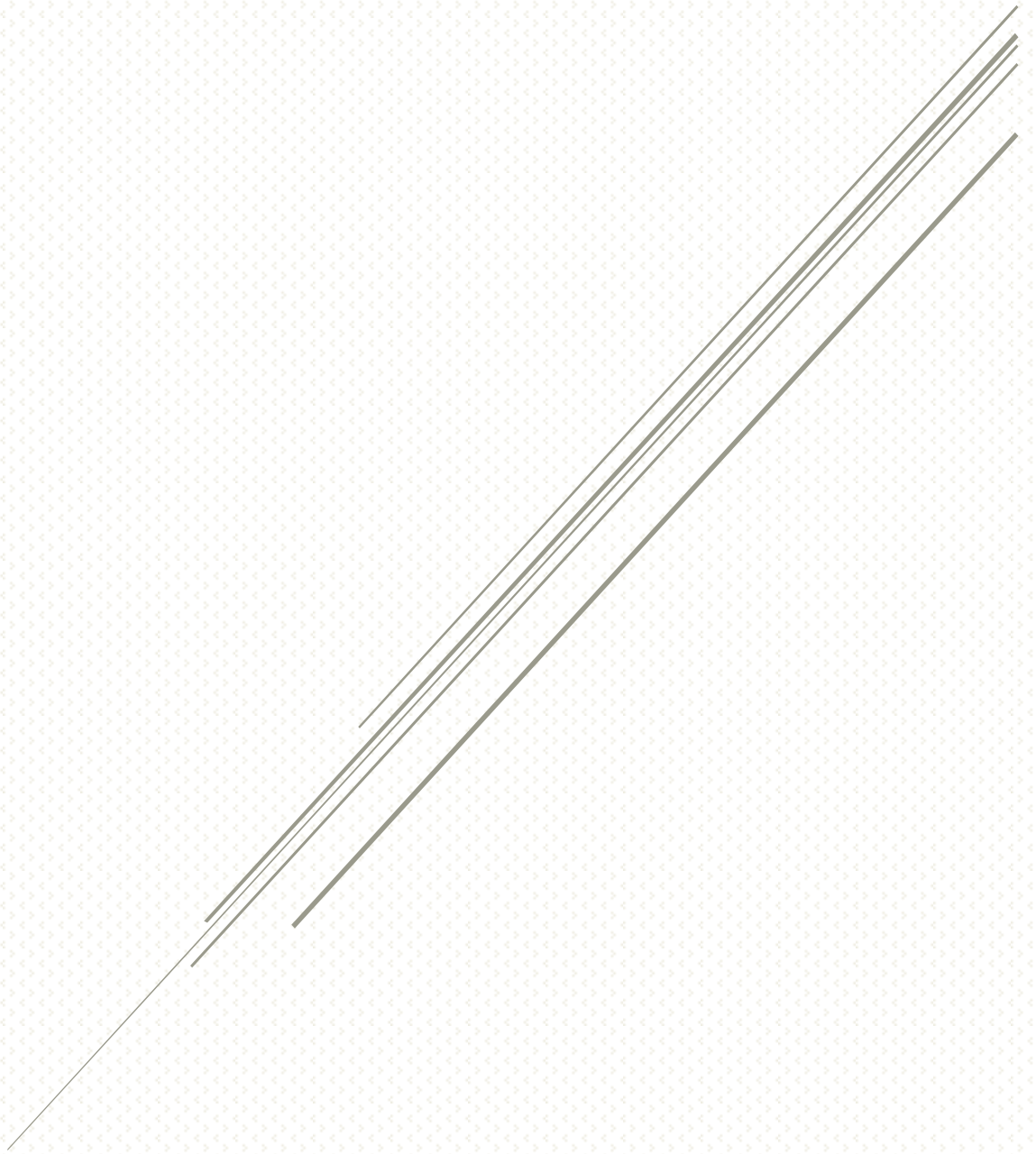


FODS QUESTION BANK SOLUTION

Jay Panchal



UNIT 4

Q1 What is Machine Learning ? What are the different types of Machine Learning Models?

Machine Learning (ML) is a subfield of artificial intelligence (AI) that focuses on developing algorithms and models that allow computers to learn from and make predictions or decisions based on data. In essence, it is a method of data analysis that automates analytical model building, enabling computers to improve their performance on a specific task over time without being explicitly programmed.

There are several different types of machine learning models, which can be categorized into three main groups:

1. Supervised Learning:

- In supervised learning, the model is trained on a labeled dataset, which means the input data is paired with the correct output. The goal is to learn a mapping from inputs to outputs.
- Types of supervised learning models include:
 - Regression: Used for predicting continuous numeric values, such as predicting house prices.
 - Classification: Used for categorizing data into predefined classes, such as spam detection or image recognition.

2. Unsupervised Learning:

- Unsupervised learning models are trained on data without explicit labels. The system tries to find patterns or structures within the data.
- Types of unsupervised learning models include:
 - Clustering: Used to group similar data points together, such as customer segmentation.
 - Dimensionality Reduction: Used to reduce the number of features in the data while retaining its essential information, such as Principal Component Analysis (PCA).

3. Reinforcement Learning:

- In reinforcement learning, an agent interacts with an environment and learns to make a sequence of decisions to maximize a reward. It's commonly used in scenarios where the model needs to learn to make a series of decisions.
- The agent takes actions, and the environment provides feedback in the form of rewards or penalties. The agent's goal is to learn a policy that maximizes the expected long-term reward.

Q2 Briefly discuss KNN algorithm with its pros and cons.

The k-Nearest Neighbors (KNN) algorithm is a simple yet effective machine learning algorithm used for both classification and regression tasks. It works by finding the k training examples (data points) in the dataset that are closest to a given test point and using their labels (in the case of classification) or values (in the case of regression) to make predictions.

How KNN Works:

The K-NN working can be explained on the basis of the below algorithm:

Step-1: Select the number K of the neighbors

Step-2: Calculate the Euclidean distance of **K number of neighbors**

Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.

Step-4: Among these k neighbors, count the number of the data points in each category.

Step-5: Assign the new data points to that category for which the number of the neighbor is maximum.

Step-6: Our model is ready.

Advantages

It is simple to implement.

It is robust to the noisy training data

It can be more effective if the training data is large.

Applications

Handwriting detection

Image Recognition

Video Recognition

Disadvantages

Always needs to determine the value of K which may be complex some time.

The computation cost is high because of calculating the distance between the data points for all the training samples.

Q3 Illustrate the steps of KNN algorithm with proper example

The K-NN working can be explained on the basis of the below algorithm:

Step-1: Select the number K of the neighbors

Step-2: Calculate the Euclidean distance of **K number of neighbors**

Step-3: Take the K nearest neighbors as per the calculated Euclidean distance.

Step-4: Among these k neighbors, count the number of the data points in each category.

Step-5: Assign the new data points to that category for which the number of the neighbor is maximum.

Step-6: Our model is ready.

KNN Example

The engineers resulting dataset is given below. Find the result of X when Maths=6 and Cs= 8. Consider value of k as 3.

Maths	CS	Result
4	3	Fail
6	7	Pass
7	8	Pass
5	5	Fail
8	8	Pass

Solution : According to Euclidean distance

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

Case 1 : 5.38

Case 2 : 1

Case 3: 1

Case 4 : 3.16

Case 5 : 2

Result for Case 2,3 and 5 is Pass . The pass result probability is greater than fail .
So it can be concluded that the result of X is Pass

Q4 Describe KNN & ANN algorithm. And also the reason why we used NN algorithm in data analysis

K-Nearest Neighbors (KNN) Algorithm:

K-Nearest Neighbors is a supervised machine learning algorithm used for classification and regression tasks. It makes predictions based on the similarity of a data point to its K nearest neighbors in a training dataset. Here's how it works:

1. Given a data point to be classified or predicted, the algorithm calculates the distance (often Euclidean distance) between that point and all other points in the training dataset.
2. It then selects the K data points with the shortest distances to the given point.
3. For classification, it assigns the class that occurs most frequently among the K neighbors as the predicted class for the data point. For regression, it calculates the average (or weighted average) of the K neighbors' values.
4. The value of K is a hyperparameter chosen by the user and impacts the algorithm's performance.

Average Nearest Neighbor Algorithm:

The Average Nearest Neighbor (ANN) is a method used to assess the spatial distribution of points in a dataset. It's often used in spatial statistics and geographic information systems (GIS). ANN measures whether points are clustered or dispersed in space and whether there is any regular pattern or randomness. Here's how it works:

1. For each point in the dataset, the algorithm calculates the distance to its nearest neighbor, forming a list of these nearest neighbor distances.
2. The algorithm then calculates the average (mean) of these nearest neighbor distances.
3. It also calculates the expected average nearest neighbor distance under the assumption of complete spatial randomness (CSR).
4. By comparing the observed average nearest neighbor distance to the expected CSR value, you can assess whether the data exhibit clustering, dispersion, or randomness.

Why We Use Nearest Neighbor Algorithms in Data Analysis:

Nearest neighbor algorithms, including K-Nearest Neighbors and Average Nearest Neighbor, are used in data analysis for several reasons:

1. Classification and Prediction: K-Nearest Neighbors is a versatile algorithm for classification and regression tasks, making it applicable to a wide range of problems.
2. Pattern Recognition: It can identify patterns and relationships in data, especially when data points are not easily separable using linear methods.

3. Spatial Analysis: Average Nearest Neighbor is valuable in spatial analysis to assess the spatial distribution and clustering of points, which is essential in various fields like ecology, epidemiology, and urban planning.

4. Simplicity: KNN is easy to understand and implement, making it a good choice for quick exploratory data analysis and benchmarking.

5. No Assumptions: KNN makes minimal assumptions about the data distribution, which can be beneficial when the data's nature is unknown or complex.

6. Versatility: Nearest neighbor algorithms can be used in both classification and regression problems, making them versatile tools in data analysis.

Q5 Describe NN analysis. Explain difference between ANN and KNN algorithm.

Nearest Neighbor Analysis" typically refers to a group of techniques used in spatial statistics and geographic information systems (GIS) to assess the spatial distribution and clustering of points in a dataset. Two common methods within this analysis are the "K-Nearest Neighbor (KNN)" and "Average Nearest Neighbor (ANN)" algorithms

Aspect	K-Nearest Neighbor (KNN)	Average Nearest Neighbor (ANN)
Purpose	Assess spatial distribution and determine if points are clustered or dispersed.	Assess spatial distribution and determine if points exhibit clustering, dispersion, or randomness.
Calculation of Neighbors	For each point, identify and consider the K nearest neighbors.	For each point, calculate the distance to its single nearest neighbor.
Result Interpretation	Provides insights into specific clustering or dispersion patterns based on the chosen value of K.	Provides a single statistic (average nearest neighbor distance) to assess overall spatial distribution.
Use Cases	Useful for analyzing and quantifying patterns in specific geographic areas with a focus on the number of neighbors considered.	Provides a general measure of spatial distribution and is often used to assess global spatial patterns.
Sensitivity to Parameters	Sensitive to the choice of K; results can vary with different values of K.	Not sensitive to a parameter; calculates the average nearest neighbor distance without the need to specify K.

Q6 Express how decision making is useful for visualization of data. explain with steps

Decision making in data visualization is a crucial process that helps transform raw data into meaningful insights and actionable information. Effective decision-making in data visualization involves a series of steps to select appropriate visualization techniques and design choices to represent and convey data in a clear and insightful manner. Here are the steps involved:

****1. Define the Objective:****

- Start by understanding the purpose of your data visualization. What message or insight are you trying to convey? Define the specific objectives of your visualization, such as showing trends, comparisons, distributions, anomalies, or relationships in the data.

****2. Know Your Audience:****

- Consider your target audience. Are they experts in the field, executives, or the general public? Tailor your visualization to match their expertise level and preferences. Knowing your audience helps determine the level of detail, complexity, and the choice of visual elements.

****3. Select the Right Data:****

- Choose the data that is most relevant to your objectives. Avoid overloading the visualization with unnecessary information. Ensure that the data is accurate, up-to-date, and well-prepared.

****4. Choose the Visualization Type:****

- Select an appropriate visualization type that best conveys the information you want to communicate. Common types include bar charts, line charts, scatter plots, heatmaps, pie charts, and more. The choice depends on the data's nature and the insights you want to highlight.

****5. Design Principles:****

- Apply fundamental design principles, such as simplicity, clarity, and consistency. Use a clean and uncluttered layout with clear labeling, proper use of colors, and easily understandable legends. Ensure that the design elements support the message.

****6. Data Mapping:****

- Map the data to visual variables. Assign data attributes to visual properties like position, size, color, and shape. For example, you might map a value to the height of a bar in a bar chart or the color of a data point in a scatter plot.

****7. Provide Context:****

- Add contextual information to help viewers understand the data. Include titles, labels, legends, and explanations where needed. Contextual information should be concise and relevant.

****8. Interactivity:****

- Consider adding interactive elements when appropriate. These might include tooltips, zooming, filtering, and highlighting. Interactivity can enhance user engagement and exploration of the data.

****9. Test and Iterate:****

- Test the visualization with a sample audience or colleagues to gather feedback. Iteratively refine the visualization based on the feedback received. Ensure that the visualization effectively conveys the intended message.

****10. Document and Share:****

- Document the visualization's design choices, data sources, and any assumptions made. Share the visualization with the intended audience through reports, presentations, dashboards, or web applications.

****11. Monitor and Update:****

- Keep the visualization up to date if the underlying data changes. Monitor the impact of the visualization on decision-making processes and be prepared to make updates as necessary.

Q7 what is regression? explain different types of regression

Regression is a supervised machine learning technique used to predict a continuous numerical value or a continuous outcome based on one or more input features. It is widely used in various fields, such as economics, finance, engineering, and data science, to model and understand the relationships between variables. In regression analysis, the goal is to fit a mathematical model to the data that can make predictions and capture the underlying patterns.

There are several types of regression, each designed for specific scenarios and types of data. Here are some of the most common types of regression:

1. **Linear Regression:**

- Linear regression is the simplest form of regression and is used to model the relationship between a dependent variable and one or more independent variables. It assumes a linear relationship between the variables. There are two main types of linear regression:

- Simple Linear Regression: Involves one dependent variable and one independent variable.
- Multiple Linear Regression: Involves one dependent variable and two or more independent variables.

2. **Polynomial Regression:**

- Polynomial regression is an extension of linear regression that allows for modeling non-linear relationships. It involves fitting a polynomial equation to the data. For example, a quadratic regression might use a model like $y = a + bx + cx^2$.

3. **Ridge Regression and Lasso Regression:**

- Ridge and Lasso regression are variations of linear regression that incorporate regularization to prevent overfitting. Ridge regression adds a penalty term for the sum of squared coefficients, while Lasso regression adds a penalty term for the absolute values of coefficients.

4. **Logistic Regression:**

- Despite the name, logistic regression is used for binary classification, not regression. It models the probability that a given input belongs to a particular class. It's commonly used in fields such as epidemiology and machine learning.

5. **Poisson Regression:**

- Poisson regression is used when the dependent variable represents count data, such as the number of events occurring within a fixed interval. It is especially useful in fields like epidemiology and social sciences.

6. **Time Series Regression:**

- Time series regression is used when the data has a temporal component. It models how the dependent variable changes over time, making it suitable for forecasting and analyzing trends in time series data.

7. **Support Vector Regression (SVR):**

- SVR is a regression technique that uses support vector machines to find a hyperplane that best fits the data. It's effective for non-linear and high-dimensional data.

8. **Decision Tree Regression:**

- Decision tree regression involves using decision trees to model the relationships between variables. It can handle non-linear data and is interpretable.

9. **Random Forest Regression:**

- Random forest regression is an ensemble method that combines multiple decision trees to improve prediction accuracy and reduce overfitting.

10. **Gradient Boosting Regression:**

- Gradient boosting regression is another ensemble method that combines the predictions of multiple weak models (typically decision trees) in a sequential manner to produce a strong predictive model.

Q8 Difference between Simple Regression and Multiple Linear Regression

Aspect	Simple Linear Regression	Multiple Linear Regression
Number of Independent Variables	1 (single predictor)	2 or more (multiple predictors)
Equation	$Y = a + bX$	$Y = a + b_1x_1 + b_2x_2 + \dots + b_nx_n$
Purpose	Modeling the relationship between a dependent variable and a single independent variable.	Modeling the relationship between a dependent variable and multiple independent variables, accounting for multiple factors.
Example	Predicting a student's exam score based on hours studied.	Predicting house sale price based on factors like bedrooms, square footage, location, and age of the house.
Complexity	Less complex, suitable for simple relationships.	More complex, suitable for cases involving multiple influencing factors.

Q9 Explain the different Evaluation Metrics used for Classification ?

Evaluation metrics are essential in assessing the performance of classification models in machine learning. These metrics help you understand how well a model is doing in terms of its ability to classify data points into different classes. Here are some common evaluation metrics used for classification:

1. **Accuracy:**

- **Formula:** $(TP + TN) / (TP + TN + FP + FN)$

- **Description:** Accuracy is a basic measure of the overall correctness of a classifier. It calculates the ratio of correctly predicted instances (true positives and true negatives) to the total number of instances. However, it may not be suitable for imbalanced datasets.

2. **Precision (Positive Predictive Value):**

- **Formula:** $TP / (TP + FP)$

- **Description:** Precision is the ratio of correctly predicted positive instances (true positives) to the total number of instances predicted as positive. It measures the model's ability to avoid false positives.

3. **Recall (Sensitivity or True Positive Rate):**

- **Formula:** $TP / (TP + FN)$

- **Description:** Recall calculates the ratio of correctly predicted positive instances (true positives) to the total number of actual positive instances. It measures the model's ability to find all positive instances.

4. **F1-Score:**

- **Formula:** $2 * (Precision * Recall) / (Precision + Recall)$

- **Description:** The F1-Score is the harmonic mean of precision and recall. It provides a balance between precision and recall, making it useful when both false positives and false negatives are important.

5. **Specificity (True Negative Rate):**

- **Formula:** $TN / (TN + FP)$

- **Description:** Specificity calculates the ratio of correctly predicted negative instances (true negatives) to the total number of actual negative instances. It is useful when you want to focus on the model's ability to identify negative instances.

6. **False Positive Rate (FPR):**

- **Formula:** $FP / (FP + TN)$

- **Description:** FPR calculates the ratio of incorrectly predicted positive instances (false positives) to the total number of actual negative instances. It complements specificity.

7. **Area Under the Receiver Operating Characteristic (ROC AUC):**

- **Description:** ROC AUC measures the area under the Receiver Operating Characteristic curve. It provides a measure of the model's ability to distinguish between positive and negative classes. A higher ROC AUC indicates a better model.

8. **Matthews Correlation Coefficient (MCC):**

- **Formula:** $(TP * TN - FP * FN) / \sqrt{((TP + FP) * (TP + FN) * (TN + FP) * (TN + FN))}$

- **Description:** MCC is a balanced measure of a model's performance, especially useful when dealing with imbalanced datasets. It takes into account all four confusion matrix values.

9. **Log Loss (Cross-Entropy Loss):**

- **Formula:** $-\sum (y_i * \log(p_i) + (1 - y_i) * \log(1 - p_i))$
- **Description:** Log Loss quantifies the uncertainty of a classifier's predictions. It penalizes incorrect predictions more than correct ones, making it a good choice for probabilistic models.

10. **Cohen's Kappa:**

- **Formula:** $(\text{Observed Agreement} - \text{Expected Agreement}) / (1 - \text{Expected Agreement})$
- **Description:** Cohen's Kappa measures the agreement between the actual and predicted classifications while accounting for agreement that may occur by chance. It is useful for evaluating inter-rater agreement and multi-class classification.

Q10 Explain the different Evaluation Metrics used for Regression ?

Evaluation metrics for regression models are used to measure the performance and accuracy of these models in predicting continuous numerical values. The choice of the appropriate metric depends on the specific problem, the characteristics of the data, and the goals of the analysis. Here are some common evaluation metrics used for regression:

1. **Mean Absolute Error (MAE):**

- **Formula:** $MAE = (1 / n) * \sum |y_i - \hat{y}_i|$

- **Description:** MAE calculates the average absolute difference between the actual values (y_i) and the predicted values (\hat{y}_i). It provides a measure of the model's accuracy, with higher values indicating greater errors.

2. **Mean Squared Error (MSE):**

- **Formula:** $MSE = (1 / n) * \sum (y_i - \hat{y}_i)^2$

- **Description:** MSE calculates the average squared difference between the actual values and the predicted values. It is more sensitive to outliers than MAE, as it penalizes larger errors more.

3. **Root Mean Squared Error (RMSE):**

- **Formula:** $RMSE = \sqrt{MSE}$

- **Description:** RMSE is the square root of the MSE. It has the same units as the target variable, making it easier to interpret. It is also more sensitive to outliers.

4. **Mean Absolute Percentage Error (MAPE):**

- **Formula:** $MAPE = (1 / n) * \sum |(y_i - \hat{y}_i) / y_i| * 100$

- **Description:** MAPE measures the percentage difference between actual and predicted values. It is useful when you want to understand the relative error. However, it can produce infinite values when the actual values are zero.

5. **Coefficient of Determination (R-squared or R^2):**

- **Formula:** $R^2 = 1 - (\sum (y_i - \hat{y}_i)^2 / \sum (y_i - \bar{y})^2)$

- **Description:** R-squared measures the proportion of the variance in the dependent variable explained by the model. It ranges from 0 to 1, with higher values indicating a better fit. It is commonly used to assess the goodness of fit of the regression model.

6. **Adjusted R-squared:**

- **Formula:** $Adjusted\ R^2 = 1 - [(1 - R^2) * (n - 1) / (n - k - 1)]$

- **Description:** Adjusted R-squared adjusts the R-squared value based on the number of predictors (k) and the sample size (n). It accounts for model complexity and prevents overfitting.

7. **Mean Bias Deviation (MBD):**

- **Formula:** $MBD = (1 / n) * \sum (y_i - \hat{y}_i)$

- **Description:** MBD measures the average bias or deviation of the model's predictions from the actual values. Positive values indicate overestimation, and negative values indicate underestimation.

8. **Median Absolute Error (MedAE):**

- **Description:** MedAE calculates the median of the absolute differences between actual and predicted values. It is robust to outliers compared to mean-based metrics.

9. **RMSLE (Root Mean Squared Logarithmic Error):**

- **Formula:** $\text{RMSLE} = \sqrt{(1/n) * \sum (\log(y_i + 1) - \log(\hat{y}_i + 1))^2}$
- **Description:** RMSLE is commonly used when the target variable or predictions have a skewed distribution. It calculates the RMSE of the logarithm of the predicted and actual values.

10. **Relative Squared Error (RSE):**

- **Formula:** $\text{RSE} = \sum (y_i - \hat{y}_i)^2 / \sum (y_i - \bar{y})^2$
- **Description:** RSE is a ratio-based metric that measures the fraction of the total variance in the data explained by the model. It ranges from 0 to 1, with higher values indicating better fit.

Q11 What is Confusion Metric ? Explain all its metrics

A confusion matrix is a table used in classification to evaluate the performance of a machine learning model. It provides a detailed breakdown of the model's predictions and actual outcomes. From the confusion matrix, several metrics can be calculated to assess the model's performance. The main metrics derived from a confusion matrix are True Positives (TP), True Negatives (TN), False Positives (FP), and False Negatives (FN). These metrics are used to calculate various performance measures, including accuracy, precision, recall, and F1-score.

Here's a breakdown of the key metrics and how they are calculated:

- **True Positives (TP):** These are the cases where the model correctly predicted the positive class (e.g., correctly identified a disease as present).
- **True Negatives (TN):** These are the cases where the model correctly predicted the negative class (e.g., correctly identified a disease as absent).
- **False Positives (FP):** These are the cases where the model incorrectly predicted the positive class when it should have predicted the negative class (e.g., a false alarm for a disease that's not present).
- **False Negatives (FN):** These are the cases where the model incorrectly predicted the negative class when it should have predicted the positive class (e.g., failing to identify a disease that's actually present).

Now, let's explore the metrics derived from these values:

1. **Accuracy:**
 - **Formula:** $(TP + TN) / (TP + TN + FP + FN)$
 - **Description:** Accuracy is the proportion of correctly predicted instances (both positive and negative) out of all instances. It provides an overall measure of model performance but may not be suitable for imbalanced datasets.
2. **Precision (Positive Predictive Value):**
 - **Formula:** $TP / (TP + FP)$
 - **Description:** Precision is the ratio of correctly predicted positive instances to all instances predicted as positive. It measures the model's ability to avoid false positives.
3. **Recall (Sensitivity or True Positive Rate):**
 - **Formula:** $TP / (TP + FN)$
 - **Description:** Recall calculates the ratio of correctly predicted positive instances to all actual positive instances. It measures the model's ability to find all positive instances.
4. **F1-Score:**
 - **Formula:** $2 * (Precision * Recall) / (Precision + Recall)$
 - **Description:** The F1-Score is the harmonic mean of precision and recall. It provides a balance between precision and recall, making it useful when both false positives and false negatives are important.

5. **Specificity (True Negative Rate):**

- **Formula:** $TN / (TN + FP)$
- **Description:** Specificity calculates the ratio of correctly predicted negative instances to all actual negative instances. It is useful when you want to focus on the model's ability to identify negative instances.

6. **False Positive Rate (FPR):**

- **Formula:** $FP / (FP + TN)$
- **Description:** FPR calculates the ratio of incorrectly predicted positive instances (false positives) to all actual negative instances. It complements specificity.

7. **Area Under the Receiver Operating Characteristic (ROC AUC):**

- **Description:** ROC AUC measures the area under the Receiver Operating Characteristic curve. It provides a measure of the model's ability to distinguish between positive and negative classes. A higher ROC AUC indicates a better model.

Q12 While predicting malignancy of tumor of a set of patients using a classification model, following are the data recorded: (a) Correct predictions – 15 malignant, 75 benign (b) Incorrect predictions – 3 malignant, 7 benign Calculate the model accuracy, error rate, sensitivity, precision, and F measure of the model.

To calculate the model accuracy, error rate, sensitivity, precision, and F-measure for predicting malignancy of tumors, you can use the values provided. Let's break down the calculations:

Given data:

- Correct predictions for malignant tumors (True Positives, TP) = 15
- Correct predictions for benign tumors (True Negatives, TN) = 75
- Incorrect predictions for malignant tumors (False Negatives, FN) = 3
- Incorrect predictions for benign tumors (False Positives, FP) = 7

Now, we can calculate the evaluation metrics:

1. ****Model Accuracy:****

- Model Accuracy = $(TP + TN) / (TP + TN + FP + FN)$
- Model Accuracy = $(15 + 75) / (15 + 75 + 7 + 3)$
- Model Accuracy = $90 / 100$
- Model Accuracy = 0.90 or 90%

2. ****Error Rate:****

- Error Rate = $1 - \text{Model Accuracy}$
- Error Rate = $1 - 0.90$
- Error Rate = 0.10 or 10%

3. ****Sensitivity (True Positive Rate or Recall):****

- Sensitivity = $TP / (TP + FN)$
- Sensitivity = $15 / (15 + 3)$
- Sensitivity = $15 / 18$
- Sensitivity ≈ 0.8333 or 83.33%

4. ****Precision (Positive Predictive Value):****

- Precision = $TP / (TP + FP)$
- Precision = $15 / (15 + 7)$
- Precision = $15 / 22$
- Precision ≈ 0.6818 or 68.18%

5. ****F-Measure (F1-Score):****

- F-Measure = $2 * (\text{Precision} * \text{Sensitivity}) / (\text{Precision} + \text{Sensitivity})$
- F-Measure = $2 * (0.6818 * 0.8333) / (0.6818 + 0.8333)$
- F-Measure ≈ 0.7511 or 75.11%

So, for the given data, the model has the following performance metrics:

- | | |
|--------------------------------|---------------------|
| - Accuracy: 90% | - Error Rate: 10% |
| - Sensitivity (Recall): 83.33% | - Precision: 68.18% |
| - F-Measure: 75.11% | |

UNIT 5

Q1 What is Clustering Algorithm ? What are its types ? List the applications of clustering algorithms.

“Clustering is the task of dividing the population or data points into a number of groups such that data points in the same groups are more similar to other data points in the same group than those in other groups”.

A way of grouping the data points into different clusters, consisting of similar data points. The objects with the possible similarities remain in a group that has less or no similarities with another group."

It does it by finding some similar patterns in the unlabelled dataset such as shape, size, color, behavior, etc., and divides them as per the presence and absence of those similar patterns.

It is an unsupervised learning method, hence no supervision is provided to the algorithm, and it deals with the unlabeled dataset.

After applying this clustering technique, each cluster or group is provided with a cluster-ID. ML system can use this id to simplify the processing of large and complex datasets.

Types of clustering Algorithms

- Partitioning Clustering
- Density-Based Clustering
- Hierarchical Clustering
- Fuzzy Clustering

Applications

The clustering technique can be widely used in various tasks. Some most common uses of this technique are:

- Market Segmentation
- Statistical data analysis
- Social network analysis
- Image segmentation
- Anomaly detection, etc.

Q2. What is Partitioning Based clustering ? What are its types ?

Partitioning-based clustering is a type of clustering algorithm that divides data points into non-overlapping clusters. The key idea is to partition the data into distinct groups where each data point belongs to only one cluster. Two common types of partitioning-based clustering algorithms are:

1. **K-Means Clustering:** K-Means divides data into K clusters, minimizing the variance within each cluster. It's widely used for its simplicity and effectiveness.
2. **K-Medoids Clustering:** Similar to K-Means but selects medoids (data points closest to the center) as cluster representatives. It's robust to outliers and noise.

These algorithms are used in applications like customer segmentation, image compression, and anomaly detection.

Q3 What is the K-Means Clustering Algorithm ? What are the steps involved in the K-Means Clustering algorithm ?

****K-Means Clustering**** is a popular partitioning-based clustering algorithm used to group data points into K clusters. Here are the steps involved:

1. ****Initialization:**** Choose K initial cluster centroids (representative points) randomly or strategically.
2. ****Assignment:**** Assign each data point to the nearest cluster centroid based on a distance measure (typically Euclidean distance).
3. ****Update:**** Recalculate the cluster centroids by computing the mean (average) of all data points assigned to each cluster.
4. ****Repeat:**** Repeat the assignment and update steps iteratively until convergence, which occurs when cluster assignments no longer change significantly or after a fixed number of iterations.

The final centroids represent the centers of the K clusters. K-Means aims to minimize the variance within clusters and maximize the separation between them. It's an iterative algorithm that can provide compact, well-separated clusters for various applications like image compression, customer segmentation, and more.

Q4 What is the K-Medoid Clustering Algorithm ? What are the steps involved in the K-Medoid Clustering algorithm ?

The **K-Medoid Clustering Algorithm** is a variant of K-Means that aims to partition data into K clusters. The key difference is that K-Medoids uses medoids, which are actual data points, as cluster representatives (centroids) instead of the mean of the points in a cluster. Here are the steps involved in K-Medoid Clustering:

1. **Initialization:** Choose K initial medoids, which can be any data points in the dataset.
2. **Assignment:** For each data point, assign it to the nearest medoid based on a distance measure (typically using the Euclidean distance).
3. **Update:** Recalculate the total dissimilarity (or cost) by considering the sum of distances between each data point and its assigned medoid.
4. **Swap:** Try swapping a non-medoid point with a medoid point and compute the new total dissimilarity. If the dissimilarity decreases, perform the swap. Repeat this step for all data points to find the optimal set of medoids.
5. **Repeat:** Iterate the assignment, update, and swap steps until convergence or a predefined stopping criterion is met.

K-Medoids aims to find medoids that minimize the total dissimilarity within clusters. This makes it robust to outliers, as medoids are less affected by extreme values compared to cluster means. K-Medoids is used in various applications, including image analysis and disease clustering.

Q5 What are the advantages and disadvantages of K Means Clustering Algorithm ?

****Advantages of K-Means Clustering Algorithm:****

1. ****Simplicity:**** K-Means is easy to understand and implement, making it suitable for quick data analysis.
2. ****Efficiency:**** It's computationally efficient and works well with large datasets, making it applicable to a wide range of data sizes.
3. ****Scalability:**** K-Means can handle high-dimensional data and works efficiently even when the number of features is large.
4. ****Clustering:**** It forms tight, well-separated clusters, which is helpful when data is well-clustered.
5. ****Versatility:**** K-Means is useful for both data exploration and preprocessing as well as part of a larger data analysis pipeline.

****Disadvantages of K-Means Clustering Algorithm:****

1. ****Sensitivity to Initial Centroids:**** The choice of initial cluster centroids can affect the outcome. Different initializations can lead to different results.
2. ****Clustering Non-Globular Shapes:**** K-Means may not perform well when clusters have complex shapes, elongated structures, or varying densities.
3. ****Number of Clusters (K):**** Selecting the optimal number of clusters (K) is challenging and often requires domain knowledge or iterative testing.
4. ****Impact of Outliers:**** K-Means is sensitive to outliers, and outliers can significantly affect the cluster centroids and results.
5. ****Equal Variance Assumption:**** It assumes that clusters have equal variance, which may not hold in real-world datasets.
6. ****Local Optima:**** K-Means can converge to local optima, leading to suboptimal solutions, especially in complex data spaces.

To mitigate some of these disadvantages, variations of K-Means, such as K-Medoids and hierarchical clustering, can be considered. The choice of clustering algorithm depends on the specific characteristics of the data and the analysis goals.

Q6 What is DBSCAN Clustering? What are the steps involved in DBSCAN Algorithm

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a density-based clustering algorithm that groups data points into clusters based on the density of their neighbors. DBSCAN is particularly effective at identifying clusters of arbitrary shapes and handling noisy data. Here are the steps involved in the DBSCAN algorithm:

1. **Parameter Selection:** Choose two key parameters: epsilon (ϵ) and minimum points (MinPts). Epsilon defines the maximum distance for a data point to be considered a neighbor of another point. MinPts specifies the minimum number of data points required to form a dense region.
2. **Core Points:** Identify core points. A data point is a core point if there are at least MinPts data points (including itself) within a distance of ϵ .
3. **Directly Density-Reachable:** Determine which data points are "directly density-reachable" from a core point. A data point, A, is directly density-reachable from another point, B, if B is a core point, and A is within ϵ distance from B.
4. **Density-Connected:** Expand clusters by defining density-connected points. Two data points, A and B, are density-connected if there is a core point, C, such that both A and B are directly density-reachable from C.
5. **Cluster Formation:** Form clusters by connecting data points based on their density-connected relationships. Each cluster is composed of all data points that are density-connected to at least one core point.
6. **Noise Points:** Assign any data points that are not core points or not part of any cluster as noise or outliers.

The DBSCAN algorithm can identify clusters of varying shapes and densities while distinguishing noise points. It's widely used in applications such as geographic information systems (GIS), image analysis, and anomaly detection.

Q7 What is Hierarchical Clustering ? Explain its types.

Hierarchical Clustering is a family of clustering algorithms that build a hierarchical representation of data points in the form of a tree-like structure, known as a dendrogram. Hierarchical clustering is commonly used to group data points into a hierarchy of clusters based on their similarity or dissimilarity. There are two primary types of hierarchical clustering:

1. **Agglomerative Hierarchical Clustering:**

- **Bottom-Up Approach:** Agglomerative clustering starts with each data point as a separate cluster and then recursively merges clusters that are closest to each other, moving up the hierarchy.
- **Steps:** It involves the following steps:
 - Start with each data point as an individual cluster.
 - Iteratively merge the two closest clusters until only one cluster remains.
 - Build a dendrogram that represents the hierarchy of cluster mergers.
- **Dendrogram Interpretation:** The dendrogram provides a visual representation of the merging process, allowing users to choose the desired number of clusters by cutting the dendrogram at a specific height.

2. **Divisive Hierarchical Clustering:**

- **Top-Down Approach:** Divisive clustering starts with all data points in a single cluster and recursively divides it into smaller clusters, moving down the hierarchy.
- **Steps:** It involves the following steps:
 - Start with all data points in a single cluster.
 - Iteratively split a cluster into two or more subclusters until each data point is in a separate cluster.
 - Build a dendrogram that represents the hierarchy of cluster divisions.
- **Dendrogram Interpretation:** The dendrogram shows the hierarchy of cluster divisions, and users can select clusters at various levels of the hierarchy.

Applications:

- **Agglomerative hierarchical clustering** is often preferred because it is more commonly used and computationally efficient.
- **Divisive hierarchical clustering** is less common but can be useful in certain scenarios, such as when you have a strong prior belief in the number of clusters you want to achieve.

Both types of hierarchical clustering provide a natural way to understand the hierarchy and relationships between clusters, making them valuable for exploratory data analysis and data visualization.

Q8 Cluster the following eight points (with (x,y) representing locations) into three clusters:

A1(2, 10), A2(2, 5), A3(8, 4), A4(5, 8), A5(7, 5), A6(6, 4), A7(1, 2), A8(4, 9)

Initial cluster centers are: A1(2, 10), A4(5, 8) and A7(1, 2).

The distance function between two points $a = (x_1, y_1)$ and $b = (x_2, y_2)$ is defined as

$P(a, b) = |x_2 - x_1| + |y_2 - y_1|$

Use K-Means Algorithm to find the three cluster centers after the second iteration.

In this K-Means clustering example, we'll perform the second iteration of the K-Means algorithm with the provided initial cluster centers and distance function. We'll assign each point to the nearest cluster center, calculate new cluster centers, and repeat the process.

Given initial cluster centers:

1. A1(2, 10)
2. A4(5, 8)
3. A7(1, 2)

Points:

- A2(2, 5)
- A3(8, 4)
- A5(7, 5)
- A6(6, 4)
- A8(4, 9)

Using the distance function $P(a, b) = |x_2 - x_1| + |y_2 - y_1|$, let's calculate the distances of each point from the initial cluster centers.

1. Distance of each point from A1(2, 10):

- A2: $|2 - 2| + |5 - 10| = 5$
- A3: $|8 - 2| + |4 - 10| = 10$
- A4: $|5 - 2| + |8 - 10| = 3$
- A5: $|7 - 2| + |5 - 10| = 10$
- A6: $|6 - 2| + |4 - 10| = 10$
- A7: $|1 - 2| + |2 - 10| = 9$
- A8: $|4 - 2| + |9 - 10| = 2$

2. Distance of each point from A4(5, 8):

- A2: $|2 - 5| + |5 - 8| = 6$

- A3: $|8 - 5| + |4 - 8| = 7$

- A4: $|5 - 5| + |8 - 8| = 0$

- A5: $|7 - 5| + |5 - 8| = 5$

- A6: $|6 - 5| + |4 - 8| = 5$

- A7: $|1 - 5| + |2 - 8| = 10$

- A8: $|4 - 5| + |9 - 8| = 2$

3. Distance of each point from A7(1, 2):

- A2: $|2 - 1| + |5 - 2| = 4$

- A3: $|8 - 1| + |4 - 2| = 9$

- A4: $|5 - 1| + |8 - 2| = 10$

- A5: $|7 - 1| + |5 - 2| = 9$

- A6: $|6 - 1| + |4 - 2| = 7$

- A7: $|1 - 1| + |2 - 2| = 0$

- A8: $|4 - 1| + |9 - 2| = 10$

Based on the distances, we can reassign the points to their nearest cluster centers:

- Cluster 1 (Center A1): A2(2, 5) and A4(5, 8)

- Cluster 2 (Center A4): A3(8, 4), A5(7, 5), A6(6, 4)

- Cluster 3 (Center A7): A7(1, 2), A8(4, 9)

Now, calculate the new cluster centers by finding the means of the points within each cluster:

New Center for Cluster 1 (A1): $(2+5)/2, (5+8)/2 = (3.5, 6.5)$

New Center for Cluster 2 (A4): $(8+7+6)/3, (4+5+4)/3 = (7, 4.33)$

New Center for Cluster 3 (A7): $(1+4)/2, (2+9)/2 = (2.5, 5.5)$

UNIT 6

Q1 What is Data Visualization ? Why is it preferred in Organisation?

Data visualization is the representation of data in a visual or graphical format to help people understand and interpret complex information more easily. It involves creating charts, graphs, maps, and other visual elements to display data patterns, trends, and insights. Data visualization is an essential part of data analysis, and it serves several crucial purposes:

Why Data Visualization is Preferred in Organizations:

- Simplifies Complex Data:** Data visualization simplifies complex datasets by representing them graphically, making it easier for individuals at all levels of an organization to understand the information.
- Enhances Decision-Making:** Visualizations provide a clear and intuitive way to present data, enabling decision-makers to quickly grasp important insights and make informed decisions.
- Identifies Trends and Patterns:** Charts and graphs reveal trends, patterns, and anomalies in data that might be challenging to detect in raw numbers or text, helping organizations respond to changing conditions.
- Improves Communication:** Visualization aids in effective communication. It allows teams to share insights and findings with colleagues, clients, or stakeholders more easily and persuasively.
- Facilitates Exploration:** Interactive data visualizations enable users to explore data, drill down into details, and ask questions to discover deeper insights, fostering a culture of curiosity and exploration.
- Supports Storytelling:** Visualizations can be used to tell compelling data-driven stories that engage and inform an audience, whether it's in a boardroom presentation or a marketing report.
- Monitors Key Performance Indicators (KPIs):** Dashboards with data visualizations make it simple to track and monitor key performance indicators and performance metrics in real-time.
- Promotes Data-Driven Culture:** By making data accessible and understandable, data visualization encourages organizations to become more data-driven and evidence-based in their decision-making.
- Cross-Functional Understanding:** Visualizations can be understood by various departments and roles within an organization, bridging the gap between technical and non-technical staff.
- Time and Cost Efficiency:** Visualizations can save time and reduce costs by making data analysis more efficient and enabling timely responses to issues or opportunities.

Q2 What are the different plots drawn using matplotlib ? Write the syntax of each

Matplotlib is a popular Python library for creating a wide range of plots and visualizations. Here are some common types of plots you can create using Matplotlib, along with simplified syntax examples:

1. **Line Plot:**

- Syntax:

```
```python
import matplotlib.pyplot as plt
plt.plot(x, y)
plt.show()
```
```

2. **Scatter Plot:**

- Syntax:

```
```python
import matplotlib.pyplot as plt
plt.scatter(x, y)
plt.show()
```
```

3. **Bar Chart:**

- Syntax:

```
```python
import matplotlib.pyplot as plt
plt.bar(x, y)
plt.show()
```
```

4. **Histogram:**

- Syntax:

```
```python
import matplotlib.pyplot as plt
plt.hist(data, bins)
plt.show()
```
```

5. **Box Plot:**

- Syntax:

```
```python
import matplotlib.pyplot as plt
plt.boxplot(data)
plt.show()
```
```

6. **Pie Chart:**

- Syntax:

```
```python
import matplotlib.pyplot as plt
plt.pie(sizes, labels)
plt.show()
```
```

7. ****Heatmap:****

- Syntax:

```
```python
import matplotlib.pyplot as plt
import seaborn as sns # Often used in combination with Matplotlib
sns.heatmap(data)
plt.show()
```
```

8. ****Violin Plot:****

- Syntax:

```
```python
import matplotlib.pyplot as plt
import seaborn as sns # Often used in combination with Matplotlib
sns.violinplot(data)
plt.show()
```
```

9. ****Area Plot:****

- Syntax:

```
```python
import matplotlib.pyplot as plt
plt.fill_between(x, y1, y2)
plt.show()
```
```

10. ****3D Plot:****

- Syntax (for a simple 3D scatter plot):

```
```python
from mpl_toolkits.mplot3d import Axes3D
import matplotlib.pyplot as plt
fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
ax.scatter(x, y, z)
plt.show()
```
```

11. ****Error Bars (Error Bar Plot):****

- Syntax

```
import matplotlib.pyplot as plt
plt.errorbar(x, y, xerr, yerr)
plt.show()
```

Q3 Write a python code to draw a pie chart using python's library.

```
import matplotlib.pyplot as plt

# Data for the pie chart
labels = ['Category A', 'Category B', 'Category C', 'Category D']
sizes = [25, 30, 15, 30] # Sizes or proportions for each category

# Explode a slice (if desired)
explode = (0, 0.1, 0, 0) # "explode" the 2nd slice (Category B)

# Create a pie chart
plt.pie(sizes, labels=labels, autopct='%1.1f%%', startangle=90, shadow=True, explode=explode)

# Add a title
plt.title("Distribution of Categories")

# Show the pie chart
plt.show()
```

Q4 Write the python code to draw histogram for following set of data using python's library.

Pop [22,55,62,45,21,22,34,42,42,4,2,8] bins = [1,10,20,30,40,50]

```
import matplotlib.pyplot as plt
```

```
# Data and bin boundaries
```

```
data = [22, 55, 62, 45, 21, 22, 34, 42, 42, 4, 2, 8]
```

```
bins = [1, 10, 20, 30, 40, 50]
```

```
# Create a histogram
```

```
plt.hist(data, bins=bins, edgecolor='k', alpha=0.7)
```

```
# Add labels and a title
```

```
plt.xlabel('Value')
```

```
plt.ylabel('Frequency')
```

```
plt.title('Histogram of Data')
```

```
# Show the histogram
```

```
plt.show()
```