

# Multiple Linear Regression

x	y

(i) Hypothesis Function

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

Error/Loss function

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

$$L(m) = (Y - mx)^2 \rightarrow \text{Gradient descent}$$

Minimize  $\rightarrow$  Random Initialization

$$m = 0.01$$

$$\alpha = 0.1$$

Derivative calculation

$$\frac{\partial L}{\partial m} = \sum_{i=1}^n (Y_i - mx_i)^2$$

$$Bx = (Bx)^T$$

$$x = \frac{1}{n} \sum x_i$$

Cost of training  $\rightarrow$  cost of training model

With  $m$  fixed  $\rightarrow$  no gradient  $\rightarrow$  A

and  $m$  into  $B$  it's

$$10.0 = 0.1$$

$$20.0 = 0.1$$

$$30.0 = 0.1$$

$$40.0 = 0.1$$

$x_1$	$x_2$	$x_3$	$y$
1	1	1	1
2	2	2	2
3	3	3	3
4	4	4	4

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

## ① Hypothesis Function

$$f = m_1 x_1 + m_2 x_2 + m_3 x_3 + c \quad (Y - \hat{Y}) = 1$$

$$= w_1 x_1 + w_2 x_2 + w_3 x_3 + w_0$$

## 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.

$$w_1 = 0.01$$

$$w_2 = 0.03$$

$$w_3 = 0.05$$

$$w_0 = 0.06$$

$$\alpha = 0.1$$

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

# Derivative calculation

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

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

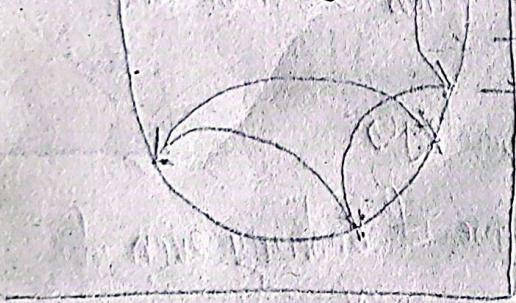
$$\frac{dL}{d\omega_3} = -2x_3(\hat{y} - y)$$

Update

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

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

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

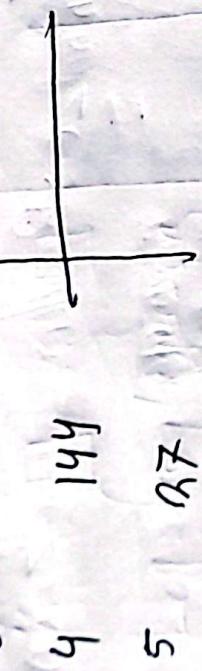


$$w_1 = 17.5$$

$$w_2 = 18.9$$

$$w_3 = 10.5$$

Loop	Error
1	107
2	102
3	122
4	144
5	137



Learning Rate:

[start from 100]

Loop	E1	E2	E3
1	50	75	98
2	8	50	97
3	48	75	96
4	52	6	95
5	12	58	94

Learning Rate: [start from 100]

$$\frac{\Delta E}{\Delta w} = \frac{E_1 - E_5}{w_1 - w_5} = \frac{107 - 137}{17.5 - 10.5} = -10 = -100$$

$$\Delta w = \alpha \cdot \Delta E = 0.05 \cdot (-100) = -5$$

$$w_2 = 17.5 - 5 = 12.5$$

$$\Delta w = \alpha \cdot \Delta E = 0.05 \cdot (-100) = -5$$

$$w_3 = 12.5 - 5 = 7.5$$

$$\Delta w = \alpha \cdot \Delta E = 0.05 \cdot (-100) = -5$$

$$w_4 = 7.5 - 5 = 2.5$$

$$\Delta w = \alpha \cdot \Delta E = 0.05 \cdot (-100) = -5$$

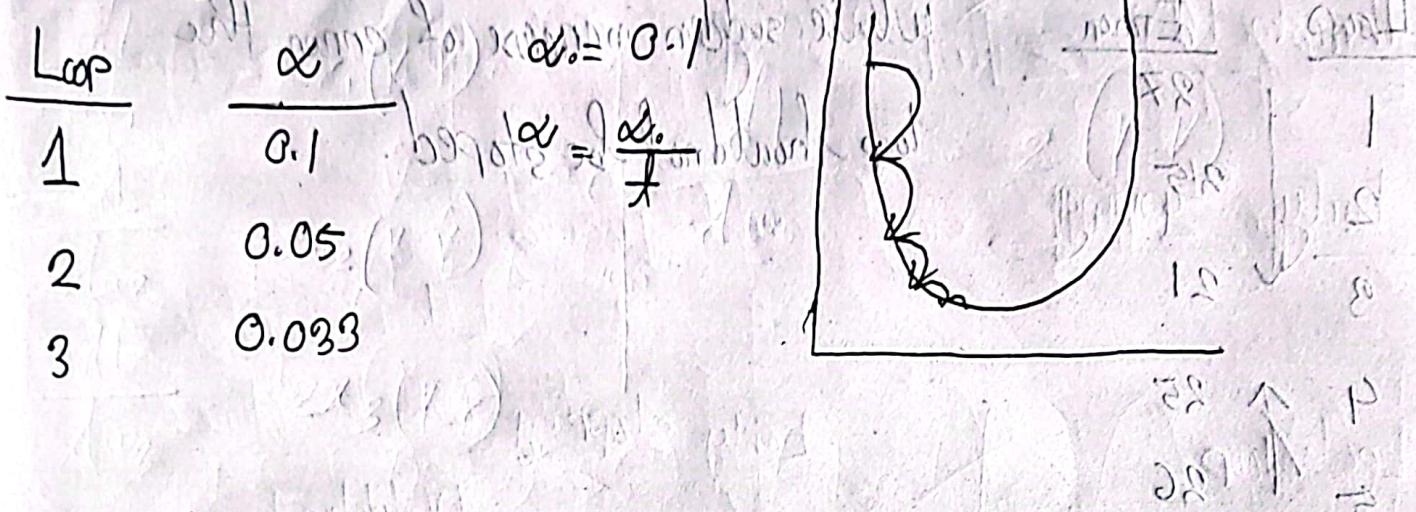
$$w_5 = 2.5 - 5 = -2.5$$

If error suddenly changes positive or negative means big  $\Delta w$

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

But if the error is small  
then update  $\Delta w$  will be small.

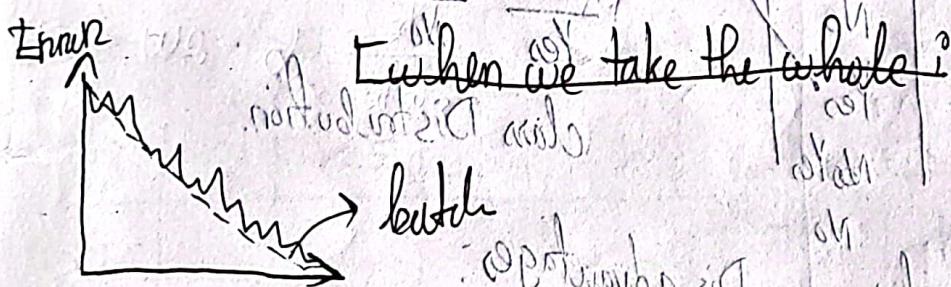
# Learning Rate Scheduling



## Gradient descent 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

taken in a loop and after several loops several avg 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

(1) 30.0

80.0

Early stopping (2.8) 20.8

Validation Loss

Softplus loss function

Categorical variable

$x_1$	$x_2$	$x_3$	$y$
			Yes
			No
			Yes
			No

Loss function for classification

$\frac{1}{n} \sum_{i=1}^n \log(p_i)$

class distribution

Categorical to Numerical

Disadvantages:

① Label Encoding: Yes > 1

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

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

$\leftarrow$  Threshold = 0.

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

$\rightarrow$  Threshold = 1.

Threshold is a value used to determine

Step 2:

$$L = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

Binary log loss function

Suppose,  $Y=0$

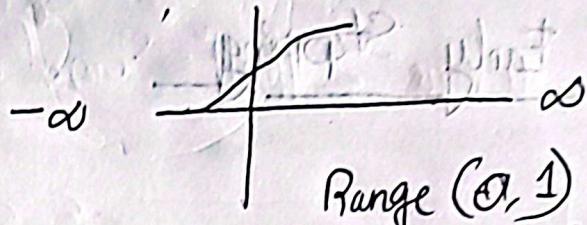
$\hat{y}=1$

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

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

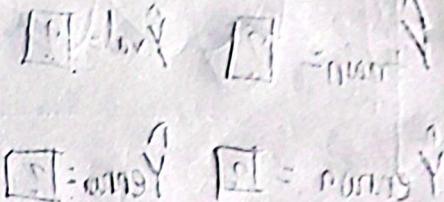
As loss function works to find the distance between actual and prediction, Thus it works the same here.

\* ~~sigmoid~~



Range (0, 1)

a decision boundary



1 = 0  
0 = 1

1 = 1  
0 = 0

- \* Sigmoid function
- \* Binary log-loss function
- \* Threshold
- \* Encoding

## Early Stopping

(E.O) epoch

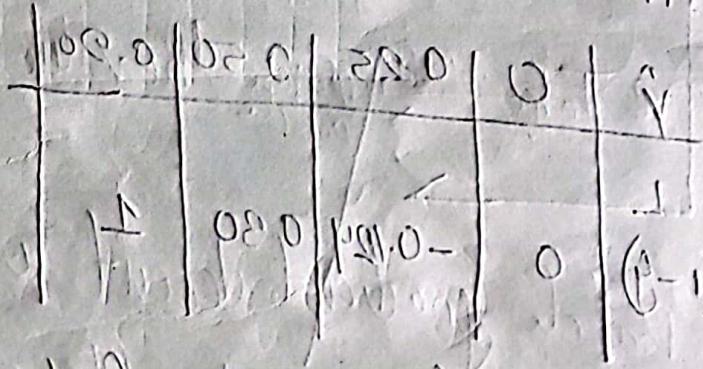
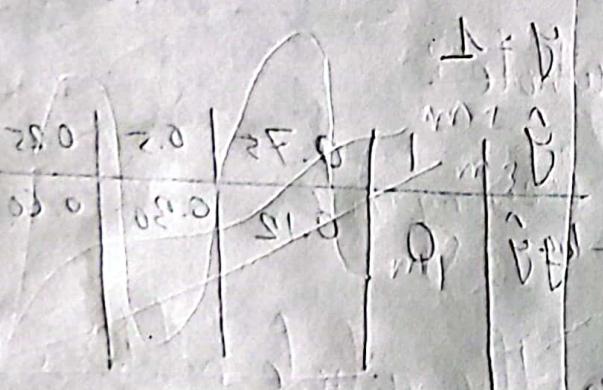
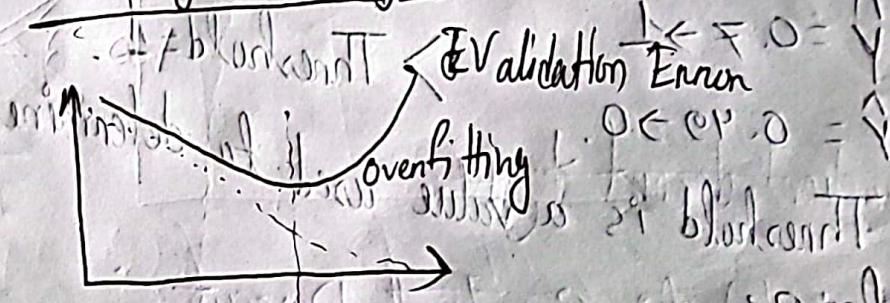
Epoch/loop

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

$$Y_{\text{train}} = ? \quad Y_{\text{val}} = ?$$

look for pattern  
not much

## Early stopping



brain tumor's record is not full but of known record not full  
 1. (V1, V2)  
 2. (V1, V2, V3)  
 3. (V1, V2, V3, V4)  
 4. (V1, V2, V3, V4, V5)

## Polynomial Regression

$$\textcircled{1} \quad f(x) = 5$$

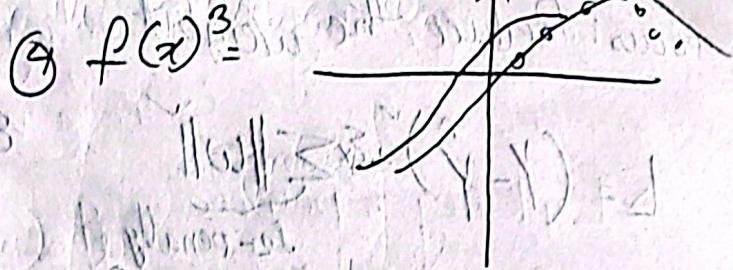
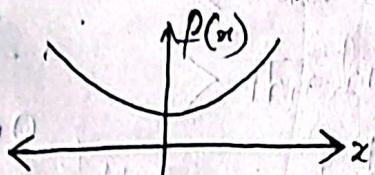
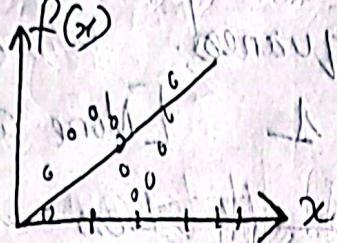
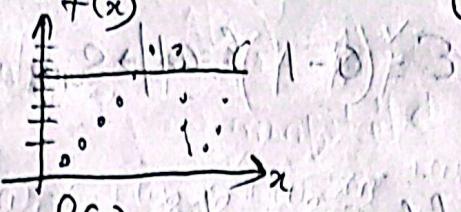
0 degree  
polynomial

$$\textcircled{2} \quad f(x) = 2x$$

One-degree  
polynomial

$$\textcircled{3} \quad f(x) = x^2$$

Two-degree  
polynomial



# The more power = more curve.

Degree of an Equation =

Highest power of that variable.

Tangent

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

2 degree  
Polynomial

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$y$

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

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

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

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

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

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

Overfitting

Underfitting

values of weight ↑

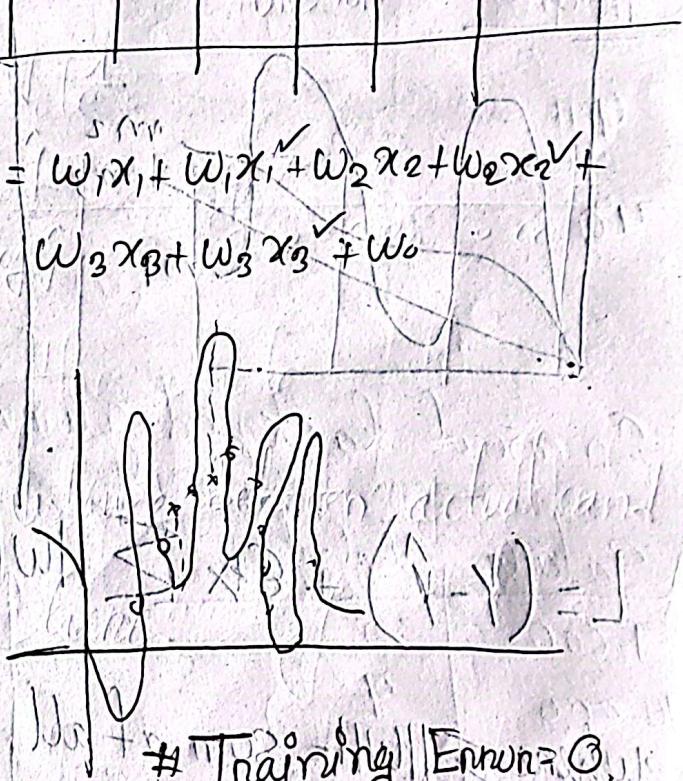
values of weight ↓

↓  
model complexity ↑

↓  
model complexity ↓

Overfitting ↑

Decrease overfitting



# Training Error = 0.

# Test Error = ∞.

Focus: Reduce the weight.

$$L = (\hat{Y} - Y)^T + \underbrace{\epsilon \sum \|w\|^2}_{L_2\text{-penalty}}$$

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

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

$$\epsilon = 1$$

$$\epsilon = 0$$

$$\epsilon = 0.1$$

$$\epsilon = 0.5$$

$$\epsilon = 1$$

$$\epsilon = 10$$

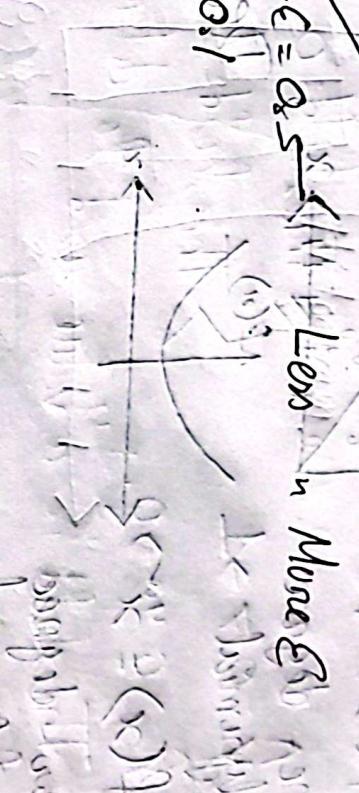
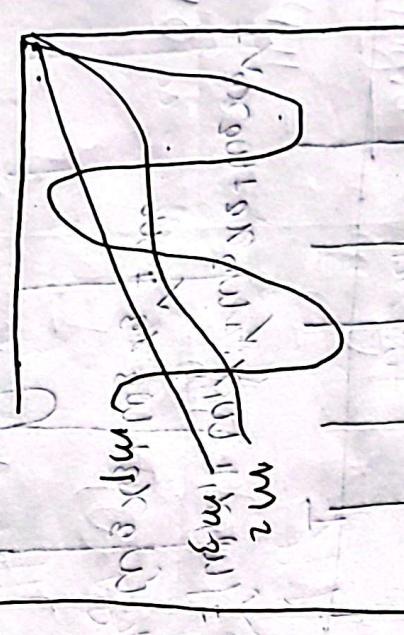
$\epsilon = 0 \rightarrow \text{No penalty} \rightarrow \text{complex}$

$\epsilon = 10 \rightarrow \text{Some } n \rightarrow \text{Reduce}(n)$

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

$$m_1 < m_2 < m_3$$

Lasso  
penalty



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

\*  $\ell_2$ -penalty  $\rightarrow$  Ridge Regression

Ridge Regression

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

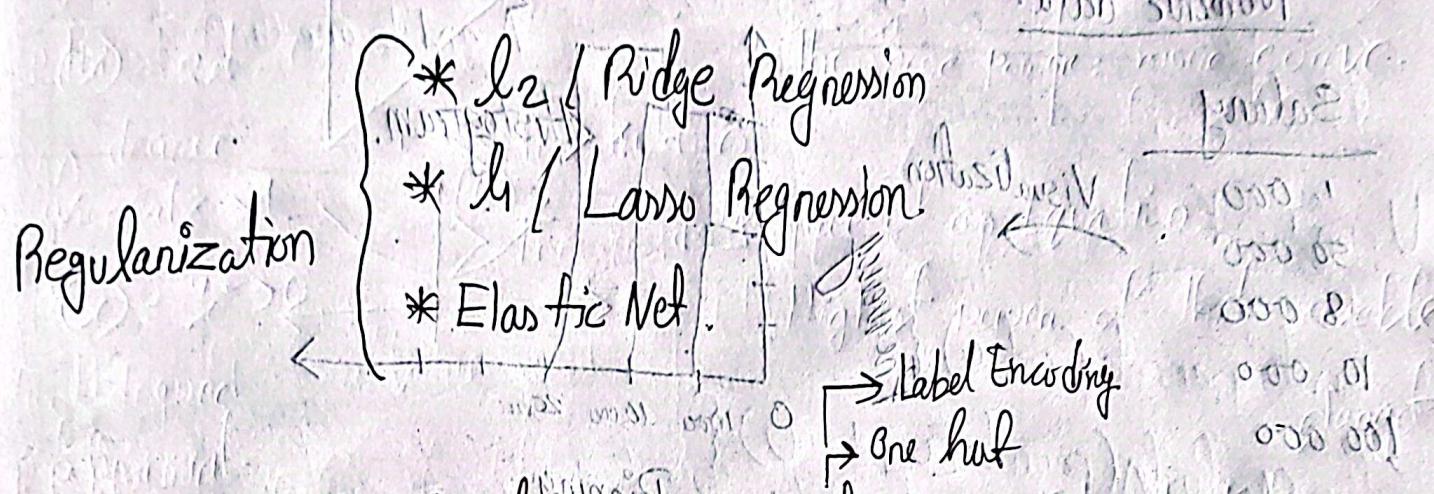
\*  $\ell_1$ -penalty  $\rightarrow$  Lasso Regression

Lasso Regression

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

Hence,  $\hat{Y}$  is given by

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



[Data scaling must for gradient descent]

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

Original	values	new values	new values	new values
Linear				
Polynomial				
Product				
Quotient				
Log				
Exponential				
Power				
Sum				
Mean				
Median				
Mode				
Product				
Quotient				
Log				
Exponential				
Power				
Sum				
Mean				
Median				
Mode				

\* Feature Engineering / Exploratory Data Analysis / Data Preprocessing

↳ Data visualization

↳ Histogram

↳ Histogram Log Transformation

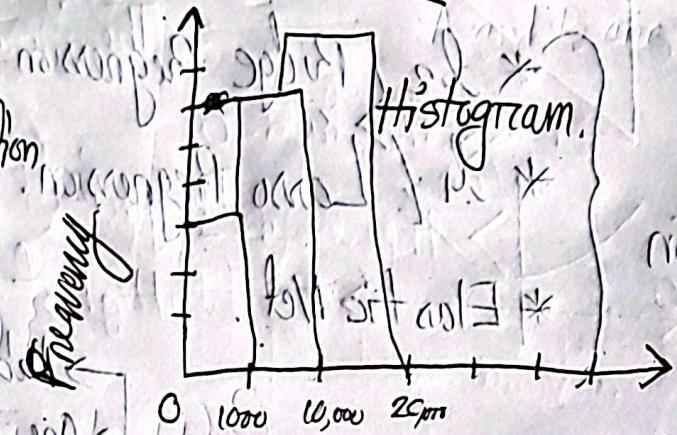
→ Exploring & preprocessing

Numeric data:

Salary

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

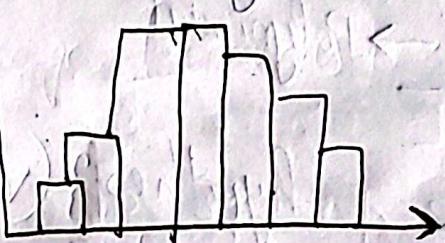
Visualization



Height

4'11  
4'8

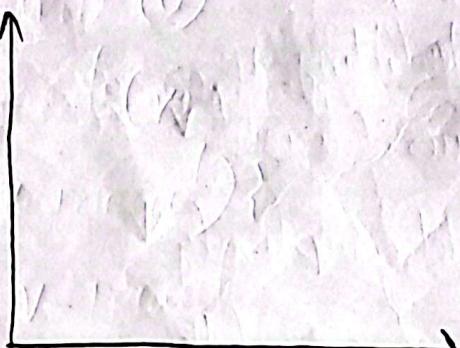
Binning



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

Salary | log(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.

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

$$= 4.$$

Min Max  
Normalizer

$$= \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Range

$$[0, 1]$$

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

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

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

int parent[6] = 

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

$$A \cap B = \emptyset$$

$$A \cup B = \{1, 2, 3, 4, 5, 6\}$$

$A$  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];
```

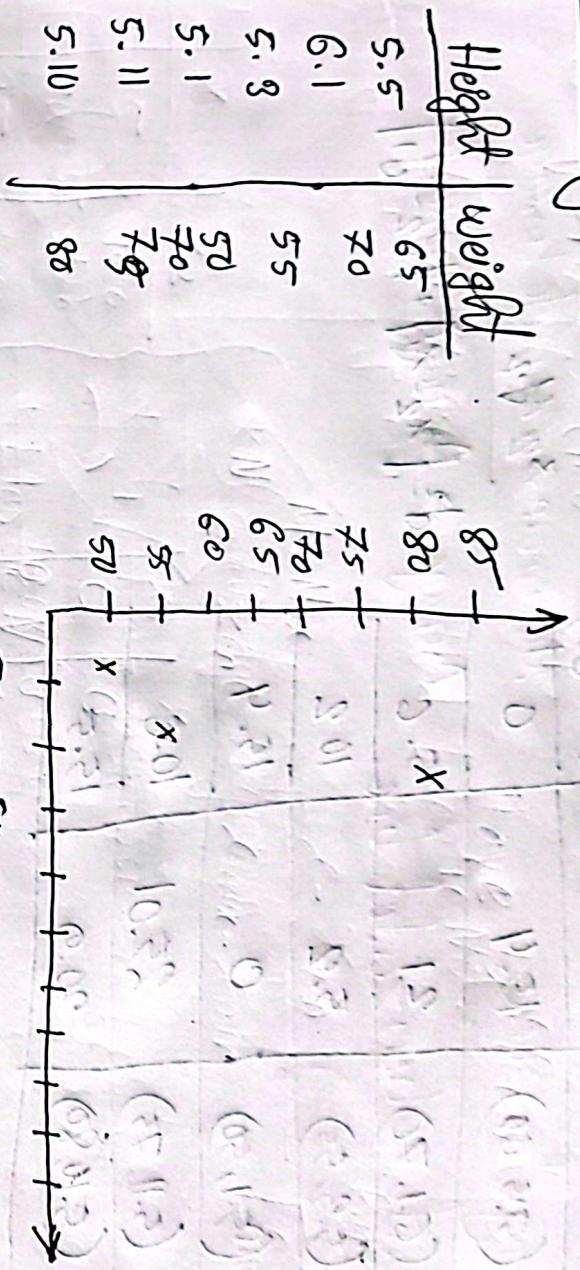
```
    if (u-parent != v-parent) {
```

```
        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]

$$C_1 = (5.1, 50)$$

$$C_2 = (5.5, 65)$$

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

	$C_1$ (5.1, 50)	$C_2$ (5.3, 65)	$x_1, y_1$
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 =  $\frac{5.1 + 5.3}{2} = 5.2$

Step 5: Updated centroid

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

$$= (5.2, 52.5)$$

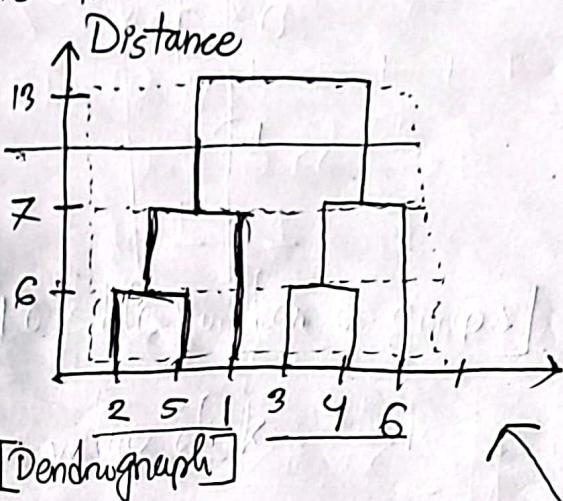
$$C_2 = \frac{5.5 + 6.1 + 5.10 + 5.11}{4}, \quad \frac{6.5 + 7.0 + 7.5 + 8}{4}$$

$$= (5.5, 72.5)$$

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



1	1, 2, 5	3, 4, 6
13	○	○
13		

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

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					

	25	34	8	6
(2,5)	0			
(3,4)	X	0		
13	26	13	0	
6	30	26	X	0

- y1. Minimum  
y2. Maximum  
y3. Average.

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

# Dendrogram

↳ Rectangular area

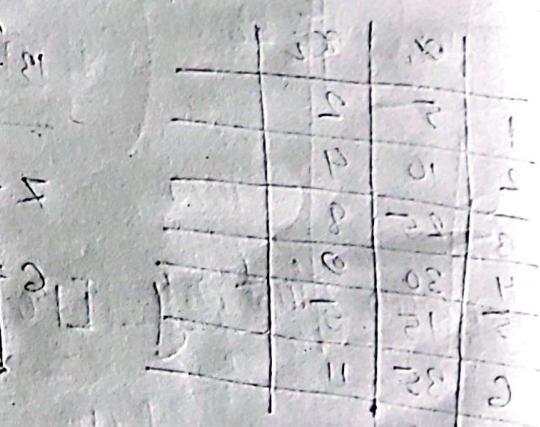
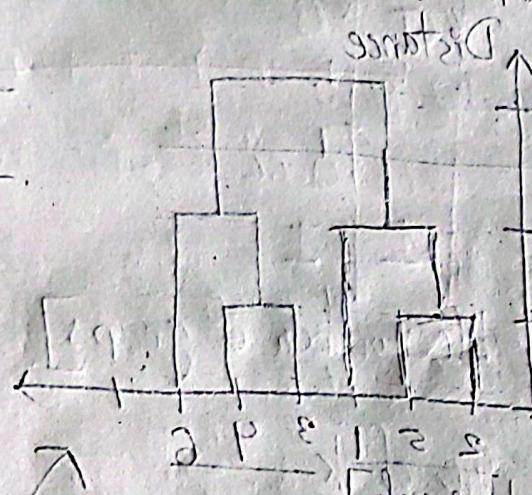
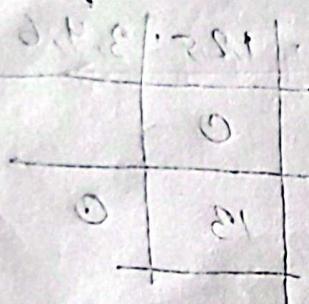
↳ but no horizontal line.

Minimum  $\Rightarrow$  Single Linkage

Maximum  $\Rightarrow$  Complete Linkage

Average  $\Rightarrow$  Average Linkage

Minimum



minimum. 1

minimum. 2

average. 3

$$\text{minimum. 1} = |1 - 2| + |2 - 3| = 2$$

