

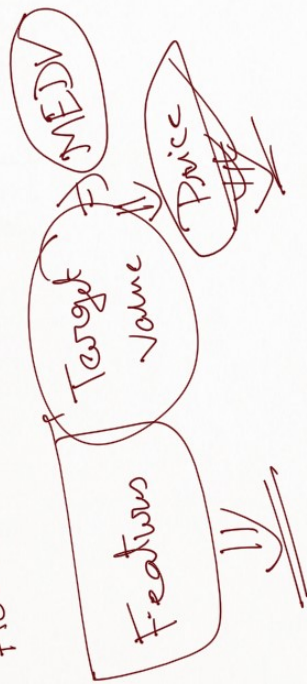
Supervised \Rightarrow Labelled Data \Rightarrow Target

Howe
 $h_{\theta}(x) =$

$$\theta^T X = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

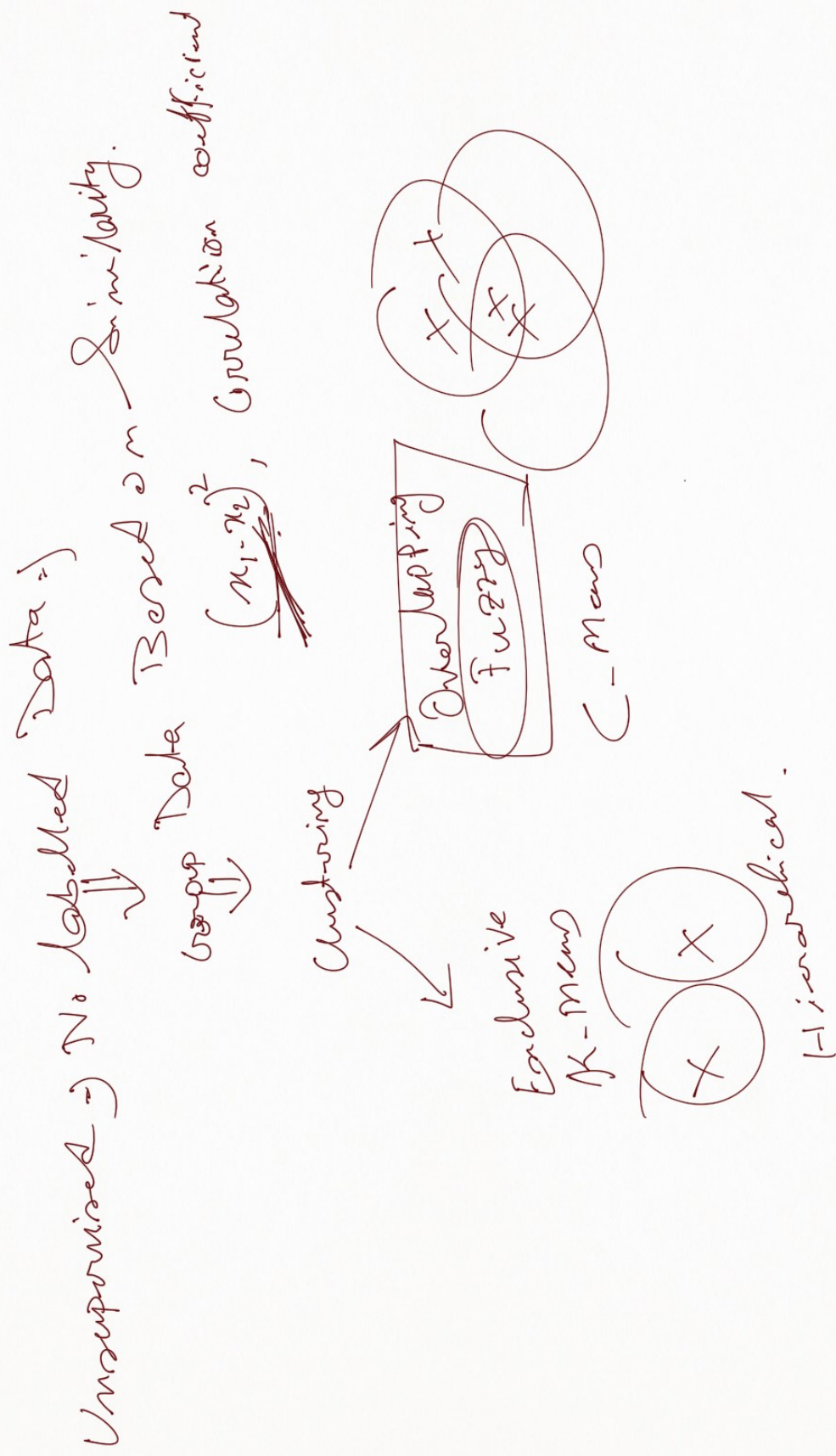
Price Prediction

Historical Data \Rightarrow Train a model \Rightarrow Trained, split
 \Downarrow
Remove Overfitting \Rightarrow Bias
 \Downarrow
Train the model
 \Downarrow

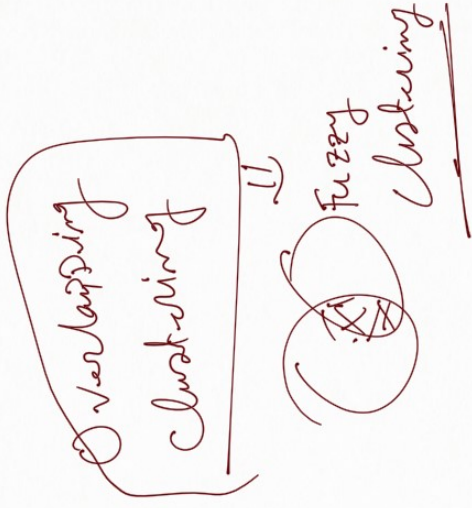
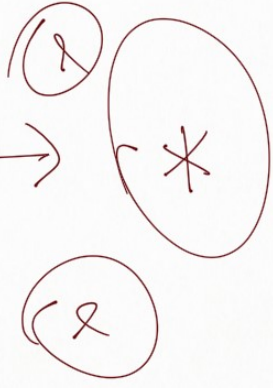


Labelled Data \Rightarrow Decurately.

Predict \Leftarrow Model \Leftarrow Train metrics based
on dependent data
Gradient Descent



clustering \Rightarrow Exclusive Clustering



Machine Learning \Rightarrow Supervised Unsupervised

Factorial $\Rightarrow y = m(x) + \epsilon$
 $n = 100, C = 20$
 $N = 5000$
 Dependent Variable

Coefficient

Source

Subpopulation
Variable

Statistical Learning

models

Training Accuracy

$$\sum_{i=0}^n \theta_i x_i \Rightarrow y = [1, x_1, x_2, \dots, x_n]$$

Dependent Variable

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$\begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 3 & 4 & 5 & 6 & 7 & 8 \\ 6 & 7 & 8 & 9 & 10 & 11 & 12 \end{bmatrix}$$

$$\begin{matrix} 3 \times 2 & 2 \times 3 \\ \hline 3 \times 3 \end{matrix}$$

$$n \times 1 \quad 1 \times n$$

$$\theta = \begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix}_{1 \times n}$$

$$h_0(x) = \theta^T X \Rightarrow$$

$$X = \begin{bmatrix} 1 & x_1 & x_2 & \dots & x_n \end{bmatrix}_{1 \times n}$$

Transpose \Rightarrow Exchange Rows & Columns

$$\underline{h_0(x)} = \theta_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$X \Rightarrow$ Features/Columns of your dataset

House Price \Rightarrow Gini, Nor, En \Rightarrow

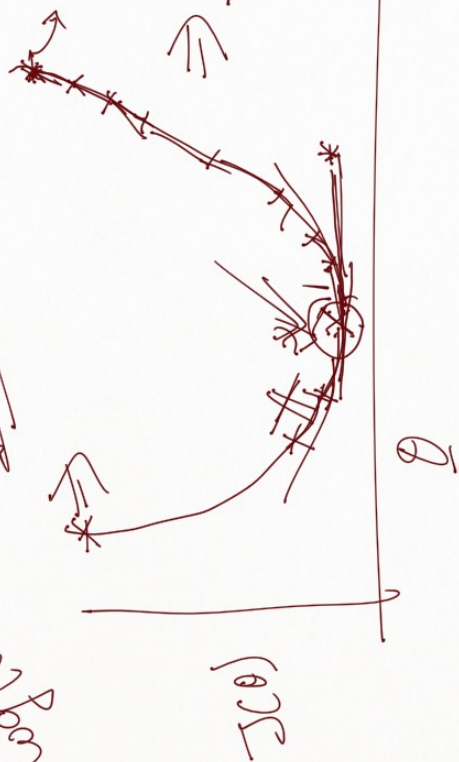
Cost Function $\Rightarrow \frac{\partial J(\theta)}{\partial \theta} = 0$

MSE \Rightarrow

$$\sum_{i=1}^n (h(\theta; x_i) - y_i)^2$$

Optimum $\theta \Rightarrow$ finding minimum

Back \Rightarrow $-(y - \hat{y}) \log(y - \hat{y})$



\Rightarrow Gradient Descent

$\theta_0 = -\frac{\partial J(\theta)}{\partial \theta}$

$\theta_0 = \theta_0 - \alpha \frac{\partial J(\theta)}{\partial \theta}$

Learning Rate \Rightarrow slope of $\frac{\partial J(\theta)}{\partial \theta}$

$J(\theta)$

Derivatives

$\frac{\partial J(\theta)}{\partial \theta} = \frac{\partial}{\partial \theta}$

Gini Index
Entropy

$$e_{h(x)} = \frac{e_{h(x)}}{1 - e_{h(x)}}$$

Rules: $1 \Rightarrow A \Rightarrow B$
 $2 \Rightarrow B$

Data Mining

Break & Butter

China Drops \Rightarrow US with $A \Rightarrow B$

China Drops join water

Ramban Forest Analysis

Market Forest Analysis

Market Forest Analysis

Rice & Purses
Fruit & Vegetable
Purses & Chocolate

$A \Rightarrow B$

Analysis Technique \Rightarrow Generate Associations/relations \Rightarrow

Retail Store \Rightarrow

Order id
1
2

Products
[- m - item]



Analysis

Market Basket

{
Rice \Rightarrow Pulad
Rice \Rightarrow Opl
Bread \Rightarrow Milk
Bread, Buttry \Rightarrow Milk
}

Apriori Algorithm of Finding out association - $\{milk, egg\} = 3$
 $S\{milk, egg\} = \frac{3}{10}$
 $A \Rightarrow B \Rightarrow 2$
 \downarrow
 Antecedent consequent

Support $C(A, B) = \frac{C(A, B)}{S(B)}$
 Confidence $\frac{S(A, B)}{S(A)}$
 Lift $\frac{C(A, B)}{S(A) \times S(B)}$
 $\frac{3}{10} \times \frac{3}{10} = \frac{9}{100}$
 $\frac{3}{10} \div \frac{3}{10} = 1$
 $\frac{3}{10} \times \frac{3}{10} = \frac{9}{100} < 1$

Frequency of a item
 Total transaction
 $\frac{F(item)}{Total} = \frac{7}{10}$
 $\frac{7}{10} \times \frac{3}{10} = \frac{21}{100}$
 $\frac{7}{10} \div \frac{3}{10} = 2.33$
 $\frac{7}{10} \times \frac{3}{10} = \frac{21}{100}$

$\{milk, egg\} = 3$
 $\{milk, egg\} = 3$
 $\{milk, egg\} = 3$

$K \Rightarrow$ Number of items

support \rightarrow threshold

Apriori Steps \Rightarrow

Reject all item whose

1. $K=1$ - Calculate Support, Confidence, Lift. \Rightarrow

2. $K=2$, Calculate Support, Confidence, Lift. \Rightarrow

3. Repeat Steps 2 by incrementing K till remaining -

we have no more item set remaining -

Threshold for support, Confidence, Lift. \Rightarrow

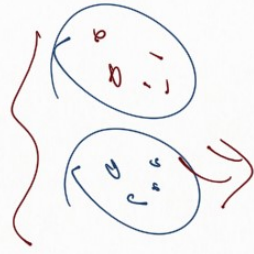
0.5 0.5 2

milk Bread Egg

milk Bread milk

Bread, milk, egg

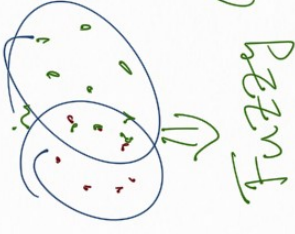
K-Mean Clustering ⇒



Exclusive Clustering



Customer Segmentation
Recommendation



Fuzzy

Clustering Overlapping Clustering.

Find out similarity of items.

K-Means \Rightarrow Hierarchical clustering \Rightarrow



$\{1, 2, 3, \dots, n\} \Rightarrow N$



\Rightarrow cluster

\Rightarrow different cluster \Rightarrow $\{c_1, c_2, c_3, \dots, c_k\}$

\Rightarrow Number of elements included c_i

$\sum_{i=1}^K |C_i| = N$ where $|C_i| \Rightarrow$ Intra cluster variance

$\Rightarrow \text{Var}(C_i) \Rightarrow$ Intra cluster variance

$\Rightarrow \text{Var}(C_i) \Rightarrow$ group similarity

$\Rightarrow W(C_i)$

minimize $\Rightarrow \sum_{i=1}^K w(i) \Rightarrow$ Euclidean Distance

$w(i) \Rightarrow$ intra cluster variance

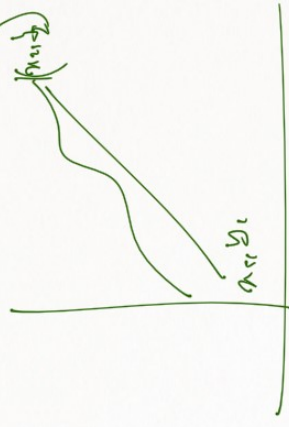
$$\begin{aligned} & \frac{(x_1, y_1) (x_2, y_2)}{\sqrt{(y_1 - y_2)^2 + (x_1 - x_2)^2}} \\ & E.D. = \sqrt{\sum_{i=1}^n (y_i - x_i)^2} \end{aligned}$$

Cost function $\Rightarrow \sum_{k=1}^K \sum_{i \in C_k} (x_{ik} - \mu_k)^2$

\Rightarrow

$$\left\{ (x_1 - x_2)^2 + (x_1 - x_3)^2 + (x_2 - x_3)^2 + \dots \right\}$$

$$\left\{ (x_4 - x_5)^2 + (x_5 - x_6)^2 \right\}$$



Similar Objects \Rightarrow

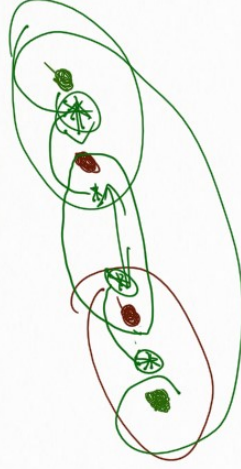
algo \Rightarrow Random Selection

k - means

mem of
All points
in a cluster

\Rightarrow Inter-Cluster Variance Cluster-Centroid

\Rightarrow Shift points betw.



$$\sum_{k=1}^K \sum_{\substack{i,j \in C_k \\ i < j}} (x_{ik} - x_{jk})^2$$

where $K \geq$ Number of clusters $\geq 2, 3, 4, \dots$

$$i, j \in C_k = \{x_1, x_2, \dots, x_K\}$$

$$\sum_{k=1}^K \left\{ \sum_{i,j \in C_k} (x_{ik} - x_{jk})^2 \right\}$$

$$i, j \in \{1, 2, \dots, N\}$$

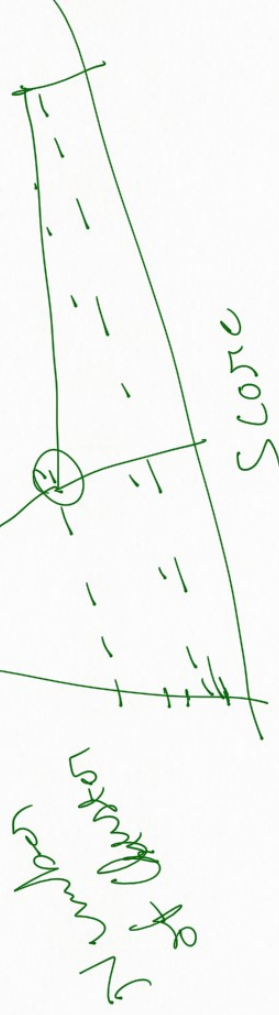
Steps of K-Means →

- 1) Guess 3-Random Centroids
- 2) Calculate variance of each point from the centroids. point To the closest centroid
- 3) Design point Centroid.
- 4) Recalculate 2-4 times no shifting occurs.
- 5) Repeat

Elbow Method

Total Variance in all clusters \Rightarrow For each $K \Rightarrow$

$$K \in R \Rightarrow K = \{1, \dots, K_{\max}\}$$



Silhouette Score = $\frac{\sqrt{a-b}}{\max(a,b)}$

from

distance of a point

cluster.

a \Rightarrow distance nearest diff.

b \Rightarrow intra cluster

b \Rightarrow Silhouette Score ≤ 1

-1 \leq

Silhouette Score ≤ 1

cluster are very far

+1 \Rightarrow cluster are very close

0 \Rightarrow cluster boundary is not clear

-1 \Rightarrow Decision

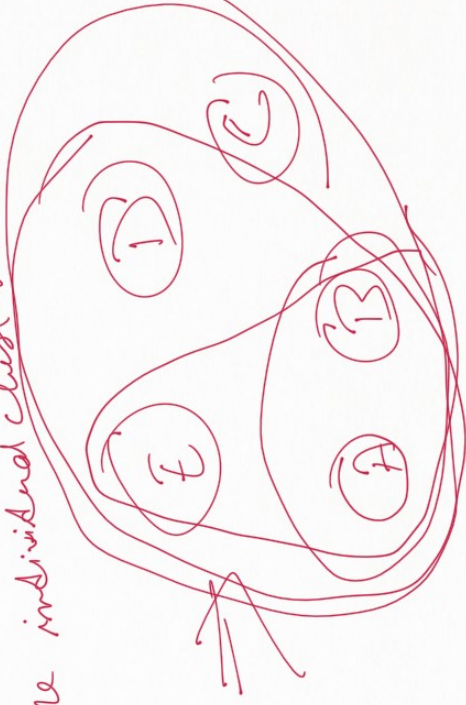
Hierarchical

cluster \Rightarrow Entire space is 1 cluster.
Individual points are individual clusters.

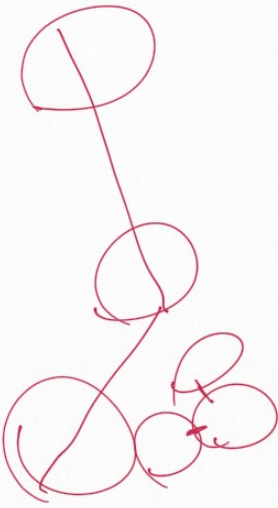
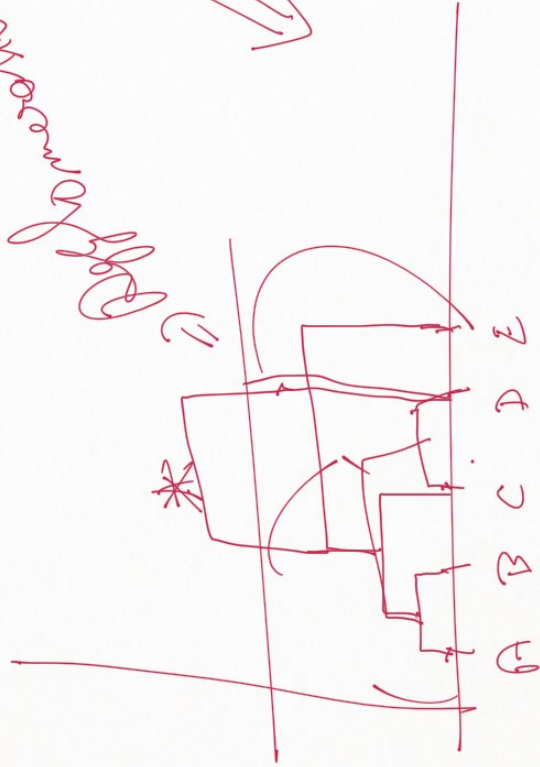
\Rightarrow Agglomerative

Bottom-up \Rightarrow Divisive

Top-down \Rightarrow Divisive



Diff. Parameter :-



1) which points to combine for one cluster
2) which cluster to form together

Linkage method :-
1) Single linkage
2) Complete linkage
3) Group Average
4) Centroid
5) Ward

Linkage method :-
1) Single linkage
2) Complete linkage
3) Group Average
4) Centroid
5) Ward

1) Single linkage
2) Complete linkage
3) Group Average
4) Centroid
5) Ward

