# Day 6 — Machine Learning Algorithms

## Unsupervised ML

① K Means Clustering
② Hierarichal clustering
③ Silhoutte Score
④ DBScan Clustering

## Unsupervised ML

| O/p | $f_1$ | $f_2$ |
|-----|-------|-------|
| | — | — |
| Clusters | — | — |
| ⇩ | — | — |
| { Similar kind of data } | — | — |

## K Means Clustering

$f_2$



cluster 1
cluster 2
$f_1$

## Custom Ensemble Technique



→ 2-3 group
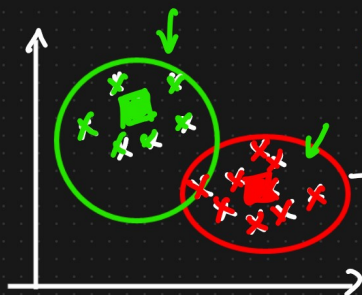
□ → | Clustering | → | For each |

K Means → K = centroids

Eucledian distance     K=2

High dimension

→ Avg

[ K=2 ]

→ centroids

① We try K values ⇒ Suitable K=2

② Initialize K number of centroids ✓

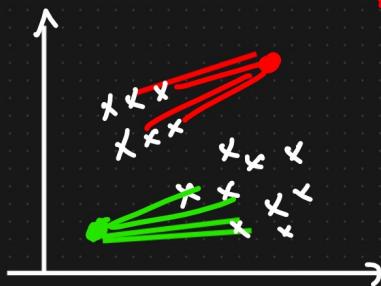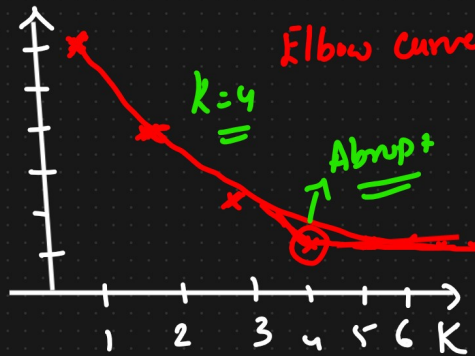③ Compute the avg to update centroids ✓

Validating

K=

Within cluster sum of square

## Elbow method (K value)

K=1

K=2 WCSS
=

for i=1,10

Elbow curve

K=4

Abrupt

1  2  3  4  5  6  K

K=4

You need find the longest vertical line
that has no horizontal line passed through it.
→ Dendogram

② Hierarchical Clustering

P1    P2

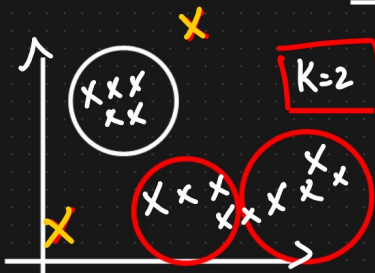P6   P3
P7  P5
       P4

4 Clusters

6
5
4
3
2
1

P1  P2  P3 P4  P5  P6  P7

Dataset is small
Dataset is large
Kmeans

Max Time   is taken by KMeans or

Hierarchical Clustering ?? ✓
    ↳ Max Time

X

K=2 ✓

XXX
XX

XXX X XX
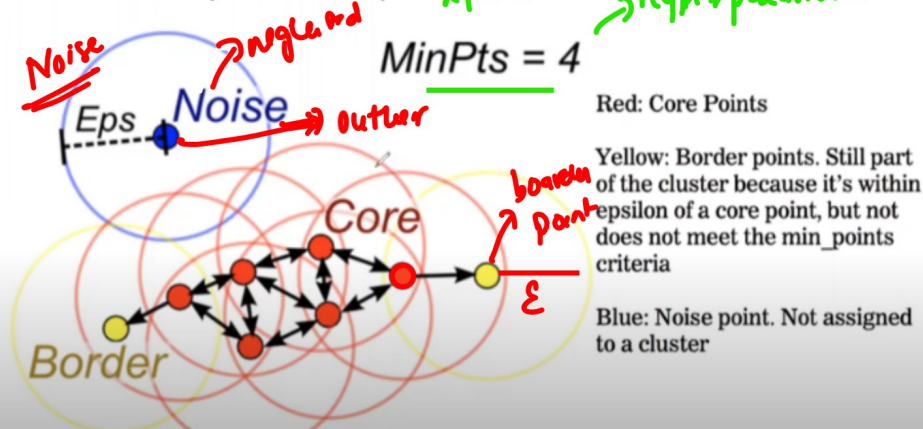
X ⇒ { KMeans r+t }

# Validate Clustering Models

## DBScan Clustering

Core point

border point

$\varepsilon$

Minpts = 4

$\varepsilon$

→ Core point

Epsilon

① Min pts

$\varepsilon$

Noise point

② Core points

③ Border points

④ Noise point

**Density-Based Spatial Clustering of Applications with Noise(DBSCAN)**

Epsilon → Hyperparameter

MinPts = 4

Noise → Neglected

Eps

Noise → Outlier

Red: Core Points

Core

border point

$\varepsilon$

Yellow: Border points. Still part of the cluster because it's within epsilon of a core point, but not does not meet the min_points criteria

Blue: Noise point. Not assigned to a cluster

Border

→ Noise Point
→ Outlier

Kmeans
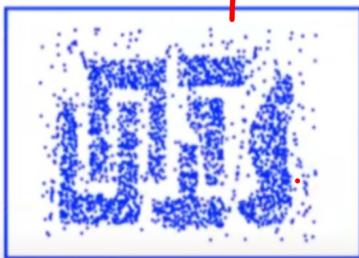
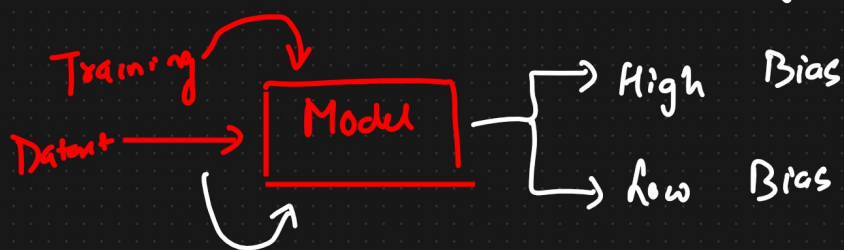KMeans

DBScan Clustering ⇓ DBScan

The left image depicts a more traditional clustering method that does not account for multi-dimensionality. Whereas the right image shows how DBSCAN can contort the data into different shapes and dimensions in order to find similar clusters.
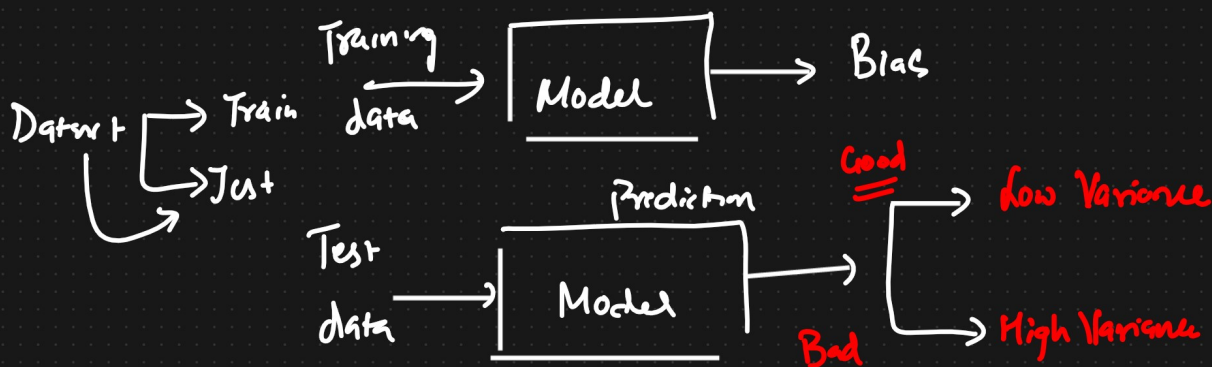
# Defn of Bias And Variance

Training Dataset = 90% ⎱
Test Dataset = 70% ⎰ ⇒ Overfitting

⇓

{ Low Bias ✓}
{ High Variance ✓}

**Bias :** It is a phenomenon that skews the result of an algorithm in favor or against an idea.

⌐→ Training dataset

Training → | Model | → High Bias
Dataset → | Model | → Low Bias

**Variance :** Variance refers to the changes in the model when using different portions of the training or test data

Training → | Model | → Bias

Dataset ⌐→ Train data
        ⌐→ Test

                    Prediction
Test
data → | Model | →

Good ═→ Low Variance

Bad → High Variance

**Model 1**

Training Acc = 90%.
Test Acc = 75%

**Model 2**

Train Acc = 60%
Test Acc = 55%.

**Model 3**

Train Acc = 90%.
Test = 92%.

$\downarrow$

$\left\{ \begin{array}{l} \text{Low Bias} \\ \text{High Variance} \end{array} \right\}$   $\left\{ \begin{array}{l} \text{High Bias} \\ \text{High Variance} \end{array} \right\}$   $\left\{ \begin{array}{l} \text{Low Bias} \\ \text{Low Variance} \end{array} \right\}$

Generalized Model

$\left\{ \begin{array}{l} \text{Low Bias} \\ \text{High Variance} \end{array} \right\}$   $\left\{ \begin{array}{l} \text{High Bias} \\ \text{High Variance} \end{array} \right\}$   $\left\{ \begin{array}{l} \text{Low Bias} \\ \text{Low Variance} \end{array} \right\}$