

I Basic Concepts

- Data Science =

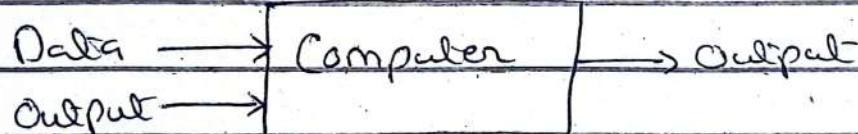
⇒ field of computer science that focuses on extracting meaningful insights, patterns and knowledge from data.

⇒ Consists of 5 stages:-

- 1) Data entry, extraction, acquisition.
- 2) Data cleansing, staging, processing, architecture
- 3) Data mining, classification, modelling, summarization
- 4) Predictive analysis, regression, text mining
- 5) Data visualization, reporting, decision making.

- Machine learning =

It is the idea that computer can learn from experience and examples.



It is like getting computers to program themselves
ML is automating the process of automation.

- ML Concepts and Algorithm =

Every ML algo has 3 components:

- 1) Representation of knowledge → decision trees, set of rules, graphical models, neural networks and others
- 2) Way to evaluate programs → accuracy, prediction posterior probability etc. [hypothesis]
- 3) Optimization of those candidate programs.

- Supervised Learning =

Data includes desired inputs and outputs. A model is prepared and through training process in which it is required to make predictions and is corrected when those predictions are wrong. This training process continues until model achieves desired outputs.

- Unsupervised Learning = Data does not include

outputs, But it is a model prepared by deducing structures present in input data

- Semi-Supervised =

Some desired outputs are provided but still model must learn the structures to organize data and make predictions.

- Deep learning =

is a ML technique that teaches computers to do what comes naturally to humans. It's similar to the neurons in the human brains.

It solves complex problems even when the data set is diverse, unstructured and inter-connected.

2

Python Concepts

- Lists =

my_list = [1, 2, 'saim', 4, 'Ali']

↓

→ It's a list that can have any type of data in it and can be accessed using indexing and slicing, etc.

→ It's mutable → can change its elements after its creation.

- Tuples =

my_tuple = (1, 2, 3, 'saim', 4)

↓

→ Tuples are same as lists but are immutable.

→ Have two functions:

i) Count → counts the repetition of certain data in the tuple.

ii) Index → tell the index position of the first occurrence of the desired/mentioned data.

- Dictionary = (mutable)

• Used keys instead of indexing.

my_Dictionary = {1: 'Ali', 2: "saim", 7: "ali"}

↓
key data.

• Functions =

→ print(my_Dictionary.get(1))

wid to get data at key '1'

→ .keys() = gives all the keys for a particular dictionary.

→ `.clear()` = empties the dictionary.

→ `.copy()` = copies one dictionary values (with keys) to another dictionary

→ `.fromkeys()` =

e.g.) `a = { 1: 5, 2: 10, 3: 15 }`.

`list_keys = [1, 2]`.

`b = a.fromkeys(list_keys, "saim")`



new dictionary will have 2 keys of 1, 2
from 'a' dictionary but both will have
value of => "saim"

→ `.items()` => shows a key and its value as
an item and all ~~the~~ items are returned

- Sets =

• Special case of lists where elements can only
be unique.

• Can't be indexed.

• Are mutable.

∴ `My_Set = {3, 1, 2, 3}` → no keys.

∴ `print(My_Set)` → prints only unique values and
no repeated ones.

- Functions =

→ `.add()` = adds a data at the end of set.

→ `.pop()` = removes the last data

→ `.discard()` = removes a specific data from
set. If that data is not in the set then it
does nothing.

3

Statistics for D.S =

- Type =

Data

Numerical

Categorical

Continuous

Discrete

Nominal

Ordinal

Interval

Ratio

↓
not True
zero

↓
True
zero

- Outliers = extreme value in data set

[55, 88, 89, 90]

↓
outlier.

- Mean = Average of data. → but greatly affected by

$$\left\{ \frac{x_1 + x_2 + x_3 + \dots + x_n}{n} \right\}$$

outliers.

- Median = Middle value

$$\left[\frac{n+1}{2} \right]^{th} \text{ value}$$

- Mode = Most occurred value in data set.

- Variance =

$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

$\therefore x_i$ = i^{th} data point

\bar{x} = mean.

n = no. of datapoint. in total.

- Standard Deviation =

$$S.D = \sqrt{\sigma^2} = \sqrt{\left(\frac{\sum (x_i - \bar{x})^2}{n-1} \right)}$$

- Coefficient of Variation =

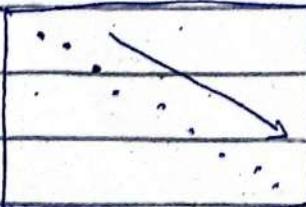
$$CV = \frac{S.D}{\text{mean}}$$

- Covariance = measures relationship b/w 2 variables

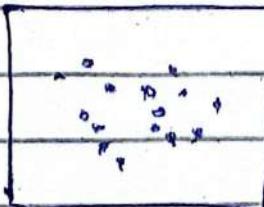
• Positive COV = If one var. increases, and also increases.

• Negative COV = If one var. increases, and var. decreases.

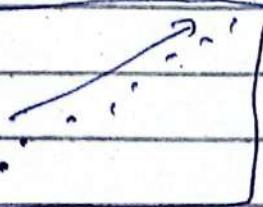
• Zero COV = no relation.



Negative COV



Zero COV



Positive COV

- Formula =

$$\text{Cov}(x,y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

\bar{x} mean of x
 \bar{y} mean of y

- Correlation tells us strength / type of relationship ^{that} 2 var. have.

$$\text{Correlation} = \frac{\text{Cov}(x,y)}{\sigma_x * \sigma_y}$$

$$\frac{\sigma_x * \sigma_y}{\text{S.D. of } x \quad \text{S.D. of } y}$$

- Inferential Statistics =

→ Distribution = It is a function that shows possible variable for a variable and the frequency of their occurrence

- Discrete Distribution

→ Uniform Dist. [11111]
 (single Dice Prob.)

→ Binomial Dist. [111111]
 (Double Dice Prob.)

- Continuous Distribution

→ Normal Dist. []

• have mean = mode = median.

• $N \sim (\mu, \sigma^2)$

Normal) \downarrow mean \downarrow Variance
 Distribution

• Standard Normal Distribution:

$$\text{mean} = 0$$

$$\text{s.D} = 1$$

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

z-score

- Central Limit Theorem:

Regardless of initial shape of population Dist., sampling Dist. will approximate to a normal Dist. As the sample size increases, sampling dist. will get narrower and more normal.

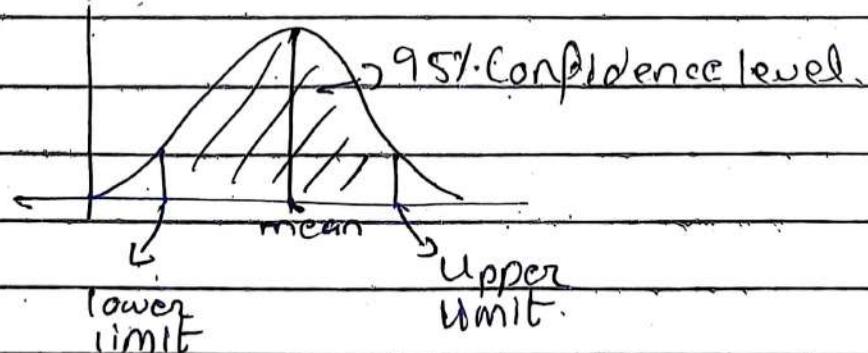
$\Rightarrow \text{CLT} \Rightarrow$

$\mu_{\bar{x}} = \mu$
$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$

→ Standard error.

• sample size ≥ 30 .

- Confidence Intervals:



H Probability

$$P(x) = \frac{\text{Preferred outcome}}{\text{Sample Space}}$$

- Expected Values: $P(x) = \frac{\text{Successful Trials}}{\text{All Trials.}}$

- Relative Frequency:

It is the number of times the event occurs ~~comes~~ divided by Total no. of trials

$$R.F = \frac{\text{Number of Success}}{\text{Total trials}}$$

- Theoretical Probability: a prob. that is expected like head has 50% chances and same goes for tails ; when tossing the coin.

- Hypothesis testing:

1) Make assumptions.

2) Take initial Position (H_0).

3) Determine alternate position (H_a)

4) Set acceptance criteria.

5) Conduct fact based tests.

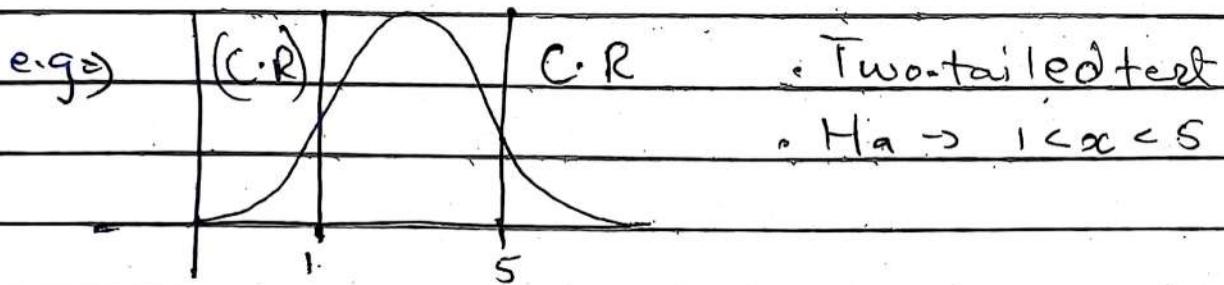
6) Evaluate Results

7) Reach one of the conclusion:

- favour alternate position

- fail to reject initial position.

- $H_0 \rightarrow$ Null Hypothesis.
- $H_a \rightarrow$ Alternate Hypothesis.
- Two-Tailed Test $\rightarrow (H_a)$ has $>$ and $<$ both.
 (H_a) has both directions.
- One-Tailed Test $\rightarrow (H_a)$ has either $>$ or $<$
 (H_a) has one direction
- Critical Region = (C.R)



5 Exploratory Data Analysis

- Unravel the data
 - To clean data
 - For feature engineering.

- Process = (Should be Quick efficient, Decisive.)

- No. of Observations. (rows)
 - No. of features (columns)
 - Data Type. (to identify type ML Problem)
 - Problem type identification (find/look for outliers).
 - Features that can be engineered.

- Find out Data Types in set =
in three groups.

① Numerical ② Categorical ③ Not sure

- Plot numerical Distributions:

- Histogram → Use it to study Distributions
 - look for: → outliers
→ Abnormal Distribution

- Plot Categorical Distributions-

- Use Bar Plots for it. \rightarrow for frequency of values.
- Box Plot as well.
- ~~use~~ use 2 ~~continuous~~ features to identify relation b/w them.

- Study Correlations-

• look for:

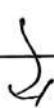
- strong relation b/w target variable and features.
- study relationship of other features with each other as well.

6 Data Cleaning =

→ Feature Scaling =

(2 Types)

- Normalized data
- Standardized Data.



Fixing range of data
b/w 2 values b/w
(0 or 1) or (-1 and 0)

Transform data to
have mean → 0 and
standard deviation → 1.

• Why do it =

→ To bring all features to a similar range.
so that no value dominates due to extreme
values.

→ To improve model performance and
training speed.

• When do it = (for algo's with rules)

→ When using ML Algo's that calculate
distance of data.

→ Neural Networks.

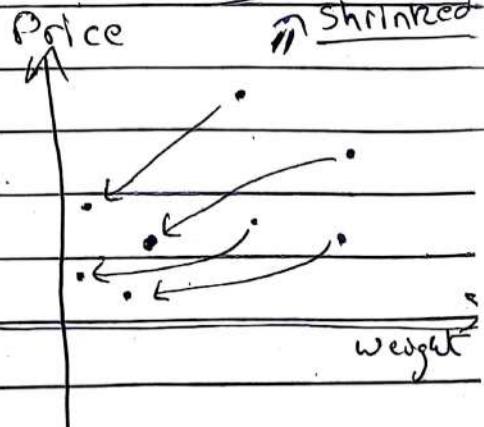
→ Knn.

→ PCA.

Values are
shrinked

• Min, Max Scalar =

$$x_{\text{new}} = \frac{x - x_{\text{min}}}{x_{\text{max}} - x_{\text{min}}}$$



- Standard Scalars

$$x_{\text{new}} = \frac{x - \mu}{\sigma}$$

→ Data Cleaning (in Videos) → [Go through the code as well.]

→ Feature Engineering:

- Study of features.
- Combining, removing or adding features is engineering.

→ We can:

- rearrange sets for clarity
- Add additional columns for more info
- Add dummy variable. to convert categorical values into numerical form

7] Linear Regression

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

\hat{y} is predicted value.

x_i is i th feature value.

n is no. of features.

θ_j is j th model parameter or j th weights.

$$\text{simplified} = \hat{y} = h_{\theta}(x) = \theta \cdot x$$

both are vectors.

It is hypothesis function using model parameter θ .

- $x \rightarrow$ feature vector
- $\theta \rightarrow$ parameter vector
- have dot product.

Cost Function

• Mean Squared Error \rightarrow MSE

$$\text{MSE}(x, \theta) = \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$

It is sum of square of all the (predicted values $-$ ^{their} actual values).

Optimize weights.

- Gradient Descent

: Optimization function.

: to optimize weights and minimize cost function.

→ Derivative of the cost function with each theta :-

$$\Rightarrow \frac{\partial}{\partial \theta_j} \text{MSE}(\theta) = \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_j^{(i)}$$

- Vectorization:

$$\nabla_{\theta} \text{MSE}(\theta) = \begin{bmatrix} \frac{\partial}{\partial \theta_0} \text{MSE}(\theta) \\ \frac{\partial}{\partial \theta_1} \text{MSE}(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n} \text{MSE}(\theta) \end{bmatrix} = \frac{1}{m} x^T (x\theta - y)$$

- Gradient Descent Step=

$$\Rightarrow \theta^{(\text{next step})} = \theta - \eta \nabla_{\theta} \text{MSE}(\theta)$$

• It is learning Rate (hyperparameter)

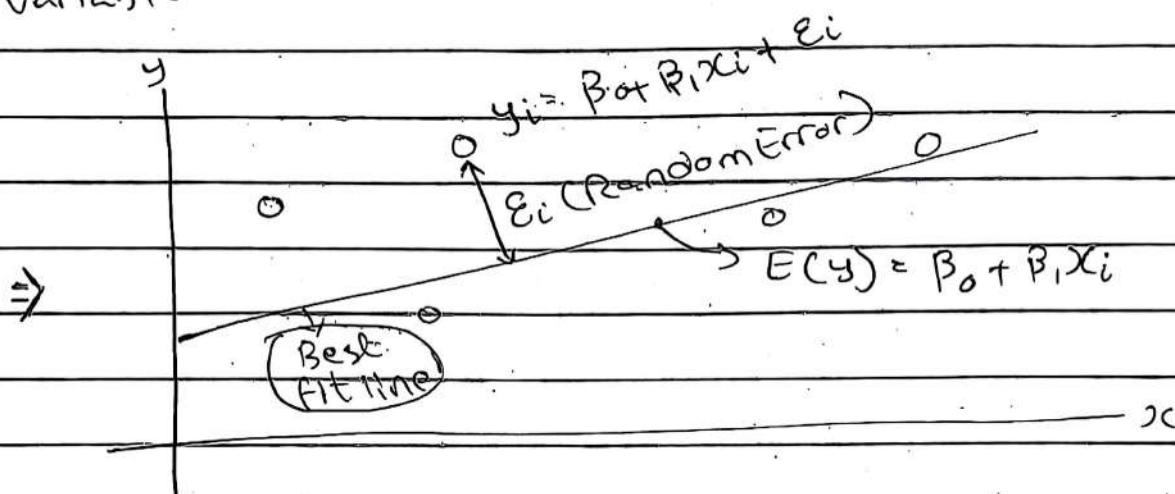
• tells size of Gradient Descent Step.

- Linear Regression Model -

→ Tells relationship b/w Variables that is a linear function.

$$\Rightarrow y_i = \beta_0 + \beta_1 x_i + \epsilon_i \rightarrow \begin{array}{l} \text{Population Slope} \\ \downarrow \\ \text{Random} \\ \text{Error} \end{array}$$

Dependent (Response) Variable Population y-intercept Independent (Explanatory) Variable.



- Least Squares = (LS)

- It helps to get the best fit line.

$$\Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n \hat{\epsilon}_i^2$$

- LS minimizes the sum of squared differences (errors) (SSE)

★ Coefficient Equations:

1) Prediction Equation:

$$\hat{y}_i = \hat{B}_0 - \hat{B}_1 x_i$$

2) Sample Slope =

$$\hat{B}_1 = \frac{SS_{xy}}{SS_{xx}} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

3) Sample y-intercept =

$$\hat{B}_0 = \bar{y} - \hat{B}_1 \bar{x}$$

$\Rightarrow \therefore$ Parameter Estimation Solution Table

x_i	y_i	x_i^2	y_i^2	$x_i y_i$
1	1	1	1	1
1	1	1	1	1
1	1	1	1	1
Σx_i	Σy_i	Σx_i^2	Σy_i^2	$\Sigma x_i y_i$

$$\therefore \hat{B}_1 = \frac{\sum_{i=1}^n x_i y_i - \left[\left(\sum_{i=1}^n x_i \right) \left(\sum_{i=1}^n y_i \right) \right]}{n}$$

$$\text{#} \sum_{i=1}^n (x_i)^2 = \left[\frac{\left(\sum_{i=1}^n x_i \right)^2}{n} \right]$$

8) \rightarrow Logistic Regression

(Prediction model)

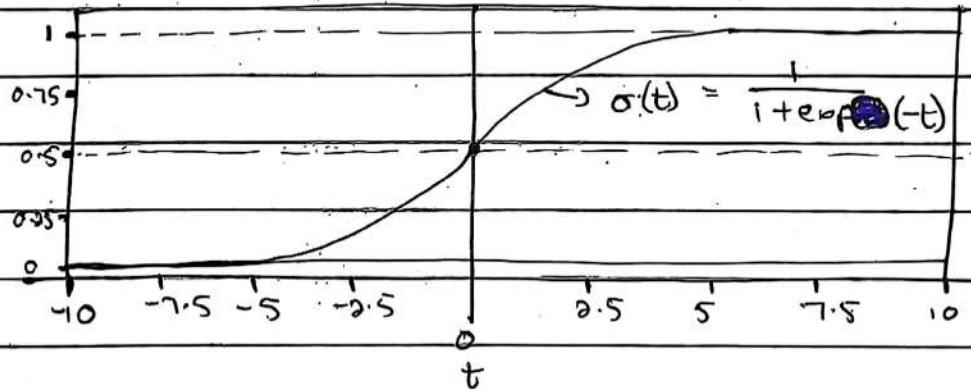
$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n.$$

$$\hat{p} = h_{\theta}(x) = \sigma(x^T \theta).$$

$$\sigma(t) = \frac{1}{1 + e^{-t}}.$$

Activation function.

• Estimating Probabilities =



$$\hat{y} = \begin{cases} 0 & \text{if } \hat{p} < 0.5 \\ 1 & \text{if } \hat{p} \geq 0.5 \end{cases}$$

• Cost function =

→ Cost function of one training instance $\Rightarrow c(\theta) = \begin{cases} -\log(\hat{p}) & \text{if } y = 1 \\ -\log(1-\hat{p}) & \text{if } y = 0. \end{cases}$

→ Cost function over

the entire training instance $\Rightarrow J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y(i) \log(\hat{p}(i)) + (1-y(i)) \cdot \log(1-\hat{p}(i))]$

Optimization Algorithms

→ Partial derivative with respect to all theta

$$\frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (\sigma(\theta^T x^{(i)}) - y^{(i)}) x_j^{(i)}$$

----->

→ K-Nearest Neighbours = (KNN) [9]

i) Parametric Models:

- Here we know which model exactly fits to the data like linear regression line.

$$y = w_1 x + e.$$

∴ Here parameters are w_1 's.

- Regression looks like a linear line.

ii) Non-Parametric Models:

- The data tells you what the regression should look like.

$$y = f(x) + e$$

(f(x) is decided by the data)

(It can be any function)

- Here no. of Parameters are decided by data
- $f(x)$ can be perfectly approximated by infinite parameter model.

⇒ EDA on IRIS Dataset

- For KNN section, Iris dataset will be used.

- It has four Attributes:

- ① sepal length in cm.
- ② sepal width in cm.
- ③ petal length in cm.
- ④ petal width in cm.

∴ sepal is outmost part
of flower, at the base
of the petals.

- For this course I will use 2 features for easier visualization:- sepal length and width.

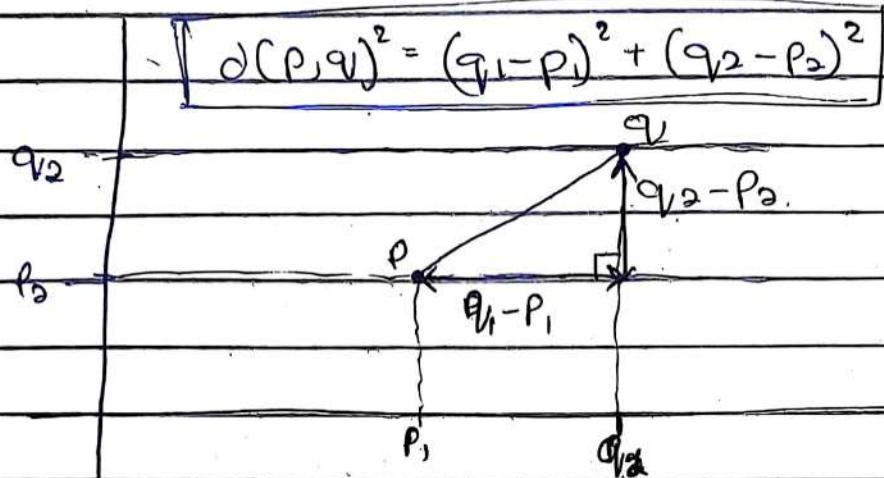
- Class for the flower can be:-

- ① Iris setosa,
- ② Iris versicolor
- ③ Iris virginica

} each has above attributes
mentioned

- KNN Algorithm Steps

① Calculate distance using ~~euclidean~~ euclidean distance:-



② Find the nearest neighbour using the python code.

③ Classify the point based ~~on~~ on the majority vote.

④ Combine all these three steps in one function to make prediction for each test point.

∴ Code for all these steps is done in the google colab.

→ Performing hyper-parameters tuning using ~~KNN~~
K-Fold cross Validation

- Never use test data for tuning → It is only for final evaluation.
- Use a validation set (or cross-validation) to tune hyper parameters.
- Dataset is usually split into train/validation /test [60%, 20%, 20%] → split]
- Cross-Validation: split data into k-groups; train on (k-1)groups; test on 1 group and repeat.
- This way every group gets used as a test set once.
- It helps in using full dataset effectively without wasting data.

→ Choosing the Distance Measurements=

- Use Euclidean distance if the dimensions are measuring similar properties (like width, height, depth)
- Use Manhattan distance if the dimensions are dissimilar (like age, weight, gender).

→ Need for normalization of data when using KNN=

- KNN uses distance-based similarity.
- So features with larger ranges would otherwise dominate
- Normalization ensures balanced influence of all features.

10

Decision Trees

→ We will do Exploratory Data Analysis on the Adult dataset

→ Objective is to find whether a person's income exceeds 50K\$/per or not. From that we will decide if that person is eligible for credit card or not.

→ Entropy =

It is used to measure impurity or disorder in the data. Tells how uncertain or mixed the data is.

$$\Rightarrow \text{Entropy } H(\pi) = - \sum \pi \log_2(\pi).$$

→ Information Gain (I_G)

$$\text{Condition Entropy } H(Y|X) = - \sum p(Y=y, X=x) \log p(Y=y, X=x) / p(X=x).$$

$$\Rightarrow I_G = H(Y) - H(Y|X)$$

[amount of uncertainty
before]

[amount of uncertainty]
after

- I_G measures how much uncertainty (impurity) is reduced after a split.

→ The Decision Tree ID3 Algorithm

- Steps =

- ① Find best feature
- ② Find best split for the best feature.
- ③ Partition classes based on step 1 and 2.

① - Partition Classes =

- We do this to separate data into smaller, purer groups, so that each group is more uniform.
- This reduces uncertainty, making decision rules simpler, and improves the accuracy of predictions.

② - Find the best split =

- 'Find the best Split' function is made in the code that find best place to split the data (on a chosen attribute) when building a decision tree.

• We do this because -

- ① Decision tree works by splitting data step by step into smaller groups.
- ② Good split makes group purer.
- ③ To measure how good a split is, we use information gain (how much uncertainty is reduced).

• Function tries all possible split values and picks the one with highest information gain.

• How it works :-

- ① Pick an attribute. → e.g "Age"
- ② check all possible values of it (like 20, 30, 40).
- ③ SPLIT dataset into left, right for each value
- ④ Calc. Information gain and choose value that gives best Info.Gain

③ - Find the best feature:-

- What is a feature:-

→ It is an attribute or column that describes each data point.

- What does function 'find best.feature' do:-

→ As we used this function in code, it looks at all the features and asks: "Which feature, if I split data on it, will give me best separation b/w classes?"

- Why do this:-

Because in the decision tree, first split is the most important.

→ If we choose wrong feature to split on, the tree will be messy and less accurate.

→ Important Hyper-Parameters in D.T =

① Criterion:-

- It's a function used to find/measure how good a split is.
- Gini → default, measure impurity.
- Entropy → uses information gain
- Log-loss → for probability.

② max_depth=

- It's the ^{no. of} max levels (depth) the tree can grow
- It's to prevent tree from growing too deep and overfitting.

③ Class-Weight =

- Assigns weight to classes (like in imbalanced sets)
- It helps the model pay more attention to minor class.

④ splitter =

- The strategy used to split nodes.
- 'best' is for choosing best split.
- 'random' is for adding randomness.
- Used for best random split to be chosen.

⑤ Min-Sample-Split =

- Min.no.of samples required to split an internal node.
→ Gini VS Entropy =

① Gini =

- Tells us how mixed up the classes are in a node
- If samples are from one class $\rightarrow Gini = 0$ (pure)
- If they're evenly split $\rightarrow Gini \approx 1$ (impure).
- Formula:-

$$Gini = 1 - \sum (P_i^2)$$

$\therefore P_i$ is proportion of samples belonging to class i .

- Range: 0 - 0.5 (for binary classification)
- Used cause it's faster to compute than entropy, work well in most cases.

② Entropy =

- Measures amount of disorder/uncertainty in data.
- Higher Entropy = More disorder.

Usage: • Provides more information-theory-based split.

- Slightly slower than Gini.

• Formula:-

$$-\sum (P_i \log_2(P_i)) = \text{Entropy}$$

→ Pruning =

- Pruning is a process of cutting down extra branches in a decision tree.
- To reduce overfitting and make tree more general and simpler.
- It is needed because D.T can become too deep and memorize training data (overfitting).
- Pruning removes unnecessary splits so the model can predict better on new data.

→ Types of Pruning:

(a) Pre-Pruning =

- Stop tree from growing too deep.
- Examples =
 - max-depth :- max-levels allowed
 - ~~min~~-sample-split :- min-samples needed to split a node.
 - min-sample-leaf :- min-samples in leafnode

(b) Post-Pruning =

- First grow a large tree, then remove weak branches.
- Uses a parameter called ccp-alpha (Complexity parameter).
- Higher ccp-alpha → more pruning → simpler tree.

→ How Is Post-Pruning Performed =

① Grow a full tree =

- First, allow the tree to grow very deep (almost no restrictions).
- This usually causes overfitting.

② Calculate "Cost-Complexity" for each subtree =

- For every possible branch (subtree), calculate :-

$$R_\alpha(T) = R(T) + \alpha \times (\text{number of leaves in } T)$$

- ∴ - $R(T)$: error of subtree.
- α : complexity penalty (ccp-alpha)
- more leaves: higher penalty.

③ Find weaker branches =

- Identify branches that don't reduce error much add complexity.
- These are the "weakest links"
- Note CC helps prune unnecessary branches to avoid overfitting.

↳ It means that model is "too smart" on the training data but "too dumb" on new data

④ Pruning them =

- Cut those weak branches (remove unnecessary splits).
- This reduces tree size without losing much accuracy.

⑤ Repeat until optimal tree is found =

- Keep pruning step by step.
- Stop when pruning starts hurting accuracy.

→ Pros of Decision Tree =

- Easy to understand and interpret
- Can work directly on raw data hence no need for data scaling/normalization.
- Handle both numerical and categorical data
- Require little data preprocessing.
- Fast training and predictions.
- Show which attributes matter most.
- Can capture non-linear relationships unlike KNN

→ Cons of Decision Tree =

- Overfitting: deep trees memorize training data and fail to generalize.
- Sensitive to [noisy data] → It can overfitting then create D.T.
- Biased with imbalanced dataset, so balance it first
- Small variation can result in different D.T.
- As branching goes deeper, you get exponentially less data.

III

Ensemble Methods + Learning

- Ensemble learning is combining many mediocre (weak) models into a one supermodel.

$$- \text{Accuracy} = 1 - \prod_{i=1}^n \text{Error}(i)$$

↓
Product of errors of base learners.

∴ note atleast one of the baselearners should be correct.

→ This formula has a problem that we are assuming all the predictions are completely independent of each other which is unrealistic.

→ Solution to this is bagging and boosting.

- Bootstrap Sampling = It is a statistical technique where you repeatedly draw samples with replacement from the original dataset to create many new "bootstrap samples" of the same size.

• Once we collect those samples, we calculate statistics on each bootstrap sample like mean, median, coefficient, regression etc

• We choose 'x' no. of samples from set and each of size 'n'. Then we perform statistic like mean for each and take average of all mean. Then decide or conclude that the avg. is the mean for whole population

- Bagging = \rightarrow [Bootstrap Aggregation]

- Its basic idea is to create different models using different bootstrapped training sets and average their prediction to produce final output of model.
- It is good for:
 - 1) Reducing variance of model. (due to replacement)
 - 2) Reducing overfitting
 - 3) Outliers are likely omitted in some of training bootstrap samples.
- Steps of bagging are:-

- 1) Generate B different bootstrap training datasets.
- 2) Train our ML Method on the b bootstrapped training set.
- 3) Final prediction:

$$\hat{f}(x) = \frac{1}{B} \cdot \sum \hat{f}^b(x)$$

- $\therefore - B$ = Total number of bootstrap models.
- $\hat{f}^b(x)$ = prediction from b^{th} model.
- Σ = sum over all B models

- Out-of-Bag Error (OOB) =

- OOB error is a way to estimate model performance in bagging without needing a separate validation set.
- When we create samples, there are high chances that some of the data points from original dataset are left out and not used in any of the samples. These are called 'out-of-bag' observations.
- So those observations are used as test set for the trained model; as they are unused and unseen by trained model.

- Random Forest=

- It is slightly modified bagging algorithm.
- It is basically collection of decorrelated DTrees
- For every tree:
 - 1) we perform bootstrapping.
 - 2) For each bootstrap, we choose attributes at random.
 - 3) We pick features randomly on each split and built tree only on them with current subsample
- Classification is done by taking a majority vote of classifications yielded by each tree in the forest after it classifies. An example

- Pros of Random Forest =

- Robust to outliers.
- Works well with non-linear data.
- Lower risk of overfitting.
- Runs efficiently on a larger dataset.

- Cons of Random Forest =

- biased while dealing with categorical variables.
- Slow Training if code is not optimized (just like our implementation took more than 3 minutes to execute the evaluation code, but for sklearn it took less than 10 seconds),
- Not suitable for linear method with a lot of sparse features.
- They are greedy algorithms.

- Boosting =

- It is a technique where we combine many weak models to build one strong model.

◦ Working =

1) Train the first model (DT)

2) Focus on the mistakes. When next tree is trained, it gives more weight to misclassified point, so it tries more to fix those mistakes.

3) Repeat the process as each tree learns from previous one and gets more better at reducing errors

4) At last, predictions are combined and finalized using classification regression.

- AdaBoost - (Adaptive Boost) :-

- One of boosting algorithms.
- It builds many weak classifiers and combine them into a strong classifier by giving more weight to the difficult-to-classify points.

→ Working =

- 1) Start with equal weights: Every training sample is given equal importance (weight).
- 2) Train the first weak learner: Fit a simple classifier (like decision stump) using these weights.
- 3) Check errors: see what point weak learner classified wrong.
- 4) Increase weight of misclassified points
- 5) Assign a weight to a learner
 - Better learners get higher weight.
 - Weaker learner still count, but less.
- 6) Repeat step 2-5.
- 7) Final Prediction = weighted majority vote (for classification) or weighted sum (for regression)

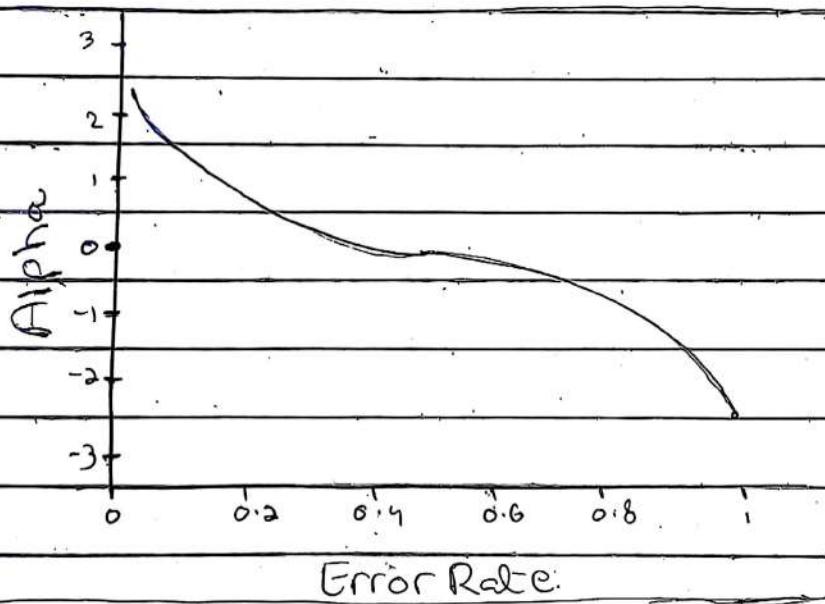
→ Decision Stumps =

- Are simplest form of decision trees.
- They have one split only (one level deep).
- Make decision based on just one feature and one threshold.
- AdaBoost: combines many stumps to form one strong model

→ AdaBoost Equation =

$$\Rightarrow \alpha_t = \frac{1}{2} \ln \left(\frac{1 - \text{Total Error}}{\text{Total Error}} \right)$$

- α is how much influence this stump would have in the final classification.
- Total Error is nothing, but the total num. of misclassifications for that training set divided by the training set size.



weight of next iteration weight of previous iteration

$$w_i = w_{i-1} \cdot e^{\pm \alpha}$$

- If a sample perform well then the sample weight decreases from what it was before; in the next iteration.
- It increases if it does not perform well

Just used to update weight of each sample after iterations.

→ AdaBoost Pseudocode

- Initially set uniform example weights.
- For each base learner:
 - 1) Train base learner ~~with~~ with a weighted sample.
 - 2) Test base learner on all data.
 - 3) Set learner weight with a weighted error.
 - 4) Set example weights based on ensemble predictions.

12 Support Vector Machine = (SVM)

- SVM =

- It is a supervised machine learning algorithm used for classification (mainly) and regression.
- It tries to find best boundary (hyperplane) that separates different classes in the data.

- How it works =

1) Plot data :- Imagine you have points of two classes (red = cats, blue = dogs)

2) Find the separating line :- SVM finds the line (in 2D), or (in 3D) plane, or hyperplane (in higher dimensions) that best separates the two classes.

3) SVM chooses the margin that leaves largest gap b/w classes.

4) Support vectors :- The point closest to the boundary are called support vectors as they define position

of the boundary.

- 5) For Non-linear Data:- If classes can't be separated by a straight line, SVM uses a kernel trick to map data into higher dimension where separation is possible

- Why to maximize this Margin?

- 1) A larger margin reduces the risk of overfitting, making the model perform well on unseen data.
- 2) With a wide margin, small variations or noise in the data are less likely to cause misclassification.
- 3) Points farther from the boundary are classified with higher certainty, improving reliability.
- 4) Margins should be maximized in order to prevent new data points from being misclassified.

- Hard vs Soft Margins

- Margin Equation:

- 1) Decision boundary:-

SVM tries to draw a line/plane to separate classes

Line is given by \Rightarrow

$$W \cdot x + b = 0$$

w =weight vector (define orientation of plane)

b =bias (shift of the plane)

x =input vector

- 2) The Margins:-

SVM creates two parallel boundaries on which the supporting vectors lie given by \Rightarrow

$$W \cdot x + b = +1 \rightarrow (\text{side of class } +1)$$

$$W \cdot x + b = -1 \rightarrow (\text{side of class } -1)$$

3) Margin Width:-

It is the distance b/w two boundaries

$$w \cdot x^+ + b = +1$$

$$w \cdot x^- + b = -1$$

$$\Rightarrow w(x^+ - x^-) = 2.$$

∴ divide both sides by $\|w\|$.

∴ $\|w\|$ is the size of the weight vector, which controls how tilted the plane is.

$$M \Rightarrow \frac{w(x^+ - x^-)}{\|w\|} = \frac{2}{\|w\|}$$

$$\Rightarrow \boxed{\frac{M}{\|w\|} = \frac{2}{\|w\|}}$$

→ To Maximize M we want to minimize $\|w\|$
so optimized eq is:

$$\Rightarrow \phi(w) = \frac{1}{2} w^T \cdot w$$

subject to :-

$$y_i(w^T x_i + b) \geq +1 \quad \text{if } y_i = +1$$

$$\text{and } y_i(w^T x_i + b) \leq -1 \quad \text{if } y_i = -1$$

→ Hard Margin:-

- It draws a line/hyperplane that perfectly separates two classes
- No mistakes are allowed - every point must be on correct side
- It requires complex kernel if some points do not fit within the boundary causing overfitting

→ Soft Margins

- More flexible.
- Separates the classes with biggest gap, but allows some points to be on wrong side.
- This makes it work better on real-world noisy data and reduces overfitting.

Now we modify equation :-

$$\Rightarrow \boxed{\frac{1}{2} \vec{w}^T \cdot \vec{w} + C \sum_{k=1}^n \epsilon_k}$$

∴ $\epsilon_k \rightarrow$ slack variable (define how much a point violates the margin)

controls overfitting. $C \rightarrow$ penalty parameter that balances

Decrease Variance margin size vs errors.

Wider Margin. \rightarrow larger $C \Rightarrow$ fewer mistakes allowed

Increase Variance Narrow margin \rightarrow smaller $C \Rightarrow$ more mistakes allowed

$$\Rightarrow \frac{1}{2} \vec{w}^T \cdot \vec{w} + C \sum_{k=1}^n \epsilon_k$$

Subject to:- $y_k(\vec{w}^T \vec{x}_k + b) \geq 1 - \epsilon_k, \epsilon_k \geq 0.$

- Kernel Trick = It is used so that SVM can handle non-linear data efficiently without going to high dimensions.

→ Formula for SVM optimization =

$$J(\alpha) = \sum_j \alpha_j - \frac{1}{2} \sum_i \sum_j y_i y_j \alpha_i \alpha_j \langle x_i^T, x_j \rangle$$

where :- x_i = data point i.

y_i = label of that data point (+1 or -1).

α_i = weight we solve for

→ $\langle x_i^T, x_j \rangle$ is dot product (similarity b/w 2 points)

$$\alpha_i \geq 0, \sum_i \alpha_i y_i = 0,$$

balances the line/hyperplane correctly b/w 2 classes so the solution can't tilt unfairly to one side

→ Now Kernel Trick =

- Imagine we have a function ' ψ ' that maps the data into another space.
- We want to optimize the dot product of two feature vectors.
- By transforming it to ψ , instead of computing $\langle x_i, x_j \rangle$, we have to compute :

$$\Rightarrow [(\psi(x_i) \cdot \psi(x_j))]$$

Here we apply transformation function to both data points and dot products to them

$$\Rightarrow K(x_i, x_j) = \varphi(x_i) \cdot \varphi(x_j)$$

↳ Kernel Function

Now we have :-

$$\min \alpha = -\sum \alpha_i + \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j K(\langle x_i, x_j \rangle).$$

- Kernel Types =

i) Polynomial Kernel =

- It is a type of Kernel Function used in SVM to allow non-linear boundaries. Instead of mapping data to higher dimension explicitly, it computes the similarity b/w 2 points as if they were transformed into a higher dimension.

- $K(x_i, x_j) = (\langle x_i, x_j \rangle + b)^d$

where - d is degree of polynomial.

- b is a constant that controls the influence of higher-order vs. lower-order terms.

ii) Gaussian Radial Basis function Kernel =

- RBF Kernel turns the similarity b/w 2 points into an exponential functions of their distance, making it possible for SVM to form very flexible, non-linear boundaries

- Kernels in SVM are basically similarity functions that tell "how close 2 data points are in meaning or position"

- Formula for similarity function is

$$\varphi(x_i, l) = \exp(\gamma|x_i - l|)$$

↓

It is feature map that maps input x_i into a new space using some landmark point l .

- Gamma HyperParameter (γ) =

It is a parameter that controls how much influence a single training point has.

→ Purpose :-

- 1) Control similarity spread:

- If γ is small, similarity curve will be wide.
Each point influences a large area. The decision boundary becomes smooth.
- If γ is large, curve is narrow. Decision boundary becomes tight and complex.

- 2) Small $\gamma \rightarrow$ high bias, low variance. Model may underfit.

- Large $\gamma \rightarrow$ low bias, high variance. Model may overfit.

- 3) Cosine Similarity Kernel = (Not important)

$$\varphi(x_i, x_j) = \frac{x_i^T x_j}{\|x_i\| \|x_j\|}$$

- To measure similarity of text documents

- SVM for Regression = (Multi-Class SVM)

We can use SVMs for regression as well instead of classification.

- The trick is to reverse the objective
- Instead of trying to fit the large possible margin b/w 2 classes while limiting violations,
- SVM regression tries to fit in as many instances as possible on the ^{margin} street while limiting margin violations.

13

K-Means

- Supervised learning:

discovers the patterns in the data that relate data attributes with a target (class) attribute.

- Unsupervised learning:

The data has no target attribute as it explores the data to find some intrinsic structure in them.

- Clustering=

- It is a technique for finding similarity groups in data called clusters.

- It is often called an unsupervised learning task as no class values denoting ~~a priori~~ grouping of the data instances are given

- Aspects of clustering=

- 1) Clustering Algorithm=

- Partitional Clustering :- Divides data into non-overlapping groups (like K-Means).

- Hierarchical Clustering :- Creates a tree of clusters, merging or splitting groups step by step.

- 2) • Distance (Similarity / Dissimilarity) Function :-

- Measures how close or far the data points are

③ Clustering Quality =

- Inter-Cluster Distance:- Should be large (clusters are far apart).
- Intra-Cluster Distance:- Should be small (points within a cluster are close).

④ Quality of clustering depend on algorithm, Distance function and the application

- K-Means Clustering =

- It is partition clustering algorithm.
 - Let $\{x_1, x_2, \dots, x_m\}$ be a set of data points (or instances) D by $\{x_i \in \mathbb{R}^r\}$ a vector in real-valued space and r is the number of attributes (dimensions) in the data.
 - The K-means algo. partitions the data (given) into K clusters. Each cluster has a ~~center~~ cluster center called centroid. K is specified by user.
- Steps for how it works =

① Decide how many groups you want the data to be divided into.

② Randomly pick K data points from dataset as the starting 'centers' of clusters.

③ For each data point, calculate its distance from every centroid and assign it to the closest one.

④ For each cluster, calculate avg. of all data points in the cluster and set that mean (avg.) as new centroid.

⑤ Keep reassigning points and updating centroids until centroids stop changing significantly.

⑥ Dataset now ~~is~~ is divided into K -clusters.

\rightarrow K-mean pseudocode = K-means(K, D).

i) Choose k data points as the initial centroids.

ii) repeat

for each data point $x \in D$ do.

compute distance from x to each centroid;

assign x to closest centroid.

end for

recompute the centroids using current cluster

~~memberships~~

until the stopping criterion is met.

\rightarrow Stopping/Convergence Criterion:

- It is no (or minimum) reassignments of data points to different clusters.
- No (or min.) change of centroids or,
- minimum decrease in sum of squared error (SSE)

$$\Rightarrow SSE = \sum_{j=1}^K \sum_{x \in C_j} \text{dist}(x, m_j)^2$$

- C_j is the j th cluster, m_j is the centroid of cluster C_j (the mean vector of all data points in C_j), and $\text{dist}(x, m_j)$ is the distance b/w data point x and centroid m_j .

\rightarrow Formula for Centroid of Cluster (Mean) = (In Euclidean Space)

$$m_j = \frac{1}{|C_j|} \sum_{x_i \in C_j} x_i$$

Diagram annotations:

- (mean) centroid
- Number of data points in cluster j
- Summation of all data points in cluster j
- Single datapoint in cluster j

- Disk Version of K-means= It is used when dataset is too large to fit into memory.

- Instead of loading whole data set, data is stored on disk.
- Small chunks of data are read from disk one at a time.
- For each batch, points are assigned to the nearest centroids and centroids are updated incrementally.
- This process continues until convergence

- Pros of K-means=

- Simple: easy to understand and implement.
- Efficient: as time complexity is $O(tKn)$, where 'n' is num. of data points, 'K' is number of clusters and 't' is num. of iterations.

- Cons of K-means=

- only applicable if mean is defined.
- For categorical data, centroid is represented by most frequent values.
- Sensitive to outliers,



→ We can solve this problem by either:

i) remove those data points that are far away from centroids.

ii) Do random sampling to reduce the chances of outliers to be selected.

- Clusters of Arbitrary Shapes:-

- Hyper-elliptical and hyper-spherical clusters are usually easy to represent.
- But irregular shaped ~~clusters~~ clusters are hard to represent hence they ~~are~~ may not be suitable/useful in some applications.
- So solution to this problem \Rightarrow Hierarchical clustering.

- Hierarchical Clustering:-

- It produces a nested sequence of clusters, a tree called Dendrogram.
- Types of Hierarchical Clustering:
 - a) Agglomerative (bottom up) Clustering:-
 - It builds tree from bottom level and merges most similar pair of clusters
 - It stops when all data points are merged into single cluster.
 - b) Divisive (top down) Clustering:- (theroot)
 - It starts with all data points in one cluster.
 - Splits root into a set of child clusters. Each child cluster is divided further recursively.
 - Stop when only singleton clusters of individual data points remain.

- Measuring the Distance b/w two clusters =

→ There are a few ways to do that:

① Single links:- $O(n^2)$.

- It measures the distance b/w two closest data points in 2 clusters, one data point from each.
- It can find arbitrary shaped clusters, but it may cause undesirable "chain effect by noisy points."

② Complete links:- $O(n^2 \log n)$.

- The distance is measured b/w the 2 furthest data points in the two clusters.
- It is sensitive to outliers because they are far away.

③ Average links:- $O(n^2 \log n)$.

- It is solution to complete-link and single-links (problems) as it finds the distances ^{b/w} all ^{data points} ~~all pairs~~. Then it take average of all those pair-wise distances.

④ Centroid Method:-

- Measures the distance b/w the centroids of 2 clusters.

- Distance Functions = (D·F)

datapoints

- Key to clustering (measure how similar/dissimilar are a.)
- There are numerous distance functions for different type of data
 - ① Numeric data
 - ② Nominal data
 - ③ Different specific ~~for~~ application.

- D·F for Numeric attributes =

- Most commonly used ones are:
 - Euclidean
 - Manhattan.
- Denoted by $\text{dist}(x_i, x_j) \rightarrow x_i, x_j$ are datapoints
- Minkowski Distance is generalized distance metric that includes above two ones as well.
It is defined as:-

$$D(x, y) = \left(\sum_{i=1}^n (x_i - y_i)^p \right)^{1/p}$$

where;

$\rightarrow x, y$ are i^{th} data points.

$\rightarrow p$ is order of norm (a parameter).

$\rightarrow p=1 \rightarrow$ Manhattan distance

$\rightarrow p=2 \rightarrow$ Euclidean distance.

$\rightarrow p=\infty \rightarrow$ Chebyshew distance

- When $p=1:$

$$\Rightarrow \text{dist}(x_i, y_i) = |x_{i1} - y_{i1}| + |x_{i2} - y_{i2}| + \dots + |x_{ir} - y_{ir}|$$

- when $p=2$:-

$$\Rightarrow \text{dist}(x_i, x_j) = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ir} - x_{jr})^2}$$

- Weighted euclidean distance :-

$$\Rightarrow \text{dist}(x_i, x_j) = \sqrt{w_1(x_{i1} - x_{j1})^2 + w_2(x_{i2} - x_{j2})^2 + \dots + w_r(x_{ir} - x_{jr})^2}$$

- Squared Euclidean Distance is another type where the under root is removed from normal Euclidean Dist. It is used cause it requires less computation than normal one and gives simillar results.

- Chebyshev Distance:-

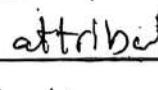
$$\text{dist}(x_i, x_j) = \max(|x_{i1} - x_{j1}|, |x_{i2} - x_{j2}|, |x_{ir} - x_{jr}|)$$

Chooses maximum distance among all the distances calculated.

Distance Function for Binary and Nominal Attributes

- Binary attributes; has two values or states but no ordering relationship, eg) Gender (Male or Female).

- We use confusion matrix to introduce the distance functions/measures.

- Confusion Matrix =  Compares 2 data points based on their binary attributes (0 or 1) to measure similarity or dissimilarity.

Data point (j)				- $a = \text{no. of attributes where both } i \text{ and } j \text{ are 1}$
		#	1	0
Data point (i)	1	a	b	- $b = \text{no. of attributes where } i = 1, j = 0$
	0	c	d	- $c = \text{no. of attributes where } i = 0, j = 1$ - $d = \text{no. of attributes where both } i \text{ and } j \text{ are 0.}$

- Once matrix is constructed, we calculate the similarity / dissimilarity b/w 2 data points.
- We use a, b, c, d from the confusion matrix to calculate Distance functions.
- Distance functions for this are :-

① → Simple matching Coefficient :-

$$\boxed{SMC = \frac{b+c}{a+b+c+d}}$$

↓

→ It is for symmetric Binary attributes where (0 and 1) states have equal importance and carry the same weight

② → Jaccard Coefficient :-

$$\boxed{JC = \frac{b+c}{a+b+c}}$$

↓

→ It is for Assymmetric Binary attributes where one of the states is more important than the other.

- By convention state 1 is more important state which is typically rare or infrequent state
- Nominal Attributes with more than two states.

$$dist(x_i|x_j) = \frac{r - q_j}{r}$$

where $\rightarrow r$ is total number of attributes

q_j is number of values that match in both x_i and x_j .

- Data Standardization =

- In the euclidean space, standardization of attributes is recommended so that all attributes can have equal impact on the computation of distances.
- It forces attributes to have a common value range.

- Interval-scaled Attributes =

- for numerical attributes
- Their values are real numbers following a linear scale.
- Eg Age of 10 and 20 have equal difference as 40-50 have.
- Two main approaches to standardize interval scaled attributes are.

① Range :-

$$\Rightarrow \text{range}(x_{if}) = \frac{x_{if} - \min(f)}{\max(f) - \min(f)}$$

where x_{if} = value of feature f for instance i.

$\min(f)$ = min. value of feature f in dataset.

$\max(f)$ = max. value of feature f in dataset.

② Z-scores :-

Transforms the attribute values so that they have a mean of zero and a mean absolute deviation of 1. The mean absolute deviation of feature f, denoted by s_f , is computed as follows:

$$\Rightarrow s_f = \sqrt{\frac{(|x_{1f} - mf| + |x_{2f} - mf| + \dots + |x_{nf} - mf|)^2}{n}}$$

where $\Rightarrow mf = \frac{1}{n} (x_{1f} + x_{2f} + \dots + x_{nf})$.

$$\Rightarrow \text{hence: } Z\text{-score} = \frac{x_{if} - m_f}{s_f}$$

- Ratio-Scaled Attributes:

- Numeric attributes, but unlike interval-scaled attributes, their scale are exponential.
- For example: Total amount of microorganisms that evolve in a time t is approximately given by Ae^{Bt} , where A and B are constants.
- Solution is log Transform: $(\log[x_{if}])$.
 $\Rightarrow \log(Ae^{Bt}) = \log(A) + Bt$.
- After log transformation, treat result as linear (interval-scaled attribute) and apply range of z-score.

- Nominal Attributes:

- We need to transform nominal attributes to numeric attributes.
- Transform nominal to binary attributes.
 - 1) The number of values of a nominal attribute is V .
 - 2) Create V binary attributes to represent them.
 - 3) If a data instance for a nominal attribute takes a particular value, the value of its binary attribute is set to '1'; otherwise '0'.

\Rightarrow like if we have Red, Blue, Green:
 if a Row(data instance) is Red then if it take Red then the value will be '1'; else if it is Green or blue in it, it will be '0' in these 2 cases.

- Ordinal Attributes:-

- It is like a nominal attribute, but its values have a numerical order like:
- Young, Middle age, and Old in Age attribute. They are ordered.
- We treat them as interval-scaled attributes.

- Mixed Attributes:-

- Our distance functions are either for all numeric or all nominal, etc.
- Practical data has different types:
 - Any subset of following:-
 - ① interval - scaled ② symmetric
 - ③ Assymmetric ④ ratio-scaled
 - ⑤ Ordinal ⑥ Nominal.
 - How to deal with it? :-
 - Convert mixed data types to one type:
 - Pick most common datatype and convert others to match it.
 - Converting to numerical:-
 - easy ones are =
 - ① Ordinal \rightarrow Numerical.
 - ② Ratio-scale \rightarrow log transform \rightarrow Numerical
 - ③ Binary \rightarrow 0 or 1.
 - Problematic =
 - Nominal (colors) \rightarrow Arbitrary numbers (Red=1, Blue=2)
 - Problem is we are creating fake ordering where it doesn't exist. Instead we can convert it to symmetric binary, which are then treated as numeric.

- Combining Individual Distances

$$\Rightarrow \text{dist}(x_i, x_j) = \frac{\sum_{f=1}^F s_{ij}^f d_{ij}^f}{\sum_{f=1}^F s_{ij}^f}$$

where $f = \text{attribute}$

$d_{ij}^f = \text{distance for attribute } f \text{ between } i \text{ and } j$

$s_{ij}^f = \text{weight for attribute } f \text{ (indicates importance)}$

\Rightarrow Allows different attributes to have importance levels and handles mixed data types by computing appropriate distance measures for each attribute type separately

- Cluster Evaluation

- The quality of clustering is a hard problem as it is hard evaluate because we do not know the correct clusters.

- Some methods used are:-

- ① Study centroids and spreads.

- ② Rules from a decision tree.

- ③ For test docs, one can read some docs in the clusters.

- Another way is Ground Truth:-

- We use some labeled data (for classification)

- Assumption: Each class is a cluster.

- After clustering, a confusion matrix is constructed.

- From matrix, we compute entropy, purity, precision, recall and F-score.

- let classes in data 'D' be $C = \{c_1, c_2, c_3, \dots, c_K\}$
- The clustering method produces K clusters, which divides 'D' into ' K ' disjoint subsets D_1, D_2, \dots, D_K
- Entropy for each cluster is measured by:

$$\text{Entropy}(D_i) = - \sum_{j=1}^K P_{D_i}(c_j) \log_2(P_{D_i}(c_j)).$$

(tells uncertainty in dataset) where • $P_{D_i}(c_j)$ is probability of outcome c_j occurring

- c_j is class in dataset.

- Purity measures the extent that a cluster contains only one class of data. It basically tell proportion of most frequent class in that cluster.

$$\text{Purity}(D_i) = \max_j(P_{D_i}(c_j)).$$

where, D_i = specific cluster.

\max = find largest value.

[14]

Principal Component Analysis (PCA)

- PCA is a dimensionality reduction technique used in ML and statistics. The main idea is:
- Data often has many features (dimensions), but some of them may be redundant or less important.
- PCA transforms the data, into a set of features called principal components which are combinations of the original features.
- These components are ordered by how much variance (information) they capture.
- Purpose of PCA:
 - ① Reduce dimensionality.
 - ② Remove noise/redundancy in data.
 - ③ Improve visualization by converting higher dimensions in 2D or 3D.
 - ④ Speed up algorithms by working with fewer features.

- PCA Drawbacks

- 1) PCA assumes that the variables are linearly correlated. If the correlations is not ~~not~~ linear, then PCA will not be efficient.
- 2) PCA performs lossy compression, which means that the information is lost when we discard insignificant components.
- 3) Since each principal component is a linear combination of original features, visualizations are not easy to interpret or relate to original features.

④ PCA results depend on the scale of data. If features are measured in different unit or have high ranges, PCA may give more weight to the larger scaled feature. Therefore, data should be standardized before applying PCA.

- Algorithm Steps:-

→ Suppose x_1, x_2, \dots, x_M are N-Dimensional vectors

→ So x is an $N \times M$ matrix.

→ N is no. of Dimensions and M is number of samples.

① Compute Data mean:-

$$\Rightarrow \bar{x} = \frac{1}{M} \sum_{i=1}^M x_i$$

② Subtract mean from each row of x (centering data):-

$$\Rightarrow \phi_i = x_i - \bar{x}$$

③ Form the matrix A :-

$$\Rightarrow A = [\phi_1, \phi_2, \dots, \phi_M] \quad (N \times M \text{ matrix})$$

④ Compute the covariance matrix C :-

$$\Rightarrow C = \frac{1}{M} \sum_{n=1}^M \phi_n \phi_n^T = A A^T \quad (N \times N \text{ matrix})$$

⑤ Compute the Eigen Values of C :-

$$\lambda_1 > \lambda_2 > \lambda_3 > \dots > \lambda_N$$

$$\Rightarrow \text{Eigen Values} \Rightarrow \det(C - \lambda I) = 0$$

⑥ Compute eigen Vector of C using those values :-

$$u_1, u_2, \dots, u_N$$

$$\Rightarrow \text{Eigen Vectors} = (C - \lambda I)(u) = 0$$

(E^N Matrix)

- Since C is symmetric, u_1, u_2, \dots, u_N form a basis
- Any vector $\phi_i = x_i - \bar{x}$ can be rewritten as a linear combination of the eigen vectors:-

$$\Rightarrow \boxed{\phi_i = b_1 u_1 + b_2 u_2 + \dots + b_N u_N = \sum_{i=1}^N b_i u_i}$$

$$\Rightarrow \therefore b = u^T \phi. \quad \therefore b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix}$$

- b contains the coefficients (weights) that tell you how much of each eigenvector is present in ϕ_i .
- So basically Principal Components = Eigen Vectors.

⑦ Dimensionality Reduction Step:- Keep only the terms corresponding to the K largest eigen values where $K < N$

$$\Rightarrow \boxed{\sum_{i=0}^K b_i u_i}$$

Note ➡ When we are done with computation of eigen values, arrange them in descending orders and reorder eigen vectors accordingly as it is ~~not~~ important for step 7, because we need most important components in the beginning.

- PCA SVD:

- In PCA Singular Value Decomposition, we directly decompose data matrix X (after mean-centering)
- Calculating covariance matrix by multiplying matrix by its transpose can cause loss of precision and numerical errors.
- Hence we use SVD technique:

$$\Rightarrow X = USV^T$$

where;

- U = left singular vectors (related to data in sample space. $(n \times n)$)
- S = Diagonal matrix of singular values (magnitude of variance in each direction). $(n \times p)$
- V = Right singular vectors, which are the principal components (directions of maximum variance). $(p \times p)$
- In short
 - 'U' gives directions of PCs.
 - 'S' gives importance of PCs.
 - 'V' gives coordinates of data in new PC space.

- Main Applications of PCA -

- 1) Image Compression.
- 2) Preprocessing Step for classification/regression/clustering tasks.
- 3) Facial Recognition: Eigen Faces
- 4) Data Visualization