

1. Uber LR and RF

Q1. What is the objective of this practical?

A: The main objective is to predict the **Uber ride fare** based on features like pickup and drop-off locations, passenger count, and time of the ride using **Machine Learning regression models** such as **Linear Regression** and **Random Forest Regression**.

Q2. What type of machine learning problem is this?

A: It's a **supervised regression problem** because the target variable (**fare_amount**) is continuous and we are training the model on labeled data.

Q3. What dataset did you use?

A: The **Uber Fares Dataset** from Kaggle, which contains information about pickup and drop-off coordinates, datetime, passenger count, and fare amount.



Data Preprocessing Questions

Q4. Why did you drop columns like `Unnamed: 0` and `key`?

A: These columns don't provide any useful information for prediction — `Unnamed: 0` is just an index, and `key` is a unique identifier. Keeping them would add unnecessary noise to the model.

Q5. How did you handle missing values?

A: Missing latitude and longitude values were replaced with their respective **column means** using:

```
data['dropoff_latitude'].fillna(data['dropoff_latitude'].mean(),  
inplace=True)
```

Q6. What did you do with zero or invalid values?

A: Zero passenger counts were replaced with the **mean passenger count**, and fares less than or equal to zero were replaced with the **mean fare**, since negative or zero fares are not realistic.

Q7. How did you extract useful features from datetime?

A: The **pickup_datetime** column was split into **month** and **hour** components to capture **time-based trends** that can affect fare prices.



Outliers and Correlation

Q8. What are outliers?

A: Outliers are data points that differ significantly from other observations in the dataset. They can skew the model's predictions if not treated properly.

Q9. How did you detect outliers?

A: Outliers were detected using **boxplots** and the **Interquartile Range (IQR)** method.

Q10. How did you treat the outliers?

A: We used the **IQR method** to cap values beyond the acceptable range using **np.clip()**:

```
data[col] = np.clip(df1[col], lb, ub)
```

This prevents extreme values from influencing the model.

Q11. What is correlation and why is it important?

A: **Correlation** measures how strongly two variables are related. It helps us understand which features most affect the target variable (fare amount). **Strong correlation with the target means that feature is useful for prediction.**

Q12. What did the correlation heatmap show?

A: The heatmap visualized relationships between numerical features. Typically, latitude/longitude values and hours show some correlation with fare amount.



Model-Related Questions

Q13. What models did you use?

A:

1. **Linear Regression** – to model a linear relationship between features and fare.
 2. **Random Forest Regression** – to capture complex, nonlinear relationships.
-

Q14. What is Linear Regression?

A: Linear Regression finds the **best-fitting straight line** (using the least squares method) that minimizes the difference between predicted and actual values.

Formula:

$$y = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

Q15. What is Random Forest Regression?

A: Random Forest Regression is an **ensemble learning** algorithm that builds multiple **decision trees** and combines their outputs (by averaging) to improve accuracy and reduce overfitting.

Q16. Why did you use both models?

A: To compare performance between a **simple linear model** (Linear Regression) and a **complex nonlinear model** (Random Forest), and determine which predicts fares more accurately.



Model Evaluation

Q17. What metrics did you use to evaluate your models?

A:

1. **MAE (Mean Absolute Error)** – Average of absolute prediction errors.
 2. **MSE (Mean Squared Error)** – Mean of squared errors.
 3. **RMSE (Root Mean Squared Error)** – Square root of MSE, gives error in fare units.
 4. **R² (Coefficient of Determination)** – Measures how much variance in the target is explained by the model.
-

Q18. What is the difference between MAE and RMSE?

A:

- **MAE** treats all errors equally.

- **RMSE** penalizes larger errors more, making it more sensitive to outliers.
-

Q19. Which model performed better and why?

A: Generally, the **Random Forest Regression** model performs better because it captures **nonlinear relationships** and is more **robust to outliers** compared to Linear Regression.

Q20. What does the R^2 value indicate?

A: R^2 indicates how well the model explains the variation in the target variable.

- **$R^2 = 1$** → Perfect prediction
 - **R^2 close to 0** → Poor prediction
-

Additional / Conceptual Questions

Q21. What is overfitting?

A: Overfitting occurs when a model learns the training data too well (including noise) and performs poorly on unseen data.

Q22. How can you prevent overfitting in Random Forest?

A: By tuning parameters like:

- `n_estimators` (number of trees)
 - `max_depth`
 - `min_samples_split`
 - `min_samples_leaf`
-

Q23. Why do we split data into training and testing sets?

A: To evaluate how well the model generalizes to new, unseen data.
Usually, **70–85%** is used for training and the rest for testing.

Q24. What is feature scaling? Was it used here?

A: Feature scaling standardizes data ranges.

It wasn't explicitly needed here because Random Forest and Linear Regression (with properly scaled geographical coordinates) handle numeric ranges acceptably.

Q25. What improvements could be made to this project?

A:

- Add more features (e.g., weather, distance, day of week).
- Use **Haversine distance** between pickup and dropoff points.
- Try **hyperparameter tuning** for Random Forest.
- Use more advanced models like **XGBoost** or **Neural Networks**.



1. What is Linear Regression?

Answer:

Linear Regression is a **supervised learning algorithm** used to predict a **continuous value** based on input features.

It assumes a **linear relationship** between the independent variables (X) and the dependent variable (Y).

Formula:

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$$
$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$$

where

- b_0 = intercept
- $b_1, b_2, \dots, b_{n-1}, b_n$ = coefficients (slopes)

It tries to find the **best-fitting straight line** that minimizes the difference between actual and predicted values (using the **least squares method**).

Example:

Predicting Uber fare based on distance, time, and passenger count.

2. What is Random Forest Regression?

Answer:

Random Forest Regression is an **ensemble learning method** that uses **multiple decision trees** to make predictions.

Each tree predicts a value, and the **average of all tree outputs** is taken as the final prediction.

It helps in:

- Capturing **nonlinear relationships** in data.
- Reducing **overfitting** compared to a single decision tree.
- Providing **high accuracy and robustness**.

Example:

Predicting Uber fares when relationships between features (like time, location, and fare) are not perfectly linear.

3. Difference Between Linear Regression and Random Forest Regression

Feature	Linear Regression	Random Forest Regression
Type of Model	Simple, linear model	Ensemble of decision trees

Relationship Captured	Only linear relationships	Handles both linear and nonlinear relationships
Complexity	Simple and fast	More complex and computationally heavier
Overfitting	Can overfit on noisy data	Less likely to overfit due to averaging across trees
Interpretability	Easy to interpret	Harder to interpret (black-box model)
Performance	Works well with simple, linear data	Performs better on complex, real-world data

⚠ 4. What are Outliers?

Answer:

Outliers are **data points that differ significantly** from other observations in the dataset.

They can occur due to **errors in data collection** or **natural variation**.

Outliers can **distort model performance**, especially for algorithms like Linear Regression.

Example:

If most Uber fares are between ₹100–₹500, but one record shows ₹5000, it's an outlier.

📊 5. What is IQR (Interquartile Range)?

Answer:

IQR stands for **Interquartile Range**, which measures the **spread of the middle 50% of data**.

It helps in **detecting outliers**.

Formula:

$$IQR = Q_3 - Q_1$$

where

- Q_1 = 25th percentile (lower quartile)
- Q_3 = 75th percentile (upper quartile)

Outlier Detection Rule:

A value is considered an outlier if:

$$\text{Value} < Q_1 - 1.5 \times IQR \text{ or } \text{Value} > Q_3 + 1.5 \times IQR$$

Example:

If $Q_1 = 100$, $Q_3 = 300$, then
 $IQR = 200 \rightarrow$ Outliers are values below $100 - 300$ or above $300 + 300$ (i.e., < -200 or > 600).

2. Email KNN and SVM

1. What is the objective of this practical?

Answer:

The objective is to classify emails as either spam or not spam using binary classification methods. We use K-Nearest Neighbors (KNN) and Support Vector Machine (SVM) algorithms to build models and analyze their performance using metrics like accuracy, precision, recall, and confusion matrix.

2. What dataset did you use and what are its features?

Answer:

We used the `emails.csv` dataset from Kaggle. It contains email content that has already been preprocessed into numeric form (features like word frequencies or TF-IDF), and a target column `Prediction` which is 1 for spam and 0 for not spam.

3. Why did you remove the 'Email No.' column from the dataset?**Answer:**

The '`Email No.`' column is just a serial number acting as an index. It does not contribute to the classification and may introduce noise, so it is dropped before training.

4. What is the purpose of splitting data into training and testing sets?**Answer:**

To evaluate the model's performance on unseen data. The training set is used to train the model, while the testing set helps assess how well the model generalizes.

5. Why did you choose `test_size = 0.20` in `train_test_split`?**Answer:**

It allocates 20% of the data for testing and 80% for training. This is a common split ratio that ensures enough data for model learning while preserving some data for evaluation.

6. What is K-Nearest Neighbors (KNN)?**Answer:**

KNN is a non-parametric, instance-based learning algorithm that classifies new data points based on the class of their nearest neighbors in the feature space. The value of `k` determines how many neighbors to consider.

7. Why did you use `n_neighbors = 2` for KNN?

Answer:

We chose **k=2** as a starting point to experiment with minimal complexity. In practice, the best **k** can be tuned via cross-validation.

8. What is Support Vector Machine (SVM)?**Answer:**

SVM is a supervised learning algorithm that finds the optimal hyperplane to separate data points of different classes. It works well, especially for high-dimensional datasets like text.

9. Why did you use a linear kernel for SVM?**Answer:**

In text classification, data is usually high-dimensional but often linearly separable. The linear kernel is faster and more efficient for such scenarios compared to non-linear kernels like RBF.

10. What are the evaluation metrics used in the project?**Answer:**

- **Accuracy:** Overall correctness of the model
 - **Precision:** How many predicted spams were actually spam
 - **Recall:** How many actual spams were correctly identified
 - **Confusion Matrix:** Shows true positives, false positives, true negatives, and false negatives
-

11. What is the difference between Precision and Recall?**Answer:**

- **Precision:** Out of all predicted spam emails, how many were correct.
- **Recall:** Out of all actual spam emails, how many were detected.

12. What is a Confusion Matrix?

Answer:

A confusion matrix is a table that summarizes the performance of a classification model by showing the true and predicted classifications for each class.

13. Which model performed better in your analysis and why?

Answer:

Typically, SVM with a linear kernel performs better than KNN for high-dimensional text data. SVM is faster, more accurate, and handles sparse data better. Performance comparison is based on accuracy and other metrics observed during evaluation.

14. What is the purpose of calculating training time for each model?

Answer:

Training time helps compare the efficiency of different models. For large datasets or real-time applications, a faster model is often preferred.

15. Why is KNN not ideal for text datasets?

Answer:

KNN computes distances between all training samples, which becomes slow and memory-intensive for high-dimensional data like text (where each word is a feature).

16. What will happen if you increase the value of `k` in KNN?

Answer:

A larger `k` smooths the decision boundary, making the model less sensitive to noise but increasing bias. Too large a `k` may lead to underfitting.

17. Why do we use `random_state` in train-test splitting and SVM?

Answer:

To ensure reproducibility. By fixing the random seed, we get consistent splits and model behavior across runs.

18. What does `classification_report` give you?**Answer:**

It provides detailed metrics including precision, recall, F1-score, and support for each class, which are essential for evaluating binary classifiers.

19. Could this classification model be used in production? What else would you check?**Answer:**

Yes, but further checks like cross-validation, parameter tuning, and robustness testing would be needed. Deployment pipeline and real-time performance should also be evaluated.

20. What preprocessing steps are usually required for email text data?**Answer:**

Common steps include tokenization, stopword removal, stemming or lemmatization, and vectorization (e.g., TF-IDF). In our dataset, this was already done before modeling.



1. Precision

Definition:

Out of all the emails that the model **predicted as spam**, how many were actually spam?

Formula:

Formula:

$$\text{Precision} = \frac{\text{True Positives (TP)}}{\text{True Positives} + \text{False Positives (TP + FP)}}$$

Interpretation:

High precision means when the model predicts spam, it's usually correct — useful when **false positives** are costly (e.g., marking important emails as spam).



2. Recall (Sensitivity or True Positive Rate)

Definition:

Out of all the actual spam emails, how many did the model correctly detect?

Formula:

Formula:

$$\text{Recall} = \frac{\text{True Positives (TP)}}{\text{True Positives + False Negatives (TP + FN)}}$$

Interpretation:

High recall means the model catches most spam — useful when **false negatives** are costly (missing spam emails).



3. F1-Score

Definition:

The **harmonic mean** of precision and recall. It balances the two metrics, especially useful when classes are imbalanced.

Formula:

Formula:

$$\text{F1 Score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

Interpretation:

High F1-score means the model has a good balance of precision and recall.



4. Support

Definition:

The number of **actual samples** in the dataset that belong to a class. For example in a spam dataset:

- Support for class "0 = Not Spam" might be 2500
- Support for class "1 = Spam" might be 1500

Interpretation:

Shows the size of each class. If one class has far more samples than the other, the dataset is **imbalanced**, and accuracy alone may be misleading.



Why These Metrics Matter?

- **Precision** cares about correctness of positive predictions.
- **Recall** cares about capturing all positive cases.
- **F1-score** balances both and is reliable for imbalanced datasets.
- **Support** shows how many actual samples belong to each class.



What is a Confusion Matrix?

A **confusion matrix** is a table used to evaluate the performance of a **classification model** by comparing its predictions with the actual labels. It helps us understand where the model is getting things right and where it is making mistakes.



Structure of a Confusion Matrix (for Binary Classification)

The matrix has **4 elements**, arranged like this:

	Predicted Not Spam (0)	Predicted Spam (1)
Actual Not Spam (0)	True Negative (TN)	False Positive (FP)
Actual Spam (1)	False Negative (FN)	True Positive (TP)

🔍 Your Confusion Matrix

```
array([
    [646,  93],
    [ 49, 247]
])
```

This can be interpreted as:

	Predicted Not Spam (0)	Predicted Spam (1)
Actual Not Spam (0)	646 (TN)	93 (FP)
Actual Spam (1)	49 (FN)	247 (TP)

✅ What Each Number Means:

1. True Negative (TN) = 646

- These are emails that were **not spam** and correctly classified as **not spam**.

2. False Positive (FP) = 93

- These are emails that were **not spam** but the model incorrectly predicted them as **spam**.
- Also called **Type I error**.

3. False Negative (FN) = 49

- These are emails that were **spam** but the model failed to detect them and classified as **not spam**.

- Also called **Type II error**.
4. **True Positive (TP) = 247**
- These are emails that were **spam** and correctly classified as **spam**.
-

⌚ Why is Confusion Matrix Useful?

It helps calculate:

- **Accuracy:** Overall correctness

$$\text{Accuracy} = \frac{TP + TN}{\text{Total samples}}$$

- **Precision:** How many predicted spams were actually spam
- **Recall:** How many actual spams were detected
- **F1-Score:** Balance between precision and recall

It gives a full picture of model performance — not just one number like accuracy.

🤖 What is KNN?

K-Nearest Neighbors (KNN) is a **supervised machine learning algorithm** used for **classification** (and also regression in some cases). In classification, it predicts the class of a new data point by looking at the **classes of its nearest neighbors** in the training dataset.

🧠 How KNN Works (Step-by-Step)

1. **Store the training data:** KNN keeps all the data points—it doesn't learn a model during training.
2. **For a new test point:**

- Calculate the **distance** between the test point and all points in the training data (usually using Euclidean distance).
 - Find the **K closest neighbors** (e.g., K = 3 means 3 nearest points).
 - Take a **majority vote** among the neighbors' classes.
 - Assign the class that appears most often.
-

Example

Suppose you're trying to classify an email as **spam (1)** or **not spam (0)**.

If **k = 3** and among the 3 nearest neighbors:

- 2 are labeled spam (1)
- 1 is labeled not spam (0)

Then the email is predicted as **spam** because most neighbors (2 out of 3) are spam.

Key Terms

- **K**: The number of neighbors to consider.
 - **Distance Metric**: A function to measure how close points are (e.g., Euclidean distance).
 - **Non-parametric**: KNN makes no assumptions about data distribution.
-

Strengths of KNN

- Simple and intuitive.
 - No training time (lazy learner).
 - Works well with small datasets.
-

Limitations of KNN

- Slow at prediction time (must compute distance to all training data).
 - Not ideal for **high-dimensional data** (like text or images).
 - Sensitive to the choice of K and distance metric.
-

Where KNN is Used

- Handwritten digit recognition (like MNIST)
 - **Recommendation systems**
 - Basic spam filtering (for comparison with better models like SVM)
-

In Your Practical

You used:

```
neigh = KNeighborsClassifier(n_neighbors=2)
```

This means the model checks the 2 nearest neighbors and classifies based on their majority class.

What is Support Vector Machine (SVM)?

Support Vector Machine (SVM) is a powerful **supervised machine learning algorithm** used for **classification** (and also regression). In classification tasks, it tries to find the best possible boundary (called a **hyperplane**) that separates the data points of different classes.

How SVM Works (In Simple Terms)

1. **Plot the data points** in space (2D, 3D, or higher dimensions).

-
2. **Find the hyperplane** (a line in 2D, a plane in 3D, etc.) that best separates the classes.
 3. The **best hyperplane** is the one that maximizes the **margin**, which is the distance between the hyperplane and the nearest data points from each class.
 4. These nearest points are called **support vectors**, and they “support” the position of the hyperplane.
-

Key Concepts

- **Hyperplane:** A decision boundary that separates data points of different classes.
 - **Margin:** The distance between the hyperplane and the closest data points from any class. SVM tries to maximize this margin.
 - **Support Vectors:** Data points that are closest to the hyperplane. These points influence its position.
-

Types of SVMs

1. **Linear SVM:**

- Used when data is linearly separable.
- Finds a straight line or flat plane to separate classes.

Example code from your project:

```
svm = SVC(kernel='linear', random_state=42)
```

2. **Non-linear SVM:**

- Used when data isn't perfectly separable by a straight line.
 - Applies a **kernel function** (like RBF or polynomial) to project data into a higher dimension where it becomes separable.
-

Advantages of SVM

- Works well with **high-dimensional data** (e.g. text classification like spam detection).
 - Effective even with a clear margin of separation.
 - Less prone to overfitting (especially with high dimensionality and proper kernel choice).
-

Limitations of SVM

- Can be slow with very large datasets.
 - Performance depends on the choice of the kernel and parameters.
 - Doesn't work well if classes overlap significantly.
-

In Your Email Spam Classification Project

You used SVM like this:

```
svm = SVC(kernel='linear', random_state=42)
svm.fit(X_train, y_train)
```

This means the model:

- Uses a **linear kernel**, good for text data.
- Trains on feature vectors extracted from emails.
- Predicts whether an email is spam (1) or not spam (0).

Applications of SVM

- Document Categorization
- Medical Diagnosis

Classifying tumors as benign or malignant.

Predicting heart disease based on clinical parameters.

- Time-Series Forecasting (SVM Regression)

Predicting weather conditions.

Forecasting energy consumption.

=====

3. Gradient Descent

Gradient Descent is an iterative optimization algorithm used to find a local minimum of a function.

At each step you:

1. Compute the gradient (derivative) of the function at the current point (this tells you which direction increases the function).
2. Move **opposite** the gradient by a small step (the step size is the learning rate).
3. Repeat until the steps become tiny or a maximum number of iterations is reached.

Write the function and its derivative.

- $f(x) = (x + 3)^2$
- $f'(x) = 2(x + 3)$

Choose:

- a starting point: $x = 2$ (given)
- a learning rate α (I used 0.1 — you can try others)
- number of iterations or a tolerance for stopping

Iterate this update rule:

- $x_{\text{new}} = x_{\text{old}} - \alpha * f'(x_{\text{old}})$

Stop when either you reach max iterations or the update is tiny.

Report final x and $f(x)$. For this problem the true minimum is $x = -3$, $f(-3) = 0$.

Q1. What is Gradient Descent?

A: Gradient Descent is an optimization algorithm used to find the minimum value of a function by iteratively moving in the direction opposite to the gradient (slope) of the function.

Q2. What is the main idea behind Gradient Descent?

A: The idea is to start from an initial point and update the variable xxx step by step toward the point where the slope (gradient) becomes zero — that is, the local minimum.

Q3. What is a gradient?

A: A gradient is the derivative (slope) of a function. It indicates the direction and rate of fastest increase of the function. In 1D, it is $f'(x)f'(x)f'(x)$; in multiple dimensions, it's a vector of partial derivatives.

Q4. Why do we subtract the gradient in the update rule?

A: We subtract the gradient because we want to move in the opposite direction of the slope — toward the function's minimum (the slope points toward increasing values).

Q5. What is the mathematical update rule used in gradient descent?

A:

$$x_{\text{new}} = x_{\text{old}} - \text{learning rate} \times f'(x)$$

Q6. What is the role of the learning rate (lr)?

A: The learning rate controls the step size of each update:

- Too high → may overshoot the minimum.
 - Too low → convergence will be very slow.
-

Q7. How do we know when the algorithm has reached the local minimum?

A: When the change in xxx becomes very small or the gradient (slope) becomes nearly zero, the algorithm has reached or is close to the local minimum.

 **Questions Based on Your Code****Q8. What function are we minimizing in this program?**

A: We are minimizing a quadratic function of the form $f(x)=ax^2+bx+c$ $f(x) = ax^2 + bx +$

$cf(x) = ax^2 + bx + c$.

Example: $f(x) = (x+3)^2 = x^2 + 6x + 9$

Q9. What is the derivative (gradient) of the function $f(x) = ax^2 + bx + c$? $cf(x) = ax^2 + bx + c$?

A: $f'(x) = 2ax + b$

Q10. What is the purpose of the loop in the program?

A: The loop repeatedly updates xxx using the gradient descent rule to approach the local minimum step by step.

Q11. What does the variable `path` store in your code?

A: `path` stores all intermediate values of xxx during each iteration — this helps to plot and visualize the path of convergence.

Q12. What will the plot show?

A:

- The **curve** of the function $f(x)$.
- Red dots ('`ro-`') showing how the gradient descent algorithm moves step by step toward the minimum point.

Q13. What is the expected local minimum for $f(x) = (x+3)^2$? $f(x) = (x + 3)^2$

A: The local minimum is at $x = -3$ with $f(-3) = 0$

Q14. What happens if the learning rate is too large?

A: The algorithm might **overshoot** the minimum and fail to converge (it may oscillate or diverge).

Q15. What happens if the learning rate is too small?

A: The convergence will be **very slow**, requiring many iterations to reach near the minimum.

Practical & Theoretical Mix

Q16. Is Gradient Descent guaranteed to find the global minimum?

A: For convex functions (like quadratics), yes — it always finds the global minimum. For non-convex functions, it might get stuck in a local minimum.

Q17. What is the stopping criterion in your code?

A: The code runs for a fixed number of iterations (`iters`). Alternatively, we could stop when the gradient becomes very small.

Q18. Why do we use `numpy` and `matplotlib` in this program?

A:

- **NumPy**: for numerical operations like arrays and generating points.
 - **Matplotlib**: for plotting the function and gradient descent path.
-

Q19. Can Gradient Descent be used for other types of functions?

A: Yes, gradient descent is a general optimization method used in many areas, especially in **machine learning** for minimizing cost functions.

Q20. What are some real-world **applications** of Gradient Descent?

A:

- **Training machine learning models** (like linear regression, neural networks)
 - **Optimization problems** in economics and engineering
 - **Curve fitting** and parameter estimation
 - **Autonomous Vehicles** - Path optimization and control.
-
-

4. Sales Data - KMeans Clustering

1. What is clustering?

Answer:

Clustering is an unsupervised machine learning technique that groups similar data points together based on their features. The goal is to ensure that data points within the same cluster are more similar to each other than to those in other clusters.

2. What is K-Means clustering?

Answer:

K-Means is an iterative algorithm that partitions data into **K clusters** by minimizing the **sum of squared distances** between each data point and its assigned cluster center (centroid).

3. What is the “K” in K-Means?

Answer:

“K” represents the **number of clusters** we want to form in the data. It is a user-defined parameter determined using methods like the **Elbow Method**.

4. What does the Elbow Method do?

Answer:

The Elbow Method helps determine the **optimal number of clusters (K)** by plotting the inertia (sum of squared distances) for different values of K. The point where the decrease in inertia slows down (forming an “elbow” shape) is the ideal number of clusters.

5. What is inertia in K-Means?

Answer:

Inertia is the **sum of squared distances between each data point and its nearest cluster centroid**. Lower inertia means tighter clusters, but too small K can lead to overfitting.

Code-Based Questions

6. Why did you use StandardScaler?

Answer:

StandardScaler standardizes data so that each feature has a mean of 0 and a standard deviation of 1.

This ensures that features like **SALES** (large values) don't dominate features like **PRICEEACH** (smaller values) during clustering.

7. Why did you use LabelEncoder for STATUS and DEALSIZE?

Answer:

Both **STATUS** and **DEALSIZE** are categorical.

LabelEncoder converts them into numerical codes (like 0, 1, 2), so K-Means can process them mathematically.

8. Why do we scale data before clustering?

Answer:

K-Means uses **Euclidean distance** to measure similarity.

If features have different ranges, one large-scale feature can dominate others, leading to incorrect clusters. Scaling fixes this.

9. What are centroids in K-Means?

Answer:

Centroids are the mean positions of all the points that belong to a particular cluster. They represent the **center of each cluster**.

10. How does K-Means work step by step?

Answer:

1. Choose the number of clusters (K).
 2. Randomly initialize K centroids.
 3. Assign each data point to the nearest centroid.
 4. Recalculate the centroids based on assigned points.
 5. Repeat steps 3–4 until centroids stabilize or a maximum number of iterations is reached.
-



Dataset & Result Questions

11. What dataset did you use and what does it contain?

Answer:

I used the `sales_data_sample.csv` dataset from Kaggle. It contains details of product sales like `QUANTITYORDERED`, `PRICEEACH`, `SALES`, `MSRP`, `DEALSIZE`, `STATUS`, etc., representing sales transactions.

12. Which features did you use for clustering?

Answer:

I selected these features:
`QUANTITYORDERED`, `PRICEEACH`, `SALES`, `MSRP`, `ORDERLINENUMBER`, `QTR_ID`, `MONTH_ID`, `YEAR_ID`, `STATUS`, and `DEALSIZE`.

13. How many clusters did you finally choose and why?

Answer:

I used the **Elbow Method** to find the point where the inertia started flattening. From the elbow curve, **3 clusters** seemed optimal for this dataset.

14. What do your clusters represent?

Answer:

Each cluster represents a group of sales records with similar characteristics — for example:

- Cluster 0: High sales, large deal size
- Cluster 1: Medium sales, medium deals
- Cluster 2: Low sales, small deals

(Interpretation may vary depending on dataset patterns.)



Visualization & Insights

15. How did you visualize the clusters?

Answer:

I used a **2D scatter plot** (e.g., **SALES** vs **PRICEEACH**) and **PCA (Principal Component Analysis)** to reduce features to 2D for visualization. Different colors represent different clusters.

16. What insights can we get from clustering this data?

Answer:

Clustering helps identify patterns in sales — such as which orders are high-value, which are small, or which quarters have similar sales behavior.

This information can help in **customer segmentation, pricing strategies, and inventory planning**.



Advanced / Conceptual

17. What are the limitations of K-Means?

Answer:

- You must predefine K (number of clusters).
 - Sensitive to **outliers**.
 - Works best for **spherical** and **equal-sized** clusters.
 - Performs poorly with categorical or non-linear data.
-

18. What is the difference between K-Means and Hierarchical Clustering?

Answer:

Feature	K-Means	Hierarchical
Cluster number	Must be predefined (K)	Dendrogram can help decide
Speed	Faster	Slower for large data
Type	Partitional	Agglomerative/Divisive
Visualization	Needs PCA or 2D plots	Dendrogram gives visual hierarchy

19. What are some alternatives to K-Means?

Answer:

- **Hierarchical Clustering**
- **DBSCAN (Density-Based Spatial Clustering)**
- **Gaussian Mixture Models (GMM)**

20. Can K-Means be used for categorical data only?

Answer:

No, K-Means works best with **numerical data**.

For categorical data, algorithms like **K-Modes** or **K-Prototypes** are better suited.