

- Data
- Python
- Database
  - MySQL
  - MongoDB
- Big data technologies
  - Hadoop
  - Spark -

- Power BI
- ML Fundamentals

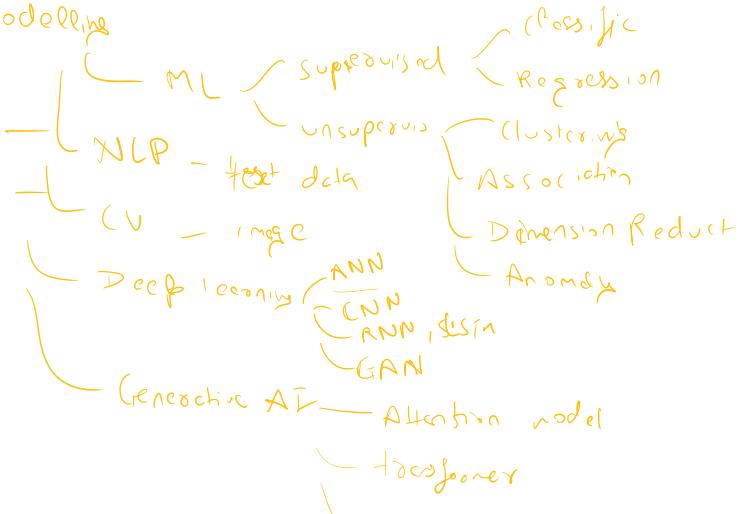
- Cloud services
  - Azure → DE → DP-201
  - AWS
- Kafka
- Airflow
- Snowflake
- Informatica
- Hive

→ Python

→ Data Preproc/ und.

→ Maths  $\begin{cases} \text{stats} \\ \text{prob} \end{cases}$  Hypothesis testing

→ Modelling



→ Deployment

Feature Engineering

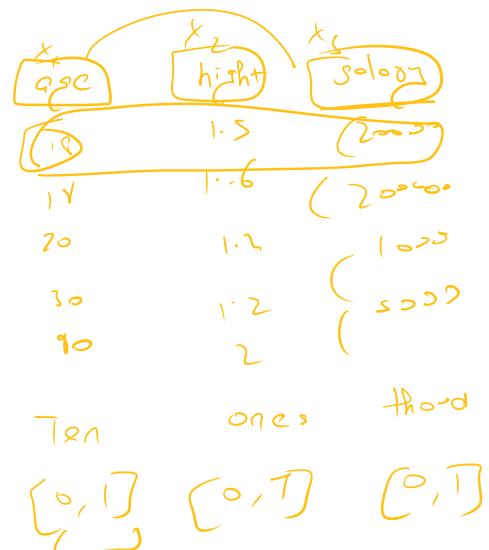
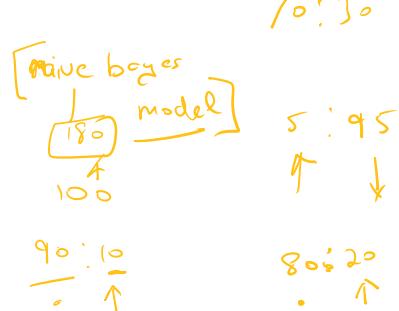
◦ Encoding —

$$y = m_1x_1 + m_2x_2 + m_3x_3 + c$$

◦ Scaling —

◦ Sampling

◦ Selection



Encoding:- Label Encoder

workclass

Private	-	0 -
Gov	-	1 -

◦ One Hot Encoder

	Priate	Gov	SE	No
1	0	0	0	0
0	1	0	0	0

Gou	-	1	-		1	0	0	0
SE	-	2	-		0	0	1	0
NW	-	3	-		0	0	0	1

→ Hi How are you  $\Rightarrow$  Hi how are you i am good how  
 → Hi, i am good how are you 1 0 0 1 0 0 0 0

## Scaling

- Standardization  $\rightarrow$  [Z-score]

• Mean of col = 0

$$\frac{18+20+25+35}{4}$$

$$Z = \frac{x - \mu}{\sigma}$$

$$\sqrt{\frac{(x-\bar{x})^2}{n-1}}$$

$$\begin{array}{c} \text{age} \\ \left[ \begin{array}{c} 18 \\ 20 \\ 25 \\ 35 \end{array} \right] \xrightarrow{\mu=24} \left[ \begin{array}{c} -0.47 \\ -0.31 \\ -0.07 \\ 0.87 \end{array} \right] \xrightarrow{\sigma=12.7} \left[ \begin{array}{c} 0 \\ 0 \\ 0 \\ 0 \end{array} \right] \end{array}$$

- Min Max Scaler  $\rightarrow$

$$[0, 1]$$

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

$$\frac{35 - 18}{35 - 18} = 1$$

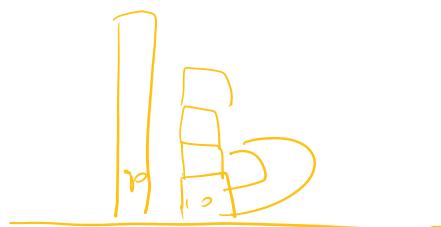
$$\frac{35 - 20}{35 - 18} =$$

$$\frac{35 - 35}{35 - 18} = 0$$

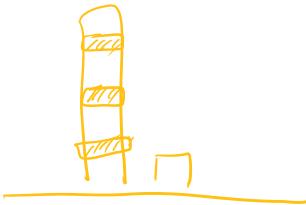
## Sampling / Balancing

### ① Random sampler

- oversampling



- under sampling



### ② SMOTE



ENN



## Feature Selection

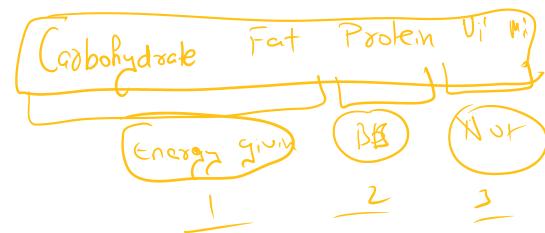
- Removing features
  - Correlation
  - Mutual-informability
- Decision

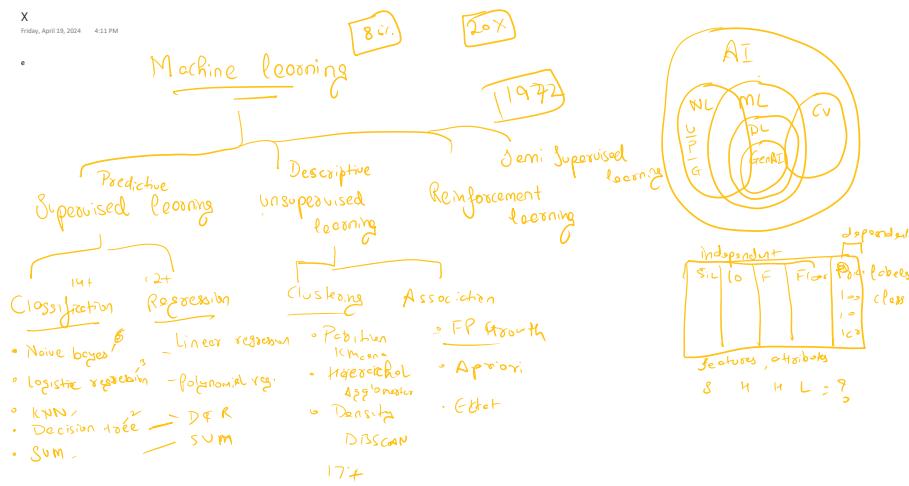
☰ 13+

"Curse of dimensionality"

- Data compression

unsupervised PCA - Principle Component Analysis  
LDA





### Noise Bayes Model

- Prob. occurrence of event / total event

$$1 \text{ Toss} = \frac{H}{T} \quad H = \frac{1}{2}$$

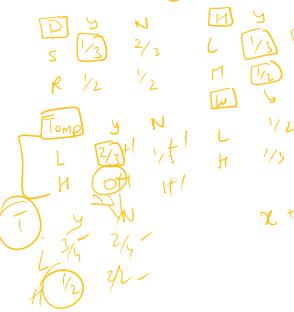
$$2 \text{ Toss} = \begin{array}{c} H \\ H \\ H \\ T \\ T \\ T \\ T \end{array}$$

- Conditionally Prob.

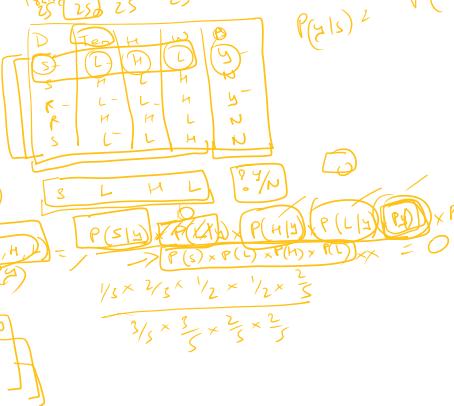
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

Training:

$$P(S) = \frac{2}{5}, \quad P(N) = \frac{3}{5}$$



P(B|A)



Assumptions:

- All feature pull equal contribution to label
- independent feature

Noise



When to Apply to

- limited data
- single fast implementation
- less dependency

- Assumption of other
- density of data
- missing data
- zero case scenario

Laplace

### Use Case

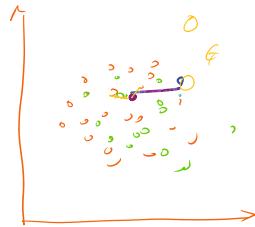
- Text classification
- Recommendation
- Medical diagnosis

Types of NB

- Discrete NB
- Continuous NB
- Naive

- Types of ND
- Gaussian NB - Continuous DTC
  - Multinomial NB / General NB
  - Bernoulli NB - Discrete
  - One vs all
  - out of core

K-nearest Neighbours =  $[K=3]$

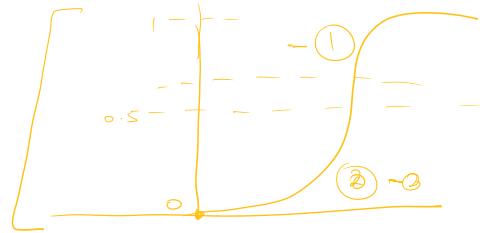


## Logistic Regression:

$[0, 1]$

A B C

$$y = \frac{1}{1 + e^{-\mu}}$$



Use Cases:

- Customer Churn
- Credit scoring
- Medical diag.

Limitations

• no outliers

• non linear

• binary classification

Assumption

- Binary target
- Independent obs.
- data linearly separable

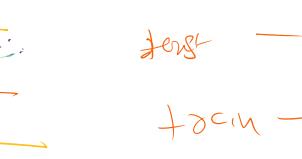


L. (lo)

when to Apply

• More amount of data

• Assumption

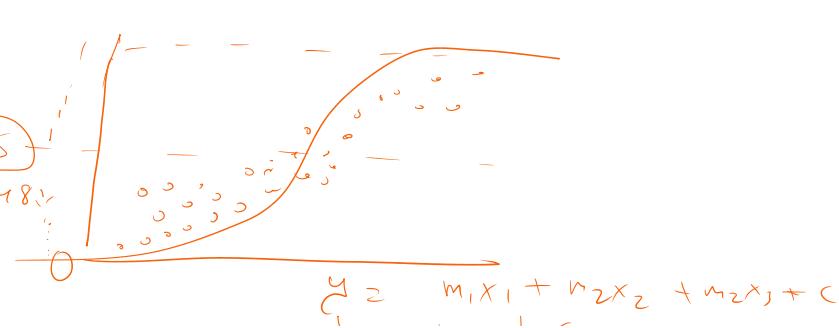


Deep learning

$$\begin{array}{c} A \\ \times \\ \times \\ \times \\ \times \end{array} = \boxed{\begin{array}{c} B \\ \times \\ \times \\ \times \\ \times \end{array}} = \boxed{\begin{array}{c} C \\ \times \\ \times \\ \times \\ \times \end{array}} = \boxed{\begin{array}{c} D \\ \times \\ \times \\ \times \\ \times \end{array}}$$

$$0.4 = \boxed{0}$$

$$0.2 = \boxed{0}$$



$$A = \begin{pmatrix} A_1 \\ A_2 \\ A_3 \\ A_4 \\ A_5 \\ A_6 \\ A_7 \end{pmatrix} \quad m \cdot A = B = \begin{pmatrix} B_1 \\ B_2 \\ B_3 \\ B_4 \\ B_5 \\ B_6 \\ B_7 \end{pmatrix}$$

$y = m_1x_1 + m_2x_2 + m_3x_3 + c$   
 $y = mx + c$

$\log e^x = m \log n + l$   
 $e^x = \frac{1}{1+e^{-x}}$        $0.5 = \Phi$   
 $0.5 \in \varnothing$

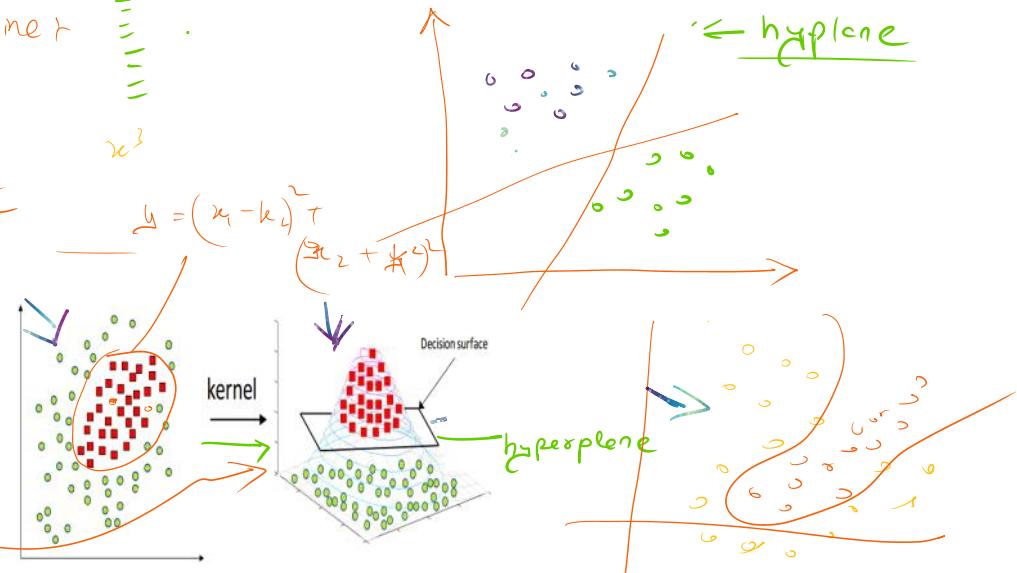
SVM

## Support Vector Machine

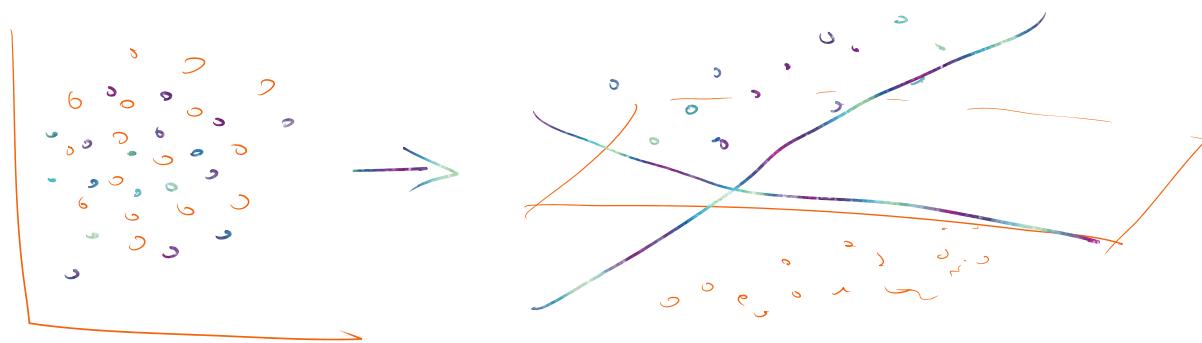
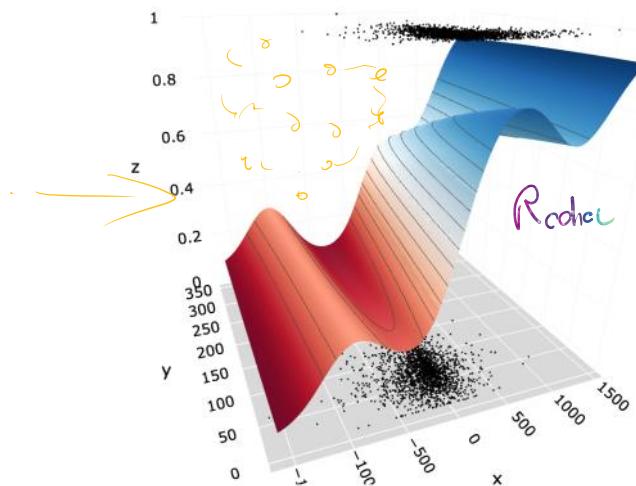
Kernel  
linear

Polynomial

Radical



Deeps



$$y = \underline{\omega^T x + b}$$

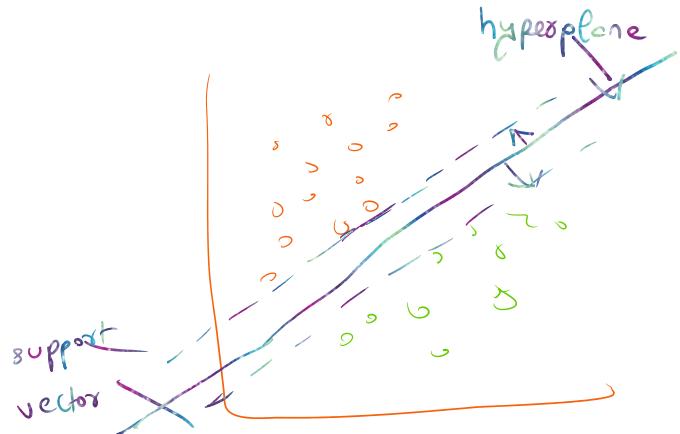
When to apply

Use Cases

Image classif.

Audio classif

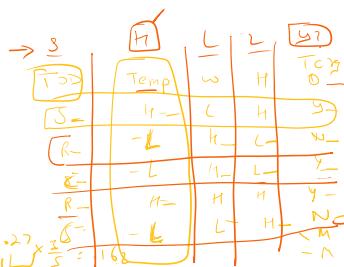
Tent. dofe



Decision Tree Model

$$IG = -\frac{P}{P+N} \log_2 \left( \frac{P}{P+N} \right) - \frac{N}{P+N} \log_2 \left( \frac{N}{P+N} \right)$$

$$TOD = \begin{cases} I(S) & S \neq R \\ I(R) & S = R \end{cases}$$



Temp : [0.62]  
wind : [0.62]  
Humid : [0.62]

$I(S) = -\frac{3}{5} \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \log_2 \left( \frac{2}{5} \right) = 0.971$   
 $I(R) = -\frac{1}{2} \log_2 \left( \frac{1}{2} \right) - \frac{1}{2} \log_2 \left( \frac{1}{2} \right) = 1 - \frac{1}{2} = \frac{1}{2} = 0.5$

Gains = Entropy of Target - Entropy of columns  
 $- TOD = 0.971 - 0.562$   
 $- Wind = " "$   
 $- Humid = " "$   
 $- Temp = 0.971 - 0.162 = 0.809$   
 - (i) no. of unique  
 - (ii) missing value

Gini index

$$1 - P(i)$$

$$I_G = 1 - \sum_i p_i^2$$

$$J_G = - \sum_i p_i \log_2(p_i)$$

Fundamental

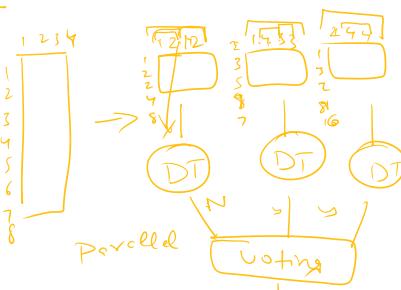
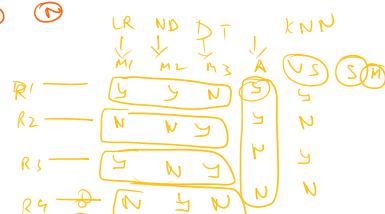
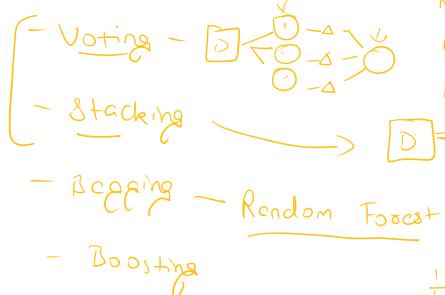
- NB -
- KNN -
- SVM -
- LR
- DT

/// //

Random Forest

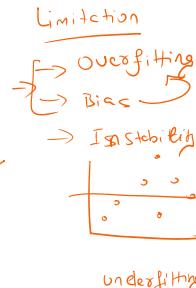


Ensemble models



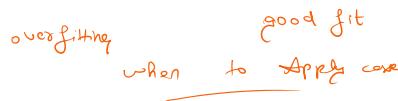
Use cases

Financial / Banks

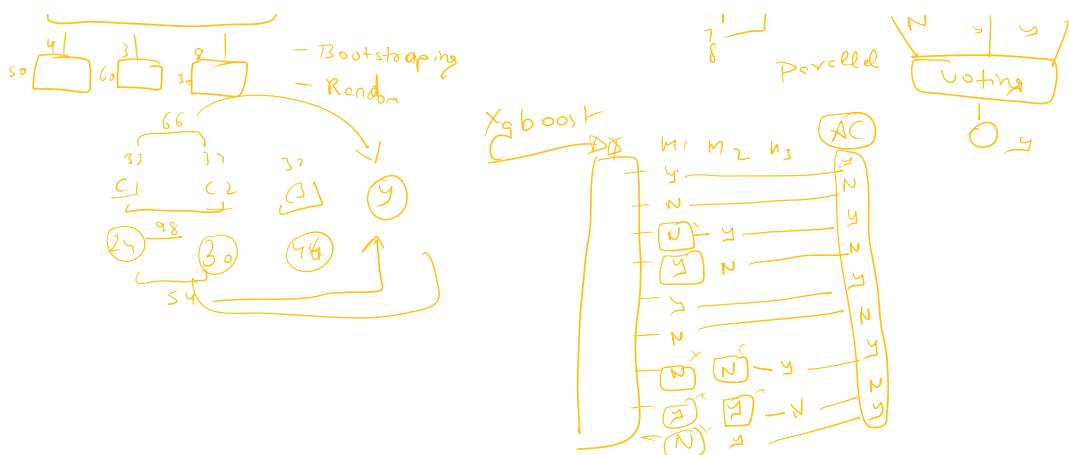


Limitation

- Overfitting
- Bias
- Instability



when to apply case



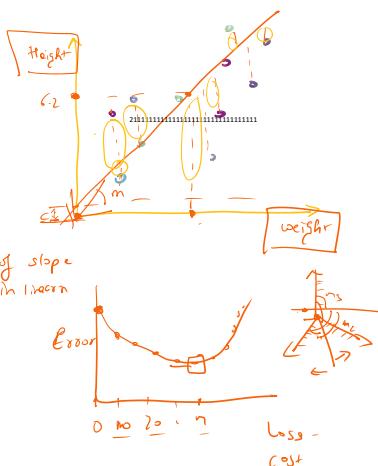
## → Regression Model

### ◦ Linear Regression

◦ Residual

$$\text{Simple Linear R: } \hat{y} = mx + c$$

$m$  = slope of line - significance of slope in linear  
 $c$  = constant



### ◦ Multiple linear regression

$$y = m_1x_1 + m_2x_2 + m_3x_3 + \dots + c$$

$$\text{◦ Loss function - total loss: } = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$\text{◦ cost function: } = \text{M}$$

## Regularization techniques

$$\text{◦ L1 - Lasso Regression: } = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \lambda \sum |m_i|$$

$$\text{◦ L2 - Ridge Regression: } = \frac{1}{2} \sum (\hat{y}_i - y_i)^2 + \lambda \sum (m_i)^2$$

$$y = M_1x_1 + M_2x_2 + M_3x_3 + \dots + c$$

Feature selection



## Use Cases

◦

Numerical

## Assumption

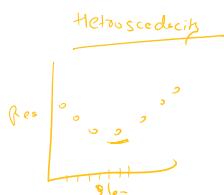
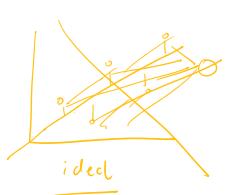
- Linearity
- Homoscedasticity
- Normality of error

## Limitation

- Assumption
- Overfitting

## Apply

- Categorical data





### Assumption:

- Cluster Shapes - Spherical
- Homogeneous Density - Almost equal datapoint
- Single Variance -

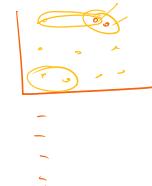
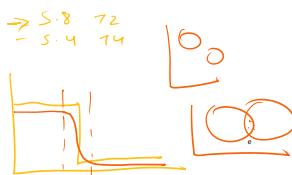
Types of Kmean

- kmean++ - default

- Mini Batch Kmean -
- Spectral Clustering -



→ 57 67 - C1  
→ 54 68 - C2  
→ 60 75 - C3



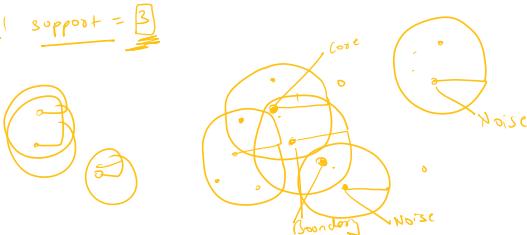
### Limitation:

- Sensitivity to initial centroid
- Sensitivity to outlier
- Number of clusters

### DB Scan

epsilon = 1   support = 3

- Core
- Boundary
- Noise



### Association Technique

- Apriori - support = 2
- FP growth

step - 1  
A - 3  
B - 3  
C - 4  
D - 4  
E - 1

1	A	B	C
2	B	C	D
3	D	A	B
4	C	D	A

step - 2

1	C	A	B
2	C	D	B
3	D	A	B
4	C	D	A

step 3



### Use Cases

### Assumption:

#### Limitations

- memory usage
- Non Specific

### Reinforcement Learning

- Interaction based
- Global based



#### Environment

- State
- Action
- Reward/Penalty
- Agent

◦ Interaction based



◦ Goal based

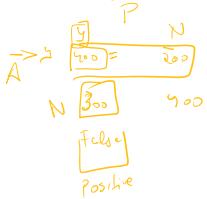
Accy  
Precision :

◦ Action

◦ Reward/Penalty

◦ Agent

◦ Policy



## Deep Learning

BB-A Neuron

◦ Activation function

- Sigmoid :  $[0, 1]$



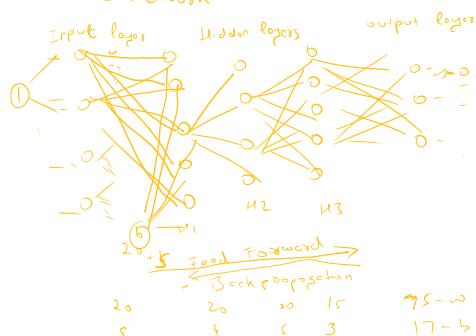
- Step :  $0, 1$



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

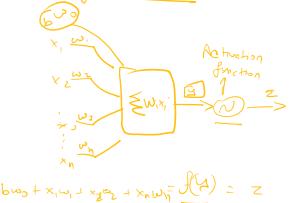


Feed Forward/Fully connected Neural Network

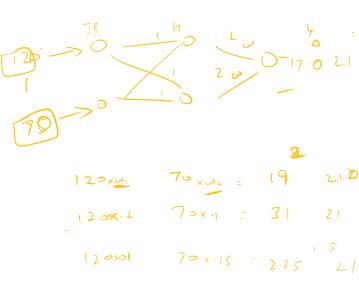


$y = mx + b$

Artificial Neuron



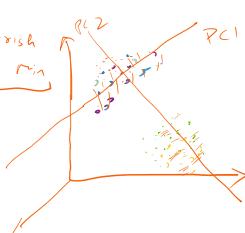
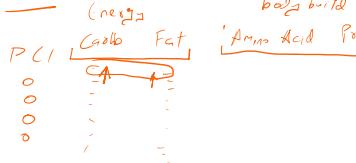
$$bias + x_1w_1 + x_2w_2 + \dots + x_nw_n \underset{f(A)}{\rightarrow} z$$



→ Convolutional NN - Images

→ Recurrent NN - Text/sequential

PCA



Natural language : - Processing

- Understanding

## Processing

~ Lowercase

• Tokenisation

• Normalization

• Stop word removal

• Stemming & lemmatization

• Cleaning & correction

## - Understanding

### - Generation

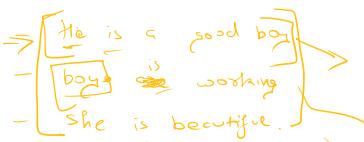
### Understanding

### Vectorization

◦ Bag of words →

TF-IDF → importance of word

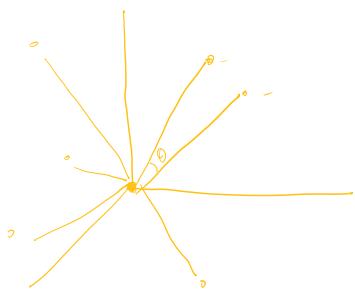
Word2Vec →



He	is	a	good	boy	working	she	-	-
1	1	1	1	1	0	0	0	0
0	1	0	0	1	1	0	0	0

- Tensor  
- tensorflow  
- nltk  
- spacy -

235 - Bank - aiwan bank - 235  
289 - Bank - Financial - 289



## Project Submission

- Sunday, 23 / June

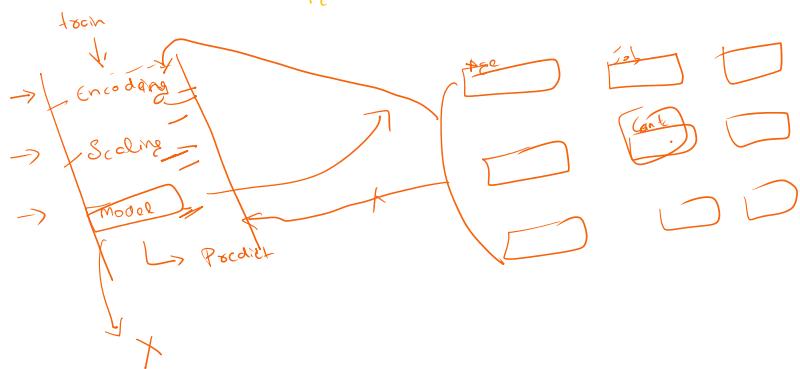
- 3 project

◦ Code

◦ Download : ◦ Business Problem [2 pages doc.]  
◦ Data Description

◦ Basic step in data processing & major challenge  
◦ Choose model, why, pros & cons

- Bank marketing dataset project



## Window functions

r - id    dep    salary | Month

select sum(salary) from emp = 470000

## Window Functions

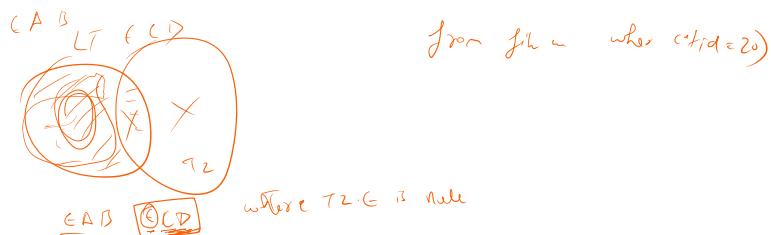
Cnpid	dept	Salary	X Total		select sum(salary) over (partition by dept) for ~
			21	23000	
1	1	50000			
2	1	60000			
3	1	65000			
4	1	55000			
5	2	110000	21		
6	2	100000	21		
7	3	20000	50		
8	3	30000	50		
9					
					490000

## Types of window function

- Aggregate
  - SUM()
  - COUNT()
  - AVG
- Rank
  - Rank()
  - Dense Rank
  - Row Number()
- Value Function
  - Lead
  - Lag
  - First Value

Cnpid	Dept id	Salary	Rank	Dense Rank	Rn	Lead
1	-1	50000	1	1	1	60000
2	-2	60000	2	2	2	70000
3	-2	70000	3	3	3	80000
4	-1	65000	4	4	4	68000
5	3	50000	5	5	5	68000
6	-1	66000	6	6	6	70000
7		70000	8			Null

Select custname from customer where title in ('saler trk')



## Fact & Dimension

