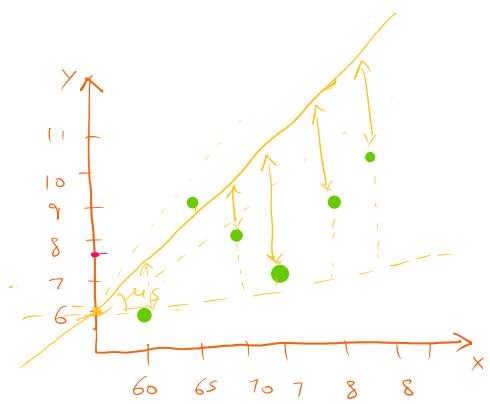


$\alpha \theta_1$	i	m	b	J
	1	45	1	100
	2	50	1	120
	3	40	1	80
	4	30	1	40
	5	20	1	10
	6	0	1	0
	7	1	2	2
	8	5	2	10
	9	8	2	18
	10	12	2	26



Assumption:

- Linearity
- Independence of errors
- Homoscedasticity
- No correlation b/w features

Limitation:

- Linearity assumption
- Outliers
- overfitting

use cases:

L1 & L2 Regularization:

$$\text{Loss} = \left[\frac{1}{2} \sum (y - \hat{y})^2 \right] + \lambda \sum |m_i| + \lambda \sum m_i^2$$

(1 - 0.005)

Decoding rate Penalty

L1 Error - Lasso Regression

L2 Error - Ridge Regression

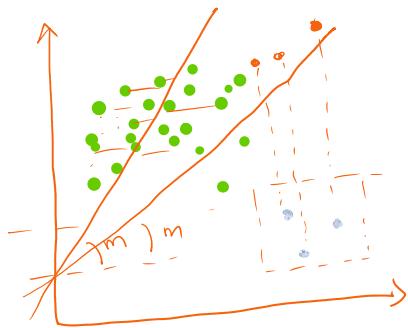


when to apply:

- Linear data
- simple & interpretable model
- moderate to high data

Trade off

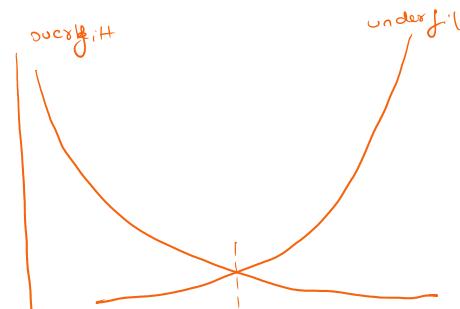
$$y = \underbrace{m_1 x_1}_{\text{SS}} + \underbrace{m_2 x_2}_{\text{SS}} + \dots + m_n x_n + m_0 - \text{Lasso}$$

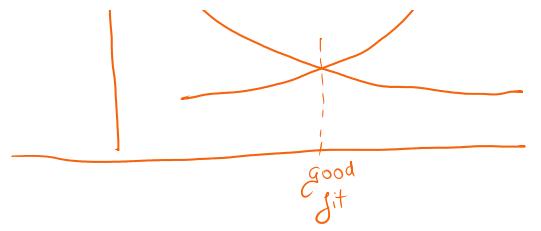


→ train error - \downarrow overfitting

→ test error - \uparrow underfitting

↑ bias & variance ↑
overfitting
underfitting





Polynomial Regression

$$\hat{y} = m_0 + m_1 x_1 + m_2 x_2^2 + m_3 x_3^3 + \dots + m_n x_n^n$$

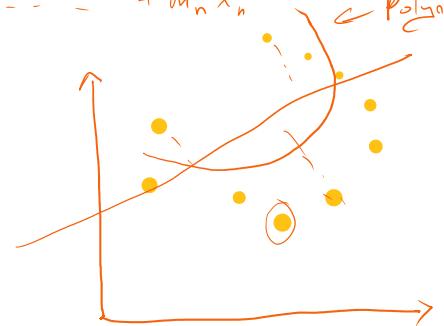
When to use

- Non linear

Assumptions :

- Linearity

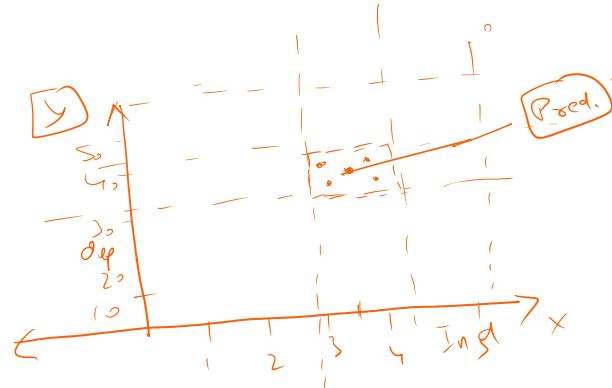
← Independence



- Decision Tree →

- SVM

- KNN



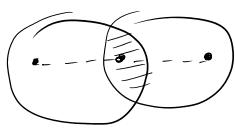
$$n = 3.5$$

$$\hat{y} = 40$$

Unsupervised learning. - (Prescriptive)

Sunday, September 11, 2024 8:06 PM

o K-Means



o Types of K-Means :-

- K-Mean ++ :-

- Mini Batch K Mean :-

- Fuzzy K Mean :-

Assumption's

- K- number of cluster
- Cluster Shape
- Homogeneous density

Limitations :-

- Sensitive to initial centroids
- Number cluster
- Sensitive to outliers
- Cluster shape assumptions

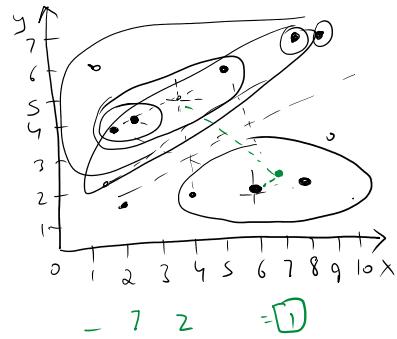
Use Cases:-

- Customer segmentation
- Document clustering
- Image compression

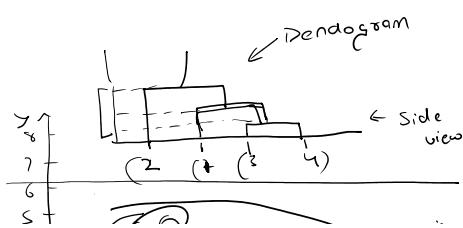
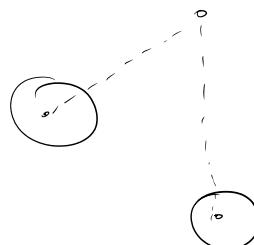
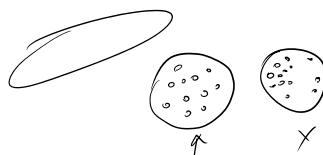
When to Apply:-

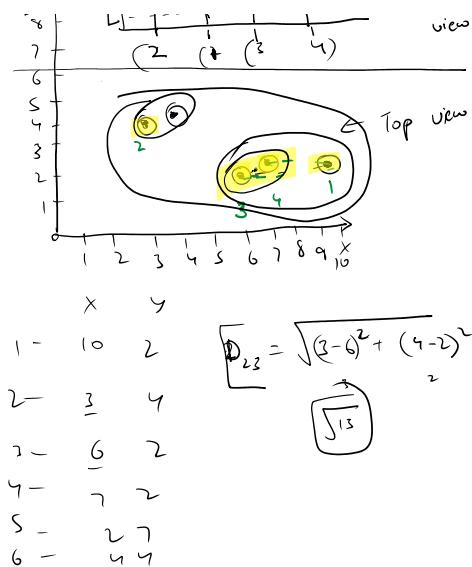
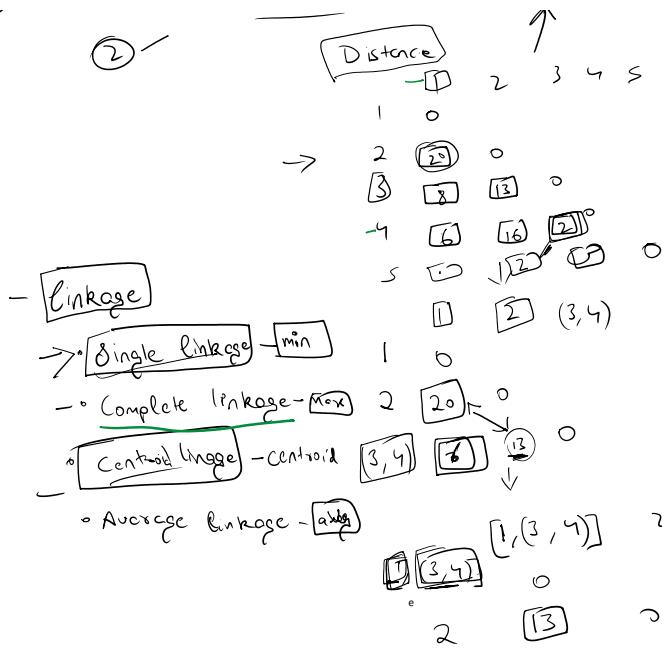
- Normal to big data
- estimation of no. of clusters
- (1), (2,3), (3,4)
(2,5)
(3,4)

Agglomerative Clustering :-



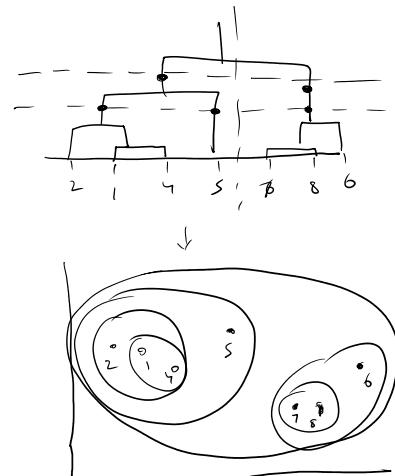
	x_i	y_i	
→ →	10	2	- 1 - No
→ →	3	4	- 2 - No
→	6	3	- 1
→	2	4	- 2
→	5	6	- 2
-	4	2	- 1





Linkage

- Single: grouping of plant
- Complete:
- Average
- Centroid/ward

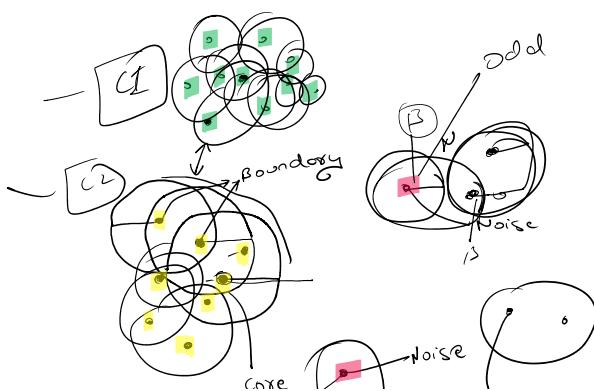


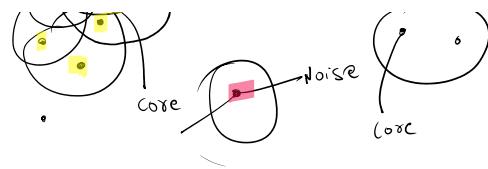
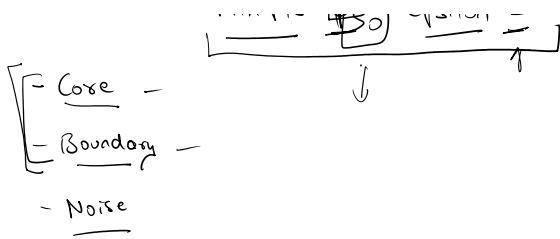
Assumptions!

- Distance / linkage choice
- Scalable
- Use cases
- When to apply
 - Explore hierarchical st. & granularity level.
 - No estimation
 - Intuitive to visualize.
- DB Scan - (Anomaly detection)

Core -

$$\text{Min Pts} = 3, \text{Epsilon} = 1$$





Limitations:

- Parameter Sensitive
- High Dimension data / very large data

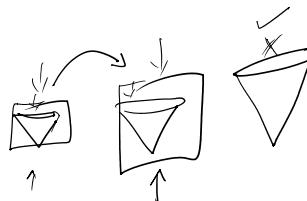
Use Cosine

- Anomaly

Association techniques

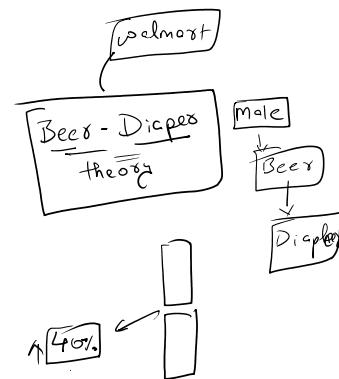
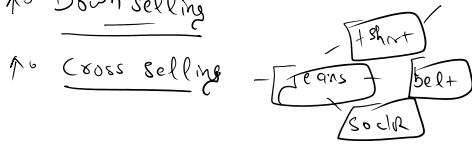
Market Basket Analysis

- FP-growth



- ↑◦ Upselling

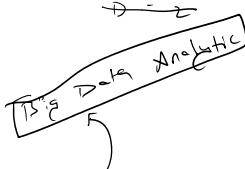
- ↑◦ Down selling



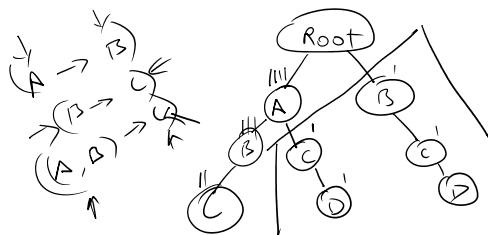
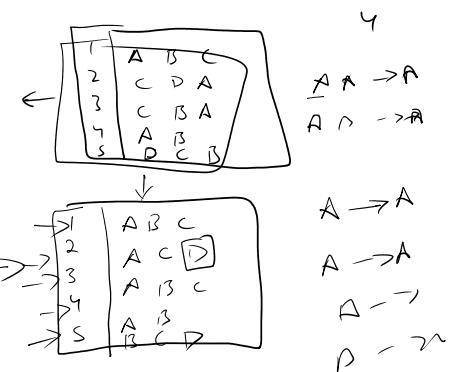
FP-growth

- threshold value = 2

$$\boxed{\text{Support}} = 2$$



Rice
Flour
Pulse
A - 4
B - 4
C - 4
D - 2



Assumption:

- structure pre-defined tabular data
- assume the existence of FP.

Use Cosine

Limitations

- Memory usage
- Handling large dataset
- Handling numeric data

Use Case

- Web usage mining.
- Biology & genomic
- Network Intrusion Detect.

When to apply:-

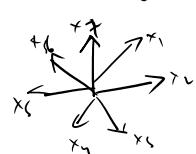
- FP
- Processed Data
- Threshold

→ ◦ Apriori \rightarrow

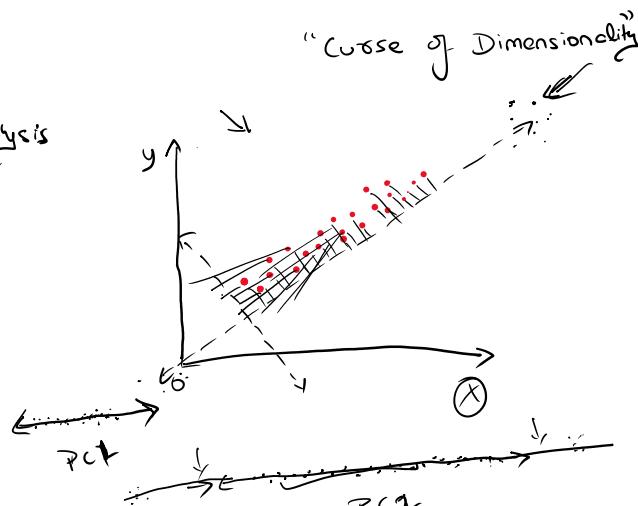
→ ◦ Eclat \rightarrow

Dimension Reduction

◦ PCA :- Principle Component Analysis



No. of PC \leq No. of features



	X	Y
F1	2.5	2.4
	0.5	0.9
	2.2	2.9
	1.9	2.2
	3.1	3.0
	2.3	2.7
	1	1.6
	1.5	1.6
	1.1	0.9

Step 1: Calc. the mean of each column:

$$\bar{x} = 1.81 \quad \bar{y} = 1.91$$

Step 2 - Calc. the covariance matrix

$$\begin{matrix} x & \bar{x} & x - \bar{x} & (x - \bar{x})(x - \bar{x}) \\ -2.5 & 1.81 & -1.41 & 0.476 \\ -0.5 & 1.81 & -1.31 & 1.7161 \\ 2.5 & 2.4 & 0.69 & 0.49 \\ 2.2 & 2.2 & 0.39 & 0.1521 \\ 1.9 & 1.9 & 0.09 & 0.0081 \\ 3.1 & 3.0 & 0.11 & 0.0121 \\ 2.3 & 2.7 & 0.41 & 0.1681 \\ 1 & 1.6 & -0.61 & 0.3721 \\ 1.5 & 1.6 & -0.11 & 0.0121 \\ 1.1 & 0.9 & -0.81 & 0.6561 \end{matrix}$$

$$\begin{matrix} x & y & x - \bar{x} & y - \bar{y} & (x - \bar{x})(y - \bar{y}) \\ 2.5 & 2.4 & 0.69 & 0.49 & 0.3381 \\ 0.5 & 0.9 & -1.41 & -1.01 & 1.4581 \\ 2.2 & 2.9 & 0.39 & 1.09 & 0.4181 \\ 1.9 & 2.2 & 0.09 & 0.31 & 0.0281 \\ 3.1 & 3.0 & 0.11 & 0.09 & 0.0111 \\ 2.3 & 2.7 & 0.41 & 0.09 & 0.0361 \\ 1 & 1.6 & -0.61 & -0.31 & 0.1921 \\ 1.5 & 1.6 & -0.11 & -0.31 & 0.0361 \\ 1.1 & 0.9 & -0.81 & -0.61 & 0.4921 \end{matrix}$$

$$C_{2 \times 2} = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{cov}(y, y) \end{bmatrix}$$

$$\text{cov}(x, x) = \frac{\sum_1^n (x - \bar{x})(x - \bar{x})}{n-1}$$

$$\text{cov}(y, y) = \frac{\sum_1^n (y - \bar{y})(y - \bar{y})}{n-1}$$

$$C_{2 \times 2} = \begin{bmatrix} 6.165 & 0.6154 \\ 0.6154 & 0.7165 \end{bmatrix}$$

$$\frac{5.538}{9}$$

Step 3: Eigen values & vector.

$$C - \lambda I = 0$$

$$\begin{bmatrix} -6165 & -6154 \\ -6154 & -7165 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = 0$$

$$\begin{bmatrix} -6165 - \lambda & -6154 \\ -6154 & -7165 - \lambda \end{bmatrix} = 0$$

$$\lambda^2 - 1.333\lambda + 0.0630 = 0$$

$$\lambda_1 = \boxed{0.490}$$

$$\lambda_2 = \boxed{1.2840}$$

$$CD = \lambda D$$

$$\begin{bmatrix} -6165 & -6154 \\ -6154 & -7165 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0.49 \begin{bmatrix} x \\ y \end{bmatrix}$$

$$-6165x_1 + -6154y_1 = 0.49x_1$$

$$6154x_1 + -7165y_1 = -0.49y_1$$

$$-6165x_1 - 0.49x_1 = -6154y_1$$

$$-5675x_1 = -6154y_1$$

$$x_1 = -1.0845y_1$$

$$\text{assuming } y_1 = 1$$

$$x_1 = -1.0845 = \frac{-1.0845}{1.7684 + 1}$$

$$= \frac{-1.0845}{\sqrt{2176.14}} = \boxed{1.475}$$

$$-1.0845 = \boxed{1.475}$$

$$\frac{PC1}{e_1} \rightarrow \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} -0.7351 \\ 0.6778 \end{bmatrix}$$

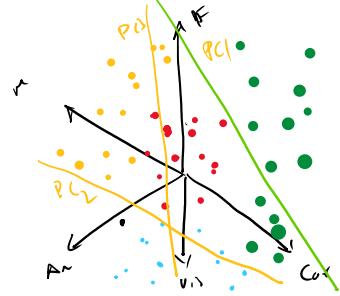
$$PC2 \rightarrow \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 0.6778 \\ 0.7351 \end{bmatrix}$$

$$PC3 \rightarrow \begin{bmatrix} x_1 \\ y_1 \end{bmatrix} = \begin{bmatrix} 1 \\ 1.475 \end{bmatrix}$$

Carbohydrates Fats Amino Acid Mineral Vitamins

Sugar
Glucose
Protein

Protein



Step 5

$$e_1^\top \begin{bmatrix} x_i & \bar{x} \\ y_i & \bar{y} \end{bmatrix}$$

$$[-0.7351 \quad 0.6778] \begin{bmatrix} 2.5 & 1.81 \\ 2.9 & 1.9 \end{bmatrix}$$

$$-0.7351 [-0.9] + 0.6778 [0.5]$$

$$= \boxed{0.166}$$

X	y	PC1
2.5	2.4	-0.166
0.5	0.9	
2.2	2.9	
1.9	2.2	
3.1	3.0	
2.3	2.7	
2	1.6	
1	1.1	
1.5	1.6	
1.1	0.9	

Assumptions

- Linearity
- mean is zero
- Uncorrelated PCA components

use cases

- Dimension Reduction
- Avoid overfitting
- Visualizing high dimension features

Limitations:

- Linearity

Limitation :-

- Linearity
 - Loss of interpretability
 - Sensitive to scaling

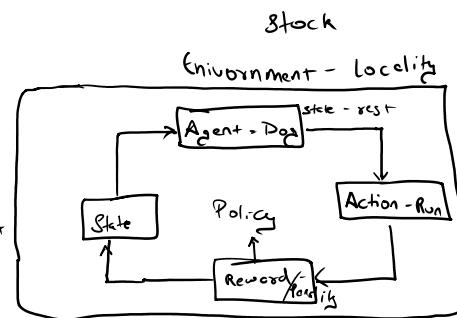
Dimension Reduction

- Kernel PCA
 - t-SNE
 - UMAP

→ Reinforcement Learning:

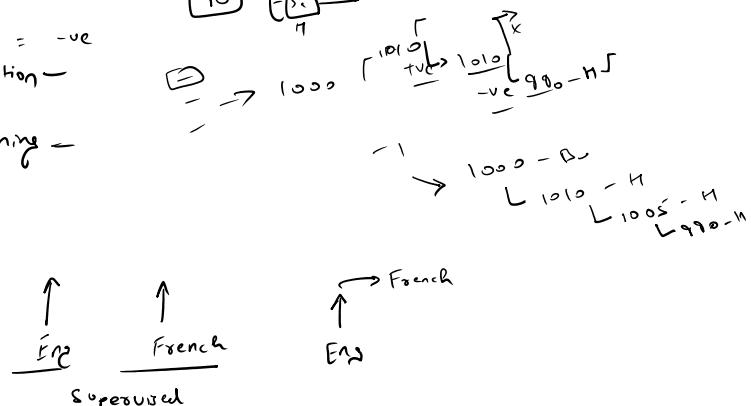
- Cement Learning:

 - Environment → Exchanged
 - Agent + → Bot
 - State → Current Position - π^*
 - Action → Buy, Sell, hold
 - Policy →
 - Types
 - Buy
 - Sell
 - Reward → Buy-Sell
 - Learning b
 - Rew
 - Pen
 - Goal orient



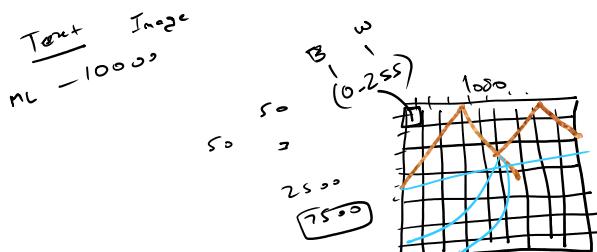
- Semi-Supervised :- language translation

- Transfer Learning



Deep learning :-

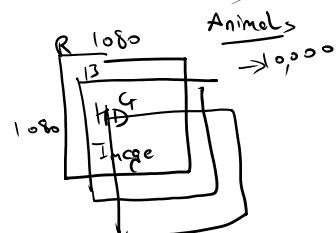
Friday, September 10, 2024 9:29 PM



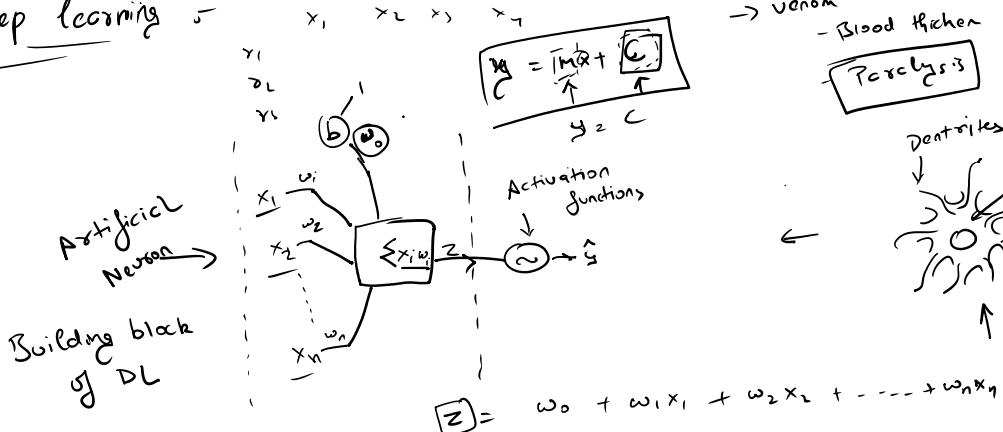
Machine learning :-

Limitation :-

- High dimensionality data :-
- Feature Selection → **Automotives**
- Image, videos, text



Deep learning :-



- x = Input

- w = weights

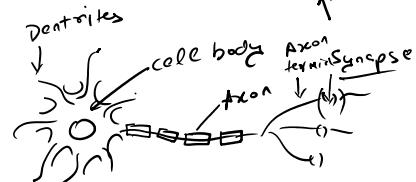
- b/w_0 = bias

- \circ = Activation function

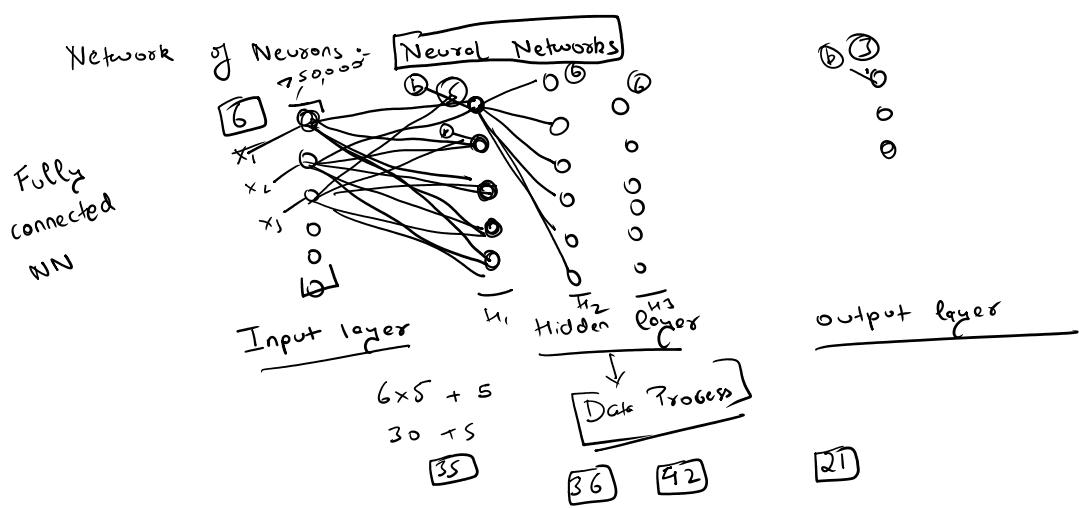
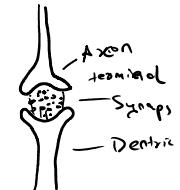
$$\text{Sigmoid} = \frac{1}{1+e^{-z}}$$

- ReLU

- tanh



Biological Neuron



Dog
Cat
Mouse

18

NN < 2 HL → Shallow Neural Network

18 NN < 2 HL \rightarrow Shallow Neural Network

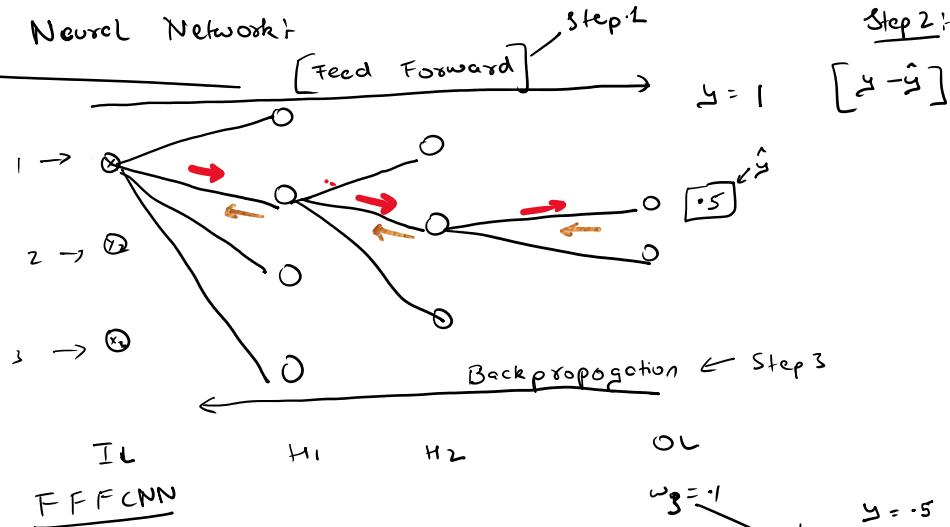
NN > 2 HL \rightarrow Deep Neural Network

1 Neuron \rightarrow Perceptron

How to train the Neural Network?

FCNN
 \downarrow

Feed Forward
Fully connected
Neural Networks



Mathematics:

Step 1:

$$z_i = \sum w_i x_i$$

$$\hat{y}_i = \frac{1}{1 + e^{-z_i}}$$

$$z_3 = (0.35 \times -1) + (0.9 \times 0.8) + b = 0.755$$

$$\hat{y}_3 = \frac{1}{1 + e^{-0.755}} = 0.68$$

$$z_4 = (0.35 \times 0.4) + (0.9 \times 0.6) + b = 0.68$$

$$\hat{y}_4 = \frac{1}{1 + e^{-0.68}} = 0.66$$

Step 3:

$$\Delta w_{ji} = \eta \cdot \delta_j \cdot o_i$$

$$\rightarrow \delta_j = O_j(1-O_j)(\epsilon_j - O_j) \quad \{j = \text{output}\}$$

$$\delta_j = O_j(1-O_j) \sum_k \delta_k w_{kj}$$

For output

$$\delta_5 = \hat{y}(1-\hat{y})(y - \hat{y}) \\ = 0.68(1-0.68)(0.5 - 0.68) = -0.18$$

$$\eta = \text{learning rate} = 1 \quad [0.1 - 0.001]$$

δ_j = error areas

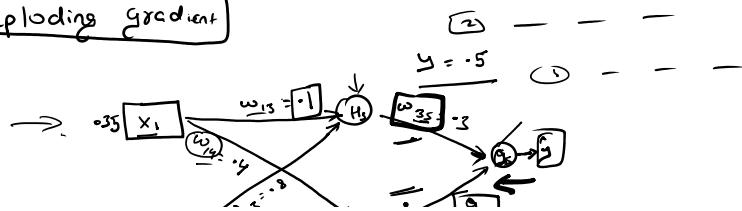
O_j = Output of neuron

ϵ_j = correct output

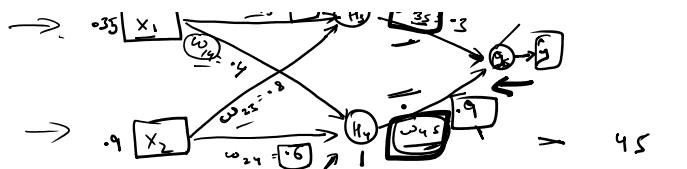
- Yes

- I - - - -

Exploding gradient



$$\delta_5 = \hat{y} (1 - \hat{y}) (y - \hat{y}) \\ = 0.687 (1 - 0.687) (0.5 - 0.687) = -0.178$$



$$\delta_3 = y_3 (1 - y_3) (\omega_{33} \times \delta_5) \\ = 0.68 (1 - 0.68) (-0.25 \times -0.178) = -0.0265$$

$$(\underbrace{-0.000001}_{\Delta \omega_{13}}) \boxed{\Delta \omega_{13}} = \underbrace{\eta \delta_3}_{-0.599999999} \underbrace{\omega_{13}}_{=0.12} = \frac{1 \times -0.178 \times -0.25}{-0.12} = \boxed{-0.11}$$

$$\boxed{\delta_4} = y_4 (1 - y_4) (\omega_{43} \times \delta_3) \\ = 0.66 (1 - 0.66) (-0.25 \times -0.0265) = -0.166$$

$$\Delta \omega_{14} = \eta \delta_4 x_1 \\ = 1 \times -0.166 \times 0.35$$

$$\omega_{43} (\text{new}) = \Delta \omega_{43} + \omega_{43} (\text{old}) \\ = -0.11 + 0.9$$

$$\omega_{14} (\text{new}) = \Delta \omega_{14} + \omega_{14} (\text{old})$$

2nd Iteration

Step 1 →

$$\text{Feed Forward} = 0.65$$

Step 2 →

$$\text{Error} = 0.5 - 0.65 = \boxed{-0.15}$$

Step 3 →

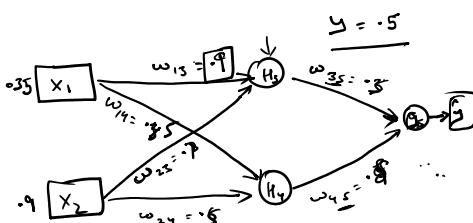
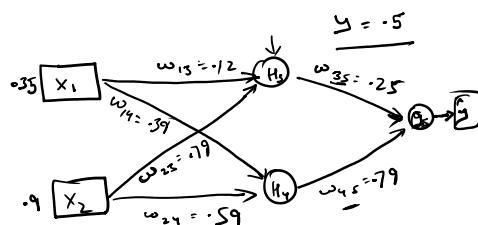
Back propagation →

3rd I

• Feed F

$$\bullet \text{Error} \rightarrow \boxed{0.01}$$

Back -



→ Types of Activation Functions

sigmoid

$$\rightarrow \underline{\text{Sigmoid}}: \quad \frac{1}{1+e^{-z}}$$

$$\text{Range: } [0, 1]$$

Use Cases:

- binary classification

Pros:

- Prob based funct., easy implement

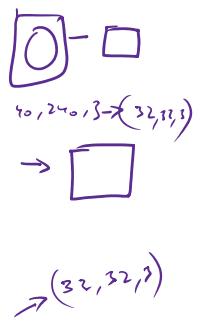
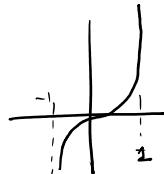
Cons:

- Vanishing Gradient

2) Hyperbolic tan : $\tanh(x) =$

$$\frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Range $[-1, 1]$



Use:

- hidden layers

Pros:

- zero centered output help optimization
- less saturation

Cons:

- vanishing Gradient

3) Relu = $\max(0, x)$ Rectified Linear Unit

Range: $[0, \infty]$

Use cases

- hidden layers

Pros:

- Computationally inexpensive
- vanishing gradient won't work

0	0
0	0
0	0
0	0
0	0

Cons:

- Dying Relu Problem

0, 25%

4) Leaky Relu

$$\alpha = (0.1) \quad \begin{cases} x & \text{if } x > 0 \\ \alpha x & \text{if } x < 0 \end{cases}$$

Range: $(-\infty, \infty)$

Regression / general

Classification

$$\sigma_{\text{softmax}} = \frac{x}{1 + e^{-x}}$$

- $\delta_{\text{sigmoid}} = \frac{x}{1 + e^{-x}}$

Range = $(-\infty, \infty)$

Use Cases

deep network

Pros: Avoid vanishing gradient

- Smooth & non-monotonic

Cons:

computationally expensive

- Softmax: $\frac{e^{x_i}}{\sum e^{x_i}}$ Range $(0, 1)$

Use Cases

multi-class classification

I close

$$\begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} - \begin{bmatrix} .1 \\ .2 \\ .7 \end{bmatrix} = \underline{1}$$

Pros: Converts logits to prob. distribution

Cons:

outlier data

- classes are mutually exclusive

Optimizer

- Gradient Descent (GD): $\Theta_{\text{new}} = \Theta_{\text{old}} - \eta \nabla \delta \Theta$

Use Cases

small datasets



Pros:

simple & effective on convex problems

Cons:

slow convergence & computationally expensive for large dataset.

- SGD: $\Theta = \Theta - \eta \nabla \delta(\Theta, x_i, y_i)$ (x, y) - individual samples

Use Cases

big datasets

Pros:

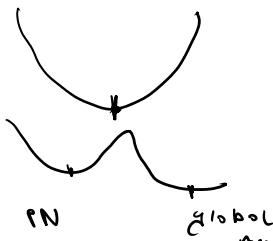
- big dataset

Pros:

- Fast processing

Cons:

Noisy updates, may not reach global minima.



o Mini batch GD = $\theta_{\text{new}} = \theta_{\text{old}} - \eta \nabla \delta(\theta, \text{batch})$

use cases:

large datasets, more stable than SGD

Pros:

- Reducing noise & variance SGD
- Faster than GD

Cons:

Very large dataset

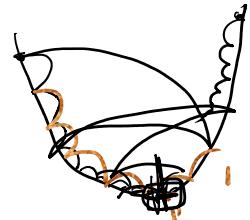
$$\begin{aligned} \theta &= \theta - \eta \nabla \delta(\theta) \\ &= \theta - \eta \cdot \underbrace{\left[\frac{1}{N} \sum_{i=1}^N \nabla \delta(\theta_i) \right]}_{\text{momentum} = \gamma} \end{aligned}$$

γ = momentum factor = .9

o Momentum Optimizer

use cases:

use in training deep network



Pros:

- Speed up the convergence
- Reducing oscillation
- Big data also

Cons:

tunable parameter γ momentum.

Adaptive optimizer

Adagrad

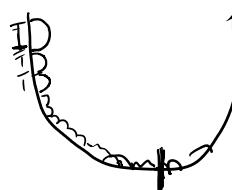
$$\theta = \theta - \frac{\eta}{\sqrt{G_t + \epsilon}} \nabla \delta(\theta)$$

G_t = sum of sq. of past gradients

ϵ = constant small value to

θ avoid zero division.

$$G_t = \sum_{i=1}^t \|\nabla \delta(\theta_i)\|^2$$



use cases:

sparse data, NLP, image

Pros:

Adaptive learning rate

Cons:

Adaptive Learning rate

Cons: learning rate keeps on decay

Wants

RMS Prop :- Root mean Sq. propagation

$$= \underline{G_t} = \beta G_t + (1-\beta) (\nabla \delta \theta)^2$$

β = decay factor = .9

$$\boxed{\theta = \theta - \frac{\eta}{\sqrt{G_t + \epsilon}} \cdot \nabla \delta(\theta)}$$

Use:-

NLP,

Pros:

Addressing adagrad cons.

Cons:

β is sensitive

Adom: Adaptive momentum estimation

- First momentum estimation:

$$m_t = \beta m_t + (1-\beta) \nabla J(\theta) \quad \beta = .9$$

- second momentum estimation

$$v_t = \beta v_t + (1-\beta) \nabla \delta(\theta)$$

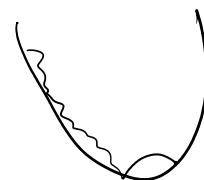
Bias correction

$$\hat{m}_t = \frac{m_t}{1-\beta}, \quad \hat{v}_t = \frac{v_t}{1-\beta}$$

$$\text{update: } \theta = \theta - \frac{\eta}{\sqrt{v_t + \epsilon}} \cdot \hat{m}_t$$

Pros:-

- Adaptive & momentum optimizer benefits
- Efficiently works with large & noisy data.
- Convergence faster.



Cons

- . Chances of overshoots
- . computationally expensive

- Chances of overshoots
- computationally expensive

$$\underline{\text{Adomax}} \quad \hat{v}_t = \max(\beta v_{t-1}, |\nabla S(\theta)|)$$

$$\theta_t = \theta - \frac{\eta}{v_t} \hat{a}_t$$