

Experiment No : 1

- * Title : Exploratory Data Analysis and visualization of Relationship in a Given data set.
- * Aim : To perform Exploratory Data Analysis (EDA) on a given dataset and visualize relationship between variables to uncover pattern trends and insights.

* Theory :

Exploratory Data Analysis (EDA) is an approach to analyzing dataset to summarize their main characteristics often using visual methods. It helps in :

- Understanding data distribution and missing values.
- Identifying correlations and relationships between features.
- Detecting anomalies and outliers.

Key Techniques in EDA :

- summary statistic : Mean, Median, standard deviation etc
- Data visualization : Histogram, scatter plots, pair plots, box plots, heatmap etc.

- Feature correlation analysis: Using correlation matrices and scatter plots.

Visualization Techniques in EDA

EDA heavily relies on visual methods to better understand data relationships and distributions.

1. Histogram: Show the distribution of a single variable.
2. Boxplots: Help detect outliers and compare distribution across categories.
3. Scatterplot: Show relationship between two numerical variables.
4. Heatmaps: Represent correlation matrices of identify strong relationships.
5. Pair plots: Provide a quick overview of relationship between multiple numerical variables.

• Importance of EDA

- 1) Helps in understanding datasets
- 2) Detect data quality issues like missing value.
- 3) Aids in feature selection by identifying relevant variables.

* Algorithm :

- Step 1 : Load and inspect the dataset
- Step 2 : Perform data cleaning (handle missing values, duplicates, outliers)
- Step 3 : Generate summary statistics (mean, median, min - max value)
- Step 4 : Visualize relationship using various plots (scatterplot, heatmap, pair plots, box plots.)
- Step 5 : Interpret the result and derive insight.

* Result :

- summary statistics of the dataset
- observation from visualization (e.g., correlation strength, presence of outliers, distribution patterns).
- key insight derived from EDA

* Conclusion :

Hence, we successfully perform the EDA and visualize the relationship for given dataset.

Experiment NO : 2

* Title : Exploratory Data Analysis
Visualization and Linear Regression
Model Training on the Housing
price dataset.

* Aim : To perform Exploratory Data
Analysis (EDA) on the Housing
price dataset visualize relationship
between variables, train a linear
Regression Model and evaluate
its performance.

* Theory :

- Exploratory Data Analysis (EDA) :
EDA is a crucial step in data
analysis that helps understand the
dataset before applying any machine
learning model. The main objectives
of EDA include:

- Identifying pattern, trends and
relationships between variables.
- Detecting missing values, outliers,
and anomalies.
- summarizing data distribution and
correlations.

- * Linear Regression :

Linear Regression is a supervised
learning algorithm used for predicting
continuous values.

The relationship between the dependent variable y (housing price) and independent variables (features like square footage, number of bedrooms, location, etc.) is modeled using the equation:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n + \epsilon$$

where,

y : Dependent variable

x : independent variable

β : coefficient

ϵ : Error term

- Key Points of Linear Regression :

1. Linearity : Relation between independent & Dependent variable is linear.
2. Independence : Observations are independent for each other.
3. Homoscedasticity : Constant variance of residual errors.
4. No multicollinearity : Independent variables should not be highly correlated.

* Algorithm :

1. Load the dataset : and perform an initial inspection.
2. Perform EDA :
 - Summary statistics and missing values analysis.
 - Visualizations (histogram, scatter plots, heatmaps).
3. Feature Engineering :
 - Handle missing data
 - Encode categorical variables.
 - Normalize /standardize numerical features if needed
4. Train - Test split : Divide the dataset into training and testing sets.
5. Train linear Regression Model on training data.
6. Evaluate Model performance using MSE, RMSE and R^2 score.

Result :

- Summary statistics of the dataset
- Insight from COO visualization
- Model performance metrics (MSE, RMSE, R² score).

Conclusion

- Key observations from the dataset.
- Find the strength of relationship between features and housing price
- Performance of linear Regression Model.
- Possible improvement

Experiment No. 1

Title : Exploratory Data Analysis and Ridge Regression Model training on a Given Dataset

Aim : To perform Exploratory Data Analysis (EDA) on a given dataset, visualize relationships between variable train a Ridge Regression model and evaluate its performance.

Theory :

Ridge Regression :

Ridge Regression is a type of linear regression that includes an L₂ regularization term to prevent overfitting by penalizing large coefficients. It modifies the cost function by adding a regularization term:

$$J(\beta) = \sum (y_i - \hat{y}_i)^2 + \lambda \sum \beta_j^2$$

where :

- 1) $J(\beta)$ is the cost function
- 2) λ (alpha) is the regularization parameter (higher λ shrink coefficient)
- 3) β represents the model coefficients

Model Evaluation Metrics:

1. Mean squared error (MSE)

MSE measures the average square difference between actual and predicted value.

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

2. Root Mean squared error (RMSE)

RMSE is simply the square root of MSE. It is used to bring the error value back to the original unit of dependent variable.

$$RMSE = \sqrt{MSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

3. R-squared (R^2) score

R^2 , also called the coefficient of determination, measures how well the independent variable

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

Algorithm :

1. Load and inspect the dataset
2. Perform EDA
 - Visualize feature distributions
 - Identify correlations
 - Detect and handle missing values or outliers.
3. Feature Engineering:
 - Encode categorical variables
 - Scale numerical features
4. Train - Test - split : split the dataset into training and testing sets.
5. Train Ridge Regression Model with different regularization strengths.
6. Evaluate Model performance using MSE, RMSE and R^2 score.

Results :

- Observation from EDA (feature distributions)
- Model performance metric (MSE, RMSE, R^2 score) for Ridge regression.

Conclusion :

Ridge Regression performance and comparison with other models. The impact of regularization on model performance.

Experiment No: 4

* Title : Exploratory Data analysis, visualization and Lasso Regression Model Training on a Given Dataset.

* Aim : To perform EDA on a given dataset, visualize relationships between variables, train a Lasso Regression model and evaluate its performance.

Theory :

- Exploratory Data Analysis (EDA)

EDA is the process of analyzing and summarizing dataset using statistical & visualization technique to uncover insights before applying machine learning models.

1- Lasso Regression :

Lasso (Least Absolute Shrinkage and selection Operator) Regression is a linear regression technique that include L1 regularization, which helps with both feature selection and preventing overfitting.

$$J(\beta) = \sum (y_i - \hat{y}_i)^2 + \alpha \sum |\beta_j|$$

Where,

' $J(\beta)$ ' is the cost function

' α ' is the regularization parameter

' β ' represent the model coefficients.

Key points of lasso regression

- Unlike Ridge regression, which shrinks coefficients close to zero, lasso can shrink some coefficients exactly to zero.
- Useful when dealing with high dimensional data.

Algorithm :

Step 1 : Load and inspect the dataset

Step 2 : Perform EDA

1. Visualize feature distribution
2. Identify correlations using heatmap.
3. Detect and handle missing values or outliers

Step 3 : Feature Engineering :

- Encode categorical variables
- scale numerical features.

Step 4 : Train - Test split Divide the dataset into training and testing sets.

Step 5 : Train Lasso Regression Model using different values of α (alpha)

Step 6 : Evaluate model performance using MSE, RMSE and R^2 score.

Results :

- Feature relationships and correlation analysis.
- Model performance metrics (MSE, RMSE, R² score) for Lasso Regression
- Important features selected by Lasso regression.

Conclusion :

The effectiveness of Lasso regression in selecting important features. Model performance evaluation and possible improvements. Future steps tuning value, trying alternative model like Ridge regression.

Experiment No : 5

i) Title : Exploratory Data Analysis, visualization and logistic Regression Model Training on the Iris dataset.

ii) Aim : To perform Exploratory Data Analysis on the Iris dataset, visualize relationship between variables, train a logistic regression model and evaluate its performance in classifying iris species.

iii) Theory :

Logistic Regression :
 Logistic Regression is a supervised machine learning algorithm used for classification tasks. It models the probability that a given input belongs to a specific class using sigmoid function.

$$P(y=1|x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n)}}$$

where, "P(y=1|x)" is the probability of an instance belonging to class 1

- 'β' represent the model coefficients
- 'x' represent the model coefficients
- The sigmoid function ensure the output is between 0 & +1.

Model Evaluation Metrics:

- Accuracy : Measures overall correctness.
- Precision, Recall and F1-score : Evaluate class-specific performance.
- Confusion Matrix : visualize correct and incorrect classifications.

Types of Logistic Regression

- Binary Logistic Regression : Two classes (e.g. Yes / No, 0 / 1)
- Multinomial Logistic Regression : More than two classes (e.g. Red, Blue, Green)
- Ordinal Logistic Regression : Ordered categories (e.g. Low, medium, high)

* Algorithm :

- P 1. Load the dataset and check its structure.
- P 2. Perform EDA :
 - Check class distribution
 - Visualize feature distributions.
 - Analyze feature relationships using pair plots & heatmaps.
- P 3. Data Preprocessing :
 - Handle missing values (if any)
 - Encode categorical variables

• Split data into training and testing sets.

4. Train a Logistic Regression Model.

5. Evaluate Model Performance using accuracy confusion matrix and classification report.

Results:

• EDA Findings:

1.) Sepal length and petal length are highly correlated.

2.) Setosa is distinctly separable from the other two species.

• Model performance:

1. Accuracy

2. Confusion Matrix

3. Classification Report.

Conclusion:

The model achieved high accuracy indicating good performance. Setosa was the easiest to classify, while Versicolor and virginica had some misclassification.

Experiment No 6:

Title : Exploratory Data Analysis , visualization and logistic Regression Model Training on The user Dataset.

Aim : To perform Exploratory Data Analysis (EDA) on the User dataset , visualize relationship between variables , train a logistic Regression model and evaluate its performance in classifying your user behaviour.

Theory :

Logistic regression is a fundamental machine learning algorithm used for classification problems . Unlike linear regression , which predicts continuous values , logistic regression predict probabilities and maps them to binary (0 or 1) or multiclass outputs.

Logistic Regression uses the sigmoid function to transform linear outputs into probabilities between 0 and 1

Sigmoid function formula :

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

where :

$$z = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

$\sigma(z)$ represents probability belong to class z

If $P(Y=1) \geq 0.5 \rightarrow$ classify as 1.
 If $P(Y=1) < 0.5 \rightarrow$ classify as 0.

The sigmoid function ensures outputs are mapped into probability range (0 to 1)

Cost Function for Logistic Regression

$$J(\theta) = -\frac{1}{n} \sum_{i=1}^n [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

where,

- y_i is actual class
- \hat{y}_i is the predicted probability for class 1
- The function penalizes incorrect predictions more if the model is highly confident but wrong

Algorithm :

1. Load and inspect the dataset
2. Perform EDA :
 - Analyze feature distributions using histograms and boxplots.
 - Identify correlation using a heatmap
 - Visualize relationships between numerical variable with pair plots.
3. Preprocess Data :
 - Handle missing values.

- Encode categorical variables
 - Scale numerical feature
 - Split data into training and testing sets.
4. Train Logistic Regression model using the preprocessed dataset.
5. Evaluate Model performance using accuracy, confusion matrix, classification report and AUC score.

iii) Results :

- EDA Findings :

1. The dataset has no missing values
2. certain features may have strong correlations affecting classification.
3. The target variable, is imbalanced.

- Model Performance :

1. Accuracy
2. Confusion Matrix
3. Precision, Recall & F1-score
4. ROC curve & AUC

iv) Conclusion :

The Logistic Regression model provided good classification accuracy. The ROC curve and AUC indicate how well the model separates classes.

Experiment No : 7

- * Title : Exploratory Data Analysis , Visualization and the Naive Bayes classification on the Wine Dataset:
- * Aim : To perform Exploratory Data Analysis (EDA) on the wine dataset , visualize relationships between variables , train a Naive Bayes classifier and evaluate its performance in classifying wine categories.

* Theory :

Exploratory Data Analysis (EDA)

EDA is a process that involves understanding dataset characteristics using summary statistics and visualizations .

- checking the distribution of variables.
- Identifying patterns , relationships and correlations between features.
- Detecting missing values and outliers.
- Preparing data for machine learning models.

* Naive Bayes Classifier :

Naive Bayes is a probabilistic classification algorithm based on Bayes Theorem , assuming Feature independence

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

where,

- $P(A|B)$ is the probability of class A given feature B.
- $P(B|A)$ is the probability of feature B given class A.
- $P(A)$ is the prior probability of class A.
- $P(B)$ is the probability of feature B occurring.

• Types of Naive Bayes classifiers:

1. Gaussian Naive Bayes : Assume normal distribution of features.

2. Multinomial Naive Bayes : Suitable for discrete feature control.

3. Bernoulli Naive Bayes : Used for binary feature data.

For wine dataset, Gaussian Naive Bayes is appropriate since features are continuous.

• Model Evaluation Metrics:

1. Accuracy
2. Precision
3. Confusion Matrix

Result :

- EDA Findings :

1. The dataset contains no missing values.
2. Certain feature show strong correlations.
3. Some classes may have overlapping features.

Conclusion :

The Naive Bayes model achieved good accuracy on the test set. Standardization of features improved model performance. The confusion matrix revealed some misclassifications, especially in overlapping feature space.

Experiment No:8

* Title : Exploratory Data Analysis
Visualization and Decision Tree
classification on the Titanic
Dataset.

* Aim : To perform Exploratory Data Analysis (EDA) on the Titanic dataset, visualize relationship between features, train a Decision Tree classifier and evaluate its performance in predicting survival.

* Theory :

Exploratory Data Analysis (EDA)

EDA helps in understanding data characteristics using statistical summaries and visualizations.

Decision Tree classifier

A decision tree is a supervised learning algorithm used for classification and regression tasks. It splits the dataset into branches based on features (conditions) following an "if-else" structure.

The Gini index and Entropy are common criteria for splitting:

- Gini Index : Measures impurity.

- Entropy : Measures information gain at each split.

A Decision Tree classification data by traversing branches from root to leaf nodes, making it interpretable but prone to overfitting.

Model Evaluation metrics

1. Accuracy : Overall correctness of prediction
2. Precision, Recall, and F1-score : Class-specific performance evaluation
3. Confusion Matrix : Breakdown of correctness and incorrect classification
4. Feature Importance : Identifies influential feature in decision making

Algorithm :

1. Load and inspect the dataset
2. Perform EDA on the dataset
3. Data Preprocessing, Handled missing values, convert categorical into numerical etc.
4. Train a Decision tree classifier using entropy / Gini criteria.
5. Evaluate Model Performance using accuracy, confusion matrix and feature importance.

Result :

EDA Findings :

- Gender : Females had a higher survival rate.
- Passenger class : Higher class passengers had a better survival rate.
- Feature correlations : Pclass, Fare and age were key predictors.

Conclusion :

The decision tree model provided good accuracy in predicting survival. Feature importance analysis revealed that gender, class and fare were crucial.

Experiment No. 9

* Title : Exploratory Data Analysis, Visualization and K - Nearest Neighbors (KNN) classification on the Digits Dataset.

* Aim : To perform exploratory Data Analysis (EDA) on the digits dataset, visualize relationships between feature, train a K-Nearest Neighbors between feature train a K-nearest neighbor classifier and evaluate its performance in recognizing handwritten digits.

* Theory :

K - Nearest Neighbors (KNN) classifier
 KNN is non-parametric, supervised learning algorithm that classifies a data point based on the majority class of its k nearest neighbors.

Key properties :

- Distance metric : Common choice include Euclidean, Manhattan and Minkowski distance.

- k -value selection : A small k may lead to overfitting while a large k may oversmooth class boundaries.

- Weighted voting: closer neighbor can be given higher importance.

Mathematically, the Euclidean distance between two points $x_1 (x_1, x_2, \dots, x_n)$ and $x_2 (y_1, y_2, \dots, y_n)$ is

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Model Evaluation metrics:

1. Accuracy
2. Precision, Recall & F1-score
3. Confusion matrix

* Algorithm:

1. Load and inspect the dataset
2. Perform EDA:
 1. visualize simple digit images
 2. Reduce dimensions using PCA for visualization
3. Data Preprocessing:
 1. Normalize pixel values.
 2. Split data into training and testing sets.

- Train a KNN classifier and optimize K value.
- evaluate model performance using accuracy, confusion matrix and classification report.

Results :

- CDA Findings :
 - 1. Each digit is represented as an 8×8 pixel image.
- PCA visualization shows clustering pattern among digits
- Feature scaling improves distance-based classification.

Conclusion :

The KNN model performed well in recognizing handwritten digits. Feature scaling was crucial for better performance.

Experiment No: 11

Title : Exploratory Data Analysis , visualization and k means clustering on the Iris Dataset.

Aim : To perform EDA and visualize the relationships among features in the iris dataset , apply the K-mean clustering algorithm and evaluate the clustering performance .

Theory :

K - Mean clustering

K Mean is an unsupervised machine learning algorithm used to group similar data points into k clusters based on feature similarity.

steps :

1. choose k (number of clusters)
2. Initialize k cluster centroid randomly
3. Assign each data point to the nearest centroid
4. Recalculate centroid based on mean of assigned points.
5. Repeat steps 3 - 4 until centroid stabilize

Mathematical Concept:
Using Euclidean distance to measure similarity between data points.

Evaluation metrics:

1. Inertia : Total distance of points from cluster center.
2. Silhouette Score : How well separated and cohesive the clusters are (-1 to 2)

3. Comparison with True Labels :

Although unsupervised, we can compare clusters with actual species labels for evaluation.

Algorithm :

1. Load the dataset and perform exploratory analysis.
2. Visualize distributions and pairwise relationships.
3. Preprocess data : scale features for clustering performance.
4. Determine optimum K using Elbow method and silhouette score.
5. Train K-means model & predict cluster.

6. Compare predicted clusters with actual species.
7. Visualize clustered data and centroids using PCA.

Result :

- cluster formed : 3 clusters corresponding closely to 3 actual species
 - silhouette score : Approximately , indicating good separation.
 - visualize : PCA plots clearly show well separated clusters.
 - comparison : K-Means performed well without supervision in approximating true species labels.
- * Conclusion : K-Means was successful in clustering the Iris data into meaningful group. EDA helped understand data distribution and relationship. PCA enables visual inspection of clustering effectiveness.

Title : Exploratory Data Analysis - Visualization and clustering with K-means on an Income Dataset

Aim :

To perform EDA and visualize relationship in the Income dataset, apply the K-means clustering algorithm to segment data based on income related pattern and evaluate clustering performance using appropriate metrics.

Theory :

K-Means Clustering

K-Means is an unsupervised learning algorithm used for clustering data into k distinct groups based on similarity.

Key steps in K-Means :

1. choose the number of clusters (k)
2. Initialize k centroids randomly.
3. Assign each point to nearest centroid.
4. Update centroid by calculating the mean of assigned points
5. Repeat 3 and 4 step until convergence

Distance Metric :

Typically uses Euclidean distance

Evaluation Metrics :

1. Inertia (within-cluster sum of squared): Lower is better

2. Silhouette Score : Measure how similar an object is to its own cluster vs other cluster.
(range -1 to 1)

3. Elbow Method : Helps determine the optimal number of clusters (k)

K-Means is not a regression technique ; it is used for clustering not for predicting a continuous target variable.

Algorithm :

1. Load & inspect the dataset

2. CDA :

- Plot distribution of income spending score , age
- Use scatter plot and pair plots

To visualize feature relationships

3. Data Preprocessing:

1. Handle missing values (if any)
2. Encode categorical variables
3. Scale /normalize the features

4. K-means clustering:

- Determine optimal k using Elbow method and/or silhouette score
- Fit the k-means model on selected features.

5. Evaluate Clustering:

- Visualize cluster using 2D or 3D scatter plots
- Print and silhouette score

Result:

- Cluster formed: K means formed 4 distinct clusters based on income, age, and spending score.
- Inertia: (displayed from code)
- Silhouette: -11-
- Visualization: PCA and scatter plot show well separated groups, indicating meaningful segmentation.

Conclusion :

individuals K Means effectively clustered data based on income-related PCA dimensional score cluster confirmed can be used for customer segmentation, targeted marketing or socio-economic analysis.