**Code-Specific Questions and Answers**

**Data Loading and Initial Exploration**

1. **What is the purpose of data = pd.read_csv('uber.csv')?**

   o **This line loads the Uber dataset into a pandas DataFrame, allowing us to manipulate and analyze the data.**

2. **Why did you use data.dropna(inplace=True) immediately after loading the data?**

   o **Dropping rows with missing values (NaNs) ensures that we only work with complete data, which helps prevent issues during model training and evaluation.**

3. **What does data.head() do?**

   o **data.head() displays the first few rows of the dataset, allowing us to quickly inspect its structure, data types, and sample values.**

4. **What is the role of column_names = data.columns?**

   o **This line retrieves the column names in the dataset and stores them in column_names, which can be useful for referencing or verifying column names in the dataset.**

5. **What is the purpose of data.info()?**

   o **data.info() provides a concise summary of the DataFrame, including data types, number of non-null values, and memory usage, which helps us understand the dataset's structure and check for missing data.**

**Dropping Unnecessary Columns**

6. **What does data.drop(columns=['Unnamed: 0', 'key'], inplace=True) achieve?**

   o **This line removes the Unnamed: 0 and key columns from the DataFrame, as they do not contain information useful for predicting fares.**

**Handling Missing Values**

7. **Why did you use data.isnull().sum() after dropping unnecessary columns?**

   o **data.isnull().sum() checks for any remaining missing values in each column, ensuring that no incomplete data remains in the dataset.**

**Date and Time Processing**

8. **Why did you convert pickup_datetime to datetime format using pd.to_datetime?**

   o **Converting pickup_datetime to datetime format allows us to easily extract components like the hour, day, and month, which may influence fare prices.**

9. **Explain data['hour'] = data['pickup_datetime'].dt.hour.**

   o **This line extracts the hour from pickup_datetime and stores it in a new column hour, which can help identify time-based patterns in fares.**

10. **Why did you drop the pickup_datetime column?**

- After extracting time-based features (hour, day, month), pickup_datetime was no longer needed, so it was dropped to reduce dataset complexity.

**Feature Scaling**

11. **What is the purpose of scaler = StandardScaler()?**

    - StandardScaler standardizes numerical features to have a mean of 0 and a standard deviation of 1, which helps improve the performance and convergence of many machine learning algorithms.

12. **Why did you select only certain columns for scaling in data[numerical_features] = scaler.fit_transform(data[numerical_features])?**

    - Only numerical features (fare_amount, pickup_longitude, pickup_latitude, dropoff_longitude, dropoff_latitude, passenger_count) were scaled, as scaling categorical data or date-related columns is generally not meaningful.

**Splitting Data into Features and Target**

13. **What does X = data.drop('fare_amount', axis=1) do?**

    - This line assigns all columns except fare_amount to X (features), as fare_amount is the target variable we aim to predict.

14. **Why is y = data['fare_amount'] defined separately from X?**

    - y is set to fare_amount, the target variable, so we can use it independently when training the model to predict this variable based on the features in X.

15. **What is the purpose of train_test_split(X, y, test_size=0.2, random_state=42)?**

    - This function splits the data into 80% training and 20% test sets to evaluate the model's performance on unseen data. Setting random_state=42 ensures the split is reproducible.

**Outlier Analysis**

16. **What does sns.boxplot(x=data['fare_amount']) visualize?**

    - This line generates a boxplot of fare_amount, helping us identify outliers in fare values, which can impact model training.

**Correlation Matrix and Heatmap**

17. **Why did you use data.corr() to create a correlation matrix?**

    - The correlation matrix reveals relationships between variables, helping us identify which features may strongly influence the target (fare_amount).

18. **Explain the purpose of sns.heatmap(corr_matrix, annot=True, cmap='coolwarm').**

    - This line visualizes the correlation matrix as a heatmap, with annotations for values. The color gradient helps quickly identify positive and negative correlations.

**Model Training and Prediction**

19. **What does linear_model = LinearRegression() achieve?**

- This line initializes a Linear Regression model, which assumes a linear relationship between features and fare_amount.

20. **What does rf_model = RandomForestRegressor(n_estimators=100, random_state=42) do?**

    - This line initializes a Random Forest model with 100 decision trees (n_estimators=100), a powerful method for capturing non-linear patterns in the data.

21. **Why is fit called on both linear_model and rf_model with X_train and y_train?**

    - The fit method trains each model on the training data (X_train and y_train), allowing them to learn patterns that will be used to predict fare_amount on new data.

## Model Evaluation

22. **What is y_pred_linear = linear_model.predict(X_test) used for?**

    - This line generates predictions from the Linear Regression model on the test data, allowing us to evaluate its performance.

23. **Explain the purpose of r2_score, mean_squared_error, and mean_absolute_error.**

    - These metrics evaluate model performance:

        - r2_score measures how well the model explains the variance in the target variable.

        - mean_squared_error calculates the average squared difference between predictions and actual values.

        - mean_absolute_error computes the average absolute error in predictions, giving a straightforward measure of prediction accuracy.

24. **Why do you print metrics for both Linear Regression and Random Forest models?**

    - Printing both models' metrics allows for a performance comparison. Generally, a higher R² and lower RMSE and MAE indicate a better model.

25. **What conclusions did you draw from comparing r2, RMSE, and MAE between the two models?**

    - By comparing these metrics, we can identify which model performed better in terms of predictive accuracy and suitability for Uber fare prediction.