

# Multiple Linear Regression

x | y      (i) Hypothesis Function

$$\hat{y} = mx + c$$

Error/Loss function

$$L = (Y - \hat{Y})^2$$

$L(m) = (Y - mx)^2 \rightarrow$  Gradient descent  
 Minimize  $\rightarrow$  Random Initialization

$$m = 0.01$$

$$\alpha = 0.1$$

↓  
 Derivative calculation

$$\frac{dL}{dm} = (Y - mx) \cdot 2x$$

$$Bx^T x \theta = (Bx)^T Y$$

$$x\theta = \frac{1}{n} \sum_{i=1}^n y_i$$

[using of known  $\theta$ ]  $\rightarrow$   $Bx^T x \theta = (Bx)^T Y$

With this we do not need  $\theta$   $\rightarrow$   $1.0 = \theta$

and we get  $y$

$$10.0 = \theta_1$$

$$20.0 = \theta_2$$

$$30.0 = \theta_3$$

$$40.0 = \theta_4$$

$x_1$	$x_2$	$x_3$	$y$
1	1	1	1
1	1	1	0
1	1	0	1
1	0	1	0
0	1	1	0

Slope  $\Rightarrow$  weight/  
 $\frac{m}{m}$  Parameter:  $(w)$

## ① Hypothesis Function

$$f = m_1 x_1 + m_2 x_2 + m_3 x_3 + c \quad \downarrow \\ = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$$

Loss function

$$L = (y - \hat{y})^2$$

$$\Rightarrow L(w_1, w_2, w_3, w_0) = (y - w_1 x_1 - w_2 x_2 - w_3 x_3 - w_0)^2 \\ = -2(y - \hat{y})$$

$$f(x) = x \quad | \quad f(x, y) = 3x + 5y^3$$

$$\frac{\partial}{\partial x} f = 3x$$

Random Initialization: [P cannot be zero]

↳ Intercept could be zero.

$$w_1 = 0.01 \\ w_2 = 0.03 \\ w_3 = 0.05 \\ w_0 = 0.06$$

As Intercept on  $c$  doesn't multiply with any other feature.

# Derivative calculation

$$\frac{dL}{d\omega_1} = -2x_1(\hat{y} - \hat{\beta})$$

$$\frac{dL}{d\omega_2} = -2x_2(\hat{y} - \hat{\beta})$$

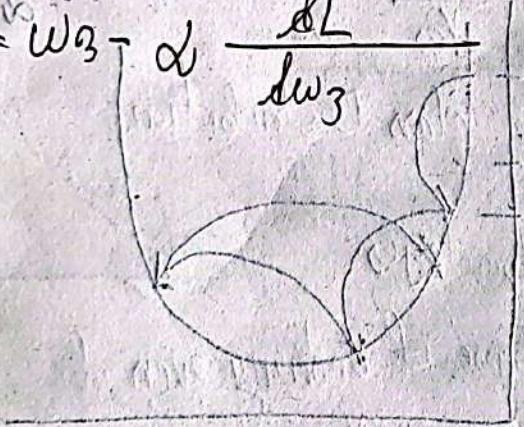
$$\frac{dL}{d\omega_3} = -2x_3(\hat{y} - \hat{\beta})$$

Update

$$\omega_1 = \omega_1 - \alpha \frac{dL}{d\omega_1}$$

$$\omega_2 = \omega_2 - \alpha \frac{dL}{d\omega_2}$$

$$\omega_3 = \omega_3 - \alpha \frac{dL}{d\omega_3}$$



$$\frac{dL}{d\omega_0} = -2(\hat{y} - \hat{\beta})$$

$w_1 = 17.5$

$w_2 = 18.9$

$w_3 = 10.5$

Loop	Error
1	107
2	192
3	122
4	144
5	37

are big

( $y - \hat{y}$ )<sup>2</sup>

are small

( $y - \hat{y}$ )<sup>2</sup>

Learning Rate:

[start from 100]

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 100$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 0.05,$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 0.005$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 0.0005$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 0.00005$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 0.000005$$

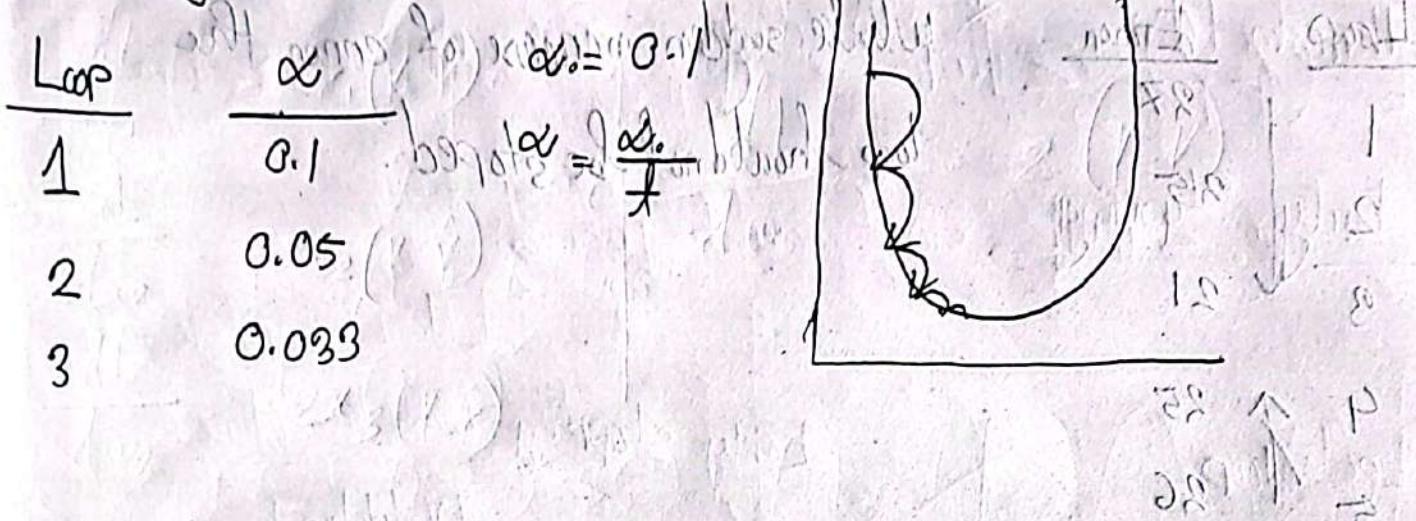
$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 0.0000005$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial w} - \frac{\partial J}{\partial w} = 0.00000005$$

If error sudden changes positive or negative means big  $\alpha$

$\alpha$  sudden change = big  $\alpha$   
 $\alpha$  change small = small  $\alpha$ .

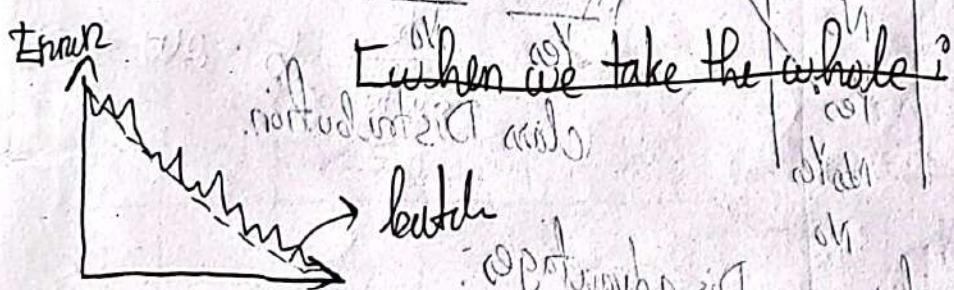
# Learning Rate Scheduling



## Gradient decent types

### ① Stochastic Gradient Descent:

→ If one data is used in every loop. Thus the overall slope is not achieved.



### ② Batch gradient descent when the whole batch is taken

In real life mini-batch is used as the memory is limited

Mini-batch gradient descent is smoother than stochastic rougher than Batch gradient.

Loop	Error
1	27
2	35
3	21
4	25
5	26

while sudden increase of error, the loop should not be stopped.

1.0

20.0

80.0

### \* Logistic Regression

$x_1$	$x_2$	$x_3$	$y$
1	2	3	1
1	2	3	0
1	2	3	1
1	2	3	0

categorical variable

class distribution

class Distribution

categorical to Numerical

Disadvantages:

① Label Encoding : Yes

② Sense of ordering even if true

Step 1: hypothesis function:  $\hat{y} = \omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3$

Step 2: Loss function is,  $L = (Y - \hat{y})^2$

$$\hat{y} = \text{Sigmoid}(\omega_0 + \omega_1 x_1 + \omega_2 x_2 + \omega_3 x_3)$$

\* Threshold = 0.5

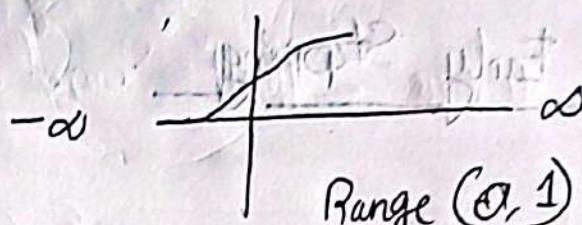
$$\hat{y} = 0.7 \rightarrow 1$$

$$\hat{y} = 0.49 \rightarrow 0.$$

$\leftarrow$  Threshold = 0.

$\rightarrow$  Threshold = 1.

Threshold is a value used to determine



Range (0, 1)

Step 2:

$$L = -y \log \hat{y} - (1-y) \log(1-\hat{y}) \quad \leftarrow \text{Binary log loss function}$$

Suppose,  $Y=0$

$\hat{y}$	0	0.25	0.50	0.75
$L$	0	-0.124	0.30	1
$-\log(1-\hat{y})$	0	-0.124	0.30	1

$\hat{y} = 1$

$\hat{y}$	1	0.75	0.5	0.25
$-\log \hat{y}$	0	0.12	0.30	0.60

As loss function works to find the distance between actual and prediction, Thus it works the same here.

\* ~~sig~~

\* Sigmoid function      \* Binary log-loss function

\* Threshold      \* Encoding

## Early Stopping

(L0) input

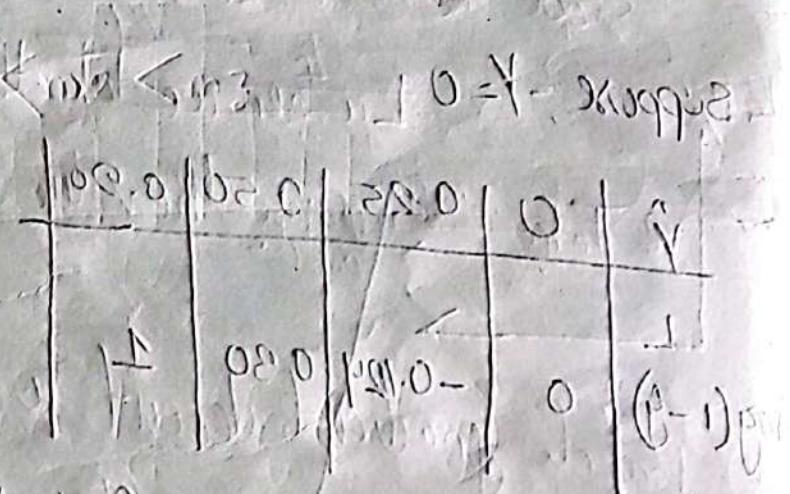
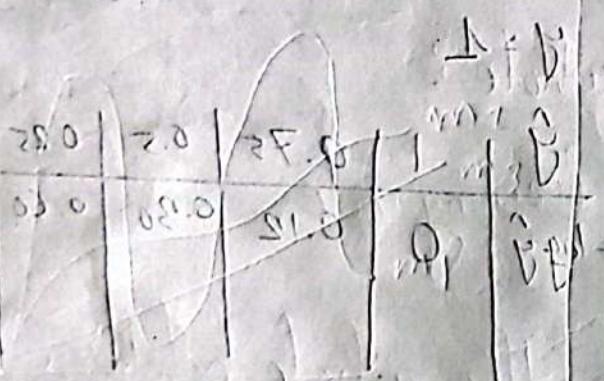
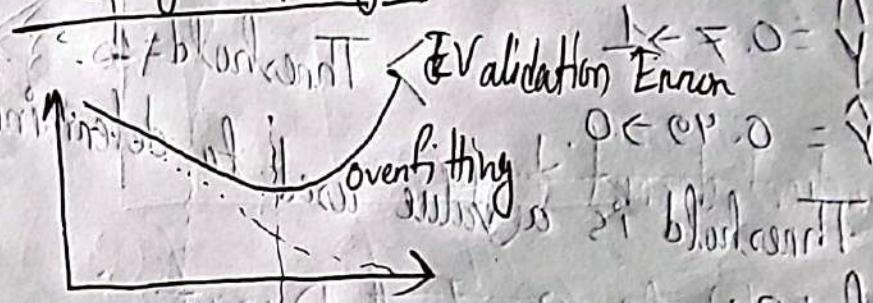
Epoch/loop

$f_{\text{train}} = ?$      $f_{\text{val}} = ?$

$Y_{\text{true}} = ?$      $\hat{Y}_{\text{true}} = ?$

look for parameter  
not trained

## Early stopping

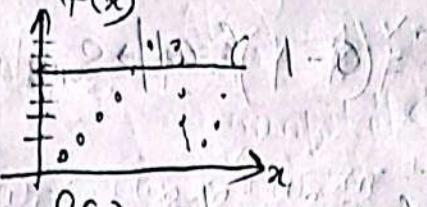


bin labels recorded from lab but bin of histogram not yet  
seen in test data set will not change

## Polynomial Regression

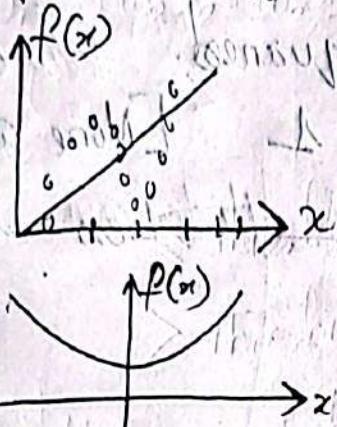
$$① f(x) = 5$$

0 degree polynomial



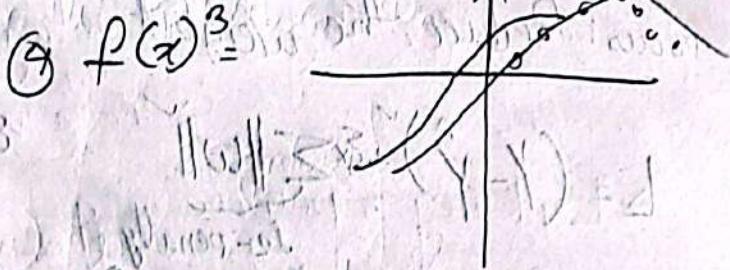
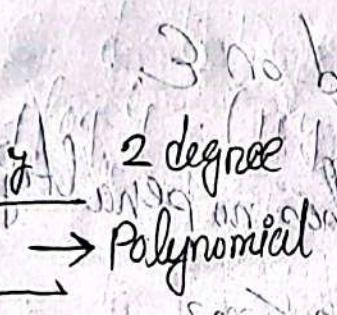
$$② f(x) = 2x$$

One-degree polynomial



$$③ f(x) = x^2$$

Two-degree polynomial



# The more power = more curve.

Degree of an Equation =

Highest power of that variable.

Tangents

$x_1$	$x_2$	$x_3$	$y$

→ 2 degree Polynomial

$x_1$	$x_2$	$x_3$	$y$

$$\text{hypothesis} = \hat{y} = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3$$

$$0^\circ = x^0$$

$$1 \text{ deg} = x$$

$$2 \text{ deg} = x + x^2$$

$$3 \text{ deg} = x + x^2 + x^3$$

$$f(x) = x^{20} + \dots$$

Overfitting

values of weight ↑



Underfitting

values of weight ↓

model complexity ↓

Decrease overfitting

Overfitting ↑

# Training Error = 0

# Test Error = ∞

Focus: Reduce the weight.

$$L = (\hat{Y} - Y)^T + \epsilon \sum \|w\|^2$$

$\ell_2$ -penalty

where,  $\sum \|w\|^2$  = Sum of weight squared.

$$\epsilon = 0 - 1 \quad \epsilon = 0, 1 \rightarrow \text{Some } n \rightarrow \text{Reduce}(n)$$

$$\epsilon = 1$$

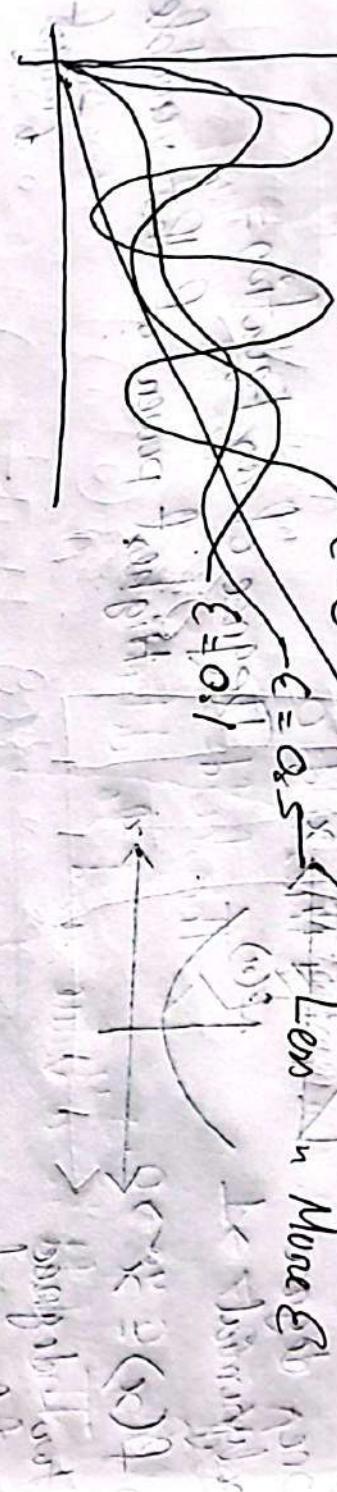
[More complex less E]

$$\epsilon = 0$$

[Less complex more E]

$$\epsilon = 0.1$$

[Less complex less E]

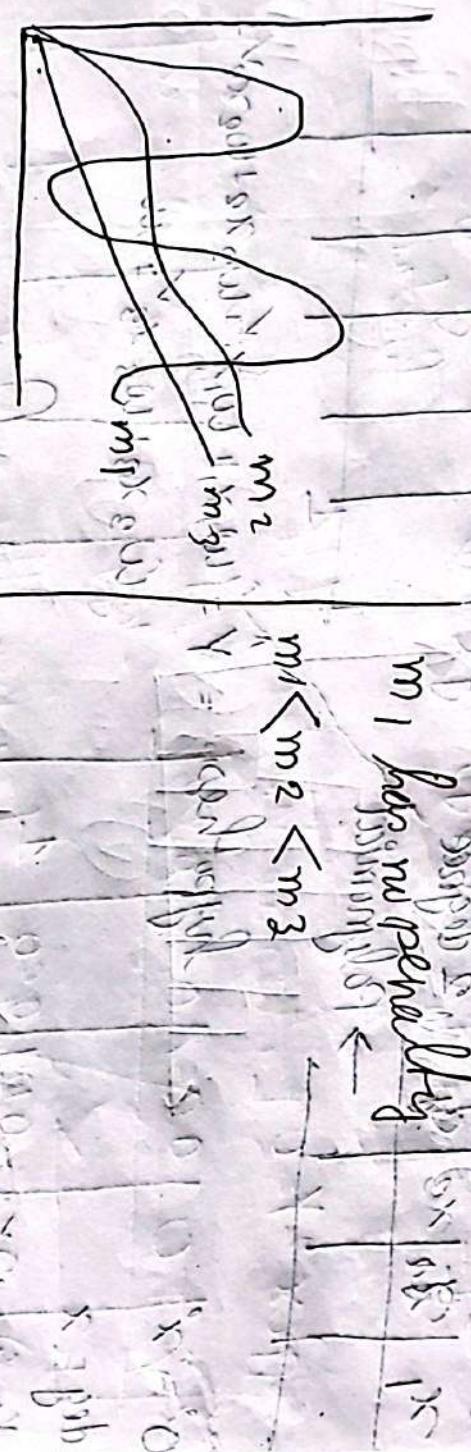


# Q Based on model change based on  $\epsilon$ .

$m_1$  has no penalty

$m_2$

$m_3$



$$L = (\hat{Y} - Y)^T + \epsilon \sum \|w\|$$

\*  $\ell_1$ -penalty  $\rightarrow$  Ridge Regression

Ridge Regression

where  $\sum \|w\|$  = sum of all weights

\*  $\ell_1$ -penalty  $\rightarrow$  Least Absolute Regression

↑ Lasso

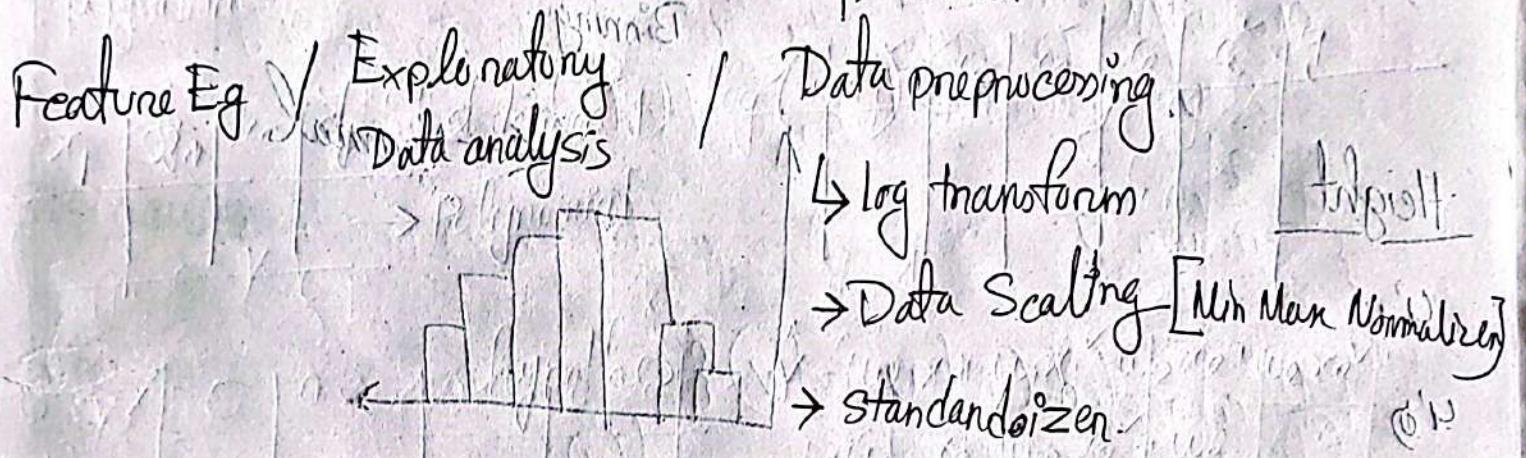
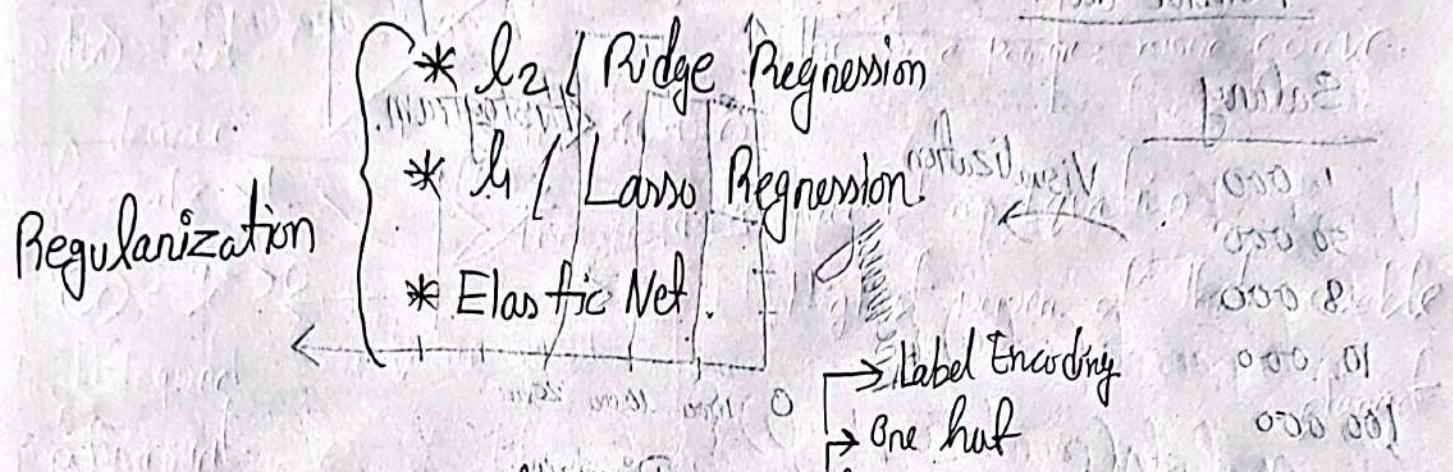
$\epsilon = 0 \rightarrow$  No penalty  $\rightarrow$  complex

$\epsilon = 0, 1 \rightarrow \text{Some } n \rightarrow \text{Reduce}(n)$

$$L = (Y - \hat{Y})^T + \frac{\lambda}{2} \sum_{j=1}^n \|w_j\|^2 + (1-\alpha) \sum_{j=1}^n \|w_j\|$$

Hence,  $\hat{Y}$

Using  $\ell_2$  &  $\ell_1$  together is called Elastic Net. It's principle is



[Data scaling must for gradient descent]

Feature Eng / Feature Composition: Creating new feature from Existing one.

Original	Transformed
1	0.0001
2	0.00001
3	0.000001
4	0.0000001
5	0.00000001

\* Feature Engineering / Exploratory Data Analysis / Data Preprocessing = 1

↳ Data visualization ↳ Histogram ↳ Histogram Transformation

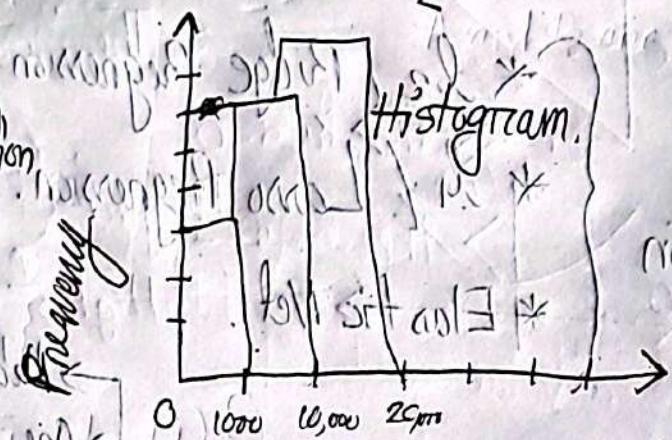
→ Exploring & preprocessing

### Numeric data:

#### Salary

1,000  
50,000  
8,000  
10,000  
100,000

Visualization



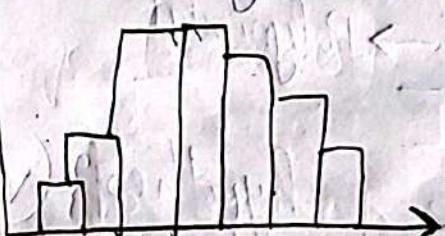
#### Height

4'11  
4'8

Measurement Unit

inches centimeters

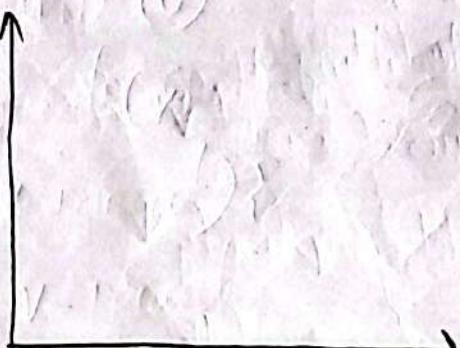
Binning



In case of ~~skewed~~ skewed column use log, log(salary)

#### Salary, $\log(\text{sal})$

Salary	$\log(\text{sal})$
1,000	3
10,000	4
2,000	3.30
1,500	3.17
2,500	3.39
12,000	4.07



$$D_{EU} = \sqrt{(4.11 - 4.0)^2 + (1000 - 50,000)^2}$$

Different Numeric Feature  
Have Different Scale/Range.

$$= \sqrt{0.04 + 49,000}$$

$$= 4.0000$$

Min Max  
Normalizer

$$= \frac{x - x_{\min}}{x_{\max} - x_{\min}} \in [0, 1]$$

	$x_1$	$x_2$
5	1000	
10	2000	
15	3000	
20	4000	

	$x_1$	$x_2$
0	0	0
0.03	0.33	
0.66	0.66	
1	1	1

Standardize

$$x_{\text{new}} = \frac{x - x_{\text{mean}}}{x_{\text{standard dev.}}}$$

Hence center is "0"

Numeration brings close to zero.

Denominator  $\rightarrow$  spread data nicely

(V) training = (U) training + V

(N) training = training + N

3 (training - V = 1 / (training - U))

training - V = [V] training

- \* Disjoint Set / Union Find Algorithm Disjoint set;  $A = \{2, 3, 5\}$
- \* Kruskal's Algorithm For identifying disjoint sets.  $B = \{1, 4, 6\}$

`int parent[6] = {0, 1, 2, 3, 4, 5}`

$$B = \{1, 4, 6\}$$

$$A \cap B = \emptyset$$

$A$  and  $B$  are disjoint sets.

### Union Find

`int parent(int v) {`

`while (parent[v] != v) {`

`v = parent[v];`

`}`

`return v;`

`int connected(int u, int v) {`

`return parent[u] == parent[v];`

`}`

`int connect(int u, int v) {`

`int u-parent = parent[u];`

`int v-parent = parent[v];`

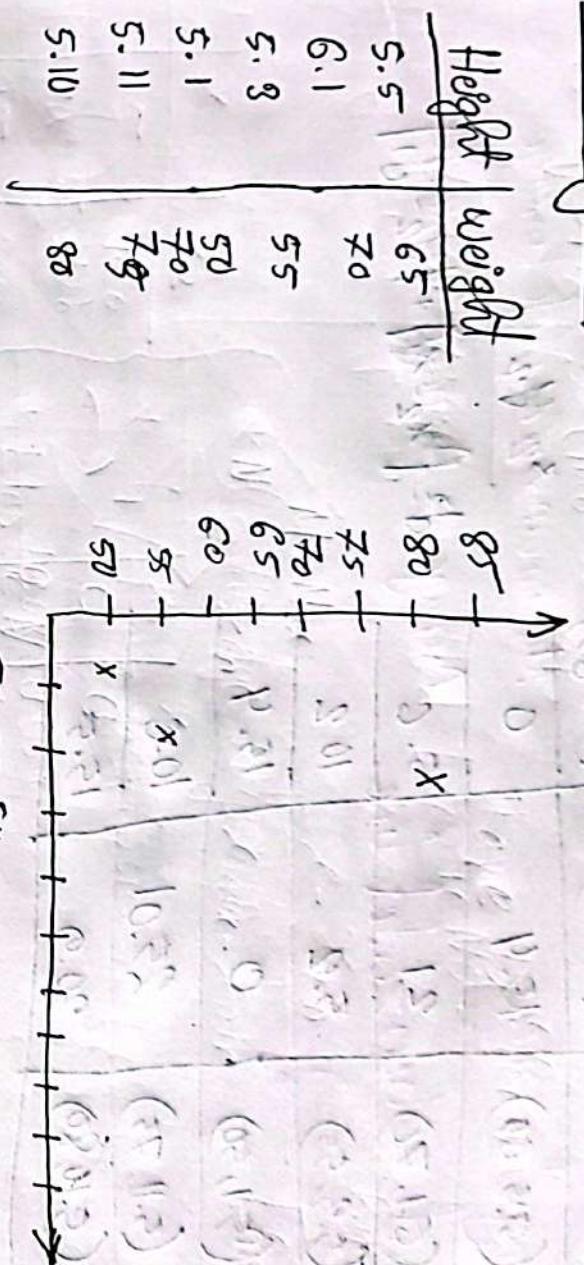
`if (u-parent != v-parent) {`

`parent[v] = u-parent;`

`}`

`}`

## Clustering



k-means clustering [Possible number of groups].

Step 1 : Select the value of "k"

$k = 2$

Step 2 : Initialize k centroids  $\leftarrow$  mean or center point of the clusters.

[start Randomly]

$$C_1 = (5.1, 70)$$

$$C_2 = (5.5, 65)$$

Step 3 : Calculate Distance from each data to each centroid.

(5.5, 65, 85)

(5.1, 70, 90)

(5.1, 75)

(5.1, 80)

(5.1, 85)

(5.1, 90)

(5.1, 95)

(5.1, 100)

(5.1, 105)

(5.1, 110)

(5.1, 115)

(5.5, 70, 85)

(5.5, 75, 85)

(5.5, 80, 85)

(5.5, 85, 85)

(5.5, 90, 85)

(5.5, 95, 85)

(5.5, 100, 85)

(5.5, 105, 85)

(5.5, 110, 85)

(5.5, 115, 85)

		$C_1$ (5.1, 50)	$C_2$ (5.3, 65)	$x_1, y_1$	$x_2, y_2$
1	(5.5, 65)	15.4	0		
2	(6.1, 70)	21	5.6		
3	(5.3, 55)	5.2	10.2		
4	(5.1, 50)	0	15.4		
5	(5.11, 75)	25.01	10.6		
6	(5.10, 80)	30.9	15.5		

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Step 4: Assign each Data to its closest centroid

1	$C_2$
2	$C_2$
3	$C_1$
4	$C_1$
5	$C_2$
6	$C_2$

$C_1$  updated  $x$  position = Average of previous positions

Step 5: Updated centroid

$$C_1 = \frac{5.1 + 5.3}{2}, \frac{50 + 55}{2} = (5.2, 52.5)$$

$$= (5.2, 52.5)$$

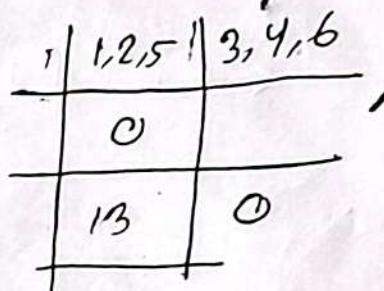
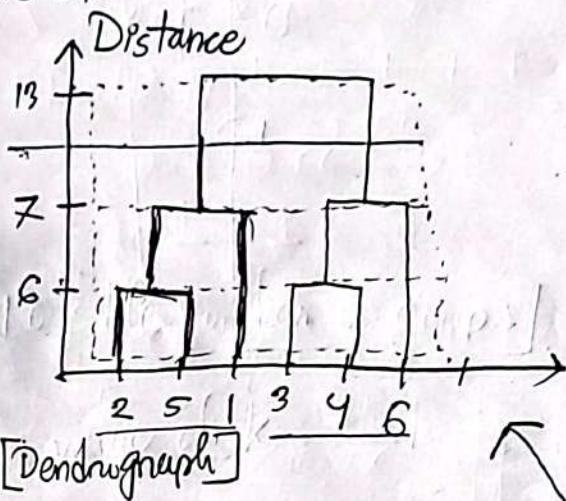
$$C_2 = \frac{5.5 + 6.1 + 5.10 + 5.11}{4}, \quad \frac{6.5 + 7.0 + 7.5 + 8}{4}$$

$$= (5.15, 72.5)$$

1, 2, 3, 4, 5, 6  
1, 2, 3, 4, 5, 6

\* Hierarchical / Agglomerative / clustering.  
Initially, each data is an individual cluster.

	$x_1$	$x_2$
1	5	2
2	10	4
3	25	8
4	30	0
5	15	5
6	35	11



Distance Matrix

	1	2	3	4	5	6
1	0	X				
2	X	0				
3	26	10	0			
4	32	25	15	0		
5	13	6	13	10	0	26
6	30	32	13	X	26	0

$$\phi = |x_2 - x_1| + |y_2 - y_1|$$

=

	1	(2,5)	3	4	6
1	0				
(2,5)	X	0			
3	26	13	0		
4	32	10	6	0	
5	30	26	13	X	0
6					

	(1,2,5)	3,4	6
(1,2,5)	0		
(3,4)	13	0	
6	26	X	0

	25	34	8	6
(2,5)	0			
(3,4)	X	0		
8	26	13	0	

# Dendrogram

↳ Rectangular area

↳ but no horizontal line.

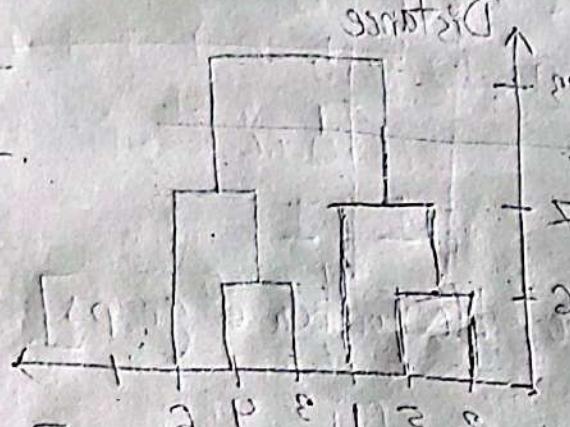
Minimum  $\Rightarrow$  Single Linkage

Maximum  $\Rightarrow$  Complete Linkage

Average  $\Rightarrow$  Average Linkage

Minimum

0	1	2	3	4	5
0	1	2	3	4	5



6	7	8	9	10	11	12	13	14	15
6	7	8	9	10	11	12	13	14	15

minimum. L

maximum. L

Average L

$$d(x, y) = |x - y|$$

$$d(x, y) = \sqrt{(x - y)^2}$$

$$d(x, y) = \sqrt{(x - y)^2}$$

0	1	2	3	4	5
0	1	2	3	4	5

0	1	2	3	4	5
0	1	2	3	4	5
0	1	2	3	4	5
0	1	2	3	4	5

6	7	8	9	10	11	12	13	14	15
6	7	8	9	10	11	12	13	14	15

0	1	2	3	4	5
0	1	2	3	4	5