

~~Ques:~~ To find the mean, median, mode and standard deviation, using their characteristic draw a box plot.

Software used: google colab, python libraries

Theory:

- ~~Central Tendency~~: The measure of central ~~tendency~~ is a single value that attempts to describe a set of data by identifying the center point within the set of data.
The mean, median, mode are all valid measures of central tendency.

- Mean: It is the sum of all values in a dataset divided by the no. of values in the dataset.

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

- Median: The median is the middle score for a set of data that has been arranged in order of magnitude.

$$\text{median} = \begin{cases} n[m/2] & , n = \text{odd} \\ n\left[\frac{m/2 + (n+1)/2}{2}\right] & , n = \text{even} \end{cases}$$

→ Mode: Most frequent score in a dataset of anything. In a histogram, the highest box is mode.

→ Variance: It is the measure of how far the set of data are dispersed out from the mean or avg. value (σ^2).

→ Standard deviation: It also represents the dispersion from mean is the square root of variance.

BOX PLOT

- Min. score: lowest score excluding outlier
- First Quartile (Q_1): 25% of score fall below the lower quartile.
- Second Quartile (Q_2): Middle score of set of data by arranging in ascending order. It is also the median.
- Third Quartile (Q_3): around 75% data is below this point
- Max. score: The highest score excluding the outlier.

CLASSTIME [$Q_1 - 1.5 * IQR, Q_3 + 1.5 * IQR$] → Range b/w Q_1 & Q_3 .

Aim: Univariate, Bivariate, Multivariate analysis using histogram, violin plot, bar graph, scatter plot, heatmap.

Software used : Python compiler, google colab

Theory :

DATA VISUALISATION

- Univariate data: Data type consist of only 1 variable.
- Bivariate data: It involves 2 different variables.
- Multivariate data: It involves 3 or more variables.
- Matplotlib: It is a comprehensive python library for creating static, animated, interactive visualisations in python.
- Histogram: It provides a visual representation of grouped freq. distribution with continuous classes.
`plt.hist(data['loan_amount'])`
- Violin plot: With the addition of a rotated kernel density plot on each side.
`plt.violinplot(data['loan_applicantIncome'])`

→ **Barplot**: Bar chart represents category of data with rectangular bars with length and height is proportional to the values which they represent.

`plt.bar(data.index, data)`

→ **Scatterplot**: Diagram where each value in the dataset is represented by a dot.

`plt.scatter(data['applicantIncome'], data['loan_amount'])`

→ **Seaborn**: is a python data visualisation library based on matplotlib that enables statistical graphics in python.

→ **Countplot**: It is used to represent the occurrence plot counts of the observations present in the categorical variable.

~~`plt.countplot(x='education', hue='loanstatus', data=data)`~~

→ **Boxplot**: also known as whisker plot, is used to display summary data values having properties like min, max, Q₁, Q₃, median and outliers.

`plt.boxplot(data['applicantIncome'])`

Experiment :

Date _____

Page No. _____

→ Heatmap : Matplotlib's heatmap represents magnitude of phenomenon to represent data distribution by differences in colors.

0 : neutral

-ve : +ve relation

+ve : -ve relation

corr = def. corr

sns. heatmap (corr, annot=True)

Observations and conclusions :

We have successfully plotted various plots of distinguishing visual representation of univariate, bivariate and multivariate analysis.

Value
Simplification

CLASSTIME

Ques: Perform EDA on a database, addressing missing values, to enhance modelling readiness.

Software used: Google colab

Theory:

Exploratory data analysis is a crucial step in the data analysis process that involves summarizing & visualizing key characteristics of the data to understand its underlying structure.

Objective of EDA -

- Understand the data
- explore relationships
- Detect patterns
- Prepare for modelling

- Handelling missing values

Missing values are a common challenge in a real world dataset. Dealing with missing value data is essential to ensure the accuracy & reliability of any subsequent data analysis or modelling, several strategies can be employed to address missing values.

Experiment :

Date _____

Page No. _____

Common Strategies:

1) Removal of rows/columns:

PROS: Simple & straight-forward

CONS: May lead to loss of valuable info.

→ suitable when missing values are random & not influential.

2) Imputation:

• Mean / Median / Mode imputation:

Filling missing values with the mean, median or mode of the variables.

PROS: Preserves the overall structure of data

CONS: May introduce bias, especially if missing data is not random.

• Interpolation or predictive imputation:

Predicting missing values based on the relationship with other variables.

PROS: Can capture more complex.

CONS: Requires a model & assumptions about relationship.

3) Advanced Techniques:

Machine learning algo such as K-nearest neighbour or multiple imputations can be used for more sophisticated imputation strategies.

Obj: To understand & implement diverse data transformation techniques in Python, including Z-score, min-max scaling, mean normalisation, max absolute scaling and robust scaling

Objectives:

- explore the theoretical concepts behind various data transformation techniques.
- Implement these techniques using python & analyse their effects in data distribution.
- Understand the importance of pre-processing in enhancing the performance of machine learning models.

Software used: Google Colab.

Theory:

Data transformation techniques:

- i) Z-score: Z-score transforms the data to have a mean of 0 and a standard deviation of 1.

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

x = original value
 μ = mean
 σ = standard deviation

Experiment :

Date _____

Page No. _____

Impact:

- centers the data around 0.
- scales data to comparable range.

2) Min-max scaling: It transforms the data to a specific range, typically $[0, 1]$.

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

Impact:

- scales data to a specific range.
- maintains the relative relationships b/w data points.

3) Mean normalisation: It transforms the data to have a mean of 0 and a range b/w -1 and 1. It is useful when the data distribution is skewed.

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

Impact:

- centers the data around 0.
- useful for data with a skewed distribution.

CLASSTIME

Experiment :

Date _____

Page No. _____

- 4) Max absolute scaling : It scales the data by dividing each value by the max. absolute value of the feature.

$$x_{\text{new}} = x / x_{\text{max}}$$

Impact :

- scales the data based on max. values.
- preserves the sign of the data.

- 5) Robust scaling : It is resistant to outliers and scales the data based on the interquartile range (IQR).

$$x_{\text{new}} = \frac{x - x_{\text{median}}}{\text{IQR or } (Q_3 - Q_1)}$$

Conclusion :

In this experiment, we have successfully explored & implemented diverse data transformation techniques.

Value
↓

CLASSTIME

Aim: To explore and implement the following distances & similarity techniques:

- i) Euclidean distance
- ii) Manhattan distance
- iii) Cosine similarity
- iv) Minkowski distance

Software Used: Google colab, python

Theory:

Distance metrics play a crucial role in various data science tasks, including clustering, classification and similarity search.

Understanding and implementing different distance metric enable data scientists to quantify the similarity or dissimilarity b/w data points, thereby facilitating effective analysis and decision making.

i) Euclidean distance:

- straight-line distance b/w points in Euclidean space.

$$d(p, q) = \sqrt{\sum_{i=1}^n (p_i - q_i)^2}$$

ii) Manhattan distance:

Experiment :

Date _____

Page No. _____

- Defⁿ: $d(p, q) = \sum_{i=1}^n |p_i - q_i|$
- Sum of absolute differences of coordinates

iii) Cosine similarity:

- $\text{cosine-sim}(p, q) = \frac{p \cdot q}{\|p\| \|q\|}$
- measures the cosine of the angle between vectors

iv) Minkowski distance:

$$\cdot d(p, q) = \left(\sum_{i=1}^n |p_i - q_i|^p \right)^{\frac{1}{p}}$$

- generalised metric including euclidean and Manhattan distances.

Conclusion: We have understood and implemented euclidean dist., manhattan dist., minkowski dist., and cosine similarity.

Final
4/5/24

CLASSTIME

Experiment : 6

Date 19/02/24
Page No. _____

Ques: Implementation of Apriori algorithm

Objective : To implement the apriori algo. using 2 different approaches.

- a) Utilizing 'mlxtend' library to apply the apriori using funcn.
- b) Manually implementing the apriori from scratch.

Software used : ~~Google Colab~~

Theory :

The apriori is a classic algo. used for frequent itemset mining in association rule learning over Transactional databases. It aims to find frequent itemsets in a dataset and generate association rules based on these frequent itemsets. It operates in 2 main steps:

i) Generate candidate itemsets

Initially, it scans the database to find frequent items. Then, it generates candidate itemsets of length ($K+1$) by joining frequent itemsets of length K .

- ii) Prune infrequent itemsets
After generating candidate itemsets, apriori prunes the candidates that contain any subset which is infrequent. This is based on the 'apriori property', which states that if an itemset is infrequent all its supersets will also be infrequent.

Implementation

- a) Using mlxtend library:

- The mlxtend library provides an easy to use implementation of the apriori funcn algo.
 - We can import the apriori funcn from the library & use it directly on our dataset to find frequent itemsets and generate association rules.
- b) Manual implementation
- Initialize : Initialize frequent itemsets of length (single items) by counting their occurrences in the dataset
 - Generate candidates : Use the frequent itemsets of length K to generate candidate itemsets of length (K+1) by joining them.

Experiment :

Date _____

Page No. _____

- Prune candidates: Remove candidates that contain subsets of K length that are infrequent.
- Repeat steps ② & ③ until no more frequent itemsets can be found.
- Generate association rule: Once frequent itemsets are found, generate association rules based on these frequent itemsets.

Conclusion: We have successfully implemented apriori algo. using both the approaches.

~~Hand~~

CLASSTIME

Aim: Implementation of FP growth algo.

Objective: We have to implement the FP growth algo. using:

- a) 'mlxtend' library
- b) mathematical programs w/o modules

Theory:

FP growth: It's a popular method for frequent itemset mining in transactional databases. It constructs a compact data structure called FP tree to represent the database and efficiently mine frequent itemsets. The algo. involved foll. steps:

i) Construct FP tree:

Initially, it scans the transactional DB to construct the FP-tree, a compact representation of frequent itemsets.

ii) Mine frequent itemset:

After constructing the FP tree, it recursively mine frequent itemsets by performing condition pattern base construction and generating conditional FP tree.

Implementation :

a) Using 'mlxtend' library

- import mlxtend
- import 'fp growth' func and apply it to the given dataset.

b) Manual implementation

Algorithm :

- 1) Construct FP-tree : Sort transactions, build tree
- 2) Mine frequent itemsets : Traverse tree, output frequent itemsets.
- 3) Insert transaction in tree :
 - a. If transaction empty, return.
 - b. Get first item, create child if not exist.
 - c. Increment count, remove item, recurse.
- 4) Mine frequent itemsets
 - a. If leaf node, return
 - b. For each child, calculate support.
 - c. If support \geq min_support, output and recurse.

Conclusion : We have successfully implemented the FP growth algo. using both methods.

- ~~Aim:~~
- i) Perform logistic regression on diabetes dataset.
 - ii) Plot the confusion matrix in the heatmap form from the model generated and print the accuracy, precision, recall and F-1 score.

~~Software used: google colab~~

Theory:

→ Logistic regression is a supervised machine learning algo used for classification tasks where the goal is to predict the prob. that an instance belongs to a given class or not.

Logistic regression is a statistical algo. which analyse the relationship b/w 2 data factors.

→ Logistic funcⁿ - ~~Sigmoid funcⁿ~~

It is used to map the predicted values to probabilities. It maps any real value into another value within a range of 0 and 1. The value of the logistic regression must be b/w 0 & 1, which

Experiment :

Date _____

Page No. _____

cannot go beyond this limit, so it forms a curve like the 'S' form. It is called the sigmoid funcⁿ or the logistic funcⁿ.

→ Types of logistic regression:

- 1) Binomial : there can be only 2 possible types of the dependent variable, such as 0 or 1, pass or fail, etc.
 - 2) Multinomial : there can be 3 or more possible unordered types of the dependent variable, such as 'low', 'high' and 'medium'.
 - 3) Ordinal : there can be 3 or more possible ordered types of dependent variables such as 'low', 'high' or 'medium'.
- Confusion matrix : It is a matrix that summarizes the performance of a machine learning model on a set of test data. It is a means of displaying the no. of accurate and inaccurate instances based on the model's predictions.

Experiment :

Date _____

Page No. _____

- True positives (TP) : occurs when the model accurately predicts a +ve datapoint.
- True negatives (TN) : occurs when the model accurately predicts a -ve data point.
- False positives (FP) : occurs when the model predicts a +ve data point incorrectly.
- False negatives (FN) : occurs when the model mispredicts a -ve data point.

ACTUAL

| | | |
|---|----------|----------|
| T | True | False |
| C | positive | positive |
| F | False | True |
| R | negative | negative |

→ Metrics based on confusion matrix data:

1) Accuracy : It is used to measure the performance of the model. It is the ratio of total correct instances to the total instances.

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

2) Precision : It is a measure of how accurate a model's +ve predictions are.

$$\text{Precision} = \frac{TP}{TP + FP}$$

Experiment :

Date _____

Page No. _____

3) Recall : It measures the effectiveness of a classification model in identifying all relevant instances from a dataset.

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

4) F-1 score : It is used to evaluate the overall performance of a classification model. It is ~~is~~ the harmonic mean of precision and recall.

$$\text{F-1 score} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Conclusion : Hence, we implemented logistic regression on diabetes dataset and plotted the corresponding confusion matrix along with accuracy, precision, recall & F-1 score.

Ques : To compare various boosting methods

- i) Gradient Boosting
- ii) XGBoost
- iii) AdaBoost
- iv) CATBoost

Software used : ~~Google Colab~~.

Theory :

→ Boosting is an ensemble modeling technique that attempts to build a strong classifier from the no. of weak classifiers.

Advantages :

- 1) Improved accuracy
- 2) Robustness to overfitting
- 3) Better handling of imbalanced data
- 4) Better interpretability.

→ Gradient Boosting - It is a boosting technique that builds a final model from the sum of several weak learning algorithms that were trained on the same dataset.

Experiment :

Date _____

Page No. _____

- XGBoost - The full form of XGBoost is Extreme Gradient Boosting algo. which is an extreme variation of previous technique. XGBoost applies a regularisation approach and outperforms a standard gradient boosting method.
- AdaBoost - It is a boosting algo. that also works on the principle of stagewise addition method where multiple weak learners are used for getting strong learners. The value of alpha parameter, in this case, will be indirectly proportional to error of the weak learner.
- CatBoost - The growth of decision trees include CatBoost is primary distinction that sets it apart from & improves upon competitors.

Conclusion : We have successfully compared various boosting methods :

- i) Gradient Boosting
- ii) XGBoost
- iii) AdaBoost
- iv) CATBoost

In all above methods, CATBoost performs the best.

CLASSTIME

10/4/19
Dhanvi

- Aim:**
- Perform classification using a Naive Bayes classifier on the given dataset.
 - Perform Regression using a regression tree on the given dataset.

Software used: google colab

Theory:

i) **Naive Bayes Classification (Diabetes dataset):** It is a probabilistic classifier based on Baye's theorem, which is expressed as:

$$P(C_k/x) = \frac{P(x/C_k) \cdot P(C_k)}{P(x)}$$

where, $P(C_k/x)$ = posterior prob. of class C_k given features x .

$P(x/C_k)$ = likelihood of features x given class C_k

$P(C_k)$ = prior probability of class C_k

$P(x)$ = probability of features x .

These classifiers are simple, yet effective for classification tasks, especially when the features are assumed to be independent given the class.

It's called 'naive' because it assumes that the presence of a particular feature in a class is unrelated to the presence of any other features.

- Regression tree (Walmart Sales dataset):
Regression tree is a decision tree based model used for regression tasks, where the target variable is continuous. Each internal node in the tree represents a feature (or an attribute), each branch represents a decision rule, each leaf node represents the outcome.
- In regression trees, the splitting of nodes is done based on a criterion that minimizes the variance of the target variable within each split. The most common criteria include mean squared error (MSE), mean absolute error (MAE), etc.
- Recursive partitioning: The process of constructing a regression tree involves recursively partitioning the data into subsets, based on the values of the

Experiment :

Date _____

Page No. _____

features. At each step, the algo. selects the feature and the split point that minimizes the chosen criterion.

- Once the tree is constructed, each instance falling into a leaf node is assigned the mean (or median) of the target variable within that leaf node as its predicted value.
- After construction, trees may be pruned to avoid overfitting, where nodes that provide little predictive power are removed.

CONCLUSION: We have successfully performed classification using Naive Bayes classifier and regression using regression tree on the following given datasets.

Vareⁱⁿ
2014

CLASSTIME

Aim: To apply different clustering methods for classification on given dataset.

Software used: Google Colab

Theory:

1) K-means clustering:

- It is a popular partitioning clustering algorithm that divides data points into K-clusters.
- It minimizes the sum of squared dist. b/w data points and their respective cluster centroids.
- The algo. iteratively assigns data points to the nearest centroid and updates centroids until it converges.

2) DB-SCAN (Density Based spatial Clustering of Application with Noise):

Unlike K-means, DB scan does not require the no. of clusters to be specified in advance. It groups together points that are closely packed, defining clusters as areas of high density separated

Experiment :

Date _____

Page No. _____

By areas of low density. It labels points as core, border or noise-based on their density and proximity to other points.

→ KNN :

KNN is a non-parametric instance-based learning algo. used for classifications & regression tasks. In clustering, KNN assigns each point to the most common class among its k nearest neighbours. It relies on the dist. metric to determine similarity b/w data points.

CONCLUSION: We have applied K-means DB-SCAN and KNN on given unlabelled dataset.