# QUESTION RELATED TO TASK 2 (AIR POLLUTION PREDICTION)

Ques 1. What challenges did you encounter while preparing the dataset for preprocessing and EDA analysis?

Ans. While preparing the dataset for preprocessing and exploratory data analysis (EDA) analysis, several challenges and issues can be identified based on the provided code and dataset:

1. Missing Data Handling:
   - The dataset contains missing values in the "pm2.5" column. Handling missing data is a common challenge in data preprocessing. In this case, the missing values are filled with the mean of the "pm2.5" column. While this is a reasonable imputation strategy, it's important to consider the potential impact on data quality and analysis.
2. Data Types and Encoding:
   - The "cbwd" column is of data type 'object,' indicating it likely contains categorical data. Depending on the analysis, you may need to encode categorical variables for machine learning models or visualization.
3. Data Visualization Challenges:
   - The code includes various data visualization commands (e.g., bar plots, scatter plots, and count plots), which can be challenging when dealing with a large dataset like this one (43,824 entries).
   - Visualizing such a large dataset can lead to cluttered and hard-to-interpret plots. Consider subsampling the data or using aggregation techniques to create meaningful visualizations.
4. Data Distribution Understanding:
   - The code generates a bar plot to visualize the distribution of data by month. Understanding the data distribution is essential for EDA, but it may require further analysis to identify trends, seasonality, or outliers.

5. Scatter Plot Challenges: The scatter plot of "TEMP" vs. "PRES" is created, but it's unclear what insights or patterns the plot is meant to reveal. Proper labeling and context are essential to interpret such plots effectively.

6. Data Scaling and Normalization: Depending on the analysis or modeling techniques you plan to use, scaling or normalizing numerical features (e.g., "TEMP," "PRES," "Iws") may be necessary.

7. Handling Outliers: Detecting and handling outliers is a crucial step in data preprocessing. Outliers can significantly affect the results of statistical analysis and machine learning models.

8. Data Exploration Objectives: The provided code shows various data visualization commands, but it's essential to have specific objectives in mind when conducting EDA. What questions or insights are you trying to gain from the data? Defining clear objectives can guide your analysis.

9. Data Volume: The dataset has a considerable volume of data (4.3+ MB). Depending on your computing resources and analysis tools, working with large datasets can be computationally intensive.

10. Interpreting and Communicating Results: EDA is not just about generating plots; it's also about interpreting and communicating the results effectively. Ensure that you can derive meaningful insights from the visualizations and communicate them to stakeholders or colleagues.

Ques 2. Describe the difference observed while modeling with linear and logistic regressor.

Ans.

| Characteristic | Linear Regression | Logistic Regression |
|---|---|---|
| Nature of Target Variable | Continuous and numeric | Binary or categorical |
| Model Output | Continuous values | Probability (0 to 1) |
| Equation and Hypothesis | Simple linear equation | Logistic (sigmoid) |
| Application | Regression tasks | Classification tasks |
| Evaluation Metrics | MAE, MSE, R-squared | Accuracy, Precision, Recall, F1-score, ROC-AUC |
| Interpretability | Coefficients indicate relationship with target variable | Coefficients indicate log-odds of the probability |

Ques 3. How did you evaluate the performance of the linear and logistic regressor?

Ans.

Here's how the performance of both types of regressors is commonly evaluated:

For Linear Regression:
  1. Regression Metrics:
  ● Mean Absolute Error (MAE): This metric measures the average absolute difference between the predicted values and the actual target values. Lower MAE indicates better performance.
  ● Mean Squared Error (MSE): MSE calculates the average of the squared differences between predicted and actual values. Lower MSE indicates better performance.

- Root Mean Squared Error (RMSE): RMSE is the square root of MSE and provides the error in the same units as the target variable.
- R-squared (R^2): R-squared measures the proportion of the variance in the target variable that is explained by the model. A higher R-squared value indicates a better fit.

2. Visualization:
    - Scatter plots: Visualize the predicted values against the actual target values to visually assess how well the model predictions align with the true values.

For Logistic Regression:
1. Classification Metrics:
    - Accuracy: Accuracy measures the proportion of correctly classified instances. It is suitable for balanced datasets.
    - Precision: Precision measures the ratio of true positive predictions to the total positive predictions. It is essential when minimizing false positives is crucial.
    - Recall (Sensitivity): Recall measures the ratio of true positive predictions to the total actual positive instances. It is important when minimizing false negatives is crucial. F1-Score: The
    - F1-Score is the harmonic mean of precision and recall, balancing both metrics.
    - Receiver Operating Characteristic Area Under the Curve (ROC-AUC): ROC-AUC measures the model's ability to distinguish between classes, especially useful in imbalanced datasets.

2. Confusion Matrix:
    - Construct and analyze the confusion matrix to understand the number of true positives, true negatives, false positives, and false negatives. It provides insights into specific types of errors made by the classifier.

3. ROC Curve and Precision-Recall Curve:

- Plot the ROC curve to visualize the trade-off between the true positive rate and the false positive rate at different probability thresholds.
- Plot the Precision-Recall curve to understand the trade-off between precision and recall at different probability thresholds.

4. Threshold Tuning:

   - Depending on the problem and the business objectives, you may need to adjust the probability threshold for classification to optimize the trade-off between precision and recall.

     The choice of evaluation metrics depends on the specific problem and the objectives of the analysis. For Linear Regression, the focus is on the accuracy of numeric predictions, while for Logistic Regression, the focus is on the accuracy of class predictions. Carefully selecting and interpreting the appropriate metrics is crucial for assessing model performance accurately.

Ques 4. Did you observe any signs of overfitting during training? If so, how did you handle it?

Ans.

Yes, based on the accuracy scores for the Ridge Regressor and Random Forest Regressor models, there are signs of overfitting observed during training. Here's the evidence of overfitting and how it's handled:

1. Evidence of Overfitting:
   - For both the Ridge Regressor and Random Forest Regressor models, the accuracy on the training dataset is noticeably higher than the accuracy on the test dataset.
   - Ridge Regressor Training Accuracy: 0.264
   - Ridge Regressor Test Accuracy: 0.277
   - Random Forest Regressor Training Accuracy: 0.371

- Random Forest Regressor Test Accuracy: 0.359

2. Handling Overfitting:
   - The code includes a check to compare the training accuracy and test accuracy for each model and prints a message indicating whether the model is overfitted or not. If the training accuracy is greater than the test accuracy, it suggests that the model is overfitting.
3. Handling Strategies:
   - The code does not explicitly include strategies to handle overfitting. However, identifying overfitting is the first step, and the next step would be to implement strategies to mitigate it.
   - Possible strategies to address overfitting may include:
   - Regularization: For Ridge Regression, increasing the regularization strength can help reduce overfitting.
   - Feature Selection: Identify and remove irrelevant or redundant features that might contribute to overfitting.
   - Cross-Validation: Optimize hyperparameters using cross-validation to ensure the model generalizes well.
   - Ensemble Methods: For Random Forest, tuning hyperparameters like max depth, minimum samples per leaf, or the number of estimators can help control overfitting.

Ques 5. How well did your model perform on the testing set compared to the training set?

Ans.

To evaluate how well the models performed on the testing set compared to the training set, we can analyze the mean absolute error (MAE), mean squared error (MSE), and root mean squared error (RMSE) for both the training and testing datasets. Here are the results:
 For Ridge Regressor:
   1. Training Data:

- MAE: 0.6119
- MSE: 0.7399
- RMSE: 0.8602

2. Testing Data:
- MAE: 0.6119
- MSE: 0.7075
- RMSE: 0.8411

For Random Forest Regressor:
1. Training Data:
- MAE: 0.5635
- MSE: 0.6321
- RMSE: 0.7950
2. Testing Data:
- MAE: 0.5672
- MSE: 0.6275
- RMSE: 0.7922

Comparing Training vs. Testing Performance:
- In general, both models exhibit slightly better performance on the training dataset compared to the testing dataset.

- This is evident from the fact that the MAE, MSE, and RMSE values for the training dataset are slightly lower than those for the testing dataset in both cases

- . Lower error metrics on the training dataset are expected because the models were trained on that data and may have learned to fit it well.