

## Agenda:

1. PCA
2. Clustering
3. Model Selection
4. Regularization
5. Boosting

PCA: - Unsupervised ML

- solve problem of 'Curse of Dimensionality'
- v.v. high # features  $\rightarrow$  Model  $\rightarrow$  Overfits
- $\therefore$  Go for PCA  $\rightarrow$  Captures variance in the I.V.
- Statistical transformation technique
- $\therefore$  accepts only numerical data & that too only X

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

PCA is applied on  $X_1, \dots, X_5$ .

- always apply PCA on scaled data.

Steps: 1. read csv 2.  $X, y$  3. Train Test Split  
↑ Scale  $X$

$\left\{ \begin{array}{l} \text{fit\_transform}(X_{\text{train}}) \\ \text{transform}(X_{\text{test}}) \end{array} \right.$   
5. PCA

$\left\{ \begin{array}{l} \text{fit\_transform}(X_{\text{train}}) \\ \text{transform}(X_{\text{test}}) \end{array} \right.$

6. Modeling

7. Evaluating

8. Deployment

Q1. What is Curse of Dim?

Q2. Describe pc.

Q3. Can PCA be used in Feature Selection? No.  
↳ is Feature Reduction

Q4. Comment on PC1.

Q5. Disadv. of Dimensionality reduction? (Original features  
lost.)

Q6. Do we scale data by PCA.

Q7. Can we use it on Large datasets?

# Clustering:

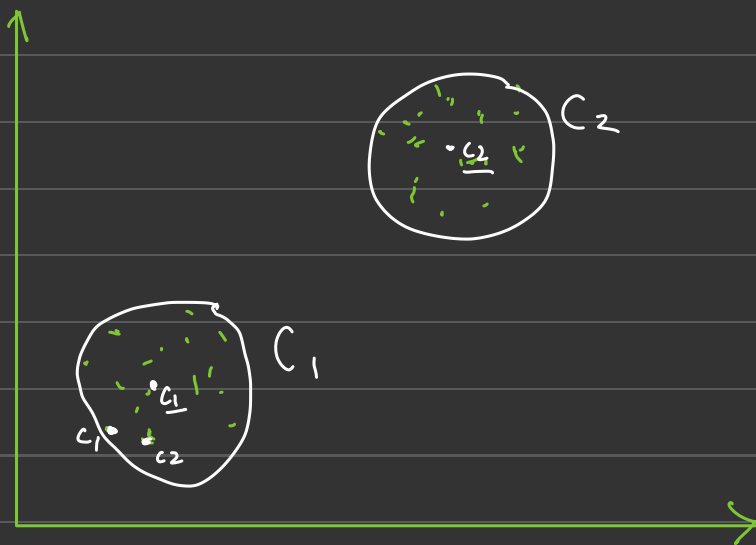
## 1. Euclidean Distance

$$(2,1) \quad d \quad (5,2)$$

$$d = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$= \sqrt{(5-2)^2 + (2-1)^2}$$

## 2. Math behind



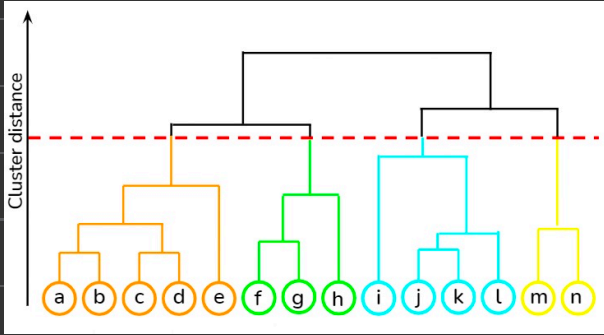
## 3. Stopping Criteria - No change in centroid values

## 4. Diff b/w kmeans & KNN

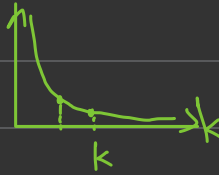
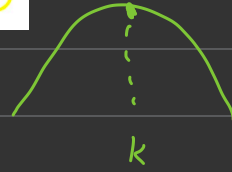
Unsupervised  
Only X

Supervised  
X & y

## 5. Hierarchical Clustering vs Kmeans Clustering



6. What is Silhouette score?
7. Significance of Hopkins test.
8. Elbow method: wcss



Choosing right value of k.

Regularization: avoid Overfitting

$$y = 15 + 1.2x_1 + 20x_2 + 39x_3$$

↓ reg

$$y = 0.9 + 1.2x_1 + 2x_2 + 5x_3$$

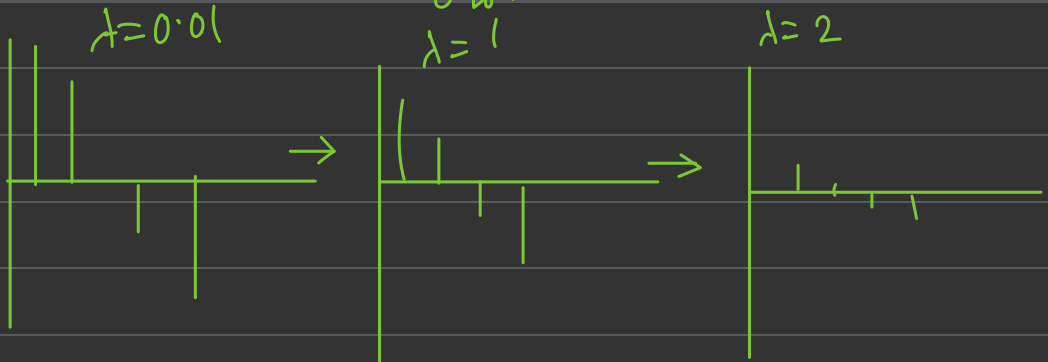
penalty

1. Ridge Reg.:  $\text{Loss} + \lambda |m|^2$   
or L2 Reg.:  $(y - \hat{y})^2$

vector of coefficient

Hyperparameter

0 to  $\infty$   
 $\lambda = 1$



$\lambda \uparrow \rightarrow$  Lesser value of coeff.

Use Ridge: When u encounter collinearity

2. Lasso Reg.:  $\text{Lasso} = \text{Loss} + \lambda |m|$   
or L1 Reg.:  $\text{Penalty}$

$$y = 15 + 1.2x_1 + 20x_2 + 39x_3$$

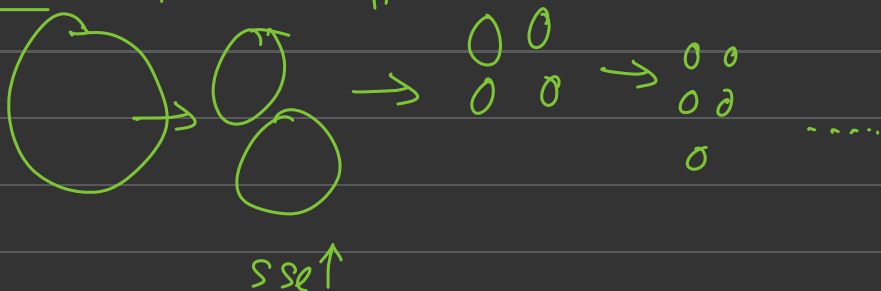
$\downarrow$  Lasso Reg.

$$y = 0.9 + \underline{0x_1 + 0x_2} + 5x_3$$

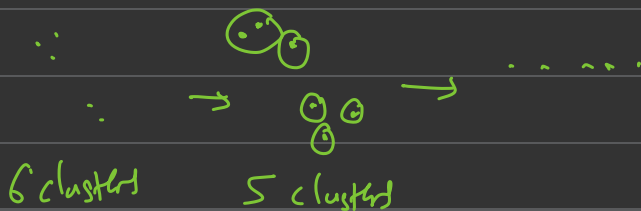
Use: Feature Selection, more robust to Outlier

## Hierarchical Cl.: Agglomerative & Divisive:

Divisive: Top-Down App.



Agglomerative: Bottom up app.



Boosting: Combine individual models into a strong learner

- Sequential learning
- Stump
- Avoid overfitting

1. Is XGBoost faster than RF?

→ XGBoost is a better performer

2. Tick adv. of XGBoost?

- a. Lot of Hyperparameters
- b. Can handle missing values
- c. Distributed computing (Parallel)

3. Disadv. of XGB:

- a. Sensitive to Outliers
- b. Manually create dummies.

4. Imp. hyperparameters of XGBoost?

<https://www.analyticsvidhya.com/blog/2016/03/complete-guide-parameter-tuning-xgboost-with-codes-python/>

↳ Only point 1 & 2

5. XGBoost vs LightGBM.

Slower

faster

powerful

relatively less powerful

6. How XGB handles miss. values?  
gblinear booster → fill miss. val. with 0.

## 7. XBoost vs AdaBoost?

1  
learns from  
errors in prediction in  
order to  
minimize error.

learns from previous stumps &  
predictions

---

## Model Selection:

1. Bias & Variance
2. Precision vs Accuracy
3. Bias Variance Tradeoff → Overfitting.
4. Cross Validation: KFold,  
Leave one out

<https://www.javatpoint.com/cross-validation-in-machine-learning>

<https://www.analyticsvidhya.com/blog/2021/05/4-ways-to-evaluate-your-machine-learning-model-cross-validation-techniques-with-python-code/>

## 5. Correlation

## 6. P-value Significance

## 7. VIF [0-5]

$$VIF_i = \frac{1}{1 - R_i^2}$$



