

. Second Project:

1. Data Exploration:

- The code starts by loading the dataset and checking its basic information using `info()` and `shape`.

2. Handling Missing Values:

- The code identifies and drops columns with missing values below a certain threshold.
- It visualizes missing values using a bar plot.

3. Data Preprocessing:

- Certain columns are filtered based on specific conditions (e.g., filtering out rows where PNEUMONIA is neither 1 nor 2).
- A new column DEATH is created based on the DATE_DIED column.
- Some columns are dropped (INTUBED, ICU, DATE_DIED) as they are deemed irrelevant for the analysis.

4. Exploratory Data Analysis (EDA):

- Several visualizations are created to understand the distribution of data across different features, medical units, gender, age, health conditions, etc.

5. Modelling:

- Several classification models (Logistic Regression, Random Forest, Gradient Boosting, SVM, KNN, Decision Tree) are trained and evaluated.
- Hyperparameter tuning is performed using GridSearchCV for each model.

- Logistic Regression:

Certainly! Logistic Regression is a statistical method commonly used for binary classification problems, where the goal is to predict the probability of an instance belonging to a particular class. Despite its name, it's used for classification rather than regression tasks.

- Random Forest:

Feature Importance: Random Forest provides a feature importance score, indicating the significance of each feature in making predictions. This can help identify the most influential

features in predicting 'seniorityAsMonths.' Handles Non-Linearity: Random Forest can model complex non-linear relationships between input features (x) and the target variable (y). It can capture interactions and patterns that may be missed by linear models.

-Gradient Boosting:

Outlier Robustness: Gradient Boosting can be less sensitive to outliers compared to some other algorithms. The ensemble nature of the model helps mitigate the impact of individual data points. Reduces Bias and Variance: Gradient Boosting reduces both bias and variance, making it less prone to overfitting compared to individual decision trees. Robust to Outliers: It's generally robust to outliers due to the nature of fitting trees to residuals. Optimized for Regression Tasks: Gradient Boosting is particularly well-suited for regression tasks, making it a suitable choice for predicting 'seniorityAsMonths.'

- SVM:

Support Vector Machine (SVM) is a supervised machine learning algorithm that can be used for both classification and regression tasks. SVM is particularly effective in high-dimensional spaces and is well-suited for tasks where there is a clear margin of separation between different classes.

-Decision Tree:

Feature Importance: Decision Trees provide a natural way to assess the importance of each feature in predicting the target variable. You can examine the tree structure to understand which features are used for splitting nodes and, consequently, have a stronger impact on predicting seniority in months.

-KNN:

is used for making predictions based on the majority class (for classification) or the mean/median (for regression) of the K-nearest neighbours of a data point.

6. Model Evaluation:

- Model accuracy is assessed, and confusion matrices are visualized for some models.
- ROC curves are plotted for Logistic Regression.

7. Conclusion:

- The code covers data loading, preprocessing, exploratory data analysis, and building and evaluating classification models.
- The choice of models allows for a comprehensive understanding of the dataset.

Remember to interpret the results critically, considering the domain and the specific goals of your analysis.