -2 \left[ (y_1 - \hat{y}_1) + (y_2 - \hat{y}_2) \right]$$

$$\hat{y}_2 = \beta_0 + \beta_1 x_{21} + \beta_2 x_{22}$$

$$= -2 \left[ \frac{(y_1 - \hat{y}_1)}{n} + \frac{(y_2 - \hat{y}_2)}{n} + \frac{(y_3 - \hat{y}_3)}{n} + \dots + \frac{(y_n - \hat{y}_n)}{n} \right]$$

$$= -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

(353)

$$= \frac{\partial L}{\partial \beta_0}$$

$y_i \rightarrow [353] \bmod$

$x_{train} \rightarrow [x_{11}, x_{12}, \dots, x_{1n}]$

$\beta_0 \rightarrow [x_{21}, x_{22}, \dots, x_{2n}]$

$\beta_1 \rightarrow [x_{11}, x_{21}, \dots, x_{n1}]$

$\beta_2 \rightarrow [x_{12}, x_{22}, \dots, x_{n2}]$

$$L = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$\frac{\partial}{\partial \beta_1} = -x_{11}$

$$L = \frac{1}{2} [(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2]$$

$$L = \frac{1}{2} [(y_1 - \beta_0 - \beta_1 x_{11} - \beta_2 x_{12})^2 + (y_2 - \beta_0 - \beta_1 x_{21} - \beta_2 x_{22})^2]$$

$$\frac{\partial L}{\partial \beta_1} = \frac{1}{2} \sum_{i=1}^n [2(y_i - \hat{y}_i)(-x_{11}) + 2(y_i - \hat{y}_i)(-x_{21})]$$

$$\frac{\partial L}{\partial \beta_1} = \frac{-2}{n} \sum_{i=1}^n [(y_i - \hat{y}_i)(x_{11}) + (y_i - \hat{y}_i)(x_{21}) + (y_i - \hat{y}_i)(x_{31}) + \dots + (y_i - \hat{y}_i)(x_{n1})]$$

$$\frac{\partial L}{\partial \beta_1} = -2 \sum_{i=1}^n (y_i - \hat{y}_i) x_{i1}$$

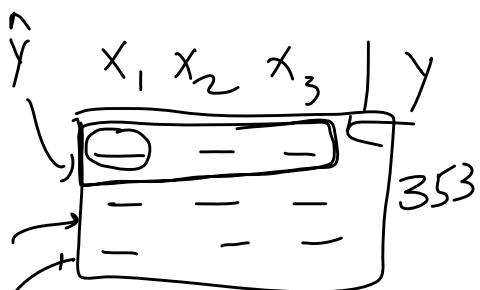
$x_{ij} \rightarrow 1 \text{ col data}$   
 $\beta_1 \rightarrow \text{values of 1 col.}$

$$\frac{\partial L}{\partial \beta_2} = -2 \sum_{i=1}^n (y_i - \hat{y}_i) x_{i2}$$

$m \text{ cols}$   
 $\beta_0 - \beta_m$

$$\frac{\partial L}{\partial \beta_m} = -2 \sum_{i=1}^n (y_i - \hat{y}_i) x_{im}$$

$\text{code}$



$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3}$$

$$\hat{y}_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3}$$

$$\beta_0 + \text{np.dot}(x_{\text{train}}, \underline{\text{coef}}) \quad \hat{y} = \beta_0 + [x_{11} \ x_{12} \ x_{13}] \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

$$\beta_0 + (353, 1) \quad \hat{y} = \text{np.dot}(\text{coef}, x_{\text{train}}) + \beta_0$$

$$\begin{array}{ccccc} x_1 & x_2 & y & \hat{y} & y - \hat{y} \\ \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} & \begin{array}{|c|c|} \hline 5 & 7 \\ \hline 7 & 8 \\ \hline \end{array} & \begin{array}{|c|c|} \hline 5 & 7 \\ \hline 6 & 8 \\ \hline \end{array} & \begin{array}{|c|c|} \hline -1 & -1 \\ \hline \end{array} & \begin{array}{|c|c|} \hline -1 & -1 \\ \hline \end{array} \end{array}$$

$$\frac{\partial L}{\partial \beta_1} = -\frac{2}{n} \sum_{i=1}^n (y_i - \hat{y}_i) x_{i1} = -\frac{2}{n} \begin{bmatrix} [-1 \ -1] & [1 \\ 3] \end{bmatrix}$$

$$\begin{aligned} \frac{\partial L}{\partial \beta_2} &= [y - \hat{y}] \begin{bmatrix} 2 \\ 4 \end{bmatrix} & [-4 + -4] && (1, 2) \\ &= [-1 \ -1] \begin{bmatrix} 2 \\ 4 \end{bmatrix} \times -\frac{2}{n} & -8 &= 8 & \begin{bmatrix} x \\ y \end{bmatrix} \\ && \begin{bmatrix} (1, 2) & (2, 2) \\ [y - \hat{y}] & [1 \ 2] \end{bmatrix} \times -\frac{2}{n} && \end{aligned}$$

$$\frac{\partial L}{\partial \beta_1} \dots \frac{\partial L}{\partial \beta_{10}} = \left[ \underline{(y_i - \hat{y}_i) \times \text{train}} \right] \times -\frac{2}{n}$$

$y_{\text{train}}$

$$y_i = 353$$

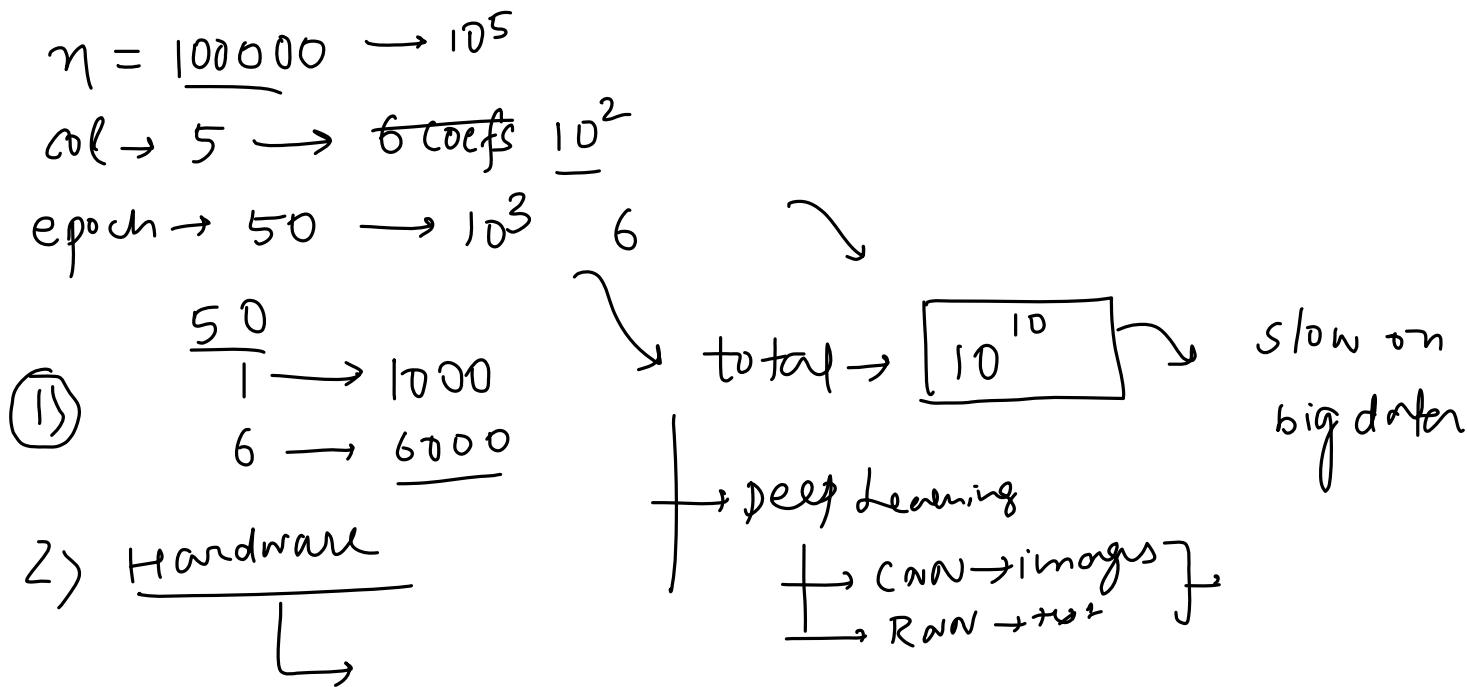
$$(353, 1) \quad (353, 10) \rightarrow (1, 353)$$

$$\begin{array}{cc} (1, 353) & (353, 10) \\ \downarrow & \downarrow \\ (1, 10) \times \begin{bmatrix} -2 \\ \frac{2}{n} \end{bmatrix} & = (1, 10) \end{array}$$

coef-clm

## The Problem with Batch GD

Tuesday, May 25, 2021 6:43 AM

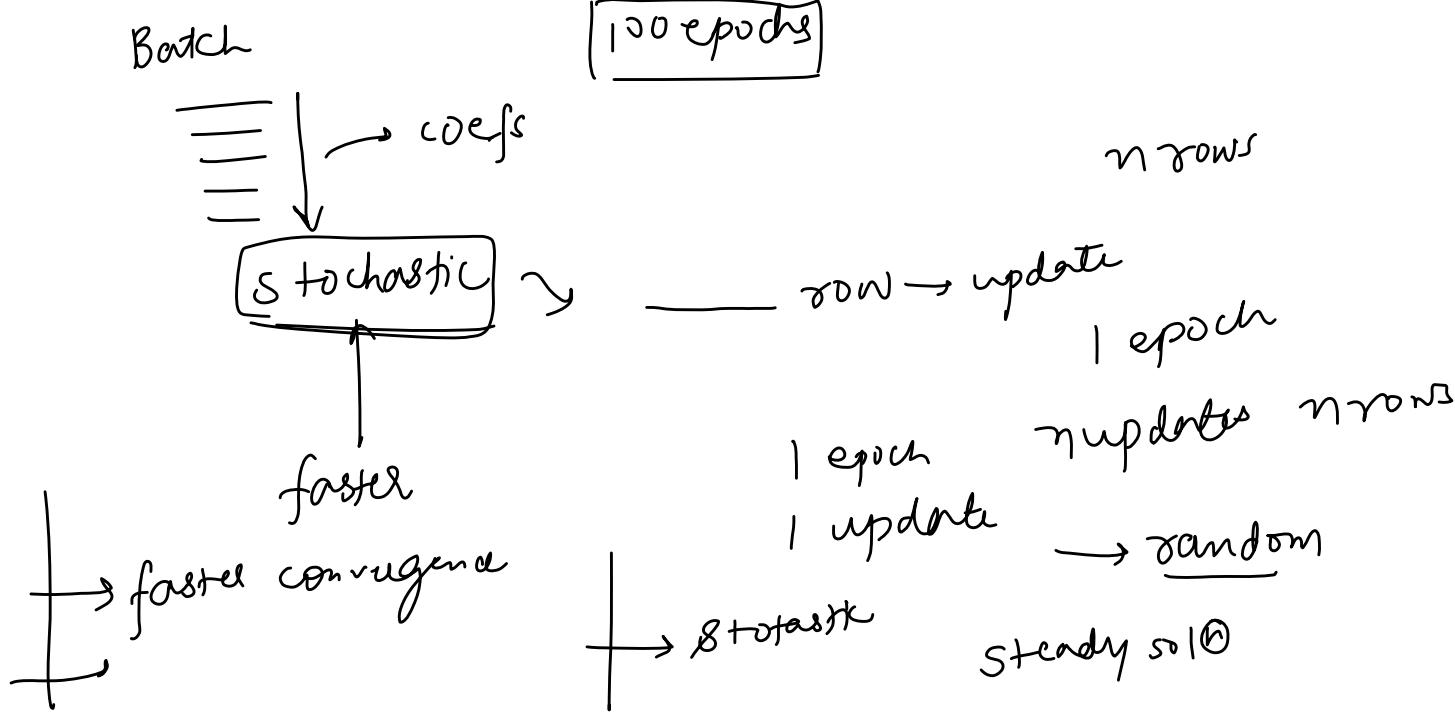


## Stochastic GD

Tuesday, May 25, 2021 6:44 AM

5, 10

100 epochs

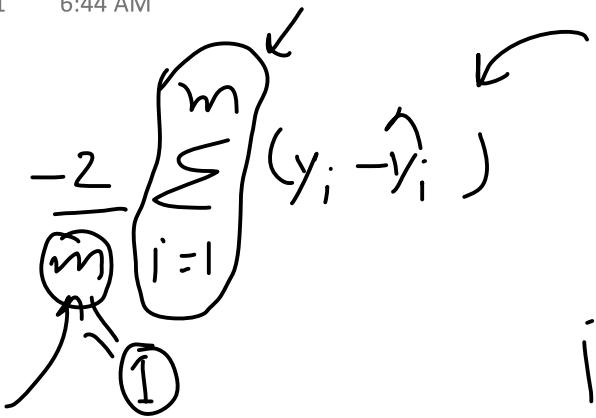


# Code

Tuesday, May 25, 2021 6:44 AM

$$\frac{\partial L}{\partial \beta_0} = -2 \sum_{i=1}^m (y_i - \hat{y}_i)$$

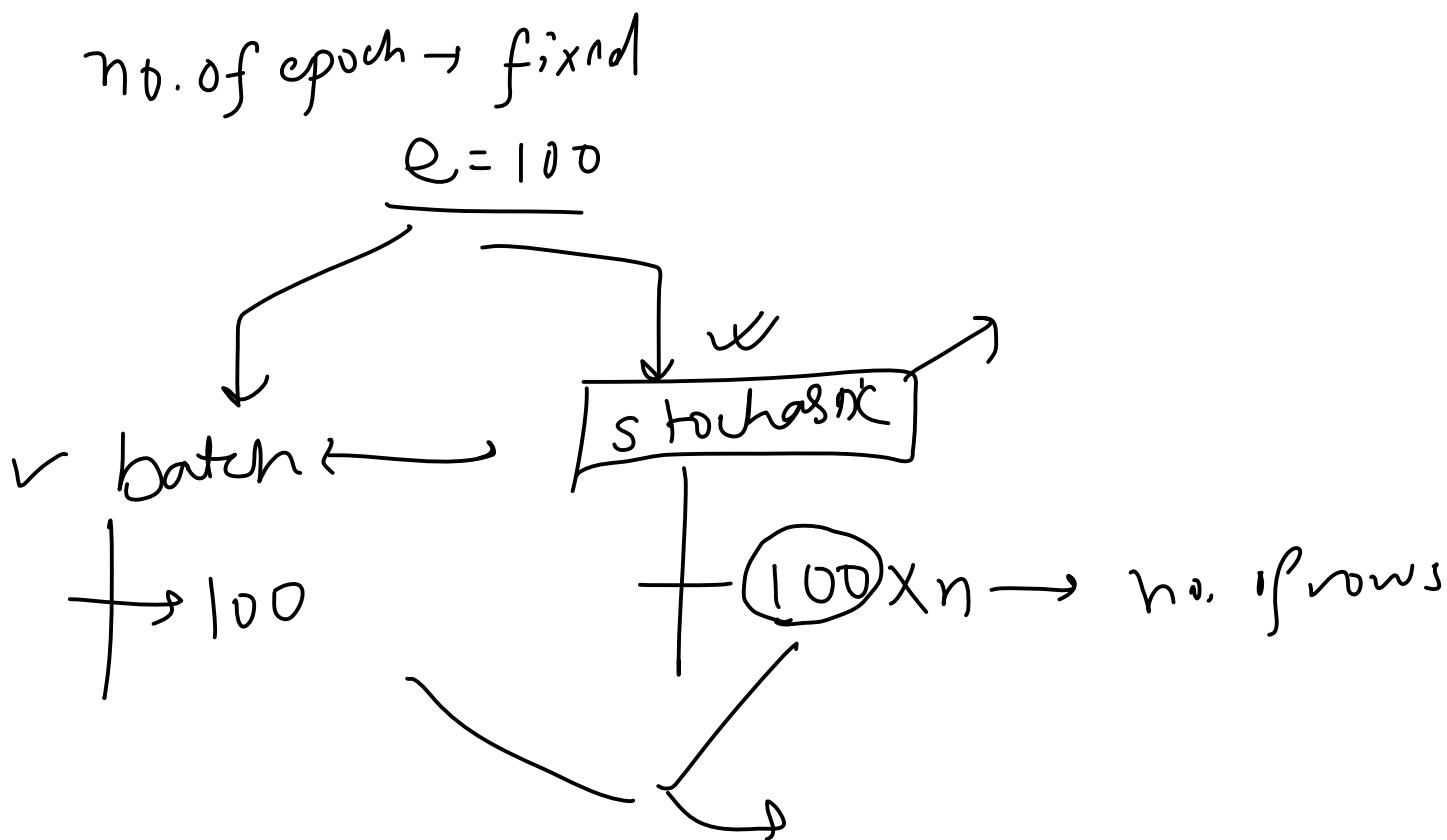
$i = id X$



$$= -2 \underline{(y_i - \hat{y}_i)}$$

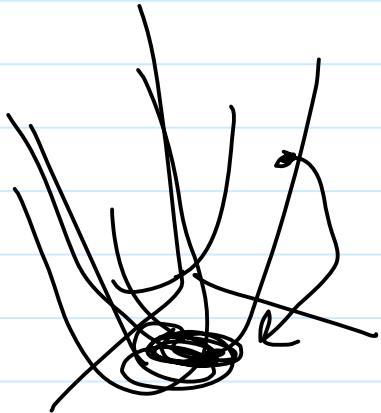
# Time Comparison

Tuesday, May 25, 2021 12:39 PM



# Visualizations

Tuesday, May 25, 2021 6:44 AM

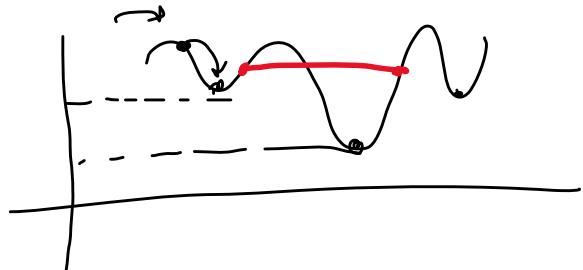


## When to use Stochastic GD

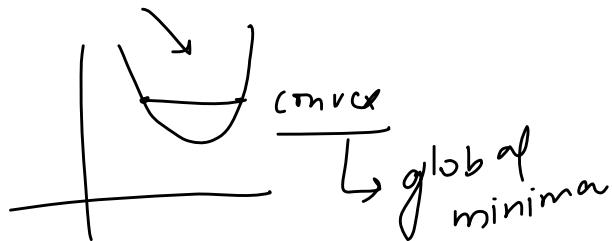
Tuesday, May 25, 2021 6:44 AM

1) Big data  $\rightarrow$  SGD

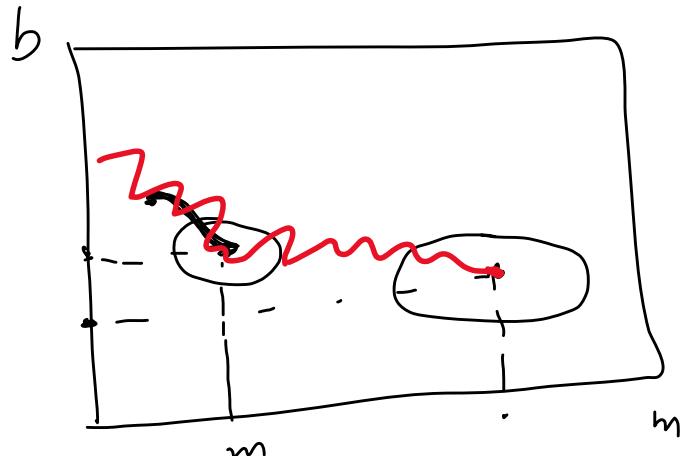
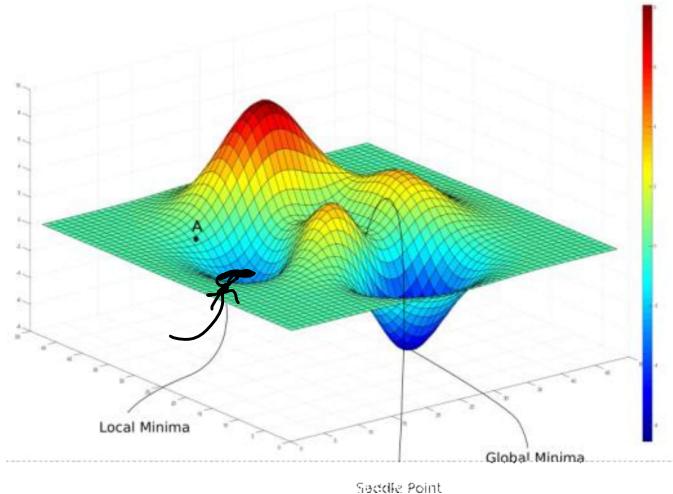
2) Non convex function



non convex



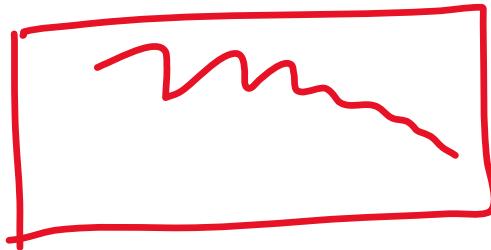
learning



# Learning Schedules

Tuesday, May 25, 2021 7:00 AM

DL



$$\eta = 100$$
$$\text{epoch} = 1$$

$$\rightarrow$$

$$lr = 0.1$$

$$lr = 0.03$$

```
t0, t1 = 5.50  
def learning_rate(t):  
    return t0/(t + t1)
```

```
for i in range(epochs):  
    for j in range(X.shape[0]):
```

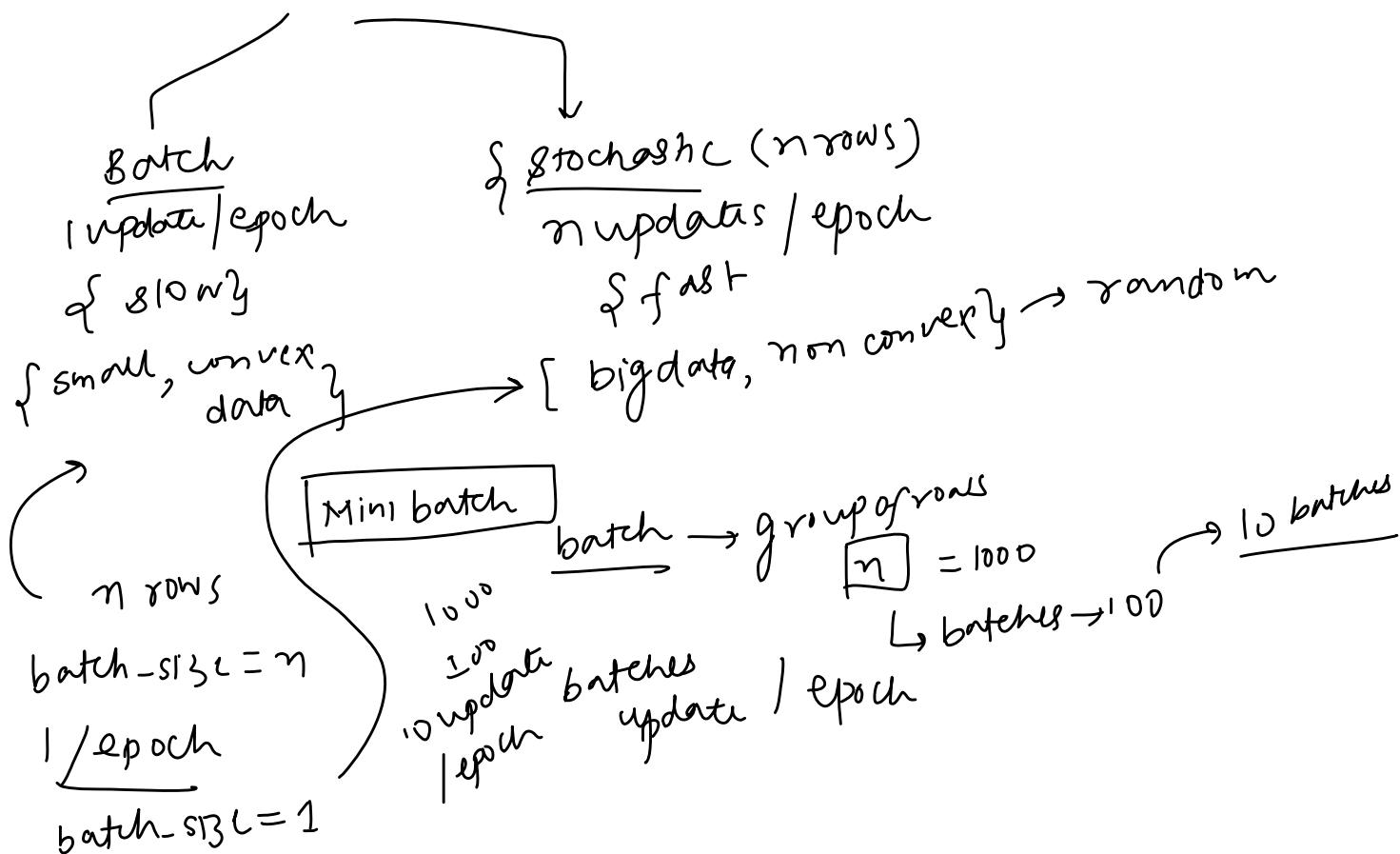
```
        lr = learning_rate(i * X.shape[0] + j)
```

# Sklearn Implementation

Tuesday, May 25, 2021 2:16 PM

## Mini-Batch Gradient Descent

Wednesday, May 26, 2021 4:47 PM



# Code

Wednesday, May 26, 2021 4:47 PM

# Visualization

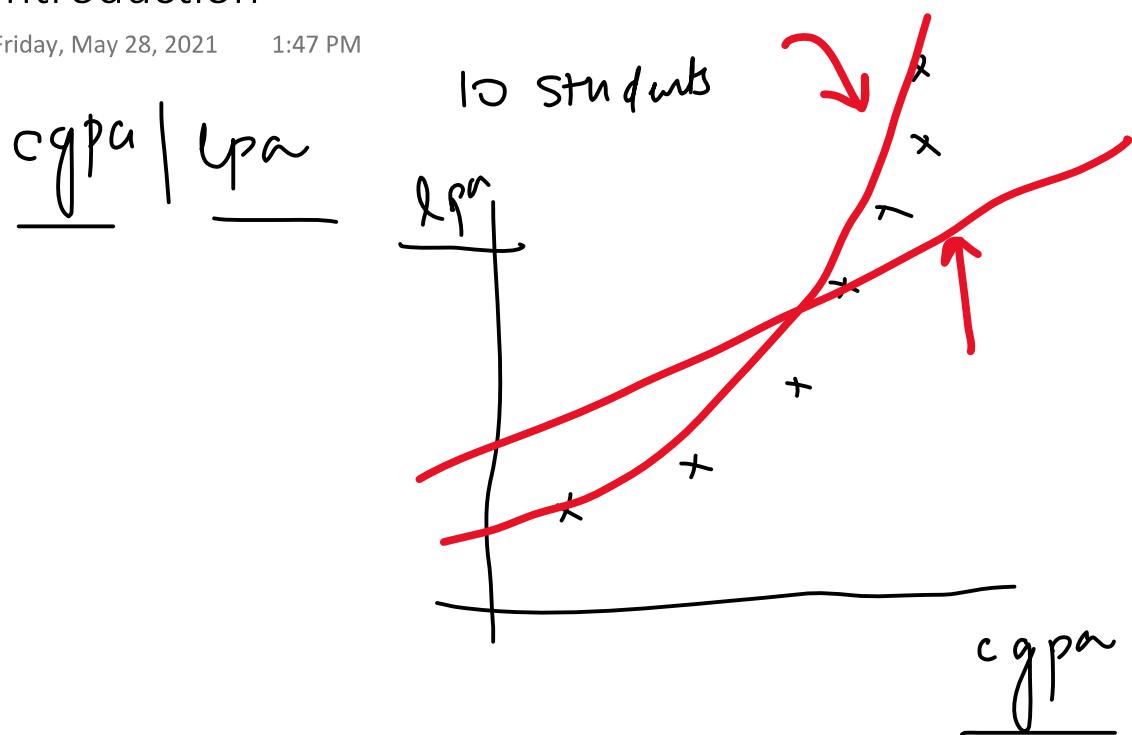
Wednesday, May 26, 2021 4:47 PM

# Sklearn Implmentation

Wednesday, May 26, 2021 4:47 PM

# Introduction

Friday, May 28, 2021 1:47 PM

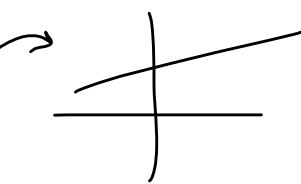


# Intuition

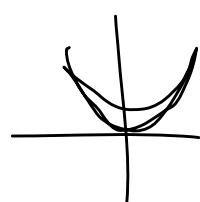
Friday, May 28, 2021 1:47 PM

polynomials

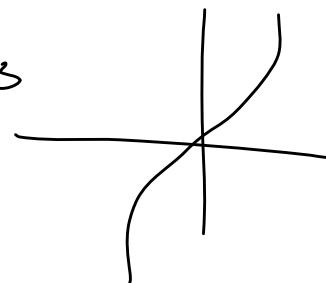
$$y = mx + b$$



$$y = x^2$$



$$y = x^3$$



$$x^3 + x^2$$

$$y = \textcircled{a}x^4 + \textcircled{b}x^3 + \textcircled{c}x^2 + dx + e$$

degree

$$\begin{array}{r} 1 \\ -x \\ \hline 2 \\ \end{array} \quad \begin{array}{r} y \\ \curvearrowright \\ \end{array} \quad \rightarrow \quad y = mx + b$$

$$\begin{array}{r} 3 \\ x^0 \quad | \quad x^1 \quad | \quad x^2 \\ \hline 1 \quad | \quad 2 \quad | \quad 4 \\ \hline 1 \quad | \quad 3 \quad | \quad 9 \end{array}$$

3

transform  
polynomial

$$y = \beta_0 + \beta_1 x^0 + \beta_2 x^2$$

# Code

Friday, May 28, 2021 1:48 PM

# Multiple Polynomial Regression

Friday, May 28, 2021 1:48 PM

# Why is it linear

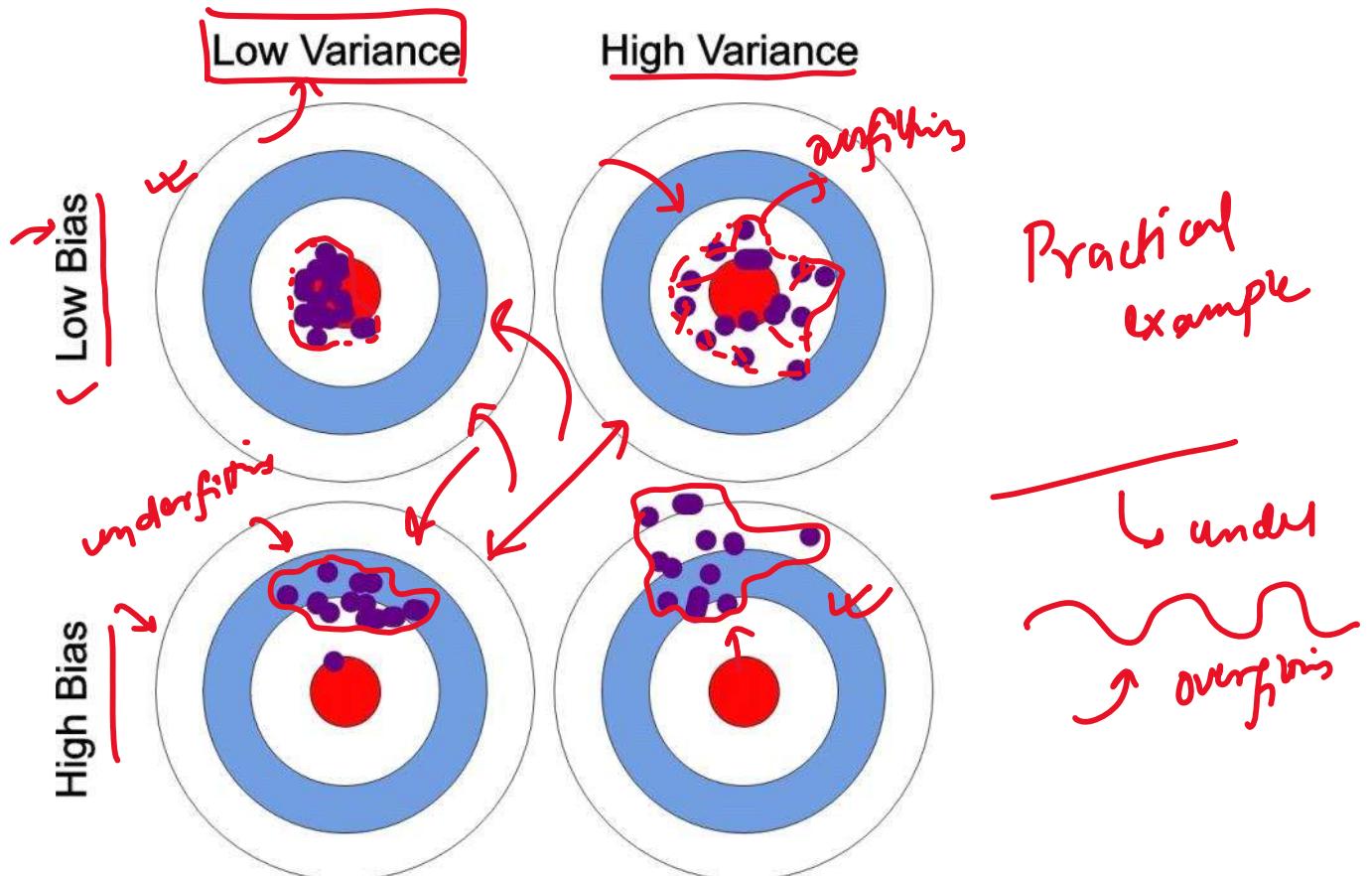
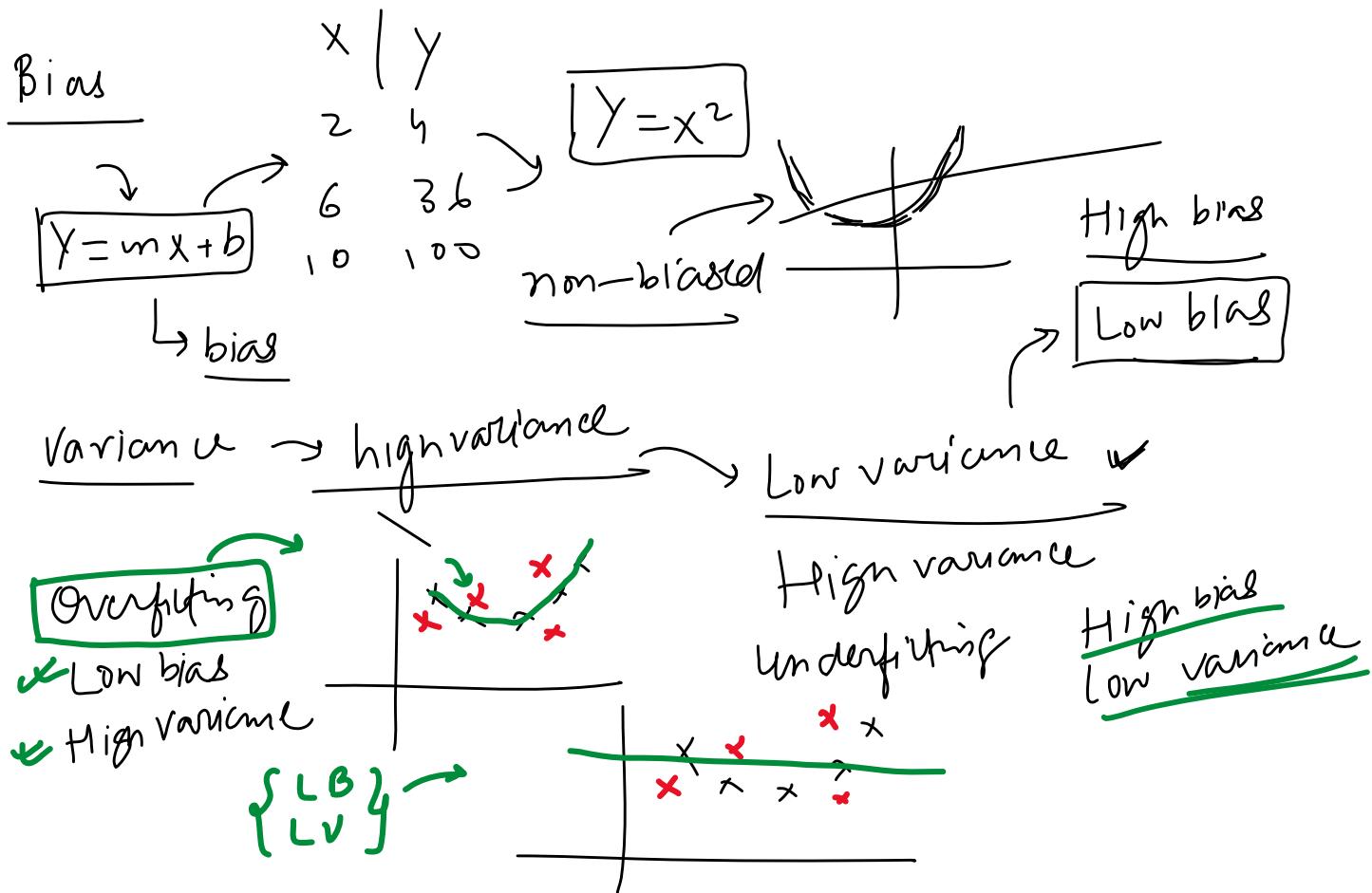
Friday, May 28, 2021 1:48 PM

# When to use Polynomial Regression

Friday, May 28, 2021 1:48 PM

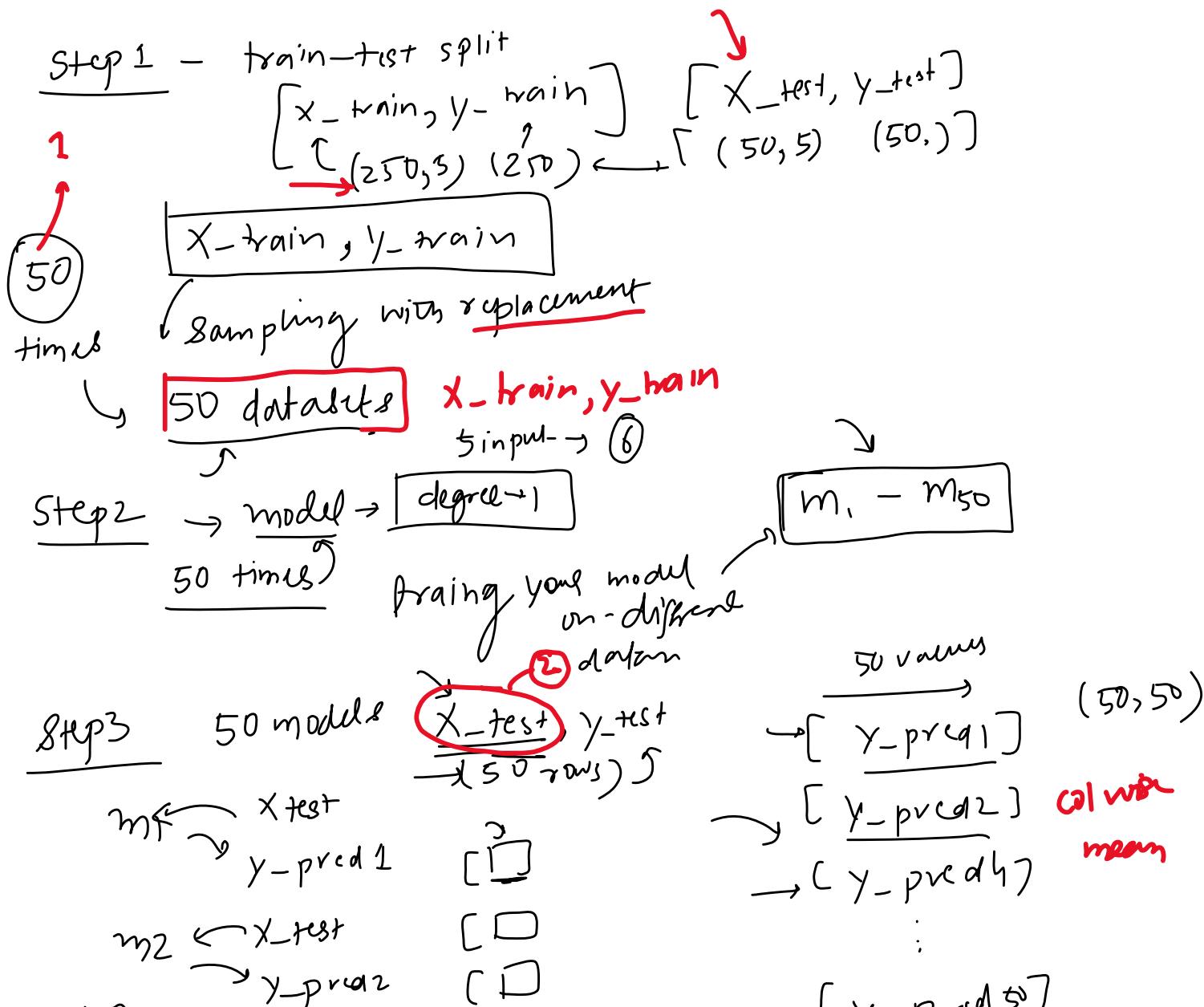
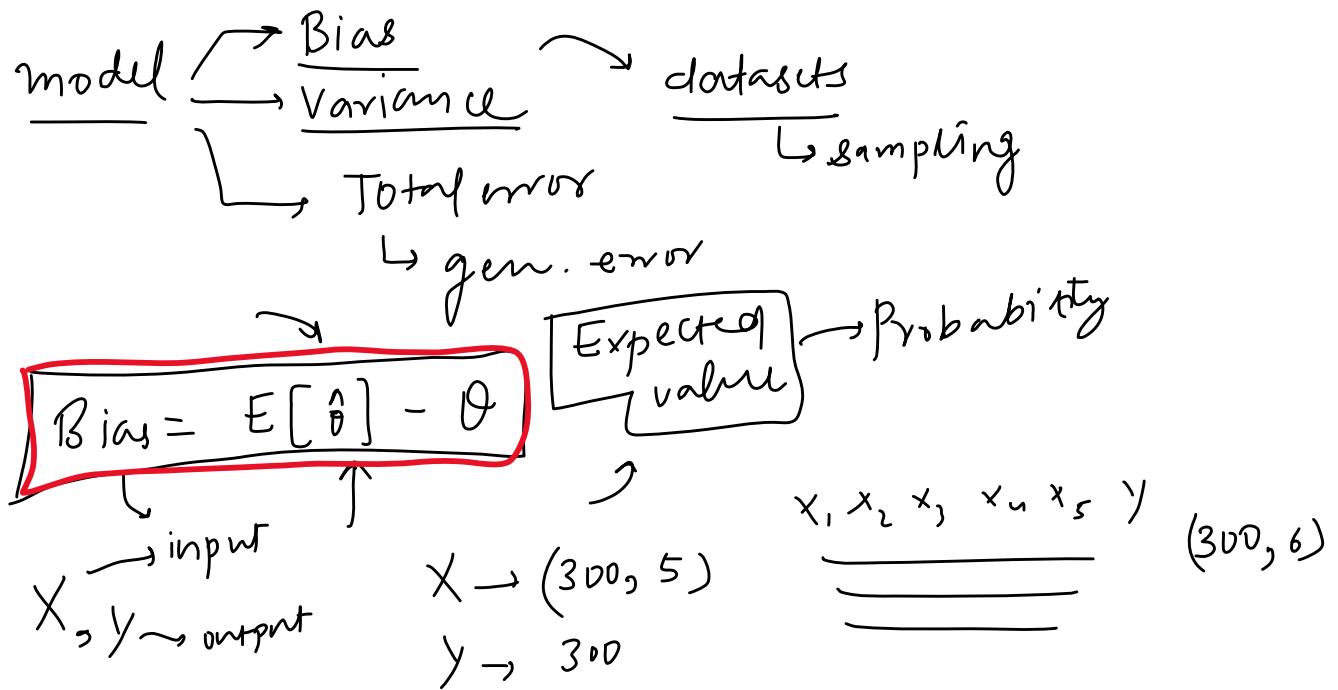
# Bias And Variance

Tuesday, June 1, 2021 12:20 PM



# How to calculate Bias and Variance

Tuesday, June 1, 2021 1:44 PM



$$\begin{aligned}
 m_2 &\leftarrow \lambda^{-1} \sigma \\
 50 &\rightarrow y_{\text{pred2}} \\
 \boxed{\begin{array}{c} 50 \\ 50 \end{array}} &\quad L \quad \boxed{\begin{array}{c} \square \\ \square \end{array}} \\
 &\quad \boxed{\begin{array}{c} \rightarrow \\ | \end{array}} \quad \boxed{\begin{array}{c} \rightarrow \\ | \end{array}} \\
 &\quad \boxed{\begin{array}{c} \theta \\ \dots \end{array}} \quad \boxed{\begin{array}{c} [y - \text{pred2}] \\ \text{mean\_prediction} \end{array}} \quad \boxed{\begin{array}{c} 50 \text{ values} \\ \dots \end{array}}
 \end{aligned}$$

$$\frac{(y_{\text{pred2}} - y_{\text{test}})^2}{50} = \boxed{\text{bias}}$$

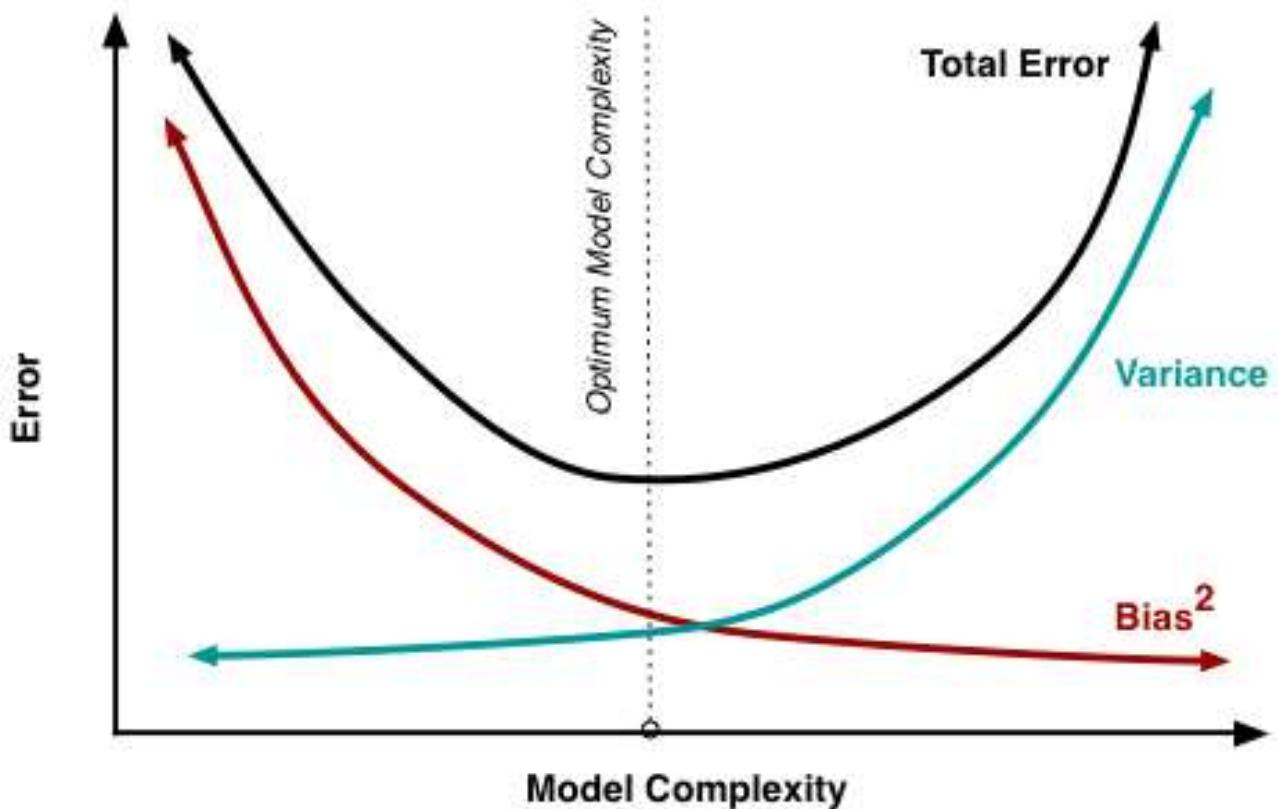
Variance  
 \(\underbrace{\dots}\_{\text{Variance}}\)

# Code

Monday, May 31, 2021 9:26 AM

# Bias Variance Decomposition for Squared Error

Monday, May 31, 2021 9:26 AM

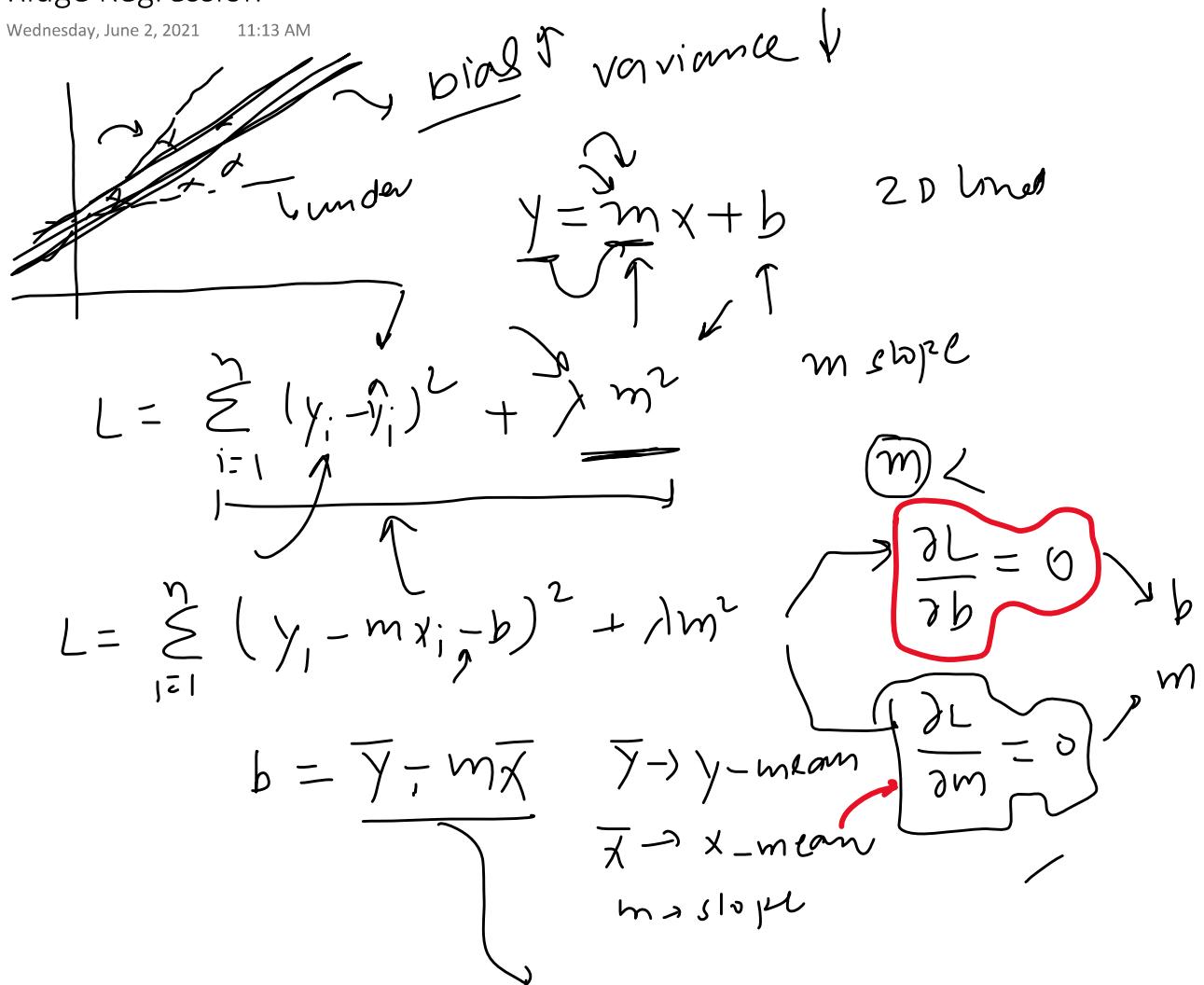


# Relationship between Bias and Variance

Monday, May 31, 2021 9:26 AM

## Ridge Regression

Wednesday, June 2, 2021 11:13 AM



$$L = \sum_{i=1}^n (y_i - mx_i - \bar{y} + \bar{mx})^2 + \lambda m^2$$

$$\frac{\partial L}{\partial m} = 2 \sum_{i=1}^n (y_i - mx_i - \bar{y} + \bar{mx}) (-x_i + \bar{x}) + 2\lambda m = 0$$

$$= -2 \sum_{i=1}^n (y_i - \bar{y} - mx_i + \bar{mx}) (x_i - \bar{x}) + 2\lambda m = 0$$

$$= \lambda n - \sum_{i=1}^n [(y_i - \bar{y}) - m(x_i - \bar{x})] (x_i - \bar{x}) = 0$$

n - ... - ... - ... = 0 - n

$$\lambda m - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) - m(x_i - \bar{x})^2 = 0$$

$$\lambda m - \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x}) + m \sum_{i=1}^n (x_i - \bar{x})^2 = 0$$

$$\lambda m + m \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})$$

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$

hyperparam.  
alpha

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$\lambda = 0 = \lambda = 10,000$

$b = \bar{y} - m\bar{x}$

# Ridge Regression for 2D data

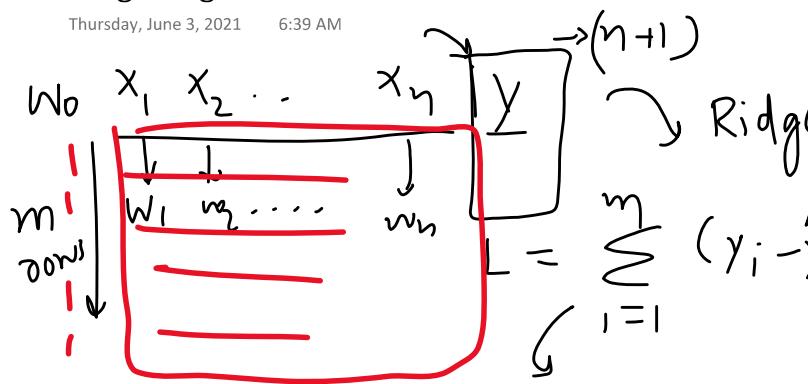
Thursday, June 3, 2021 6:38 AM

# Code

Thursday, June 3, 2021 6:39 AM

## Ridge Regression for nD data

Thursday, June 3, 2021 6:39 AM



$$= (\mathbf{x}\mathbf{w} - \mathbf{y})^\top (\mathbf{x}\mathbf{w} - \mathbf{y})$$

$$\begin{aligned} & \text{m vals} \\ \mathbf{y} &= \begin{bmatrix} y \\ \vdots \\ y \end{bmatrix} \quad (\mathbf{n} \times 1) \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \quad \mathbf{x} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \end{aligned}$$

Normal LR  $\rightarrow$  Ridge

$$L = (\mathbf{x}\mathbf{w} - \mathbf{y})^\top (\mathbf{x}\mathbf{w} - \mathbf{y}) + \lambda \|\mathbf{w}\|^2$$

$$[\mathbf{w}_0 \mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n] \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \underbrace{\mathbf{w}^\top \mathbf{w}}_{\lambda(w_0^2 + w_1^2 + w_2^2 + \dots + w_n^2)}$$

$$(\mathbf{a} - \mathbf{b})^\top = \mathbf{a}^\top - \mathbf{b}^\top$$

$$L = \underbrace{(\mathbf{x}\mathbf{w} - \mathbf{y})^\top (\mathbf{x}\mathbf{w} - \mathbf{y})}_{\frac{dL}{dw}} + \lambda \mathbf{w}^\top \mathbf{w}$$

$$L = [(\mathbf{x}\mathbf{w})^\top - (\mathbf{y})^\top] (\mathbf{x}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^\top \mathbf{w}$$

$$= (\mathbf{w}^\top \mathbf{x}^\top - \mathbf{y}^\top) (\mathbf{x}\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^\top \mathbf{w}$$

$$= \mathbf{w}^\top \mathbf{x}^\top \mathbf{x}\mathbf{w} - \mathbf{w}^\top \mathbf{x}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{x}\mathbf{w} + \mathbf{y}^\top \mathbf{y} + \lambda \mathbf{w}^\top \mathbf{w}$$

$$L = \underbrace{w^T X^T X w}_{\text{red bracket}} - 2 \underbrace{w^T X^T y}_{\text{red bracket}} + \underbrace{y^T y}_{\text{red bracket}} + \lambda \underbrace{w^T w}_{\text{red bracket}}$$

$$\frac{dL}{dw} = \cancel{\rho} X^T X w - \cancel{\lambda} X^T y + 0 + \cancel{\lambda} \lambda w = 0$$

$$X^T X w + \cancel{\lambda} w = X^T y$$

$$(X^T X + \lambda I) w = X^T y$$

3 (4x4)  
(n x 1, n x 1)

$$\begin{bmatrix} \ddots & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & 0 & \ddots \end{bmatrix} \boxed{w = (X^T X + \lambda I)^{-1} X^T y}$$

$$w = (X^T X)^{-1} X^T y \quad \boxed{[ \quad ]}$$

# Code

Thursday, June 3, 2021 6:39 AM

## Ridge Regression using Gradient Descent

Friday, June 4, 2021 2:17 PM

Vector form loss

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \xrightarrow{\text{brace}} \quad L = (Xw - y)^T (Xw - y) + \lambda \|w\|^2$$

$\underbrace{L = (Xw - y)^T (Xw - y) + \lambda w^T w}_{\substack{\text{m rows} \\ X \\ Y \\ w \\ (n+1)}}$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & x_{23} & \dots & x_{2n} \\ 1 & x_{m1} & x_{m2} & x_{m3} & \dots & x_{mn} \end{bmatrix} \quad \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad \begin{bmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$w_0, w_1, \dots, w_n$  (parameters)

$$w_0 = w_0 - \eta \frac{\partial L}{\partial w_0} \quad w_1 = w_1 - \eta \frac{\partial L}{\partial w_1} \quad \dots \quad w_n = w_n - \eta \frac{\partial L}{\partial w_n}$$

$$w_{\text{new}} = \underline{w_{\text{old}} - \eta \left[ \frac{\Delta L}{\Delta w} \right]} \rightarrow \text{gradient} \left[ \begin{array}{c} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{array} \right]$$

$$\begin{aligned} L &= \frac{1}{2} (Xw - y)^T (Xw - y) + \frac{1}{2} \lambda w^T w \\ &= \frac{1}{2} (w^T X^T - y^T) (Xw - y) + \frac{1}{2} \lambda w^T w \\ &= \frac{1}{2} \left[ w^T X^T X w - \cancel{w^T X^T y} - \cancel{y^T w X} + y^T y \right] + \frac{1}{2} \lambda w^T w \end{aligned}$$

$$= \frac{1}{2} L^{\text{new}} \lambda' \lambda' \quad \boxed{L' \quad \boxed{\lambda'}}$$

$$= \frac{1}{2} \left[ \underbrace{w^T x^T x w}_{2w^T x^T y} - \cancel{2y^T w x} + y^T y \right] + \frac{1}{2} \lambda \underline{w^T w}$$

$$\frac{dL}{dw} = \frac{1}{2} \left[ \cancel{2x^T x w} - \cancel{2y^T x} \right] + \frac{1}{2} \cancel{2\lambda w}$$

$$= \boxed{x^T x w - y^T x + \lambda w} = \frac{dL}{dw} \left( \frac{\Delta L}{\Delta w} \right)$$

$w = \begin{bmatrix} w_0 & w_1 & \dots & w_n \\ 0 & 1 & \dots & 1_m \end{bmatrix}$  starting

Epochs  $w = w - \eta \frac{dL}{dw}$

$w \rightarrow$  final answer

epoch times

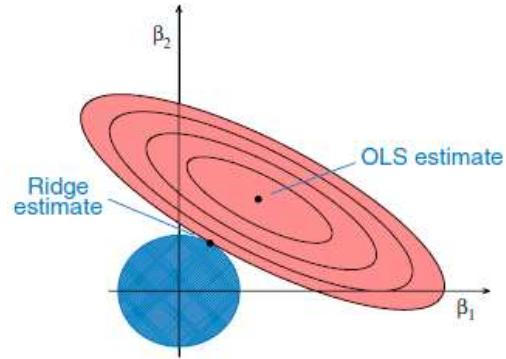
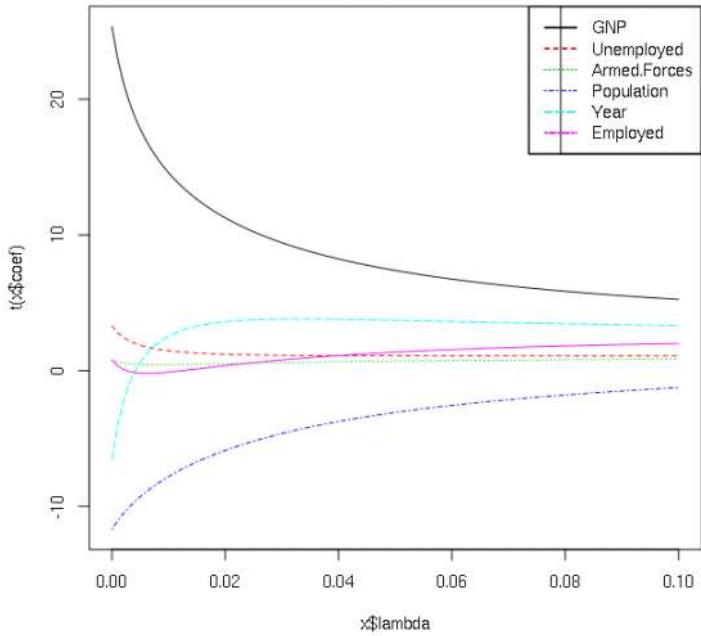
$$\boxed{w = w - n \frac{dL}{dw}}$$

$$\boxed{\frac{dL}{dw} = \underline{x^T x w} - \underline{x^T y} + \cancel{\lambda w}}$$

# Notes

Friday, June 4, 2021 4:47 PM

Why is it called ridge



## 5 Key Understandings

Saturday, June 5, 2021 4:20 PM

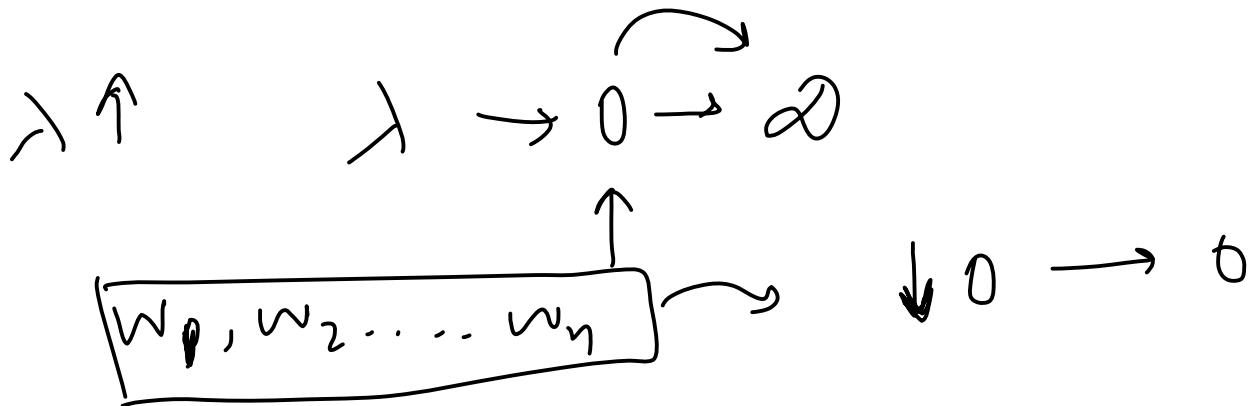
$$\sum_{i=1}^n (y_i - \hat{y}_i)^2 + [\lambda \|w\|^2]$$

$\lambda (w_1^2 + w_2^2 + \dots + w_n^2)^2$

Shrinkage      Coef      Overfitting ↗

# 1. How the coefficients get affected?

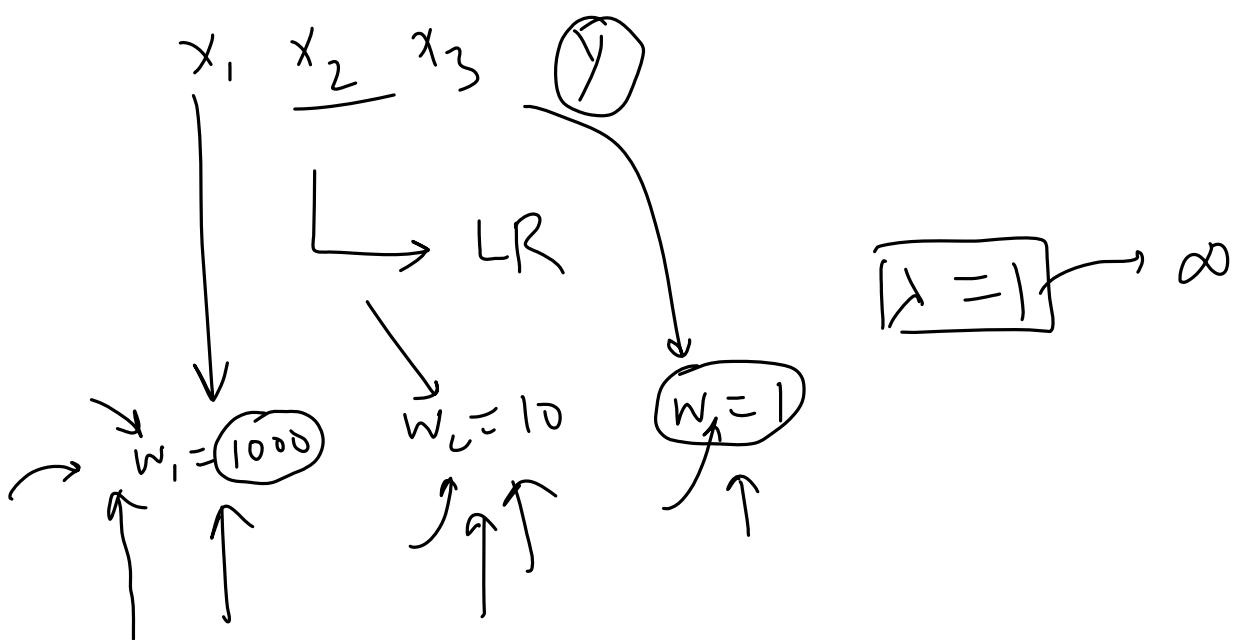
Saturday, June 5, 2021 4:20 PM



## 2. Higher Values are impacted more

Saturday, June 5, 2021 4:21 PM

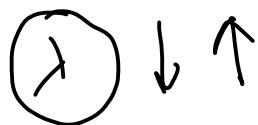
Never reaches 0



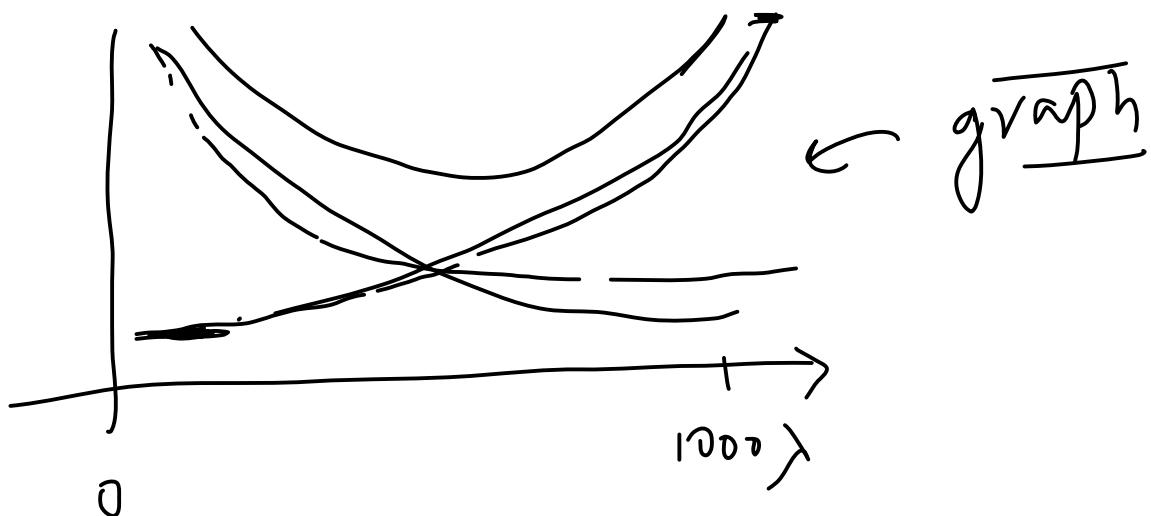
### 3. Bias Variance Tradeoff

Saturday, June 5, 2021 4:21 PM

Bias Variance

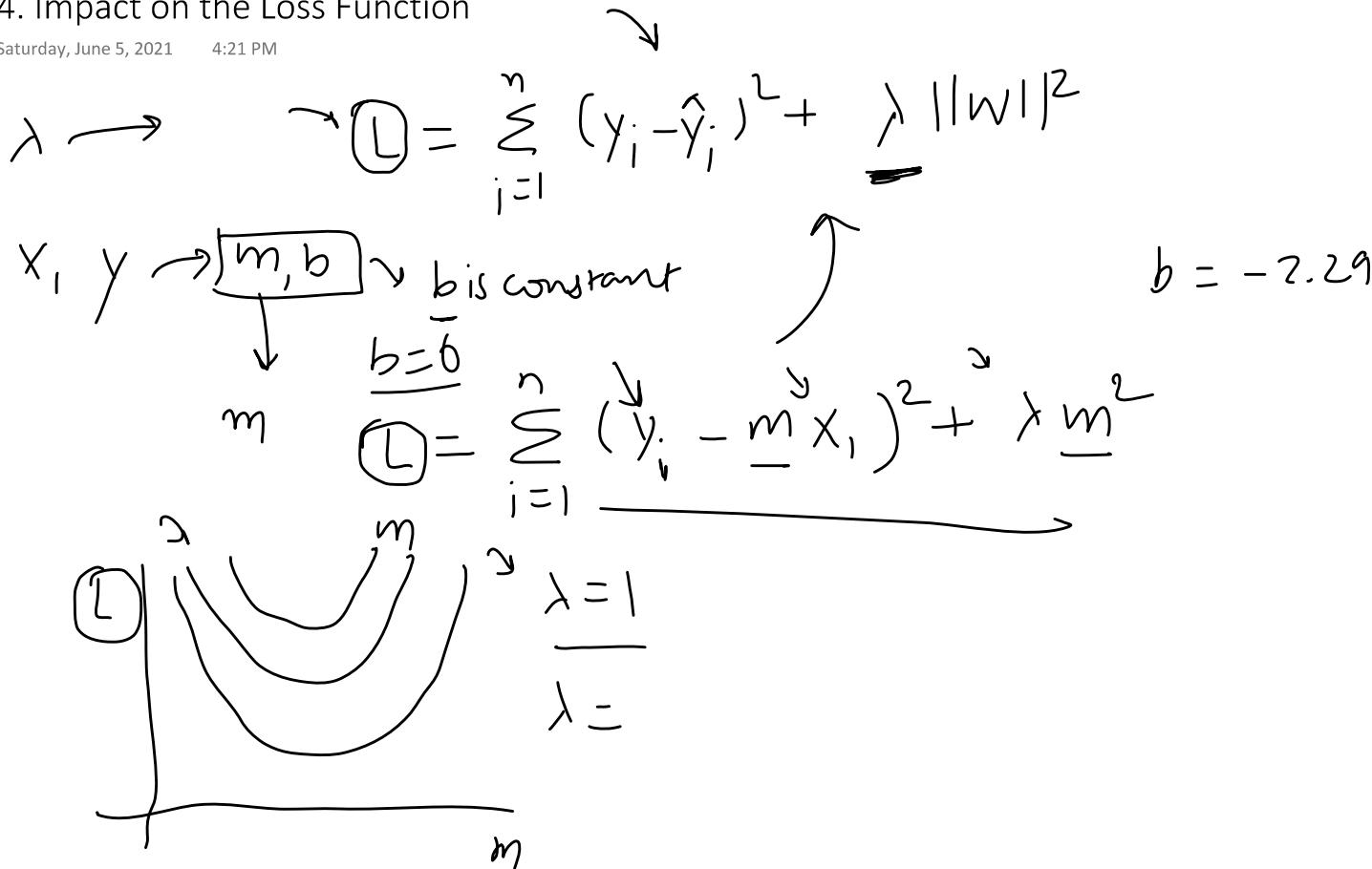


Bias ↓ overfit Variance ↑  
Bias ↑ underfitting Variance ↓



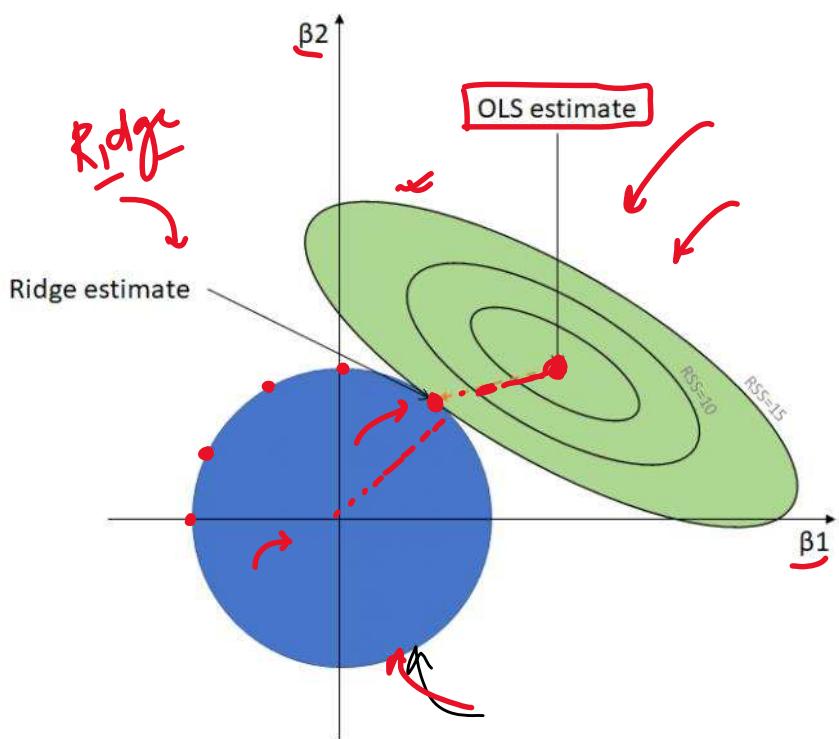
#### 4. Impact on the Loss Function

Saturday, June 5, 2021 4:21 PM



## 5. Why called Ridge

Saturday, June 5, 2021 4:22 PM



Hard constraint  
Ridge constraint

$$2 \text{ coef} \quad \beta_1 \quad \beta_2 \quad \beta_0$$

$$L = \text{MSE} + \lambda \|w\|^2$$

contour

$$(y_i - \hat{y}_i)^2$$

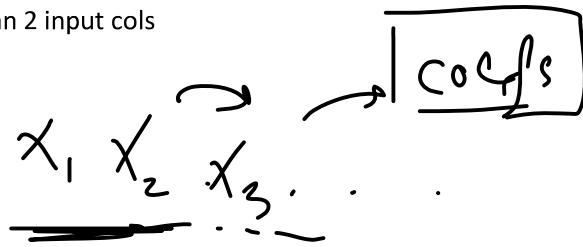
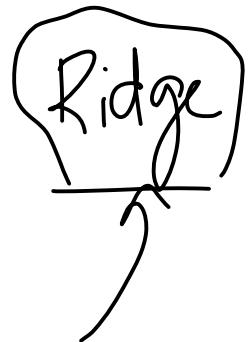
$$\sum_{i=1}^n \frac{(y_i - (\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}))^2}{\text{MSE}}$$

$$\boxed{\lambda (\beta_1^2 + \beta_2^2)}$$

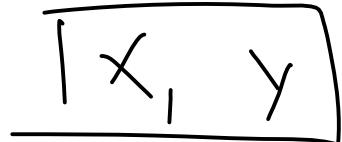
# Practical Tip

Monday, June 7, 2021 1:20 PM

Use ridge when there are more than 2 input cols



$\geq 2$



## Lasso Regression

Thursday, June 10, 2021 6:42 AM

L1 Regularization

overfitting

$$y = mx + b \quad \lambda \uparrow \quad \hat{Y} = \underline{\text{under}}$$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \|\mathbf{w}\|_1$$

$\lambda (w_1^2 + w_2^2 + \dots + w_n^2)$

$$\boxed{\lambda > 0}$$

$$\rightarrow 0 \quad \underline{w_1 \rightarrow w_n} \rightarrow \underline{\text{coeff}}$$

overfitting | under

alpha

Lasso

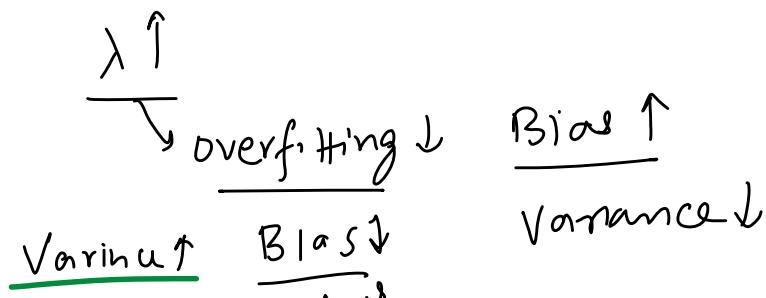
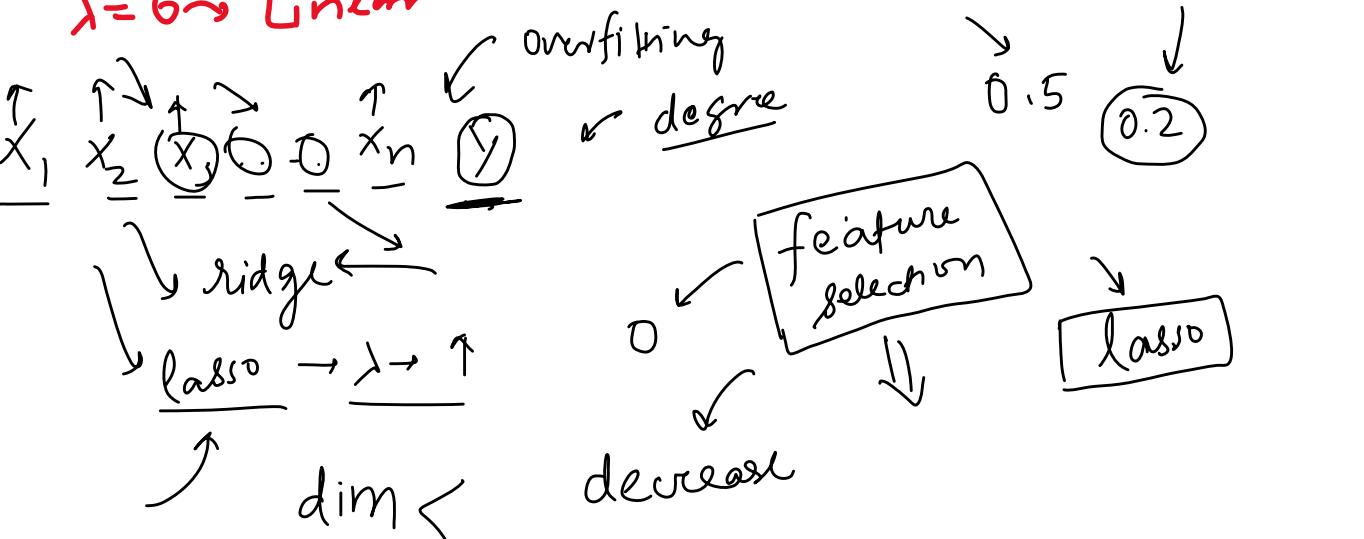
L1 norm

$$\|\mathbf{w}\|_1$$

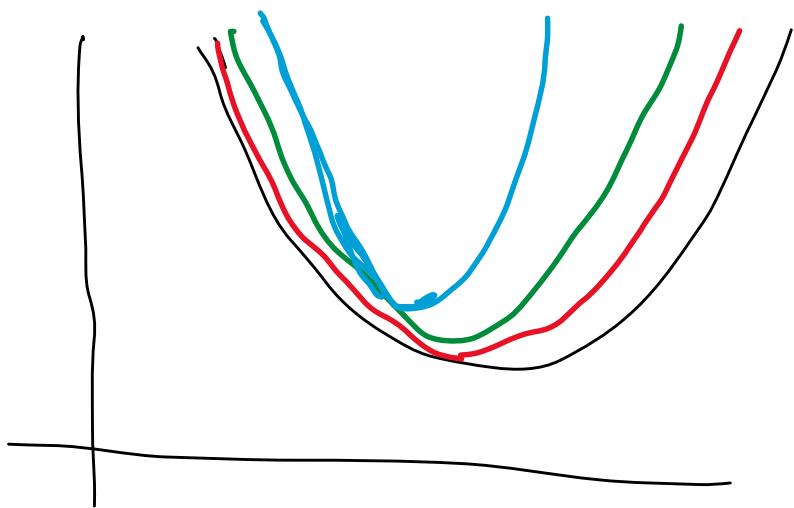
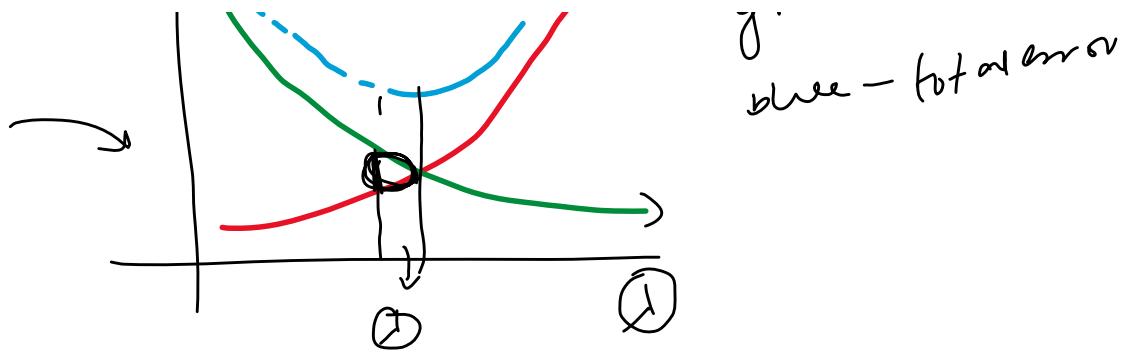
$$|w_1| + |w_2| + |w_3| + \dots + |w_n|$$

underfitting

$\lambda = 0 \rightarrow \text{Linear}$



Red → bias  
green → variance  
blue → total error



## Understanding Sparsity

Friday, June 11, 2021 11:21 AM

alpha	age	sex	bmi	bp	s1	s2	s3	s4	s5	s6
0.0000	-9.160885	-205.462260	516.684624	340.627341	-895.543609	561.214533	153.884786	126.734316	861.121400	52.419828
0.0001	-9.118336	-205.337133	516.880570	340.556792	-883.415291	551.553259	148.578680	125.355917	856.480254	52.467627
0.0010	-8.763583	-204.321125	518.371729	339.975385	-787.690766	475.274718	106.786540	114.632063	819.739542	52.872100
0.0100	-6.401088	-198.669767	522.048548	336.348363	-383.709187	152.663678	-66.060583	75.611090	659.869402	55.828128
0.1000	6.642753	-172.242166	485.523872	314.682122	-72.939323	-80.590053	-174.466515	83.616653	484.363285	73.584154
1.0000	42.242217	-57.305508	282.170831	198.061386	14.363544	-22.551274	-136.930053	102.023193	260.104308	98.552274
10.0000	21.174004	1.659796	63.659772	48.493240	18.421492	12.875448	-38.915435	38.842464	61.612405	35.505355
100.0000	2.858979	0.629452	7.540604	5.849997	2.710879	2.142134	-4.834047	5.108223	7.448466	4.576129
1000.0000	0.295726	0.069290	0.769004	0.597829	0.282900	0.225936	-0.495607	0.527031	0.761497	0.471029
10000.0000	0.029674	0.006995	0.077054	0.059915	0.028412	0.022715	-0.049686	0.052870	0.076321	0.047241

Ridge

Lasso  
sparsity

$\lambda \uparrow \quad w \rightarrow 0$

single  $x | y \rightarrow$

alpha	age	sex	bmi	bp	s1	s2	s3	s4	s5	s6
0.0000	-9.160885	-205.462260	516.684624	340.627341	-895.543596	561.214523	153.884780	126.734314	861.121395	52.419828
0.0001	-9.071288	-205.337332	516.780313	340.539730	-888.652320	555.952271	150.585260	125.453044	858.639860	52.379002
0.0010	-8.264924	-204.213177	517.641106	339.751339	-826.653342	508.609613	120.899583	113.924518	836.314382	52.011583
0.0100	-1.361404	-192.944226	526.348511	332.649058	-430.205495	191.277876	-44.048113	68.990747	688.384976	47.939528
0.1000	0.000000	-113.976046	526.737112	292.635423	-82.691928	-0.000000	-152.691332	0.000000	551.077200	7.169852
1.0000	0.000000	0.000000	363.882636	27.278420	0.000000	0.000000	-0.000000	0.000000	336.135971	0.000000
10.0000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	-0.000000	0.000000	0.000000	0.000000
100.0000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	-0.000000	0.000000	0.000000	0.000000
1000.0000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	-0.000000	0.000000	0.000000	0.000000
10000.0000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	-0.000000	0.000000	0.000000	0.000000

feature selection

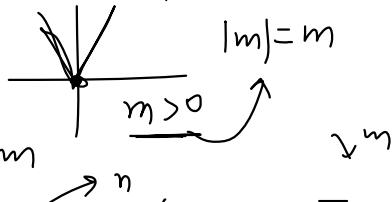
$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2 + \lambda}$$

simple  $x | y \rightarrow y = m x + b$

$$b = \bar{y} - m \bar{x}$$

$\bar{y} \rightarrow \text{mean}(y)$   
 $\bar{x} \rightarrow \text{mean}(x)$

$$m = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$



$$b = \bar{y} - m \bar{x}$$

$$m = ?$$

$$0$$

$$|m| = -m$$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \pm \lambda |m|$$

$$\frac{d}{dm} \sum_{i=1}^n \frac{(y_i - mx_i - \bar{y} + m\bar{x})^2}{2} + 2\lambda m$$

$$\frac{d}{dm} L = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})^2 + 2\lambda m = \sum_{i=1}^n (y_i - mx_i - \bar{y} + m\bar{x})(-x_i + \bar{x}) + 2\lambda m = 0$$

$$m \in (x_i - \bar{x})^2 = \sum (y_i - \bar{y})(x_i - \bar{x}) - \lambda$$

$$-L \geq \sum (y_i - \hat{y})^2 + \lambda \sum |x_i|$$

$$\begin{aligned} & -\sum [(y_i - \bar{y})(x_i - \bar{x})] - m(x_i - \bar{x})^2 + \lambda = 0 \\ & -\sum (y_i - \bar{y})(x_i - \bar{x}) + m \sum (x_i - \bar{x})^2 + \lambda = 0 \end{aligned}$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) - \lambda}{\sum (x_i - \bar{x})^2}$$

Lasso Coeff Sparsity

$$\boxed{\text{for } m > 0} \quad \lambda > 0$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) - \lambda}{\sum (x_i - \bar{x})^2}$$

$$m = \frac{(YX) - \lambda}{X^2}$$

$$\left\{ \begin{array}{l} YX = 100 \\ X^2 = 50 \end{array} \right.$$

$$\lambda \uparrow \quad \downarrow$$

$$\lambda = 0 \quad m = 2 \quad \boxed{m=0} \quad \lambda = 100 \quad m = -1 \quad \boxed{m=-1}$$

$$\lambda = 50 \quad m = 1$$

$$\boxed{\text{for } m=0}$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2}$$

$$\boxed{\text{for } m < 0}$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}$$

$$m = \frac{YX + \lambda}{X^2} = \frac{100 + \lambda}{50}$$

$$= \frac{100 + 150}{50} =$$

$$(m = \bar{5})$$

$$\boxed{m < 0} \quad \boxed{m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x}) + \lambda}{\sum (x_i - \bar{x})^2}}$$

$$\lambda > 0$$

$$m = -\frac{100 - \lambda}{50} \leftarrow$$

$$m = -\frac{100 + \lambda}{50}$$

$$\lambda = 0 \quad m = -2$$

$$\lambda = 50 \quad m = -1$$

$$\lambda = 100 \quad m = 0 \rightarrow 1$$

$$\lambda = 150 \quad m = -5$$

$$1) \rightarrow 0 \rightarrow$$

$$2) \rightarrow \text{stop}$$

$$m = \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2 + \lambda}$$

$$0 \quad \lambda \rightarrow \text{numerical}$$

$$\lambda = 100000000$$

$$\lambda \rightarrow \text{Denominator}$$

Lasso  $\lambda \rightarrow \text{numerical}$

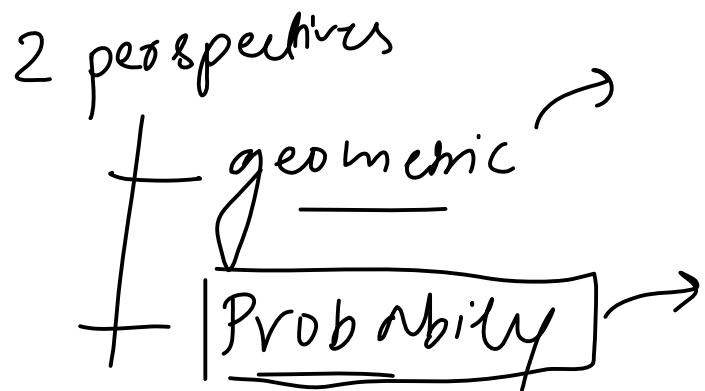
## ElasticNet Regression

Saturday, June 12, 2021 1:04 PM

$\text{Ridge} \quad \lambda(w_1^2 + w_2^2 + \dots + w_n^2)$	$\text{Lasso} \quad \lambda( w_1  +  w_2  + \dots +  w_n )$
$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \ w\ _2^2$ <div style="display: flex; justify-content: space-between; align-items: center;"> <span style="margin-right: 20px;"><math>\underbrace{\text{mse}}</math></span> <span><math>\overbrace{\text{overfitting}}</math></span> </div>	$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \ w\ _1$ <div style="display: flex; justify-content: space-between; align-items: center;"> <span style="margin-right: 20px;"><math>\underbrace{\text{mse}}</math></span> <span><math>\overbrace{\text{feature selection}}</math></span> </div>
$\lambda \uparrow \quad w \downarrow 0$	$\lambda \uparrow \quad w \rightarrow 0$
$\downarrow$	$\downarrow$
$\text{EN Reg}$	$L = \sum (y_i - \hat{y}_i)^2 + \underline{a} \ w\ ^2 + \underline{b} \ w\ $
$\frac{\lambda=1}{a=0.5} \quad l1\_ratio = 0.5$	$\left\{ \begin{array}{l} \lambda, l1\_ratio \\ \boxed{l1, \lambda} \end{array} \right.$
$\uparrow$	$\uparrow$
$l1\_ratio > 0.9$	$q_0 \times \text{mid} \approx \text{lambda}$
$x_1 \quad   \quad x_2$	$\left\{ \begin{array}{l} \lambda = a + b \\ l1\_ratio = \frac{a}{a+b} \\ l1 = \frac{a}{\lambda} \\ a = l1 \times \lambda \\ b = \lambda - a \end{array} \right.$

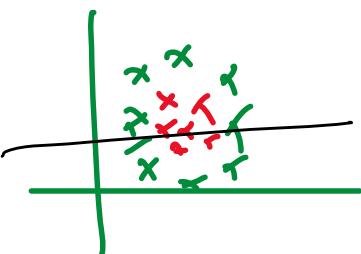
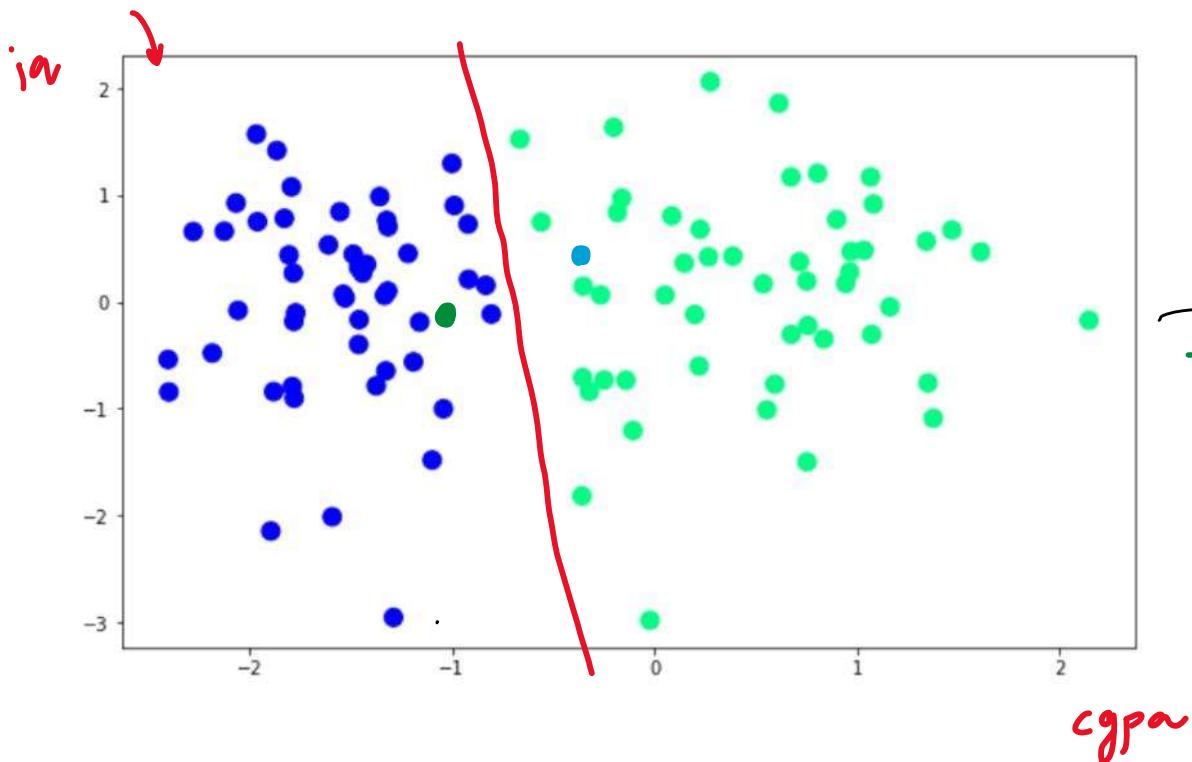
# Introduction

Tuesday, June 15, 2021 12:07 PM



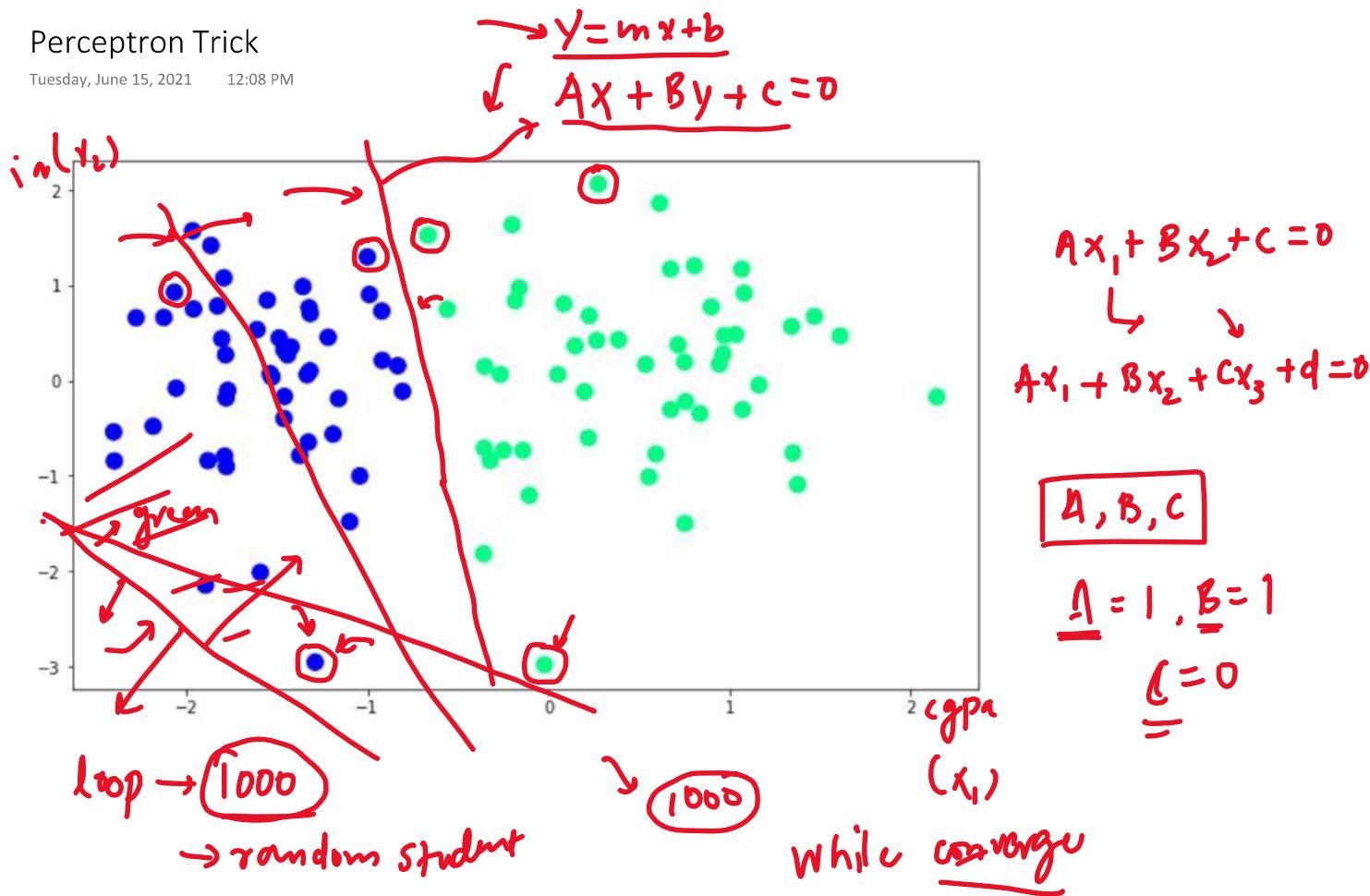
# Requirement

Tuesday, June 15, 2021 12:07 PM



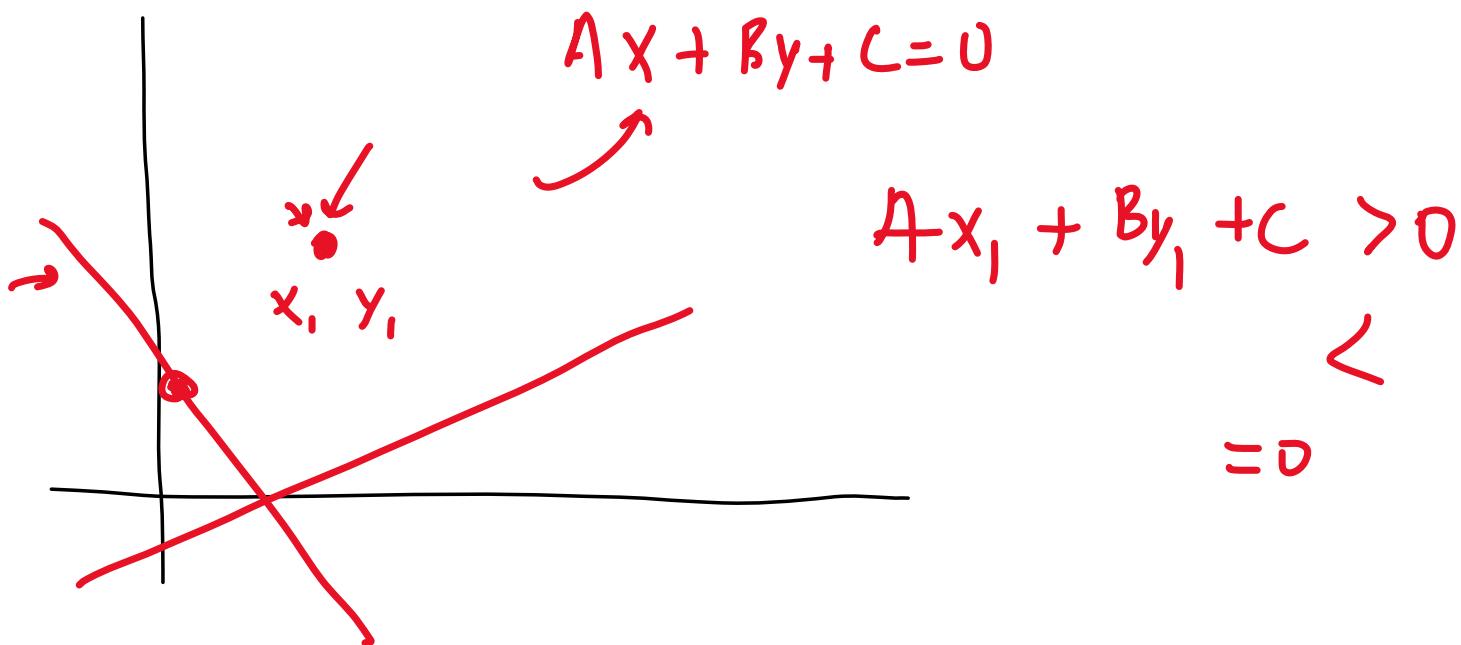
## Perceptron Trick

Tuesday, June 15, 2021 12:08 PM



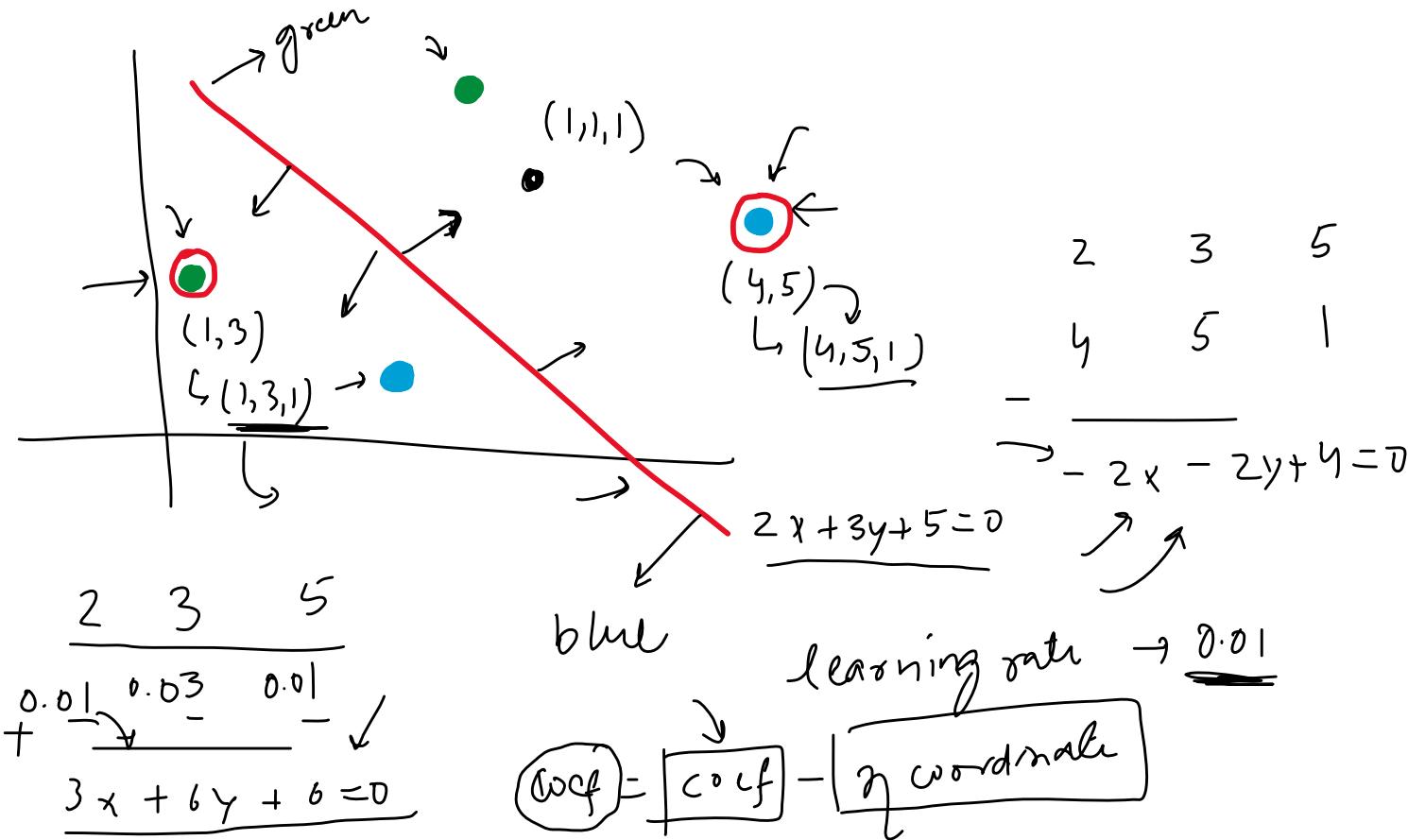
## How to label regions?

Tuesday, June 15, 2021 1:12 PM



## Transformations

Tuesday, June 15, 2021 1:31 PM



## Algorithm

Tuesday, June 15, 2021 2:31 PM  
 $x_0 (x_1)$   $x_2$   $y$   
 cgpa placid

	$x_1$	$x_2$	$y$
1   7.5	61		①
1   8.9	109		1
1   7.0	81		0

$Ax + By + C = 0$   
 $w_0 + w_1 x_1 + w_2 x_2 = 0$   
 $w_0 = C, w_1 = A, w_2 = B$   
 $w_0 x_0 + w_1 x_1 + w_2 x_2 = 0$   
 $w_0 x_1 + w_1 x_2 + w_2 x_3 = 0$   
 $\sum_{i=0}^2 w_i x_i = 0$   
 $[w_0 \ w_1 \ w_2] \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}$

$= \geq 0 \rightarrow ①$   
 $< 0 \rightarrow 0$   
 $-ve \quad +ve \quad L_f \ 0$   
 $x_i \in P$   
 $y \cdot x_i \in N$

epoch  $\rightarrow 1000, \eta = 0.01$   
 for i in range (epochs):

randomly select a student  
if  $x_i \in N$  and  $\sum_{i=0}^2 w_i x_i \geq 0$  {  
 $w_{new} = w_{old} - \eta x_i$   
if  $x_i \in P$  and  $\sum_{i=0}^2 w_i x_i < 0$   
 $w_{new} = w_{old} + \eta x_i$

$x_i \in N \quad w_{old} < 0$   
 $x_i \in P \quad w_{old} > 0$

## Simplified Algo

Tuesday, June 15, 2021 2:44 PM

```

if  $x_i \in N$  and  $\sum w_i x_i \geq 0$ 
     $w_n = w_0 - \eta x_i$ 
else if  $x_i \in P$  and  $\sum w_i x_i < 0$ 
     $w_n = w_0 + \eta x_i$ 

```

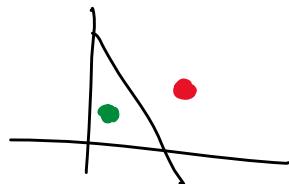
for  $i$  in 1000  
 random student  $\leftarrow 0$

$$w_n = w_0 + \eta(y_i - \hat{y}_i)x_i$$

$$\begin{cases} w_n = w_0 \\ w_n = w_0 + \eta x_i \end{cases}$$

$$w_n = w_0 - \eta x_i$$

$x_i$	$\hat{y}_i$	$y_i - \hat{y}_i$
1	0	0
-1	0	0
1	1	0
-1	1	-1



for  $i$  in range(epochs):  
 select a random student ( $i$ )

$$w_n = w_0 + \eta(x_i - \hat{y}_i)x_i$$

$$Ax + By + C = 0$$

$$y = mx + b$$

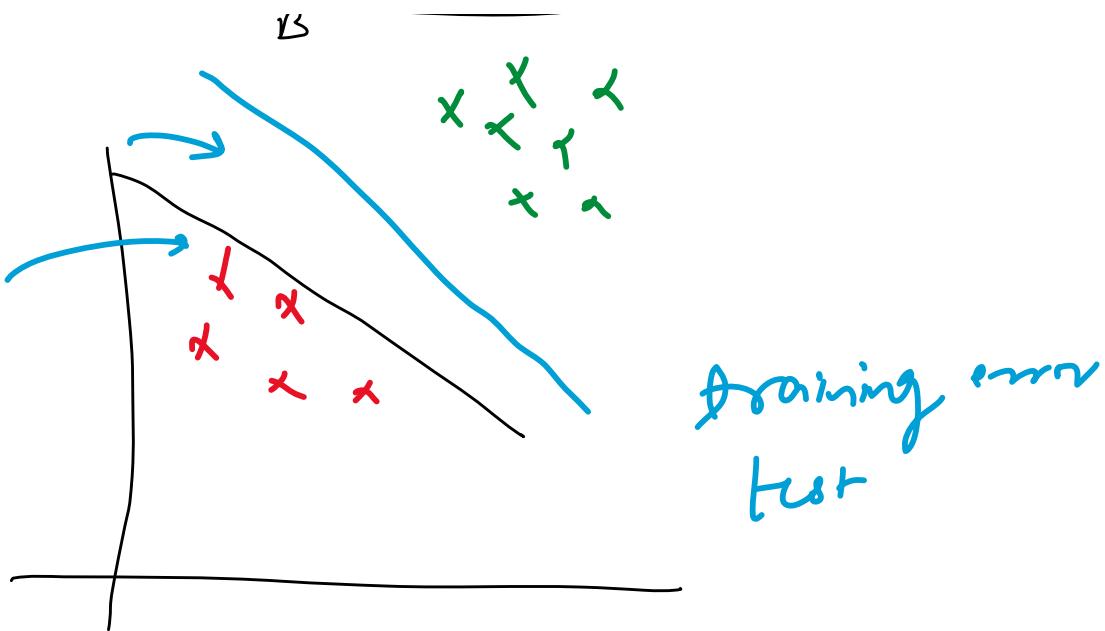
$$m = -\frac{A}{B}$$

$$C = -\frac{C}{B}$$

$A, B, C$

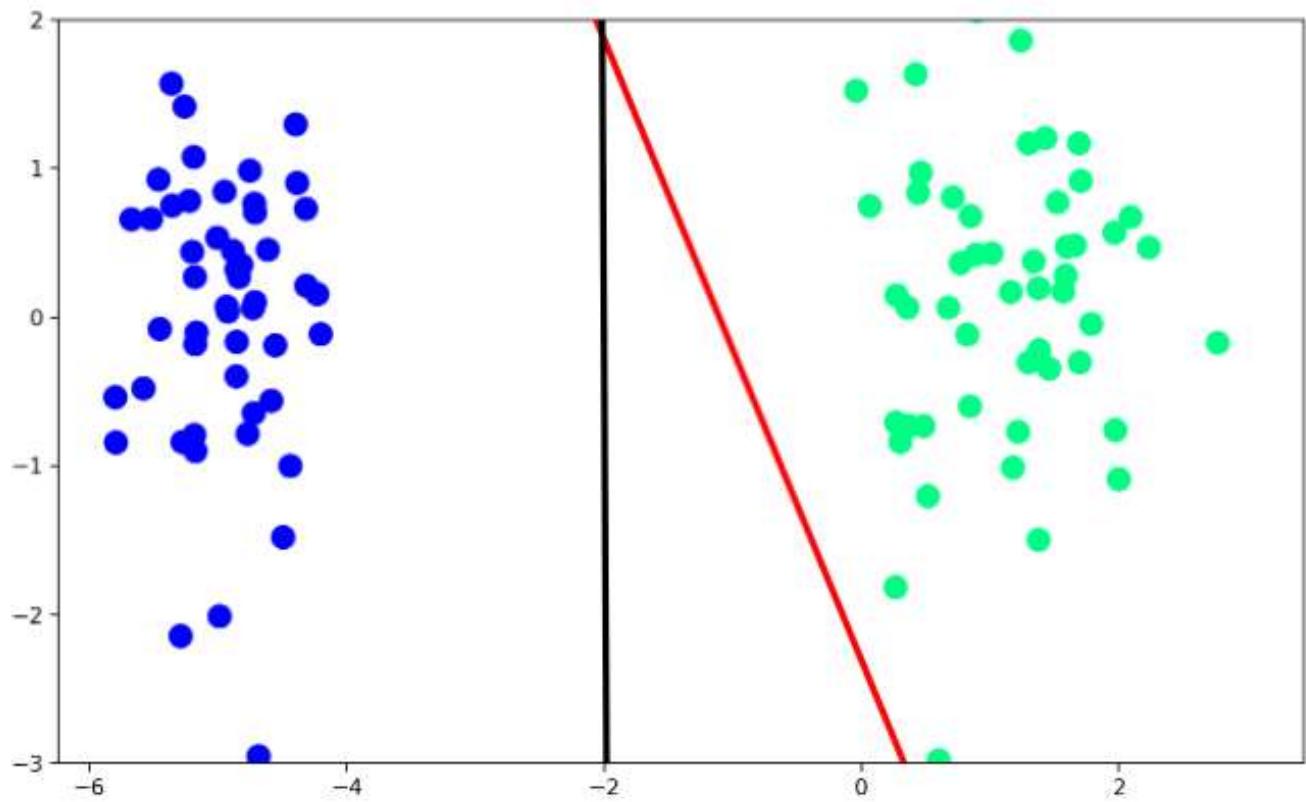
if  $[ - ]$   $\rightarrow [ ]$

$\checkmark \times \times$



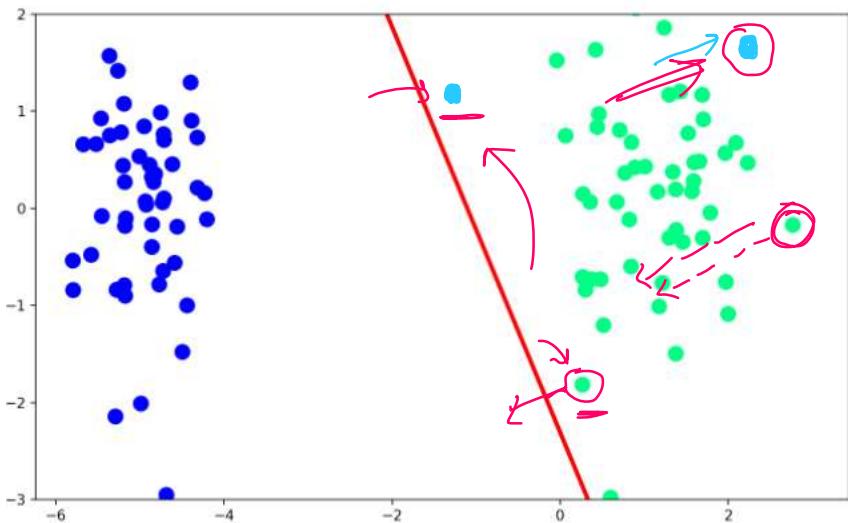
# Problem with Perceptron

Thursday, June 17, 2021 12:18 PM



## Possible Solution?

Thursday, June 17, 2021 12:18 PM



$$w_{\eta} = w_0 + \eta(y_i - \hat{y}_i)x_i$$

→ misclassified  
line - pull

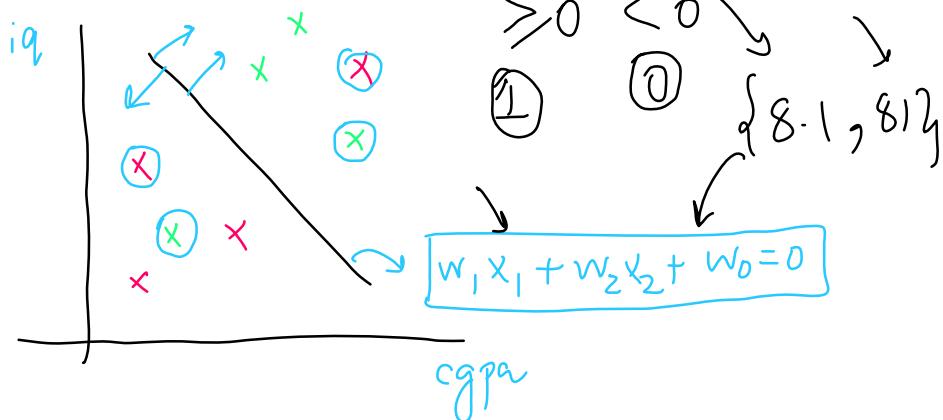
→ correctly  
line push  
wash out

$$w_{\eta} = w_0 + \eta \sum_{i=1}^n (y_i - \hat{y}_i)x_i$$

$y_i$	$\hat{y}_i$	$y_i - \hat{y}_i$
1	1	0
0	0	0
1	0	1
0	1	-1

$$\rightarrow (y_i - \hat{y}_i) \neq 0 \quad \text{model predict} \\ \sum w_i x_i = [0, 1] \\ \rightarrow (y_i - \hat{y}_i)$$

$$w_1 \times 8.1 + w_2 \times 81 + w_0 =$$



cgpa	iq	(xi) placed
9	91	0
8.8	78	1
8.1	102	1
7.9	98	1

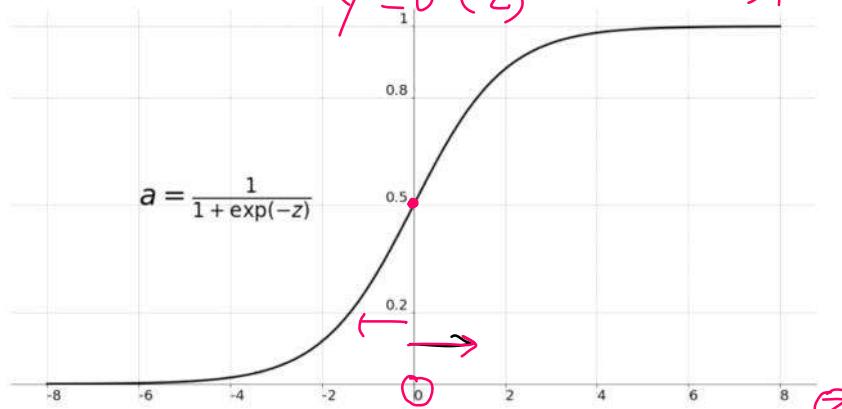
## The Sigmoid Function

Thursday, June 17, 2021 12:19 PM

### Sigmoid Function

$$y = \sigma(z) \rightarrow 1$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

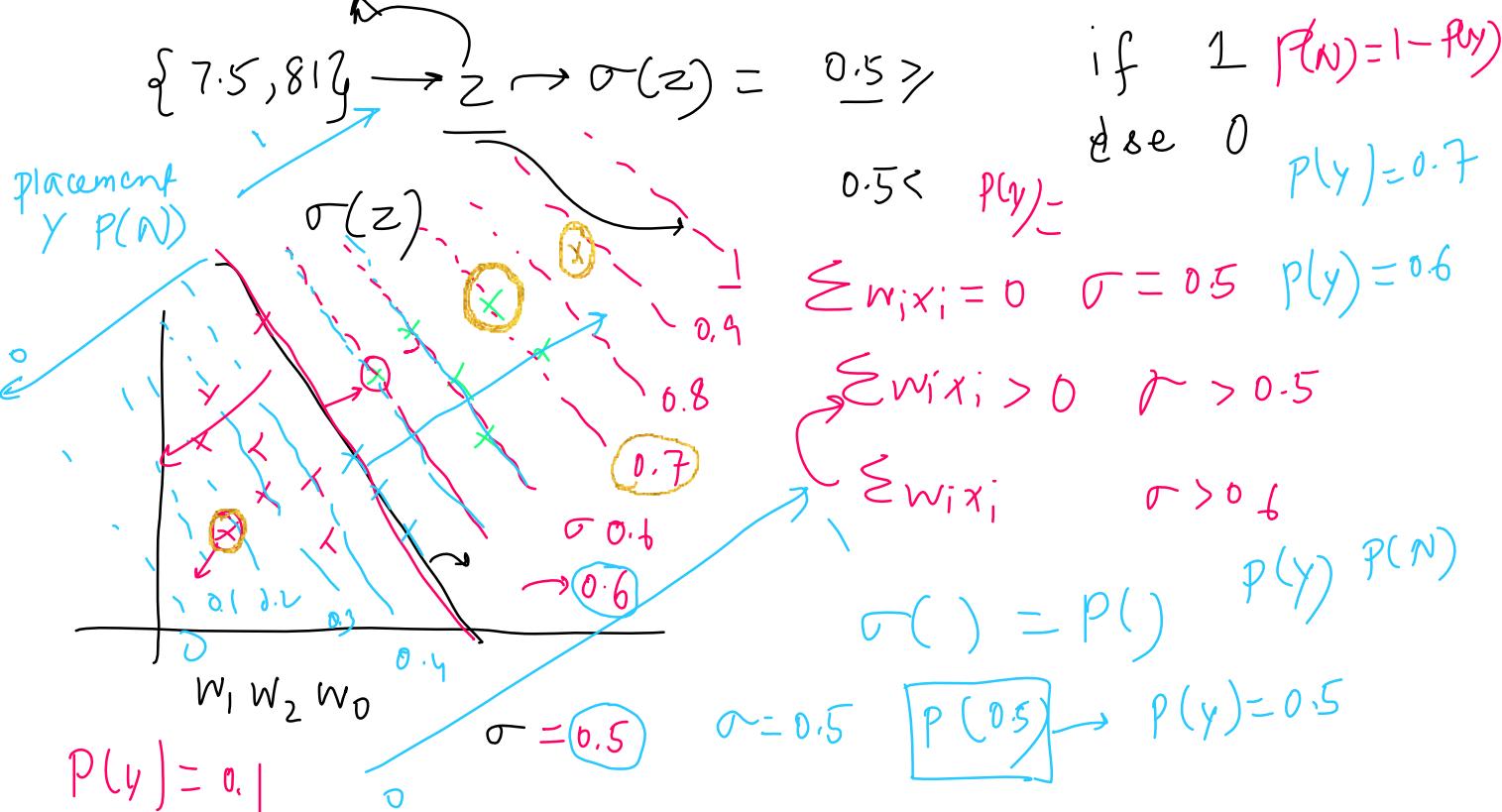


$$\begin{cases} -\infty < z < \infty \\ 0 < y < 1 \end{cases}$$

$$y_i = \frac{w_1 \times 7.5 + w_2 \times 81 + w_0}{\sum w_i x_i} = \sigma(z)$$

$\sigma > 0.5 \rightarrow z \geq 0 \rightarrow 1$

$\sigma < 0.5 \rightarrow z < 0 \rightarrow 0$



## Impact of Sigmoid

Thursday, June 17, 2021 3:27 PM

$$w_n = w_0 + \eta(y_i - \hat{y}_i) x_i$$

$$\hat{y}_i = \sigma(z)$$

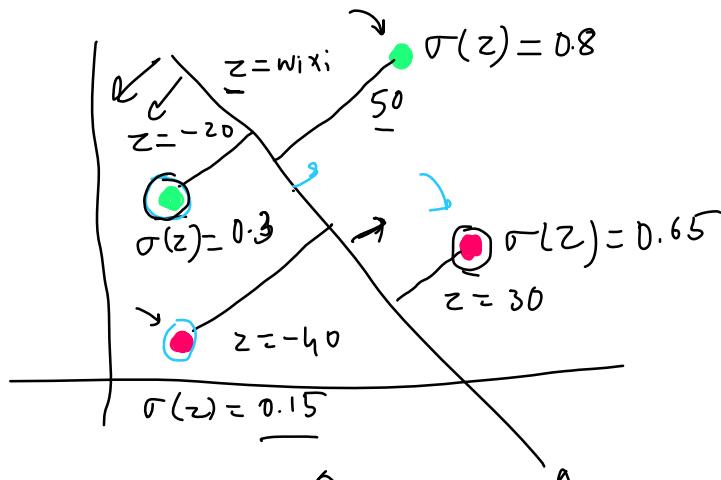
$$\text{where } z = \sum w_i x_i$$

$$w_n = w_0 + \eta \times 0.2 \times x_i$$

$$w_n = w_0 - \eta \times 0.65 \times x_i$$

$$w_n = w_0 + \eta \times 0.7 \times x_i$$

$$w_n = w_0 - \eta \times 0.15 \times x_i$$



$y_i$	$\hat{y}_i$	$y_i - \hat{y}_i$
1	0.8	0.2
0	0.65	-0.65
1	0.3	0.7
0	0.15	-0.15

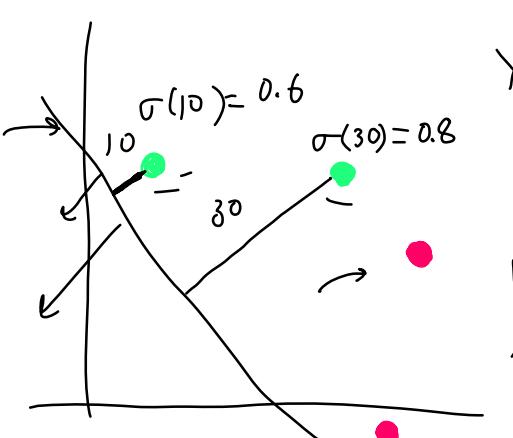
$$w_n = w_0 + \eta(y_i - \hat{y}_i) x_i$$

$$\hat{y}_i = \sigma(z)$$

$$\text{where } z = \sum w_i x_i$$

$$w_n = w_0 + [\eta \times 0.4 \times x_i] = x_1$$

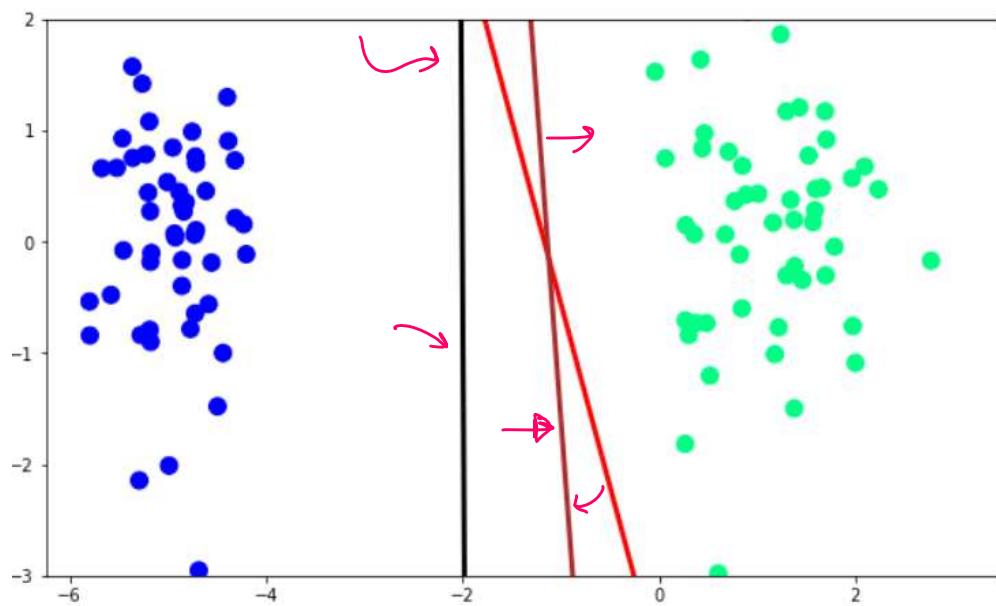
$$w_n = w_0 + [\eta \times 0.2 \times x_i] = x_2$$



$x_1 > x_2 \rightarrow \underline{\text{code implement}}$

## The Problem

Friday, June 18, 2021 2:20 PM



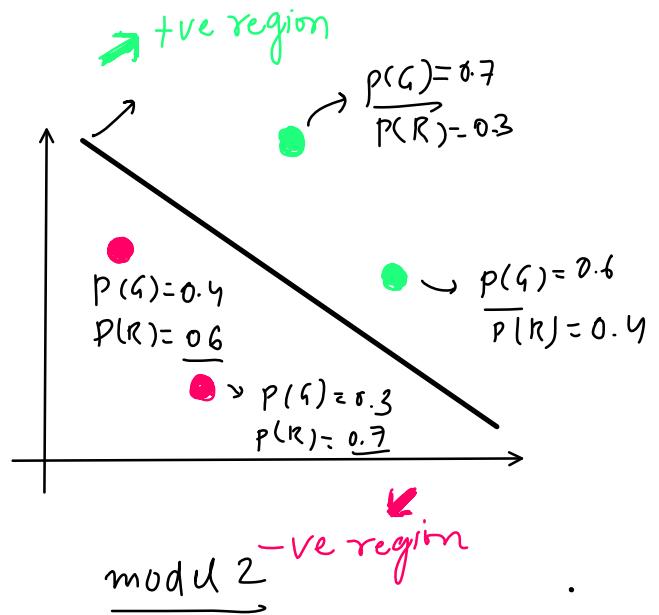
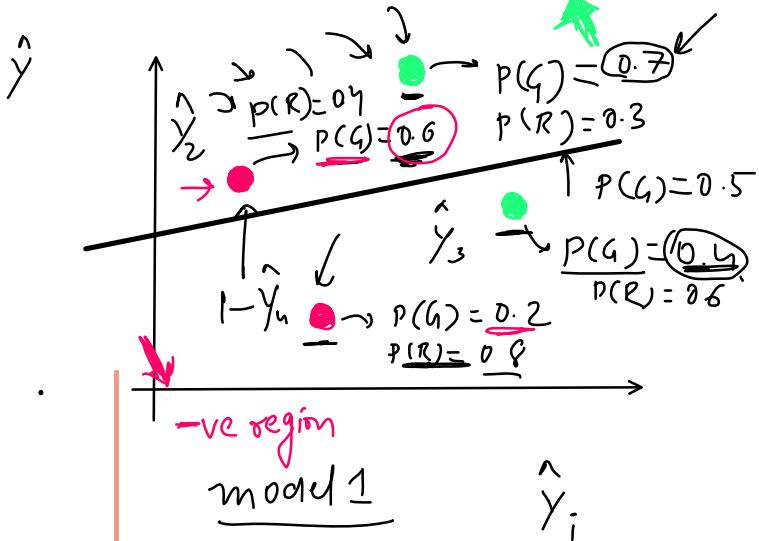
red → step  
brown → Sigmoid  
black → Sigmoid learn

for i in 1000  
random

Loss function  
error function

## Maximum Likelihood and Cross-Entropy

Thursday, June 17, 2021 12:19 PM



$$\text{model 2} \rightarrow 0.7 \times 0.6 \times 0.6 \times 0.7 \\ = 0.176$$

$$\text{model 1} \rightarrow \frac{0.7 \times 0.4 \times 0.4 \times 0.8}{0.089}$$

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

$$\log(\max) = -\log(0.7) - \log(0.4) - \log(0.4) - \log(0.8)$$

$0-1 = -ve$

Cross entropy

minimize

No

$$-\log(\hat{y}_1) - \log(\hat{y}_2) - \log(\hat{y}_3) - \log(\hat{y}_4)$$

$(1 - \hat{y}_i)$

$$\underline{y_i = 1}$$

$$-y_i \log(\hat{y}_i) - \frac{(1-y_i) \log(1-\hat{y}_i)}{\boxed{t}}$$

$$-\underline{y_i \log(\hat{y}_i)} = -\log(\hat{y}_i)$$

$$= -\log(0.7)$$

$$y_2 = 0 \quad y_3 = 1$$

$$y_2 = 0 \quad y_3 = 1 \quad = -\log(0.7)$$

$$\frac{-y_2 \log(\hat{y}_2) - (1-y_2) \log(1-\hat{y}_2)}{y_3 \log(\hat{y}_3)} - \log(1-\hat{y}_3) = -\log(0.4)$$

$$-y_3 \log(\hat{y}_3)$$

$$-\log(y_3) = -\log(0.4)$$

$$y_4 = 0$$

$$-(1-y_4) \log(1-\hat{y}_4)$$

$$-\log(1-\hat{y}_4)$$

$$-\log(0.8)$$

$$L = \sum_{i=1}^n -y_i \log(\hat{y}_i) - (1-y_i) \log(1-\hat{y}_i)$$

MSE

$$L = -\frac{1}{n} \sum_{i=1}^n y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

Closed form gradient descent

$\min_{w_1 w_2 w_0}$

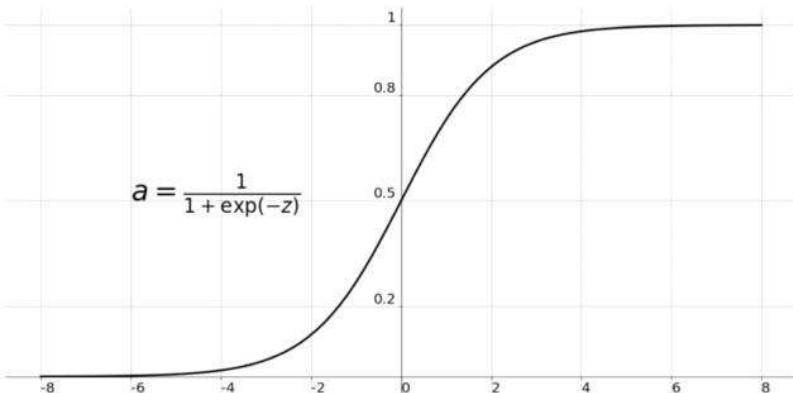
log-loss error

binary cross entropy

## Derivative of Sigmoid

Thursday, June 17, 2021 12:20 PM

## Sigmoid Function



$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma'(x) = \frac{d}{dx} \left( \frac{1}{1 + e^{-x}} \right)$$

$$\frac{d}{dx} \left( \frac{1}{x} \right) = \frac{d}{dx} (x)^{-1}$$

$$= -x^{-2} = -\frac{1}{x^2}$$

$$\frac{d}{dx} \left[ \frac{1}{1 + e^{-x}} \right] = \frac{d}{dx} \left[ (1 + e^{-x})^{-1} \right] = -\frac{1}{(1 + e^{-x})^2} \frac{d}{dx} (1 + e^{-x})$$

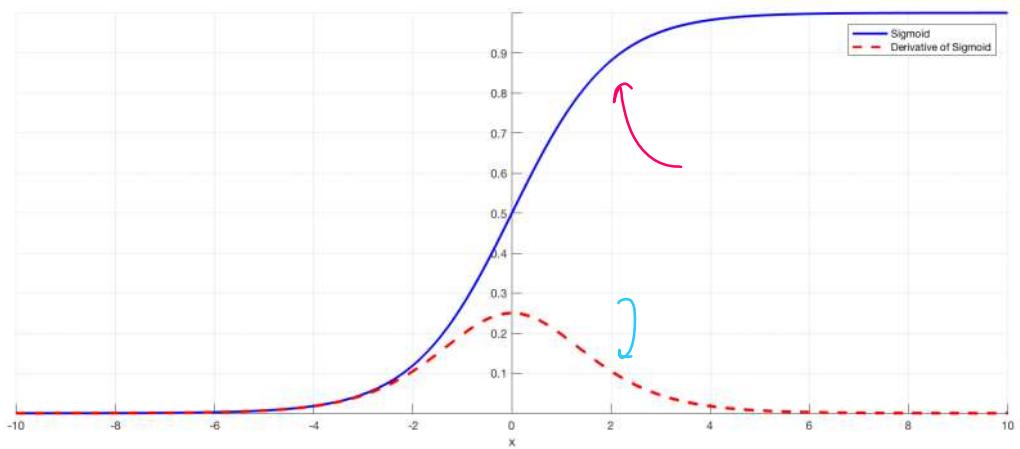
$e^{-x} = e^{-x}$        $-x = -1$

$$= -\frac{1}{(1 + e^{-x})^2} \frac{d}{dx} (e^{-x}) = -\frac{e^{-x}}{(1 + e^{-x})^2} \frac{d}{dx} (-x) = \frac{e^{-x}}{(1 + e^{-x})^2}$$

$$\frac{1 \cdot e^{-x}}{(1 + e^{-x})(1 + e^{-x})} = \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} = \sigma(x) \left[ \frac{e^{-x}}{1 + e^{-x}} \right]$$

$$= \sigma(x) \left[ \frac{1 + e^{-x} - 1}{1 + e^{-x}} \right] = \sigma(x) \left[ \frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right]$$

$$\sigma(x) [1 - \sigma(x)] \Rightarrow \sigma'(x) = \boxed{\sigma(x) [1 - \sigma(x)]}$$



## Gradient Descent

Monday, June 21, 2021 11:50 AM

Classification  $\{x_1, x_2\} \rightarrow y$

$$y = \sigma(w_0 + w_1 x_1 + w_2 x_2)$$

$$\hat{y}_i = \sigma(z) = \sigma(w_0 + w_1 x_1 + w_2 x_2 + \dots)$$

$$L = -\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

$$L(w_0, w_1, w_2) = \arg \min_{w_0, w_1, w_2} L(w_0, w_1, w_2)$$

$$\text{gradient descent}$$

$$\text{for } i \text{ in epochs}$$

$$w_{new} = w_{old} - \eta \frac{\partial L}{\partial w_{old}}$$

$$w_0 = w_0 - \eta \frac{\partial L}{\partial w_0}, \quad w_1 = w_1 - \eta \frac{\partial L}{\partial w_1}, \quad w_2 = w_2 - \eta \frac{\partial L}{\partial w_2}$$

n calc → n+1 derivatives

$$w_j = w_j - \eta \frac{\partial L}{\partial w_j} \quad j=0, 1, 2, \dots, n$$

$$L = -\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

$$\frac{\partial L}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \left[ \frac{\partial L}{\partial w_j} (y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)) \right]$$

$$\hat{y}_i = \sigma(z) = \sigma(w_0 + w_1 x_1 + w_2 x_2 + \dots)$$

$$\frac{\partial L}{\partial w_j} y_i \log(\hat{y}_i) = \frac{\partial L}{\partial w_j} y_i \log(\sigma(z)) - y_i \frac{\partial L}{\partial z} \sigma'(\sum_{j=0}^n w_j x_j)$$

$$\frac{\partial L}{\partial w_j} = \frac{y_i \log(\sigma(\sum_{j=0}^n w_j x_j)) - \hat{y}_i}{\hat{y}_i} = y_i \frac{\partial L}{\partial w_j} \sigma(\sum_{j=0}^n w_j x_j)$$

$$\begin{aligned} & \frac{\partial L}{\partial w_j} = y_i \log(\sigma(\sum_{j=0}^n w_j x_j)) \\ & = y_i \frac{\partial L}{\partial w_j} \underbrace{\log(\sigma(\sum_{j=0}^n w_j x_j))}_{\hat{y}_i(1-\hat{y}_i)} = y_i \frac{\hat{y}_i(1-\hat{y}_i)}{\hat{y}_i} \frac{\partial L}{\partial w_j} \sum_{j=0}^n w_j x_j \end{aligned}$$

$$y_i(1-\hat{y}_i) \sum_{j=0}^n \frac{\partial L}{\partial w_j} w_j x_j = \boxed{y_i(1-\hat{y}_i) \sum_{j=0}^n x_j}$$

$$\frac{\partial L}{\partial w_j} \frac{(1-y_i) \log(1-\hat{y}_i)}{\sigma(\sum_{j=0}^n w_j x_j)} \Rightarrow (1-y_i) \frac{\partial L}{\partial w_j} \log(1-\hat{y}_i) \quad \hat{y}_i = \sigma(z)$$

$$\frac{(1-y_i)}{(1-\hat{y}_i)} \frac{\partial L}{\partial w_j} \frac{(1-\hat{y}_i)}{\hat{y}_i} \Rightarrow -\frac{(1-y_i)}{(1-\hat{y}_i)} \frac{\partial L}{\partial w_j} \hat{y}_i \Rightarrow -\frac{(1-y_i)}{(1-\hat{y}_i)} \frac{\partial L}{\partial w_j} \sigma(z)$$

$$\Rightarrow -\frac{(1-y_i)}{(1-\hat{y}_i)} \sigma(z) \frac{(1-\sigma(z))}{\sigma(z)} \frac{\partial L}{\partial w_j} \sigma(z) = -\frac{(1-y_i)}{(1-\hat{y}_i)} \hat{y}_i(1-\hat{y}_i) \frac{\partial L}{\partial w_j} \sigma(z)$$

$$\Rightarrow -\hat{y}_i(1-y_i) \frac{\partial L}{\partial w_j} \sum_{j=0}^n w_j x_j \Rightarrow -\hat{y}_i(1-y_i) \sum_{j=0}^n \frac{\partial L}{\partial w_j} (w_j x_j)$$

$$\Rightarrow \boxed{-\hat{y}_i(1-y_i) \sum_{j=0}^n x_j}$$

$$\frac{\partial L}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \left[ y_i(1-\hat{y}_i) \sum_{j=0}^n x_j - \hat{y}_i(1-y_i) \sum_{j=0}^n x_j \right]$$

$$\approx \boxed{-\frac{1}{m} \sum_{i=1}^m \left[ y_i(1-\hat{y}_i) - \hat{y}_i(1-y_i) \right] \sum_{j=0}^n x_j}$$

$$\frac{\partial L}{\partial w_j} = -\frac{1}{m} \sum_{i=1}^m \left[ y_i (1 - \hat{y}_i) - \hat{y}_i (1 - y_i) \right] \sum_{j=0}^n x_j$$

$$= -\frac{1}{m} \sum_{j=1}^m \left[ y_i - y_i \cancel{\hat{y}_i} - \hat{y}_i + y_i \cancel{\hat{y}_i} \right] \sum_{j=0}^n x_j$$

$$\boxed{\frac{\partial L}{\partial w_j} = -\frac{1}{m} \sum_{j=1}^m (y_i - \hat{y}_i) \sum_{j=0}^n x_j}$$

$$\frac{\partial L}{\partial w_0} = -\frac{1}{m} [1+0] [1+1]$$

$$= -\frac{1}{2} [1][1] = -1$$

$$\begin{array}{cccc} x_1 & x_2 & y_i & \hat{y}_i \\ 1 & 2 & 1 & 0 \\ 3 & 4 & 0 & 0 \end{array}$$

$$= -\frac{1}{2} [1][0] \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

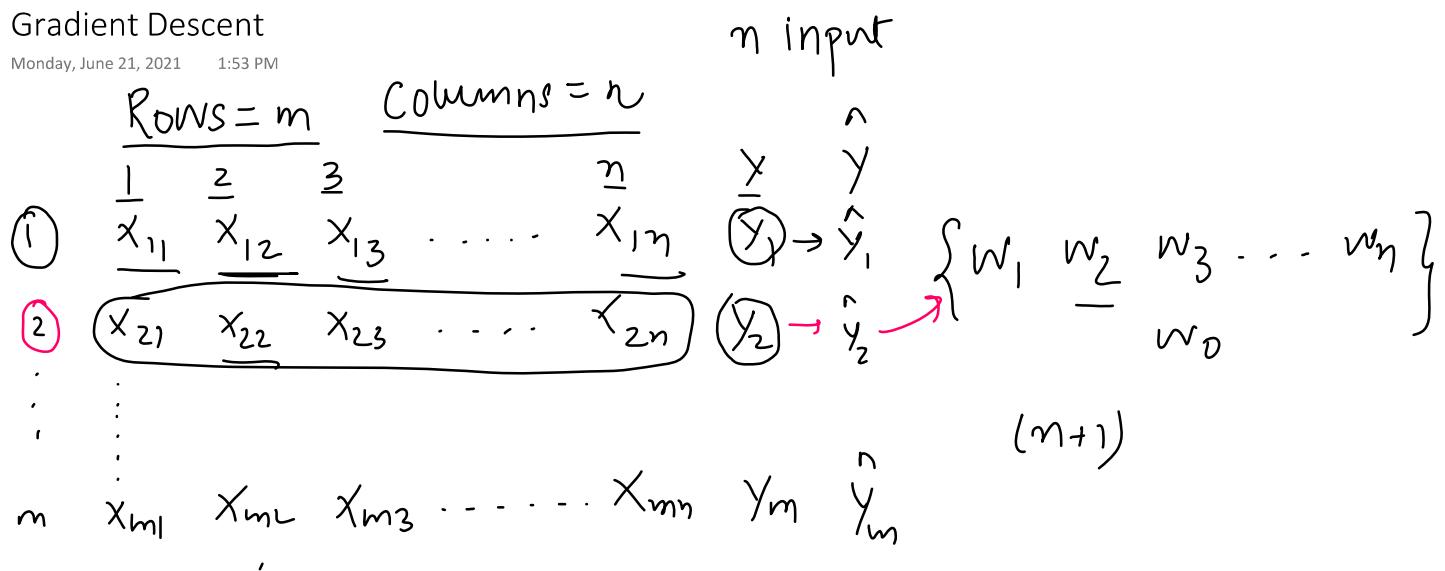
$$= -\frac{1}{2} [1] = -\frac{1}{2}$$

no. of rows = m = 2

no. of cols = n = 2

## Gradient Descent

Monday, June 21, 2021 1:53 PM



$$\sigma(w_0 + w_1 x_{11} + w_2 x_{12} + w_3 x_{13} + \dots + w_n x_{1n}) = \hat{y}_1$$

$$\sigma(w_0 + w_1 x_{21} + w_2 x_{22} + w_3 x_{23} + \dots + w_n x_{2n}) = \hat{y}_2$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} = \begin{bmatrix} \sigma(w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n}) \\ \vdots \\ \sigma(w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n}) \\ \vdots \\ \sigma(w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn}) \end{bmatrix}$$

$$\hat{Y} = \sigma \left( \begin{bmatrix} c \\ w_0 + w_1 x_{11} + w_2 x_{12} + \dots + w_n x_{1n} \\ w_0 + w_1 x_{21} + w_2 x_{22} + \dots + w_n x_{2n} \\ \vdots \\ w_0 + w_1 x_{m1} + w_2 x_{m2} + \dots + w_n x_{mn} \end{bmatrix} \right)$$

$$\hat{Y} = \sigma \left( \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & & & & \end{bmatrix} \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix} \right)$$

$$\hat{y} = \sigma \left( \begin{bmatrix} 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & & & & \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \right)$$

$$\hat{y} = \sigma(xw)$$

$$L = -\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

$$L = -\frac{1}{m} \left[ \sum_{i=1}^m y_i \log(\hat{y}_i) + \sum_{i=1}^m (1-y_i) \log(1-\hat{y}_i) \right]$$

(1-y) log  
-(1-xw)

$$\sum_{i=1}^m y_i \log(\hat{y}_i) = y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3 + \dots + y_m \log \hat{y}_m$$

$$\begin{bmatrix} y_1 & y_2 & y_3 & \dots & y_m \end{bmatrix} \begin{bmatrix} \log \hat{y}_1 \\ \log \hat{y}_2 \\ \vdots \\ \log \hat{y}_m \end{bmatrix}$$

$$\begin{bmatrix} y_1 & y_2 & y_3 & \dots & y_m \end{bmatrix} \log \left( \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} \right)$$

$$y \log \hat{y} = y \log(\sigma(xw))$$

$$L = -\frac{1}{m} \left[ y \log \hat{y} + (1-y) \log(1-\hat{y}) \right]$$

min

where  $\hat{y} = \sigma(xw)$

↑ GD [w] find

where  $\hat{y} = \sigma(xw)$   $L(GD) [WJ] \dots$

Loss function in Matrix form

$$L = -\frac{1}{m} \left[ y \log(\sigma(wx)) + (1-y) \log(1 - \sigma(wx)) \right]$$

minimum

for i in epochs:

$$w = w - \eta \frac{\Delta L}{\Delta w}$$

gradient descent

$$w = [ ]$$

$$\rightarrow \rightarrow \rightarrow \rightarrow \rightarrow$$

$$w_0 \rightarrow w_n$$

$$(n+1)$$

$$\left\{ \frac{\Delta L}{\Delta w} \right\}$$

$$\frac{\Delta L}{\Delta w} = \left[ \frac{\partial L}{\partial w_0}, \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial w_2}, \dots, \frac{\partial L}{\partial w_n} \right]$$

$$\frac{\Delta L}{\Delta w} \quad L = -\frac{1}{m} \left[ y \log \hat{y} + (1-y) \log (1-\hat{y}) \right]$$

$$\frac{dL}{dw} =$$

$$\frac{d}{dw} y \log \hat{y} \Rightarrow y \frac{d}{dw} \log \hat{y} \Rightarrow \frac{y}{\hat{y}} \frac{d}{dL} (\hat{y})$$

$$\Rightarrow \frac{y}{\hat{y}} \frac{d}{dL} \sigma(wx) \Rightarrow \frac{y}{\hat{y}} \sigma(wx) [1 - \sigma(wx)] \frac{d}{dw} (wx)$$

$$= \frac{y}{\hat{y}} \hat{y} (1-\hat{y}) X = \boxed{Y(1-\hat{y})X} \quad \hat{y} = \sigma(wx)$$

$$= \frac{1}{\hat{y}} \cdot 1 - \hat{y} \cdot \frac{\partial}{\partial w} \log(1 - \hat{y})$$

$$= -\frac{(1-y)}{(1-\hat{y})} \frac{\partial}{\partial w} \sigma(wx) \Rightarrow -\frac{(1-y)}{(1-\hat{y})} \left[ \sigma(wx) \left[ 1 - \sigma(wx) \right] \right]$$

$$\Rightarrow -\frac{(1-y)}{(1-\hat{y})} \hat{y} (1-\hat{y}) X = \boxed{-\hat{y} (1-\hat{y}) X}$$

$$\begin{aligned} \frac{\partial L}{\partial w} &= -\frac{1}{m} \left[ y(1-\hat{y})X - \hat{y}(1-y)X \right] \\ &= -\frac{1}{m} \left[ y(1-\hat{y}) - \hat{y}(1-y) \right] X \\ &= -\frac{1}{m} \left[ y - y/\hat{y} - \hat{y} + y/\hat{y} \right] X \end{aligned}$$

$$\boxed{\frac{\Delta L}{\Delta w} = -\frac{1}{m} (y - \hat{y}) X}$$

gd



$$\boxed{w = \underline{w} + \eta \frac{1}{m} (y - \hat{y}) X}$$

$$\underline{w} = \underline{w} + \eta \frac{1}{m} (\underline{y} - \hat{\underline{y}}) \wedge$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_n \end{bmatrix}^{(n+1), 1}$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1n} \\ 1 & x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & & & & \\ 1 & x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}^{m, n+1}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}^{1, m}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix}^{1, m, 1}$$

$$\underline{w} = \underline{w} + \left[ \frac{\eta}{m} \right] (\underline{y} - \hat{\underline{y}}) \underline{X}$$

$$\underline{\underline{w}} = \frac{\underline{w}}{(n+1, 1)} \quad (1, m) \quad \boxed{m, (n+1)} \rightarrow (1, \underline{n+1})$$

$\downarrow$   
 $(n+1), 1$

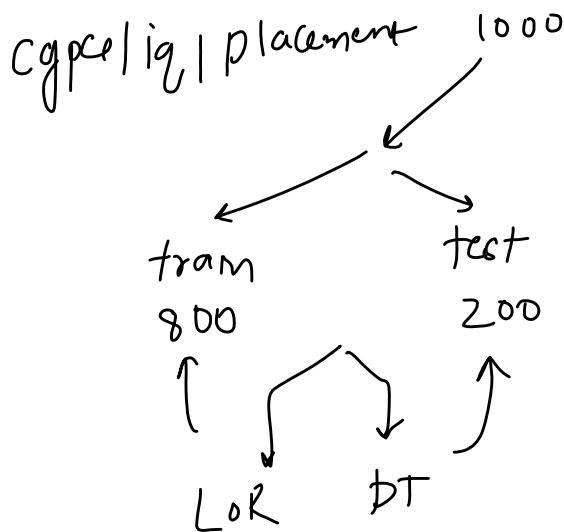
# Accuracy

Tuesday, June 22, 2021 1:12 PM

binary

✓

Actual Label	Logistic Regression Prediction	Decision Tree Prediction
1	✓ 1	✓ 1
0	✗ 1	✗ 1
0	✓ 0	✓ 0
0	✓ 0	✓ 0
1	✓ 1	✓ 1
1	✓ 1	✓ 1
0	✗ 1	✓ 0
0	✓ 0	✓ 0
0	✓ 0	✓ 0
1	✓ 1	✓ 1



$$\frac{8}{10} = 0.8 \rightarrow 80\%$$

$\text{Accuracy} = \frac{\text{no. of } \checkmark}{\text{total predictions}}$

$$\frac{9}{10} = 90\%$$

## Accuracy of multi-classification problem

Wednesday, June 23, 2021 8:44 AM



Actual Label	<u>Logistic Regression Prediction</u>	<u>Decision Tree Prediction</u>
0	✓ 0	0
0	✓ 0	0
0	✓ 0	0
2	✓ 2	2
0	✓ 0	0
2	✓ 2	2
0	✓ 0	0
2	✓ 2	2
1	✓ 1	1
1	✓ 1	1

iris  
setosa, virginica / versicolor  
0 1 2

$$\text{accuracy} = \frac{\# \text{ correct predictions}}{\# \text{ total}}$$

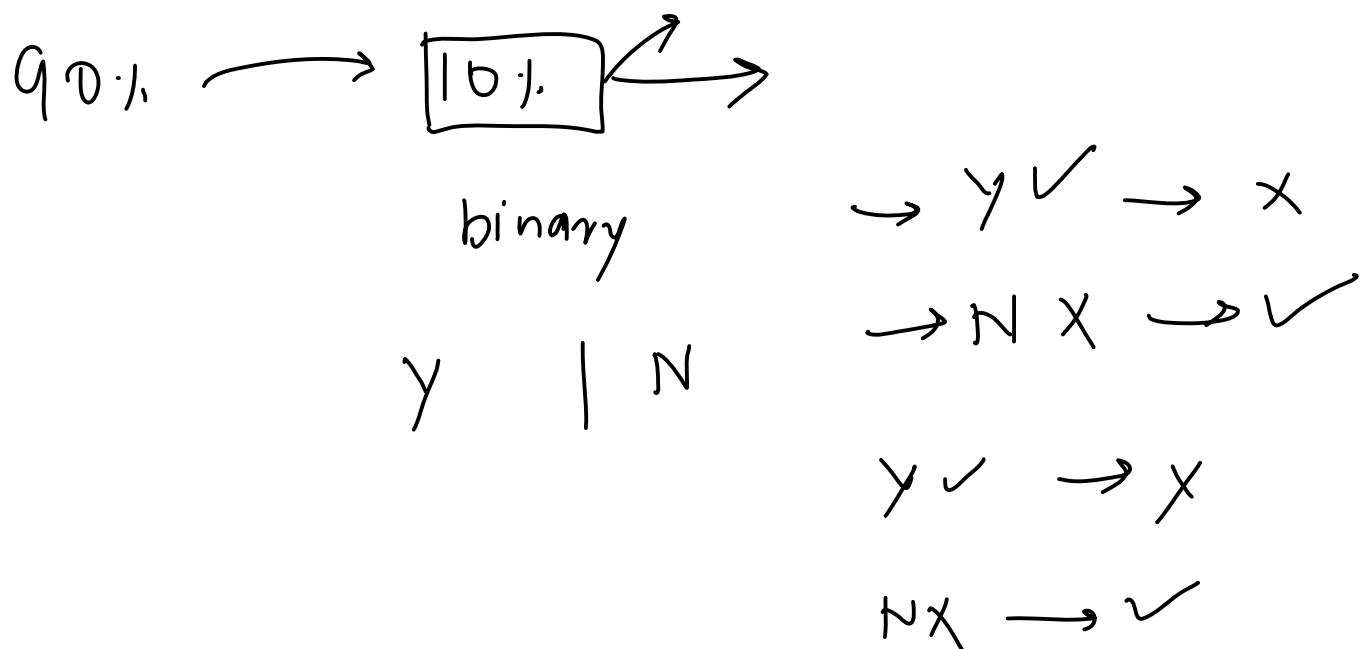
$$= \frac{10}{10} = 1 = 100\%$$

# How much accuracy is good?

Wednesday, June 23, 2021 11:14 AM

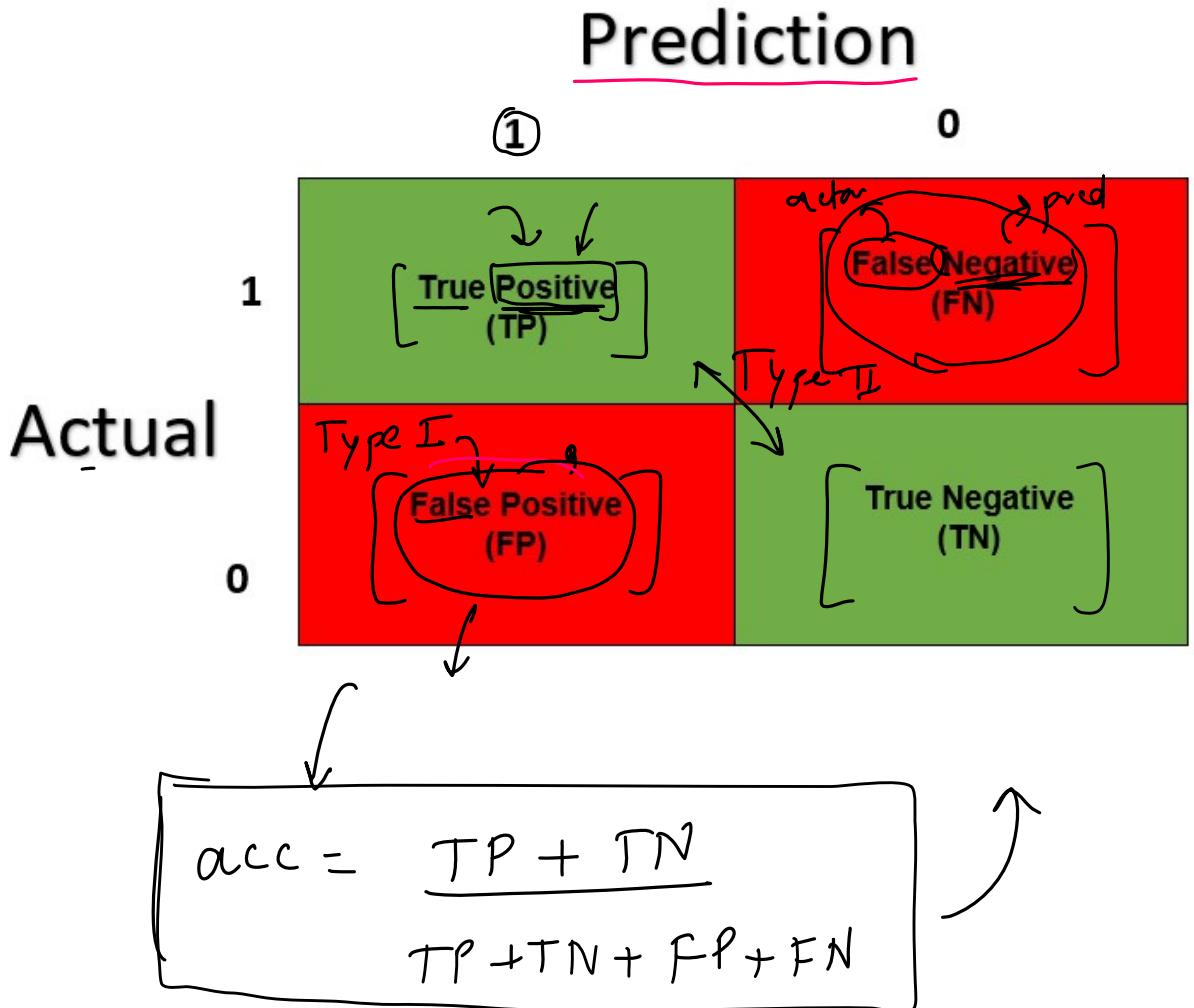
# The Problem with Accuracy

Wednesday, June 23, 2021 11:14 AM



## Confusion Matrix

Tuesday, June 22, 2021 1:12 PM



Extended  
↳ { echo }  
          { or }  
          { not }

duplicate  
↳ echo

$$1 + 2 + 3 = 6 \rightarrow \text{outcome}$$
$$\begin{array}{r} xy \\ \hline zy \end{array}$$

# Type 1 Error

Wednesday, June 23, 2021 8:45 AM

# Type 2 Error

Wednesday, June 23, 2021 8:45 AM

# Confusion Matrix for Multi-classification Problem

Wednesday, June 23, 2021 8:45 AM

0, 1, 2 →

10 x 10

		Predicted				
		0	1	2		
Actual	0	7	0	5		
	1	2	21	6		
	2	9	0	13		

3 x 3  
binary  
2 x 2

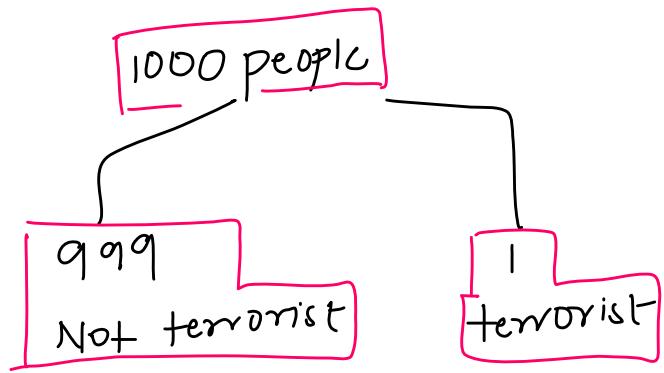
True pred / total = accuracy

3 x 3

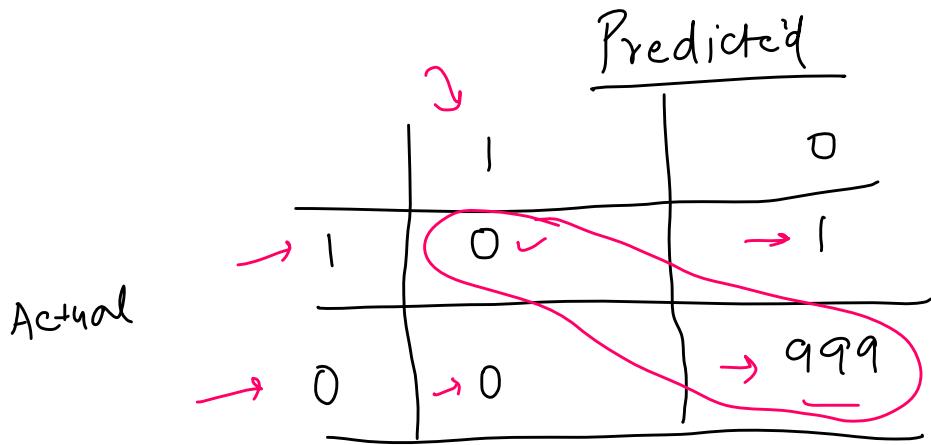
When accuracy is misleading?

Wednesday, June 23, 2021 8:45 AM

## Imbalanced Dataset



model → No one is terrorist



$$\text{Accuracy} = \frac{999}{999 + 1}$$
$$= 99.9\%$$

## Precision

Wednesday, June 23, 2021 8:46 AM

{ spam: 1, not spam: 0 }

Actual

①(A)

Predicted

②(B)

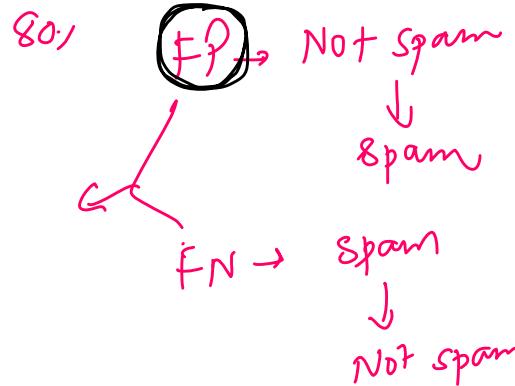
	Sent to Spam	Not sent to spam
Spam	100	170 FN
Not Spam	30 FP	700

	Sent to Spam	Not sent to spam
Spam	100	190
Not Spam	10	700

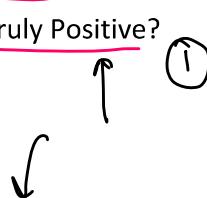
$$P_A = \frac{100}{100 + 30}$$

$$P_B = \frac{100}{100 + 10}$$

$$P_A < P_B$$



What proportion of predicted Positives is truly Positive?



0

Binary  
tvs

	Sent to Spam	Not sent to spam
Spam	True Positive	False Negative
Not Spam	False Positive	True Negative

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

## Recall

Wednesday, June 23, 2021 8:46 AM

has cancer: 1, no cancer: 0

Predicted

Actual

	Detected Cancer	Not Detected
Has Cancer	1000	200 <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">FN</span>
No Cancer	800 <span style="border: 1px solid red; border-radius: 50%; padding: 2px;">FP</span>	8000

A 90%

$$\text{Recall}_A = \frac{1000}{1200}$$

$R_A > R_B$

B 90%

$$R_B = \frac{1000}{1500}$$

$R_B < R_A$

What proportion of actual Positives is correctly classified?

	Detected Cancer	Not detected Cancer
Has Cancer	<u>True Positive</u>	<u>False Negative</u>
No Cancer	False Positive	True Negative

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

## F1 Score

Wednesday, June 23, 2021 8:46 AM

$$F1 \text{ score} = \frac{2PR}{P+R}$$

$$P = 0 \quad R = 100$$

$$F1 = 0 \quad F1 \text{ score} = 50$$

$$F1 \text{ score} = \frac{P+R}{2}$$

(A)

$$P = R = 80$$

$$\frac{80 + 80}{2} = 80 \quad F1 = 80$$

$$\frac{2 \times 80 \times 80}{160} = 80$$

(B)  $\frac{75}{100} = 75$

$$P = 60 \quad R = 100$$

$$\frac{100 + 60}{2} = 80$$

$$\frac{2 \times 60 \times 100}{160} = 75$$

## Multi-class Precision and Recall

Wednesday, June 23, 2021 8:46 AM

Positive
Binary  
① ↗  
Yes  
No
Multi  
② class  
③
Dog Cat Rabbit

$$\frac{2 \times 0.86 + 0.66}{0.86 + 0.66}$$

Predicted

	Dog	Cat	Rabbit	Total	Recall
Dog	25	5	10	40	0.62
Cat	0	30	4	34	0.88
Rabbit	4	10	20	34	0.58
Total	29	45	34		
Precision	0.86	0.66	0.58		

Actual

$$F1_D = \frac{2 P_D R_D}{P_D + R_D}$$

$$F1_C = \frac{2 P_C R_C}{P_C + R_C}$$

$$F1_R = \frac{2 P_R R_R}{P_R + R_R}$$

$$R_D = \frac{25}{40} = 0.62, \quad R_C = \frac{30}{34} = 0.88, \quad R_R = \frac{20}{34} = 0.58$$

+ macro recall ↗
+ weighted f1

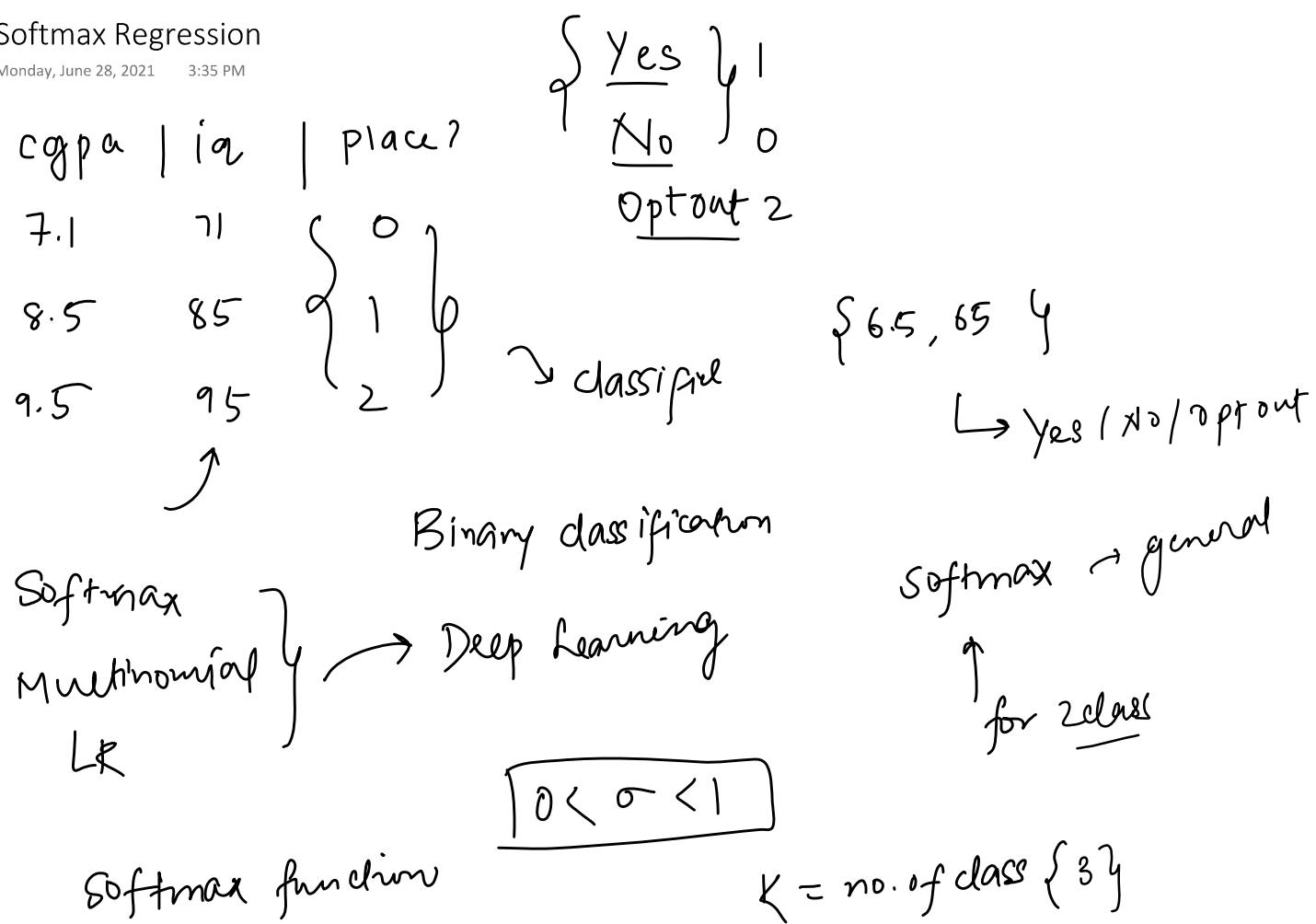
+ weight recall ↗
+ weighted f1

# Multi-class F1 Score

Thursday, June 24, 2021 11:14 AM

## Softmax Regression

Monday, June 28, 2021 3:35 PM



yes → 1

no → 2

opt → 3

(Yes)

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

$$\sigma(z)_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\sigma(z)_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

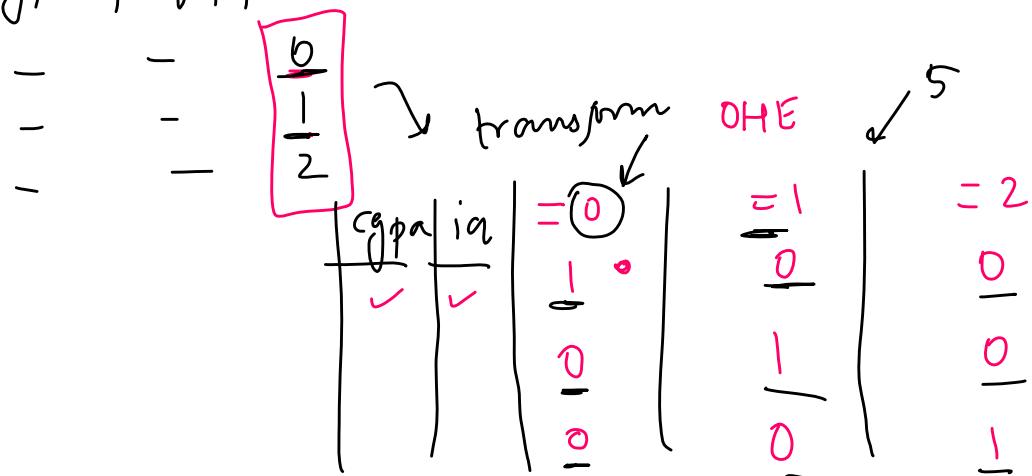
$$\sigma(z)_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

# Training Intuition

Monday, June 28, 2021 3:35 PM

$$\{y_0, y_1, y_2\}$$

cgpa | iq | place?



3 coeff

$$w_1, w_2, w_0$$

$$D \rightarrow D_1$$

$$D \rightarrow D_2$$

$$D \rightarrow D_3$$

$$cgpa | iq | = 0$$

$$m_1$$

$$w_1^{(0)}, w_2^{(0)}, w_0^{(0)}$$

{  
loss function  
gradient descent}

$$cgpa | iq | = 1$$

$$m_2$$

$$w_1^{(1)}, w_2^{(1)}, w_0^{(1)}$$

$$w_1^{(2)}, w_2^{(2)}, w_0^{(2)}$$

$$cgpa | iq | = 2$$

$$m_3$$

## Loss Function

Monday, June 28, 2021 3:47 PM

$$L = -\frac{1}{m} \sum_{i=1}^m y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)$$

$\leftarrow \{1, 2, 3\} \rightarrow i=1$

$$L = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(\hat{y}_k^{(i)})$$

$x_1 \quad x_2 \quad y \quad y_{k=1} \quad y_{k=2} \quad y_{k=3}$

$x_{11}$	$x_{12}$	1	1	0	0
$x_{21}$	$x_{22}$	2	0	1	0
$x_{31}$	$x_{32}$	3	0	0	1

$w_1^{(1)} \quad w_2^{(1)} \quad w_0^{(1)}$

$$y_1^{(1)} \log(\hat{y}_1^{(1)}) + y_2^{(1)} \log(\hat{y}_2^{(1)}) + y_3^{(1)} \log(\hat{y}_3^{(1)}) +$$

$$+ y_1^{(2)} \log(\hat{y}_1^{(2)}) + y_2^{(2)} \log(\hat{y}_2^{(2)}) + y_3^{(2)} \log(\hat{y}_3^{(2)}) +$$

$$+ y_1^{(3)} \log(\hat{y}_1^{(3)}) + y_2^{(3)} \log(\hat{y}_2^{(3)}) + y_3^{(3)} \log(\hat{y}_3^{(3)}) +$$

$$L = y_1^{(1)} \log(\hat{y}_1^{(1)}) + y_2^{(2)} \log(\hat{y}_2^{(2)}) + y_3^{(3)} \log(\hat{y}_3^{(3)})$$

$$\hat{y}_1^{(1)}, \hat{y}_2^{(2)}, \hat{y}_3^{(3)}$$

softmax

$$\hat{y}_1^{(1)} = \sigma(w_1^{(1)}x_{11} + w_2^{(1)}x_{12} + w_0^{(1)})$$

$$y_2^{(2)} = \sigma(w_1^{(2)}x_{21} + w_2^{(2)}x_{22} + w_0^{(2)})$$

$$y_3^{(3)} = \sigma(w_1^{(3)}x_{31} + w_2^{(3)}x_{32} + w_0^{(3)})$$

$$\begin{bmatrix} w_1^{(1)} & w_2^{(1)} & w_0^{(1)} \\ w_1^{(2)} & w_2^{(2)} & w_0^{(2)} \\ w_1^{(3)} & w_2^{(3)} & w_0^{(3)} \end{bmatrix}$$

(L)

$$\frac{\partial L}{\partial w_1^{(1)}}, \frac{\partial L}{\partial w_2^{(1)}}, \frac{\partial L}{\partial w_0^{(0)}} \dots \quad q \text{ due}$$

gradient

*q-values init=1*

$$\left[ \begin{array}{c} \quad \\ \quad \\ \end{array} \right] =$$

loop  $\rightarrow$  1000 epochs

$$w_1^{(1)} = w_1^{(1)} - \eta \frac{\partial L}{\partial w_1^{(1)}}$$

$$w_2^{(1)} = w_2^{(1)} - \eta \frac{\partial L}{\partial w_2^{(1)}}$$

;

# Prediction

Monday, June 28, 2021 3:35 PM

$$\begin{aligned}
 S_x = \{7, 70\} &\Rightarrow \overline{Y, N, OPR} \\
 \downarrow & \downarrow \\
 m_1 & \text{Yes} \\
 w_1^{(1)} & w_2^{(1)} w_0^{(1)} \\
 \underline{z}_1 = 7 \times w_1^{(1)} + 70 \times w_2^{(1)} + w_0^{(1)} & \quad \downarrow \\
 \downarrow & \downarrow \\
 m_2 & \text{No} \\
 w_1^{(2)} & w_2^{(2)} w_0^{(2)} \\
 \underline{z}_2 = 7 \times w_1^{(2)} + 70 \times w_2^{(2)} + w_0^{(2)} & \quad \downarrow \\
 \downarrow & \downarrow \\
 m_3 & \text{Opt Out} \\
 w_1^{(3)} & w_2^{(3)} w_0^{(3)} \\
 \underline{z}_3 = 7 \times w_1^{(3)} + 70 \times w_2^{(3)} + w_0^{(3)} & \quad \downarrow \\
 \downarrow & \downarrow \\
 \sigma(y) = \frac{e^{\underline{z}_1}}{e^{\underline{z}_1} + e^{\underline{z}_2} + e^{\underline{z}_3}} & \quad \downarrow \\
 \underline{\underline{\sigma(y)}} = 0.40 & \quad \downarrow \\
 \sigma(N) = \frac{e^{\underline{z}_2}}{e^{\underline{z}_1} + e^{\underline{z}_2} + e^{\underline{z}_3}} & \quad \downarrow \\
 \underline{\underline{\sigma(N)}} = 0.35 & \quad \downarrow \\
 \sigma(O) = \frac{e^{\underline{z}_3}}{e^{\underline{z}_1} + e^{\underline{z}_2} + e^{\underline{z}_3}} & \quad \downarrow \\
 \underline{\underline{\sigma(O)}} = 0.25 &
 \end{aligned}$$

# Sigmoid Vs Softmax

Monday, June 28, 2021 3:47 PM

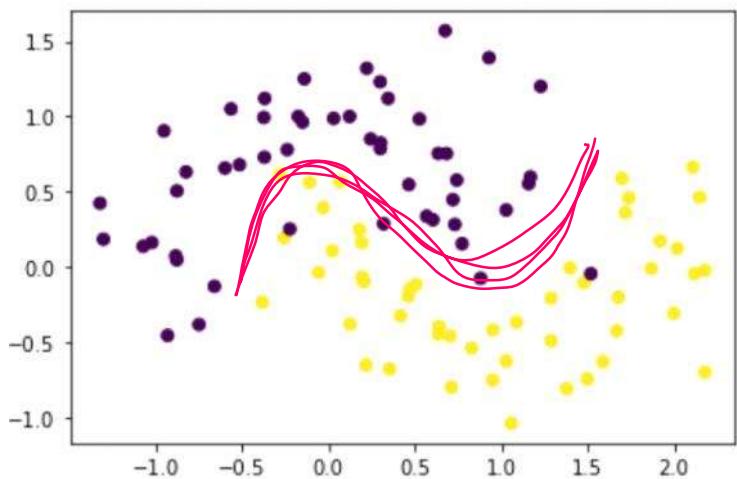
# Code Sample

Monday, June 28, 2021 3:36 PM

## Polynomial Logistic Regression

Monday, June 28, 2021 6:43 PM

$$\underline{x_1} \underline{x_2} \underline{y} \rightarrow \{0, 1\}$$



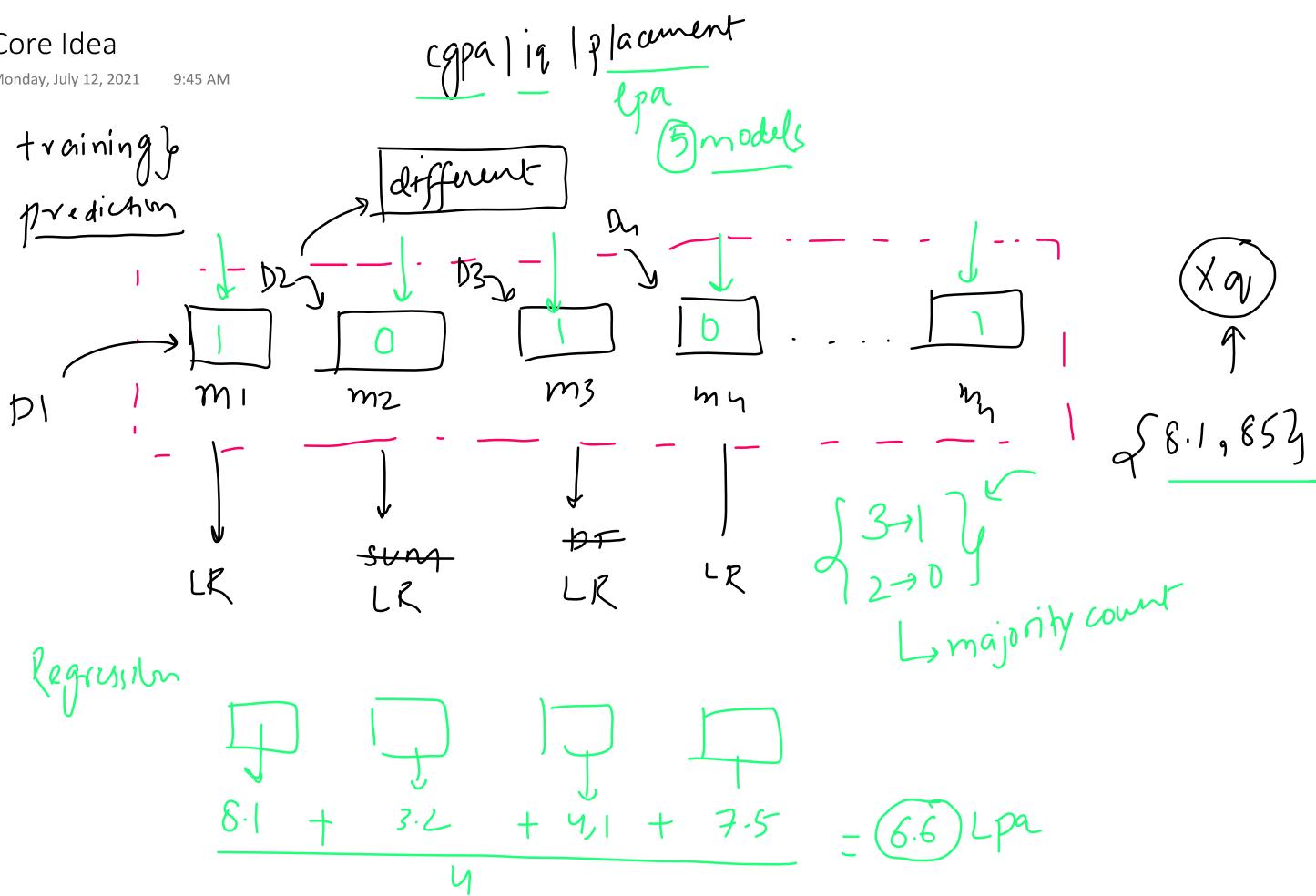
Linear       $x^0 x^1 x^2 x^3$   
 polynomial  
 degree = 2  
 $\begin{bmatrix} x_1 & x_2 \end{bmatrix} \rightarrow 2 \text{ cols}$   
 degree  
 $w_0 \rightarrow$   
 $\begin{bmatrix} x_1^0 & x_1^1 & x_1^2 & x_2^0 & x_2^1 & x_2^2 \end{bmatrix} \rightarrow 6 \text{ cols}$

# Wisdom of the Crowd

Monday, July 12, 2021 9:56 AM

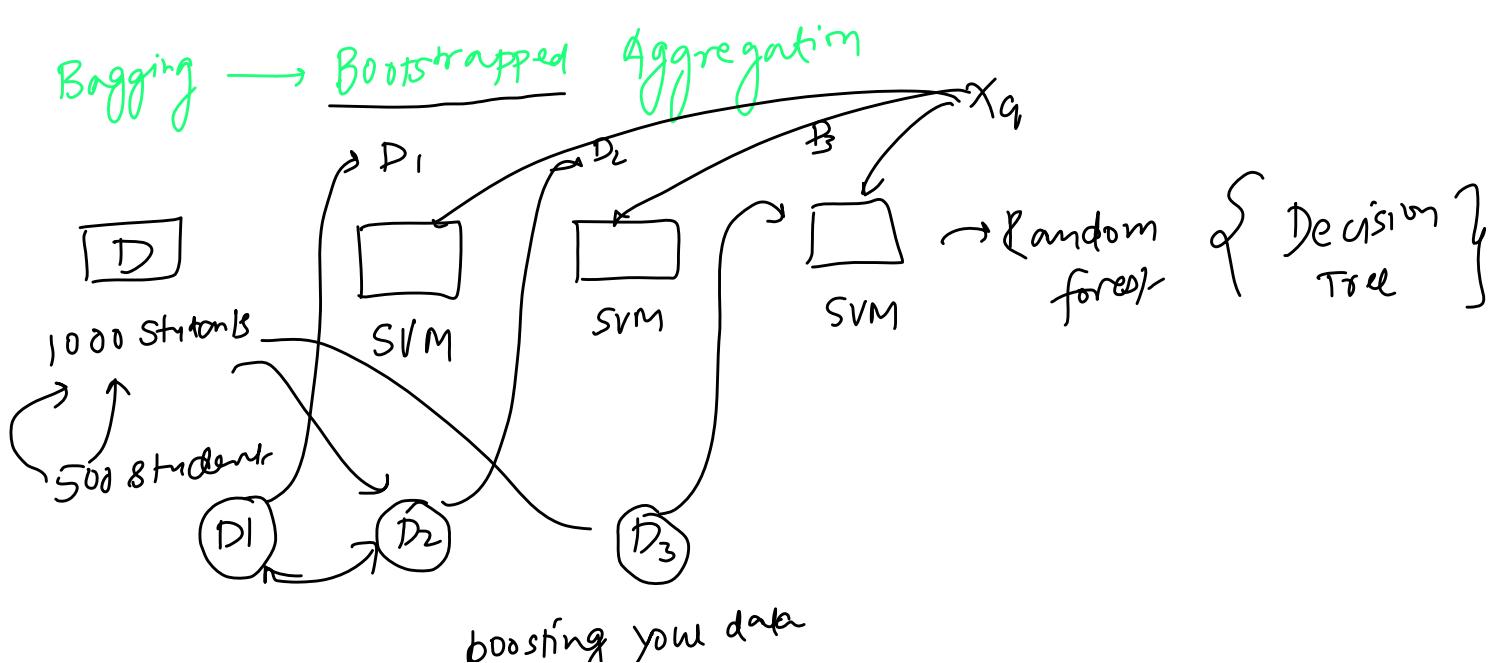
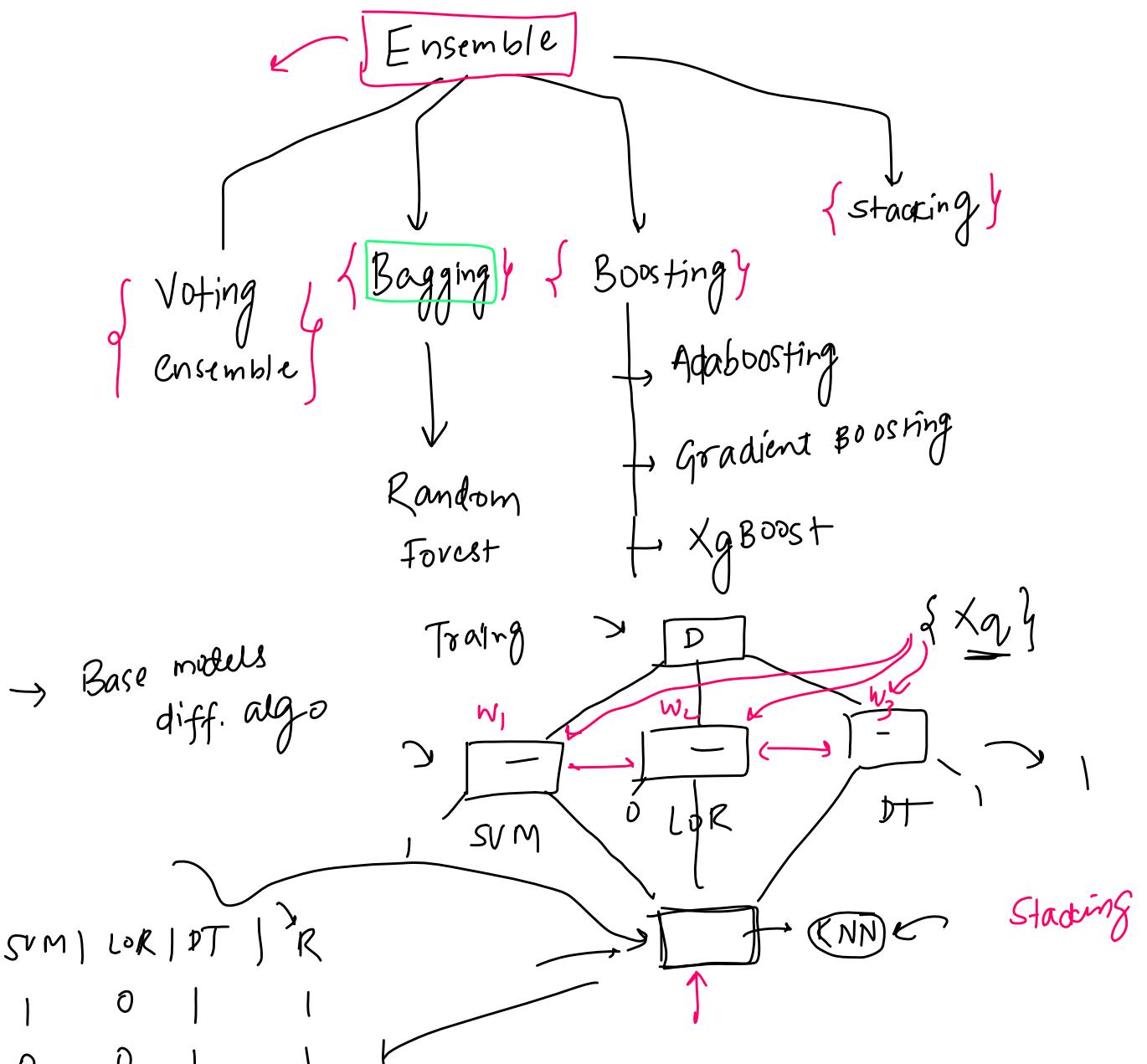
## Core Idea

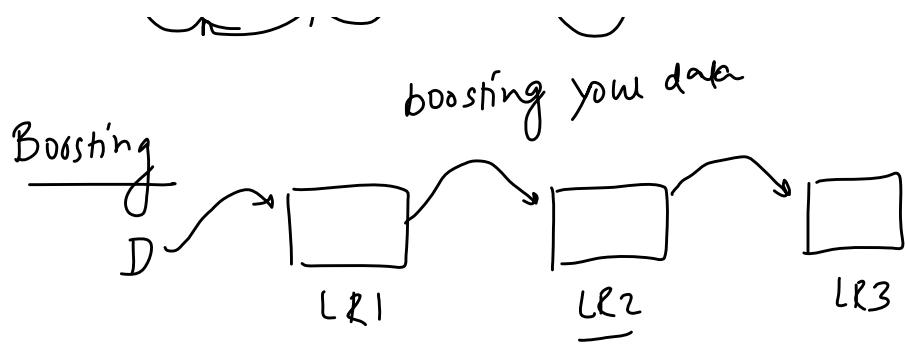
Monday, July 12, 2021 9:45 AM



# Type of Ensemble Learning

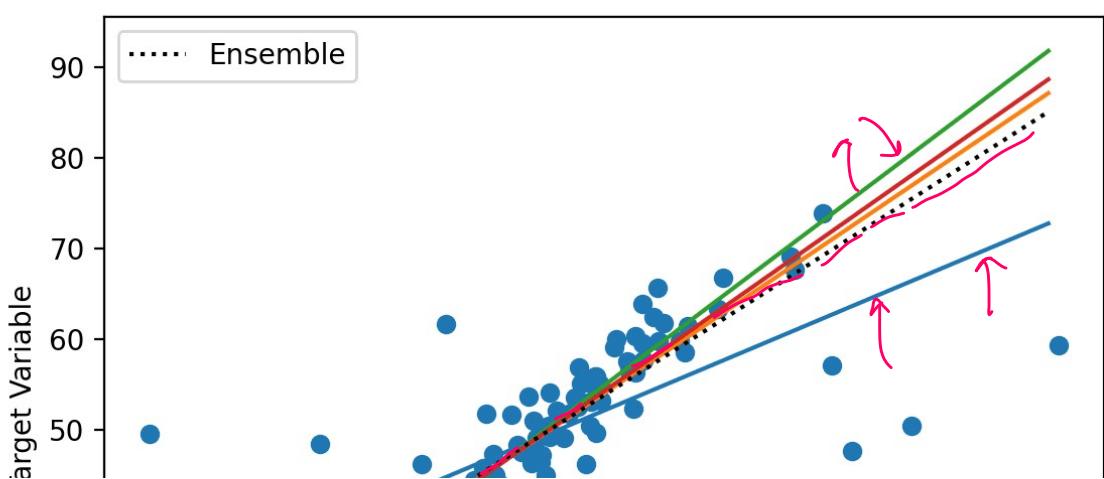
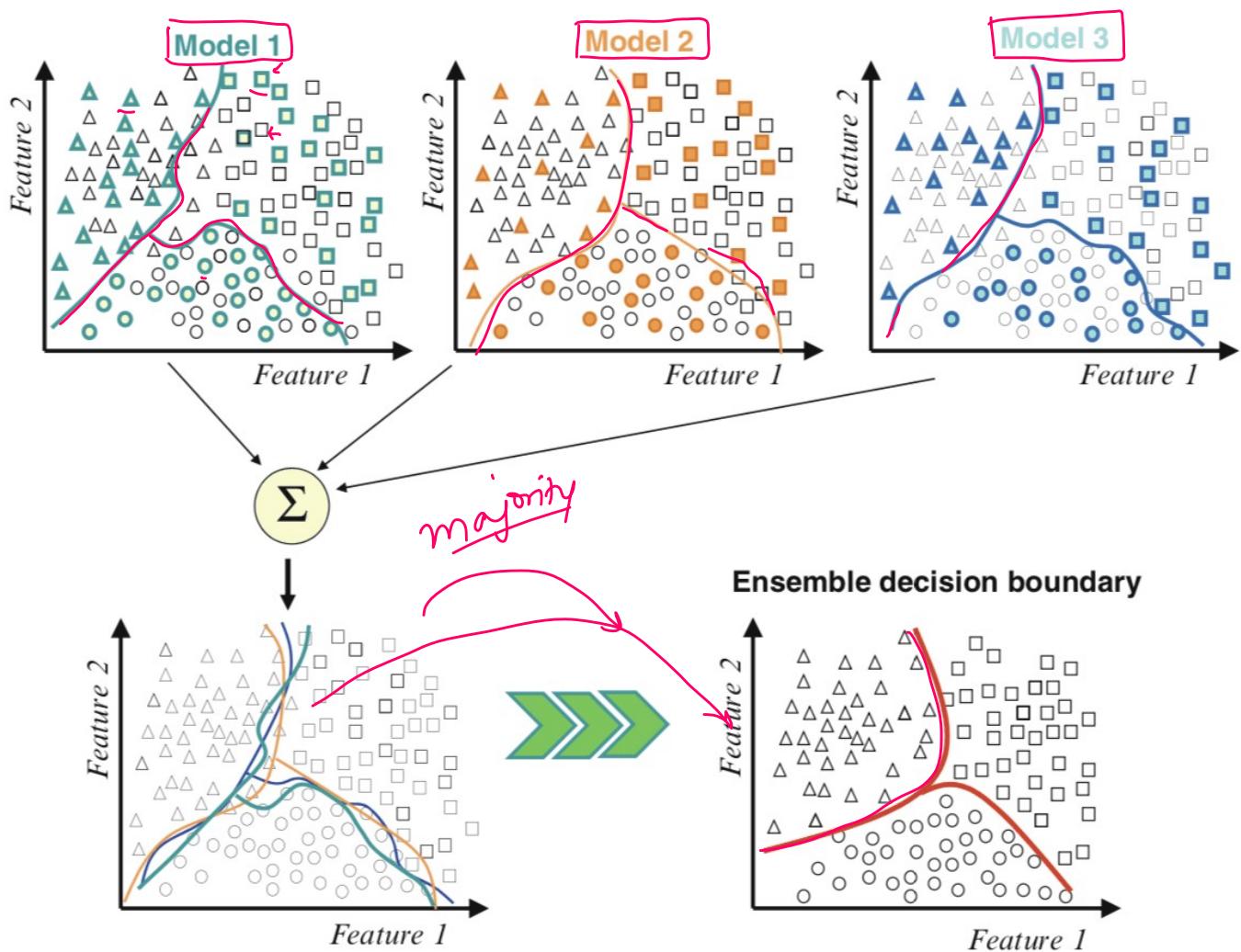
Monday, July 12, 2021 9:45 AM

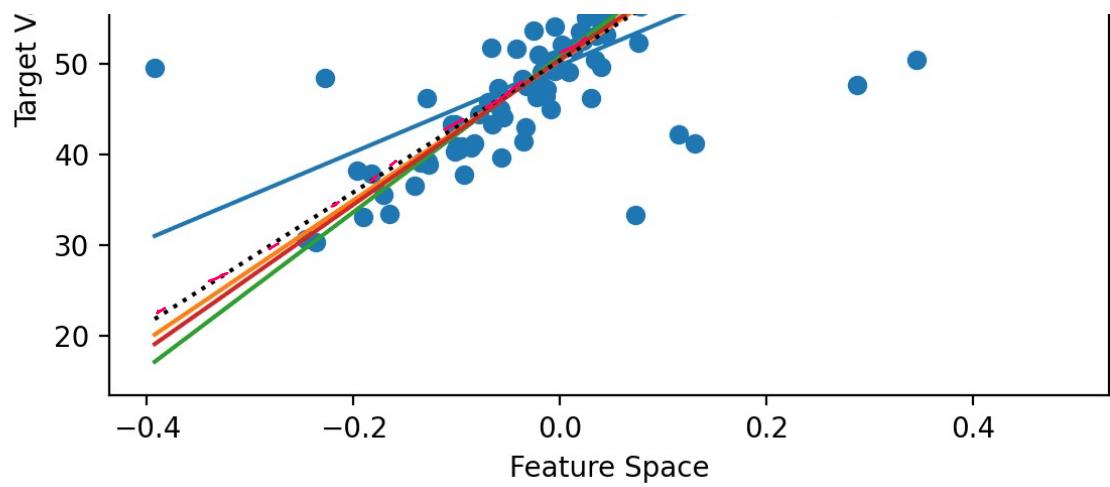




## Why it works?

Monday, July 12, 2021 9:45 AM





## Benefits

Monday, July 12, 2021 9:45 AM

1) Improvement in performance

2) Bias Variance



Low Bias + Low Variance

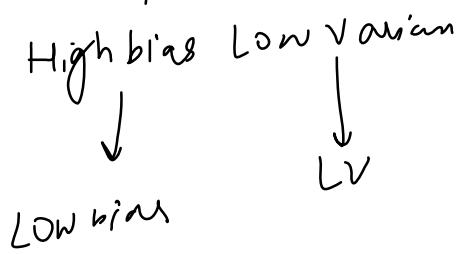
Low Bias      High Variance



→ Ensemble

Low Bias → LV

3) Robustness



# When to use?

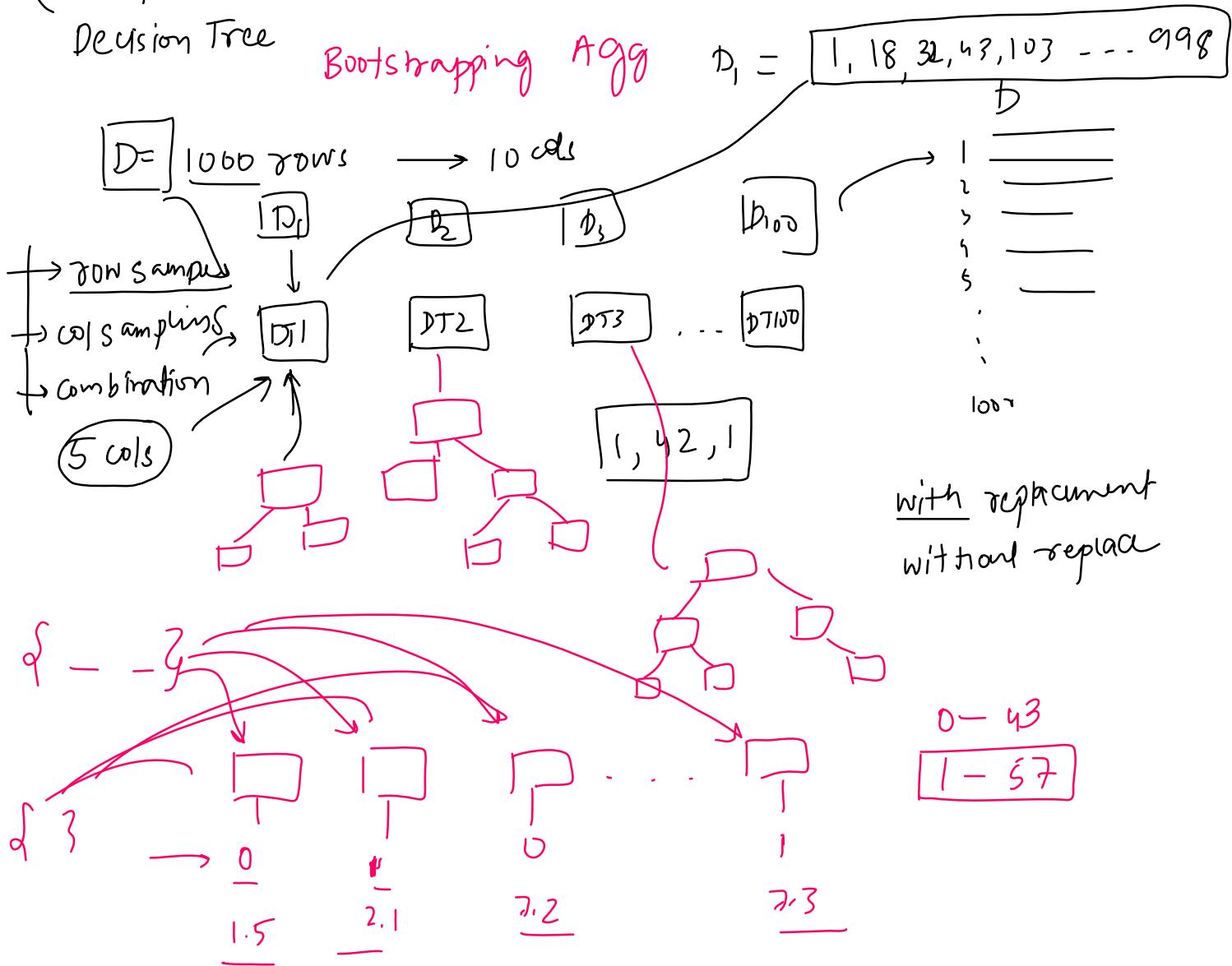
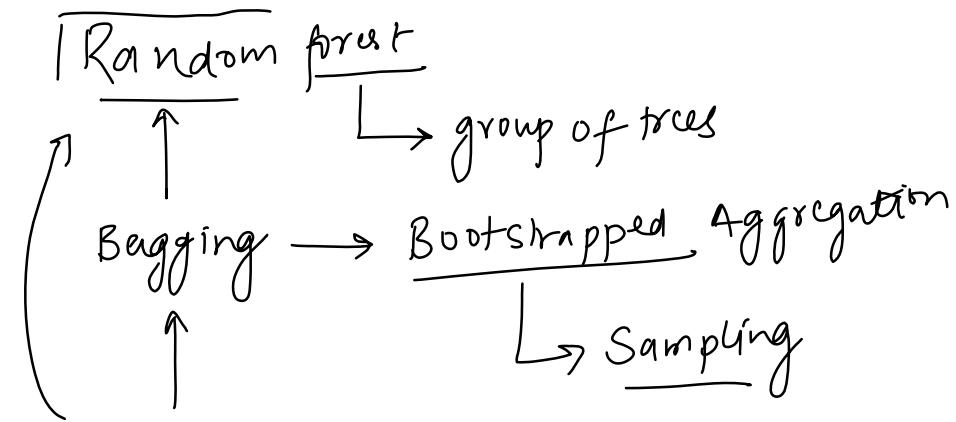
Monday, July 12, 2021 9:46 AM

Always



# Intuition

Monday, July 19, 2021 4:53 PM

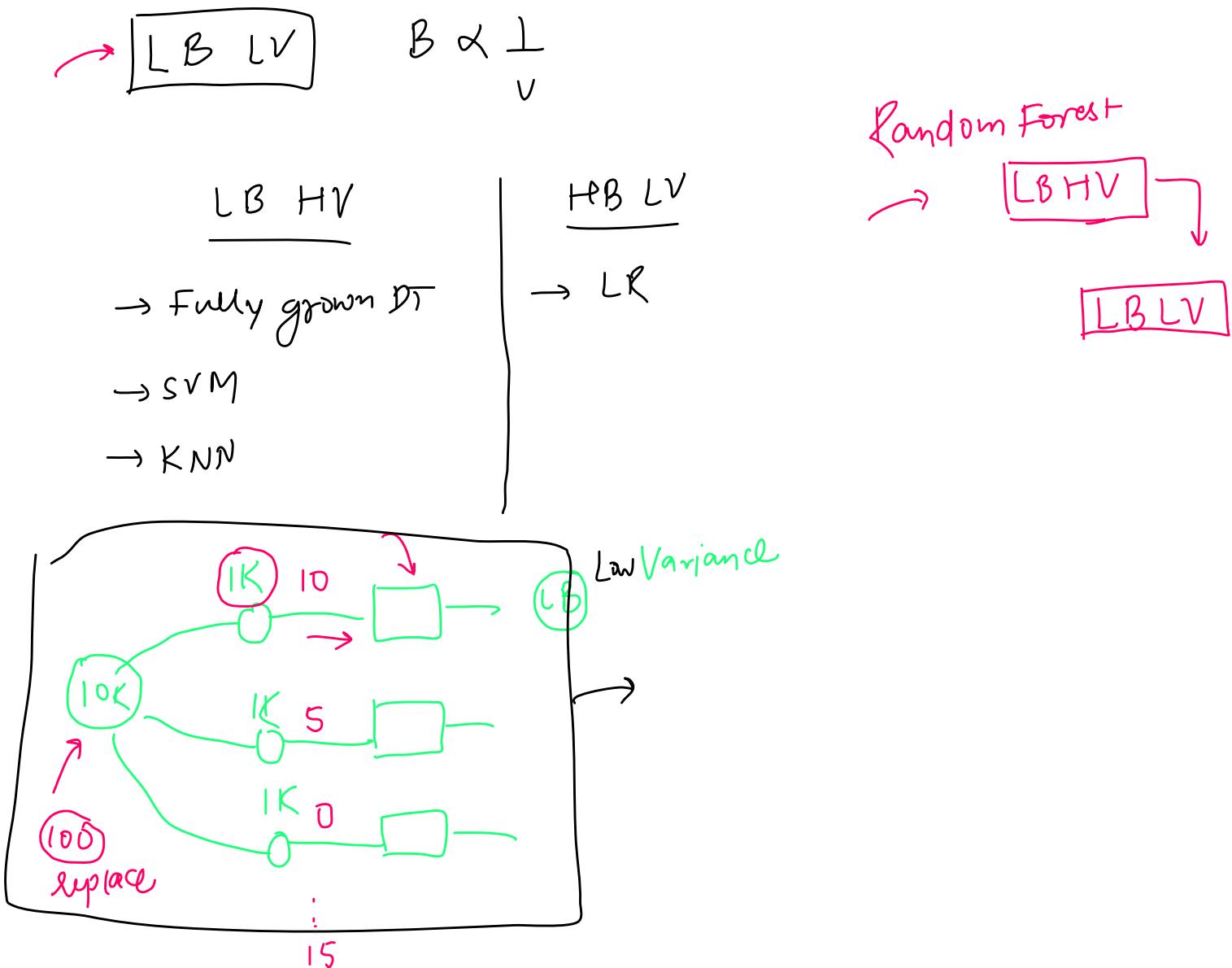


# Demo

Monday, July 19, 2021 5:05 PM

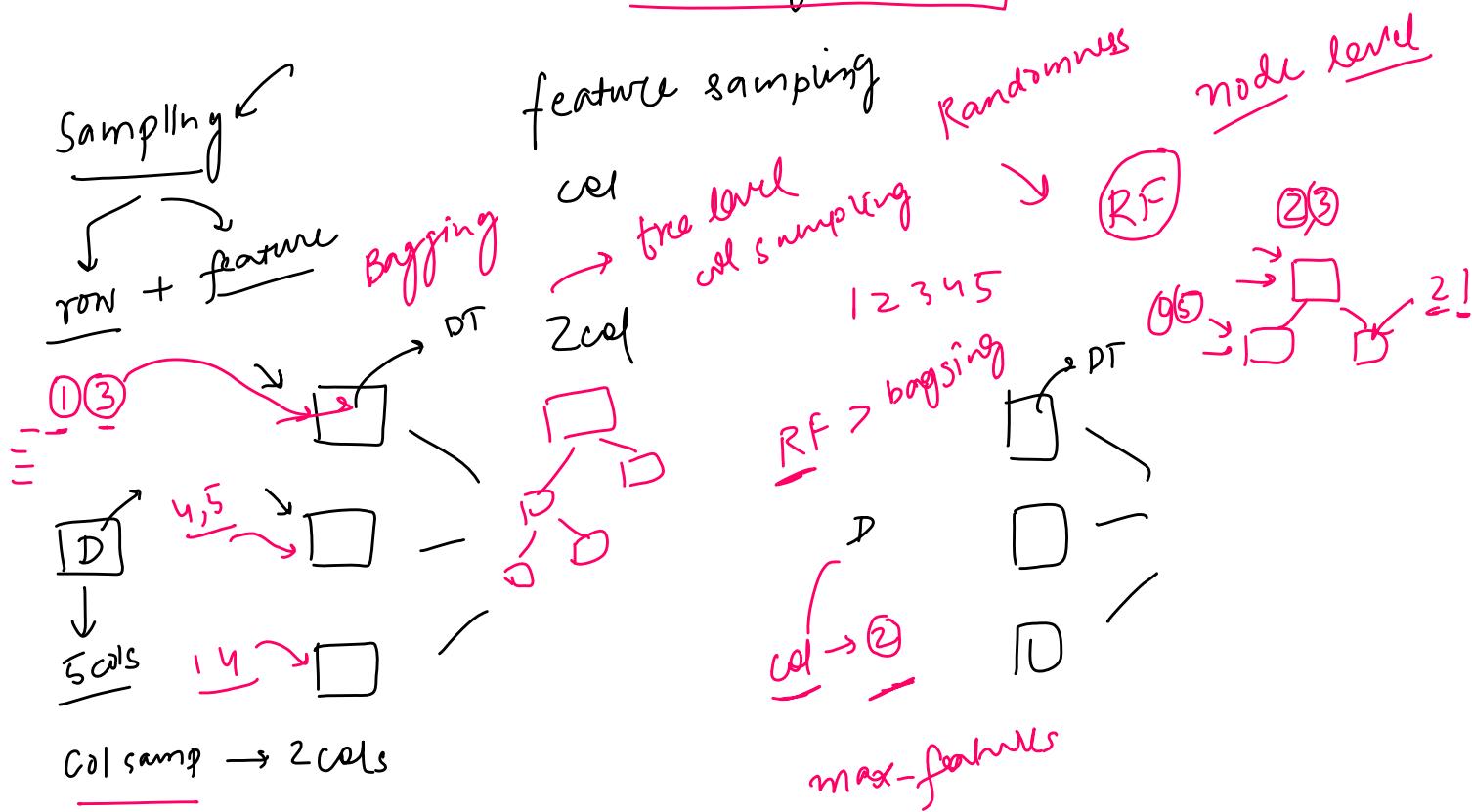
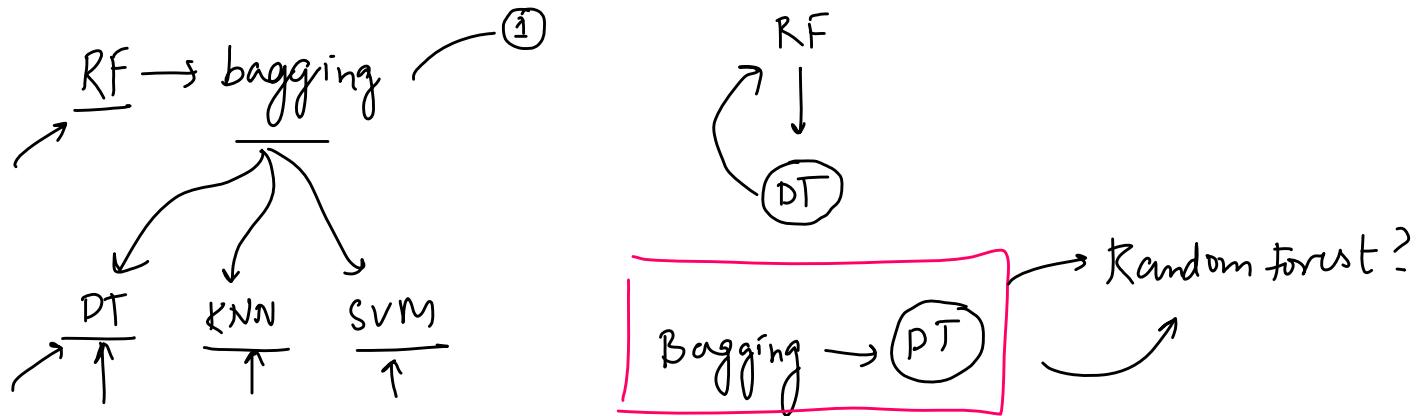
# Bias Variance and Random Forest

Tuesday, July 20, 2021 10:52 AM



# Bagging Vs Random Forest

Monday, July 19, 2021 4:53 PM



# Hyperparameters

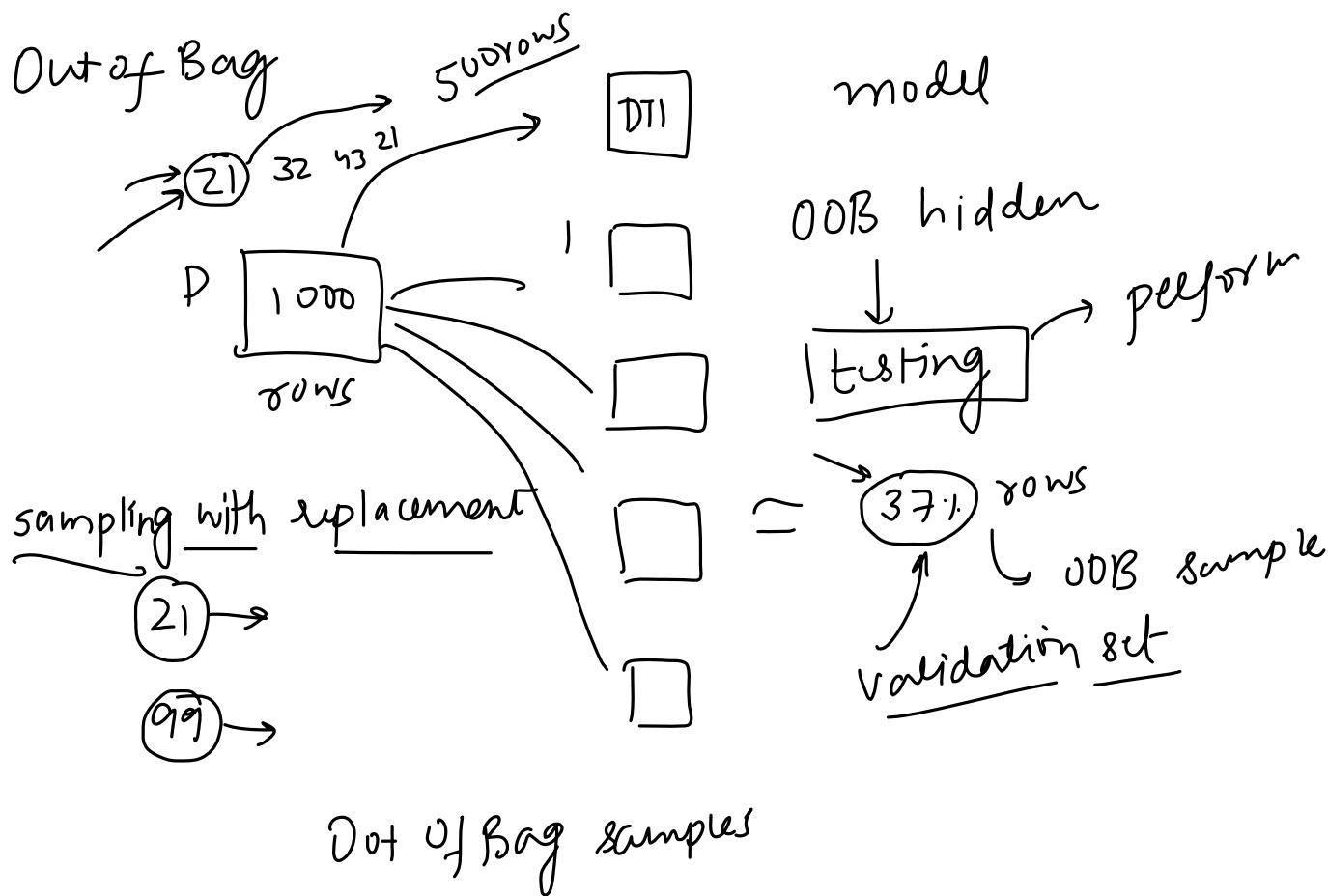
Monday, July 19, 2021 4:56 PM

# Code Demo

Monday, July 19, 2021 5:05 PM

# OOB Evaluation

Monday, July 19, 2021 5:14 PM

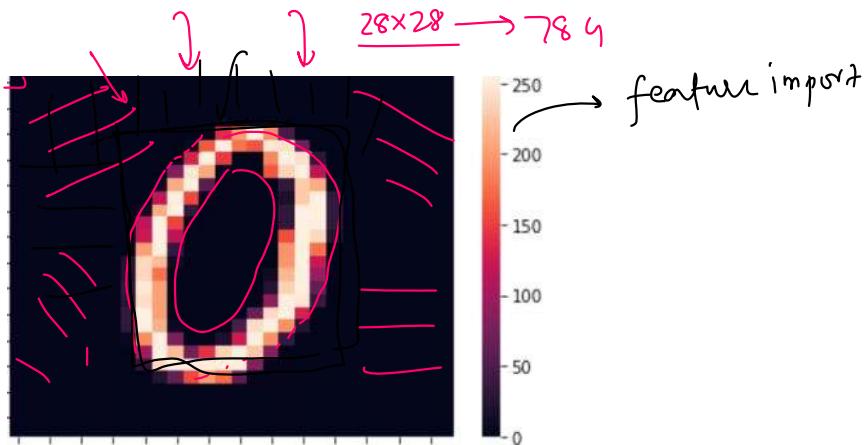


## Feature Importance

Monday, July 19, 2021 4:56 PM

MNIST

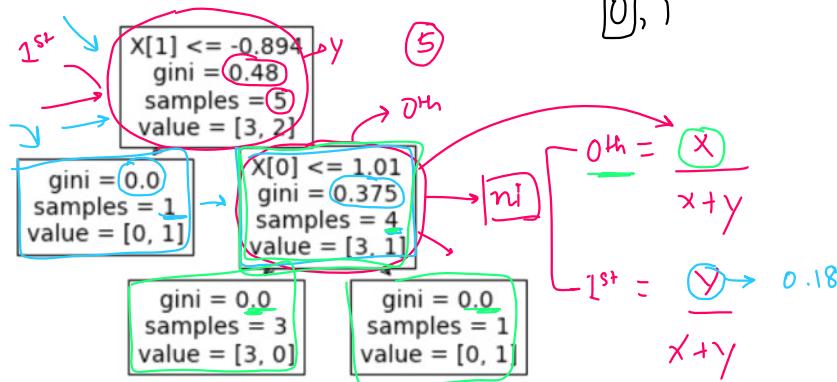
label	pixel0	pixel1	pixel2	pixel3	pixel4	pixel5	pixel6	pixel7	pixel8	...	pixel774	pixel775	pixel776	pixel777	pixel778
0	1	0	0	0	0	0	0	0	0	...	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0
2	1	0	0	0	0	0	0	0	0	...	0	0	0	0	0
3	4	0	0	0	0	0	0	0	0	...	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0



formula

$$ni = \frac{N-t}{N} \left[ \text{impurity} - \left( \frac{N-t-\text{right\_samples}}{N-t} \times \text{right\_impurity} \right) - \left( \frac{N-t-L}{N-t} \times \text{left\_impurity} \right) \right]$$

$$f_{ik} = \frac{\sum_{j \in \text{node split on feature } k} ni}{\sum_{j \in \text{all nodes}} ni}$$



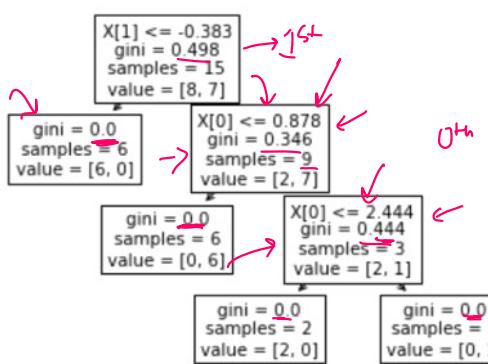
$$= \frac{5}{5} \left[ 0.48 - \frac{4}{5} \times 0.375 - \frac{1}{5} \times 0. \right] = 0.18$$

$$\frac{4}{5} \left[ 0.375 \right] = x = 0.8 \times 0.375 = 0.30 = x$$

$$0^{th} = \frac{0.3}{0.3 + 0.18} = 0.625$$

$$0^{th} = \frac{0.3}{0.3 + 0.18} = 0.625$$

$$1^{st} = \frac{0.18}{0.3 + 0.18} = 0.375$$



0<sup>th</sup> Node

$$\frac{15}{15} \left[ 0.49 - \frac{9}{15} \times 0.34 \right]$$

$$= 0.290$$

↓

$$\frac{9}{15} \left[ 0.34 - \frac{3}{15} \times 0.44 \right]$$

$$= 0.118 \leftarrow$$

$$\frac{3}{15} \left[ 0.44 \right] \leftarrow$$

$$= 0.088$$

$$f_i[1] = \frac{0.29}{0.48} = 0.60$$

$$f_i[0] = \frac{0.11 + 0.08}{0.29 + 0.11 + 0.08} = 0.48$$

$$= 0.39$$

# Extra Trees

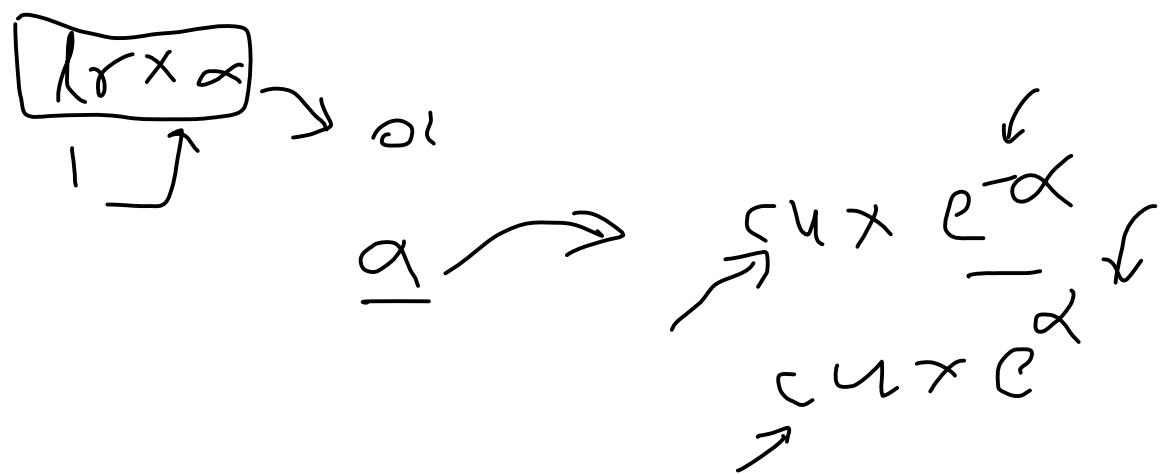
Monday, July 19, 2021 4:57 PM

# Before Starting

Monday, August 9, 2021 10:42 AM

1) Weak classifiers ( $> 50\%$ )

$$\alpha = \frac{1}{2} \log \left( \frac{1-p}{p} \right)$$



# The Big Idea

Monday, August 9, 2021 10:43 AM

# Step by Step Breakdown

Monday, August 9, 2021 10:43 AM

# Points to remember

Monday, August 9, 2021 10:43 AM

# Code Walkthrough

Monday, August 9, 2021 10:43 AM

# Example and Hyperparameters

Monday, August 9, 2021 10:43 AM

## Algorithm

Monday, September 20, 2021 8:23 AM

$$\rightarrow \{(x_i, y_i)\}_{i=1}^n \quad n=3 \quad x_i \ y_i$$

Input: training set  $\{(x_i, y_i)\}_{i=1}^n$ , a differentiable loss function  $L(y, F(x))$ , number of iterations  $M$ .

$\rightarrow$  1. Initialize  $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$ .

$\rightarrow$  2. For  $m = 1$  to  $M$ :

$\rightarrow$  (a) For  $i = 1, 2, \dots, N$  compute

$$r_{im} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}.$$

*residual / pseudo-residual*

(b) Fit a regression tree to the targets  $r_{im}$  giving terminal regions  $R_{jm}$ ,  $j = 1, 2, \dots, J_m$ .

(c) For  $j = 1, 2, \dots, J_m$  compute

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma).$$

(d) Update  $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$ .

3. Output  $\hat{f}(x) = f_M(x)$ .

$$\oplus \quad f_1(x) = f_0(x) + dT$$

$$f_2(x) = f_1(x) + dT$$

$$f_2(x) = f_1(x) + dT$$

$$f_1(x) + dT$$

$$f_0(x) + dT$$

$\rightarrow$  reursion

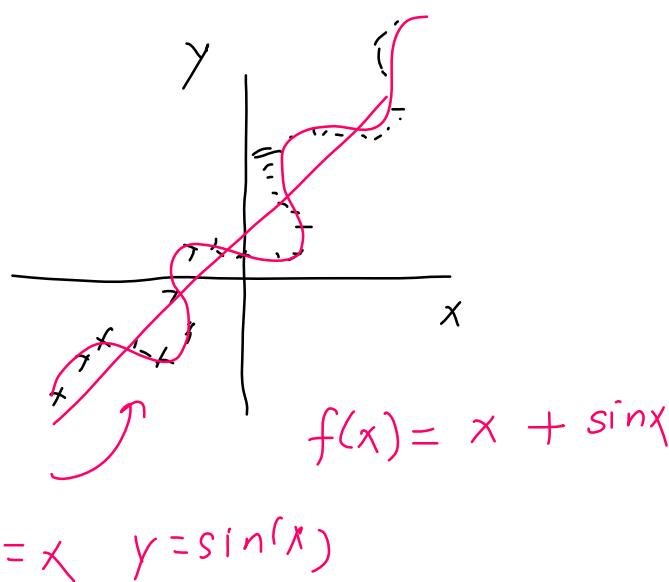
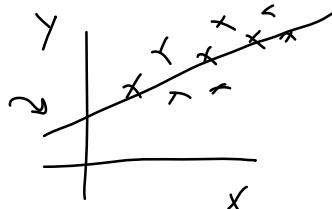
$$f_4(x) = f_0(x) + \dots$$

$$f_1(x) \ f_2(x)$$

## Additive Modelling

Monday, September 20, 2021 8:28 AM

$$\begin{array}{c} \text{X} \\ \downarrow \\ \text{Y} \rightarrow f(\cdot) \Rightarrow Y = f(x) \\ \text{x} \nearrow \\ y = f(x_1, x_2, x_3) \\ x_1, x_2, x_3 | y \end{array}$$



additive

$$F(x) = f_0(x) + f_1(x) + f_2(x) + \dots$$

$\uparrow$        $\uparrow$   
DI      PT

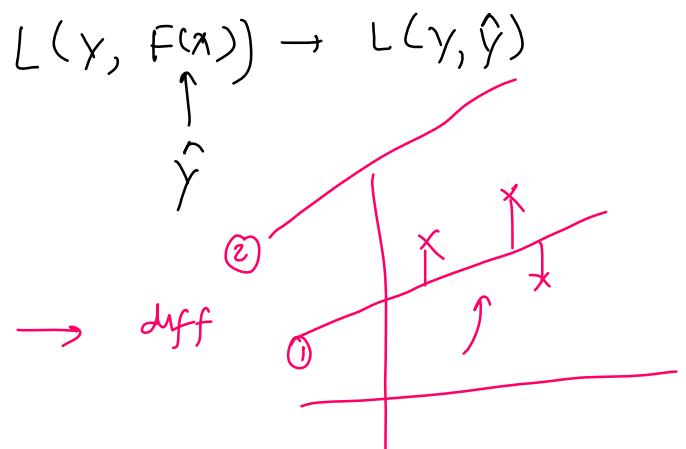
## Explanation

Monday, September 20, 2021 8:25 AM

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

actual      pred

$L = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$



$$y = f(x) \rightsquigarrow$$

$$f(x) = f_0(x) + \underbrace{f_1(x) + f_2(x) + \dots + f_n(x)}$$

$$f_0(x) = \underset{\gamma}{\operatorname{argmin}} \sum_{i=1}^n L(y_i, \underline{\gamma})$$

$$L = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$f_0(x) = \underset{\gamma}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n (y_i - \underline{\gamma})^2$$

$$\frac{d f_0(x)}{d \gamma} = \frac{d}{d \gamma} \frac{1}{2} \sum_{i=1}^n (y_i - \underline{\gamma})^2 = \frac{1}{2} \sum_{i=1}^n \frac{d}{d \gamma} (y_i - \underline{\gamma})^2$$

$$\sum_{i=1}^n (y_i - \underline{\gamma}) \frac{d}{d \gamma} (y_i - \underline{\gamma}) = - \sum_{i=1}^n (y_i - \underline{\gamma}) = 0$$

$$\sum_{i=1}^n (\underline{\gamma} - y_i) = 0$$

$$\sum_{i=1}^3 (\underline{\gamma} - y_i) = 0 \Rightarrow (\underline{\gamma} - 192) + (\underline{\gamma} - 144) + (\underline{\gamma} - 91) = 0$$

$$\sum_{i=1}^3 (y_i - \bar{y}) = 0 \Rightarrow (\bar{y} - 192) + (\bar{y} - 144) + (\bar{y} - 9) = 0$$

$$3\bar{y} = 192 + 144 + 9$$

mean  
 $F_m(x) = f_0(x) + f_1(x) + f_2(x) + \dots + f_m(x)$   
mean of output

$$\bar{y} = \frac{192 + 144 + 9}{3}$$

$$F(x) = \underbrace{f_0(x)}_{\text{mean}} + \underbrace{f_1(x)}_{\text{leaf}} + \underbrace{f_2(x)}_{\text{leaf}} + \dots + \underbrace{f_m(x)}_{\text{leaf}}$$

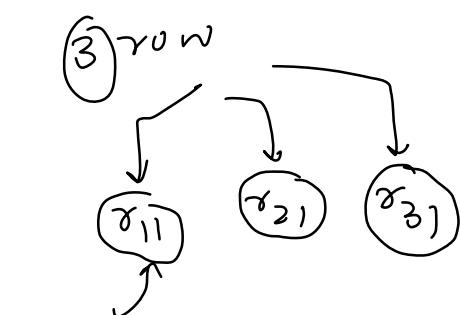
$$m = 1$$

$$\hat{\sigma}_{im} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$$

$$\hat{\sigma}_{ii} = - \left[ \frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_0}$$

$$\hat{y}_i = f(x_i)$$

$$\hat{\sigma}_{ii} = - \left[ \frac{\partial}{\partial \hat{y}_i} L(y_i, \hat{y}_i) \right]_{f=f_0}$$



$$L = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\hat{\sigma}_{ii} = - \left[ \frac{\partial}{\partial \hat{y}_i} \frac{1}{2} (\hat{y}_i - \hat{y}_i)^2 \right]_{f=f_0}$$

$$= \left[ (y_i - \hat{y}_i) \right]_{f=f_0} = \left[ (y_i - f(x_i)) \right]_{f=f_0}$$

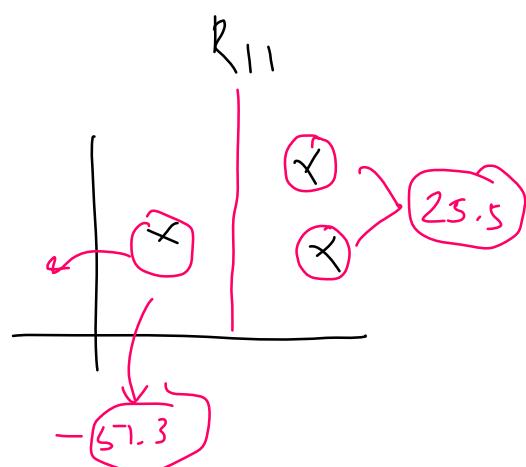
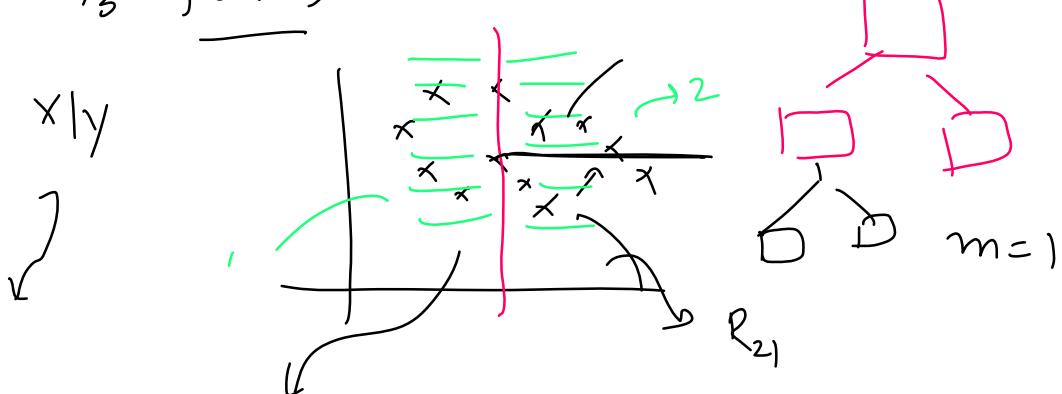
$$= \lfloor (y_i - \hat{y}_i) \rfloor_{f=f_0} = \lfloor (y_i - f(x_i)) \rfloor_{f=f_0}$$

$$\varepsilon_{ij} = \underline{(y_i - f_0(x_i))}$$

$$\varepsilon_{11} = y_1 - \underline{f_0(x_1)} = 192 - 142 =$$

$$\varepsilon_{21} = y_2 - \underline{f_0(x_2)} = 144 - 142 =$$

$$\varepsilon_{31} = y_3 - \underline{f_0(x_3)} = 91 - 142 =$$



$$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma)$$

$$\gamma_{j1} \rightarrow \gamma_{11} = \arg \min_{\gamma} \sum_{x_i \in R_{11}} L(y_i, f_{m-1}(x_i) + \gamma)$$

$$\begin{array}{c}
 \downarrow \quad \downarrow \\
 \gamma_{11} \quad \gamma_{21} \\
 \curvearrowleft \quad \curvearrowright \quad \curvearrowup \quad \curvearrowdown
 \end{array}$$

$$\gamma_{11} = \arg \min_{\gamma} \frac{1}{2} \sum (y_i - (f_0(x_i) + \gamma))^2$$

$$\frac{\partial L}{\partial \gamma} = \frac{1}{2} \times 2(y_i - (f_0(x) + \gamma)) \frac{d}{d\gamma} (\underbrace{y_i - f_0(x)}_{=0} - \gamma) = 0$$

$$= \underbrace{(y_i - f_0(x) - \gamma)}_{=0} = 0$$

$$= y_i - f_0(x) - \gamma = 0$$

$$\gamma_{11} = q_1 - 142 - \gamma = 0$$

$$\boxed{\gamma = q_1 - 142 = -51}$$

$$\begin{aligned}
 \gamma_{21} &= \arg \min_{\gamma} \sum_{x_i \in R_{21}} L(y_i, f_0(x_i) + \gamma) \\
 &= \arg \min_{\gamma} \frac{1}{2} \sum_{i=1}^2 (y_i - (f_0(x_i) + \gamma))^2 \\
 &= - \sum_{i=1}^2 (y_i - f_0(x_i) - \gamma) = 0 \\
 &= \sum_{i=1}^2 (y_i - f_0(x_i) - \gamma) = 0 \\
 &= y_1 - f_0(x_1) - \gamma + y_2 - f_0(x_2) - \gamma = 0 \\
 &= 142 - 142 - \gamma + 144 - 142 - \gamma = 0
 \end{aligned}$$

336  
284

$$336 - 284$$

$$52 - 2\gamma = 0$$

$$\gamma = \frac{52}{2} - 26$$