

## EDA:

- 1) Stats-based Analysis
- 2) Graph-Based Analysis
- 3) Pre-processing:

## Model training

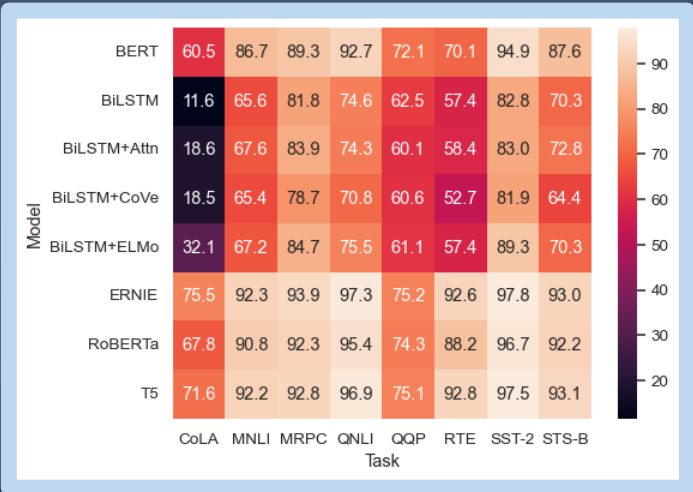
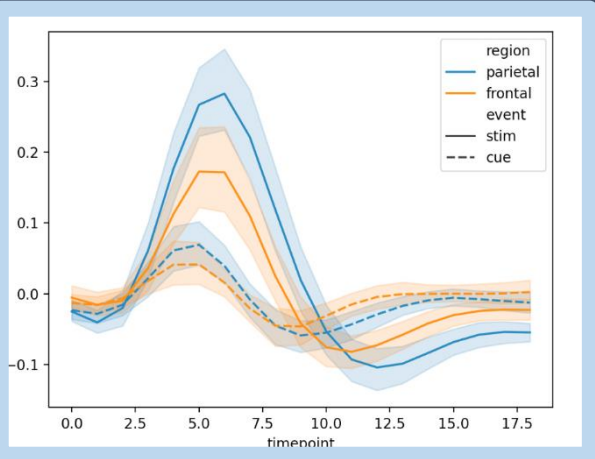
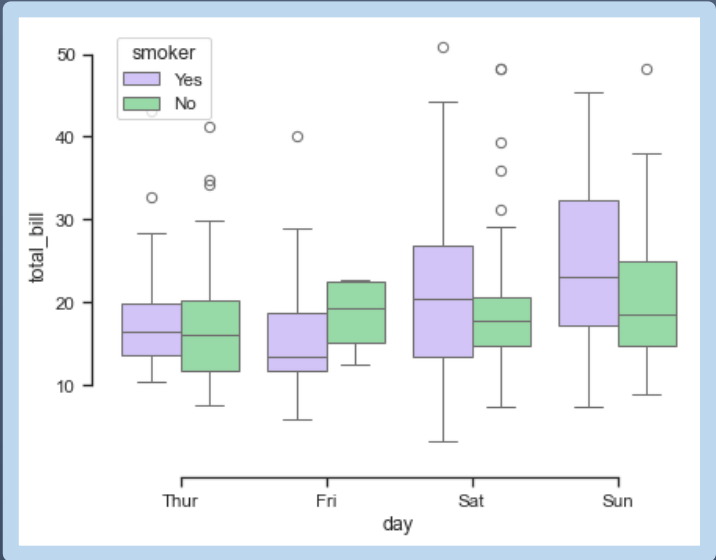
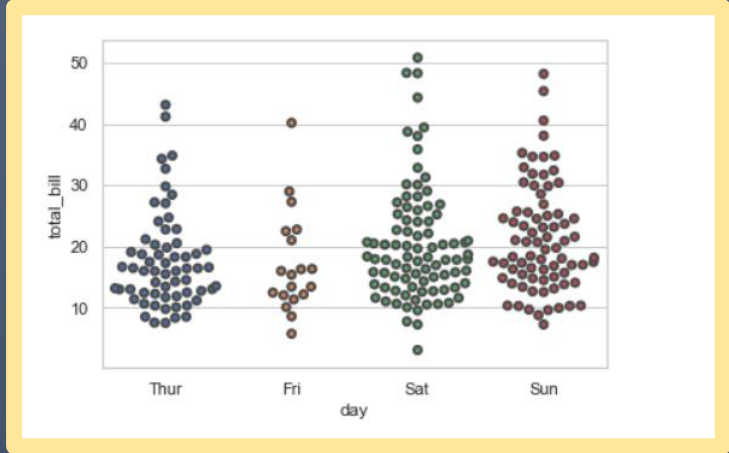
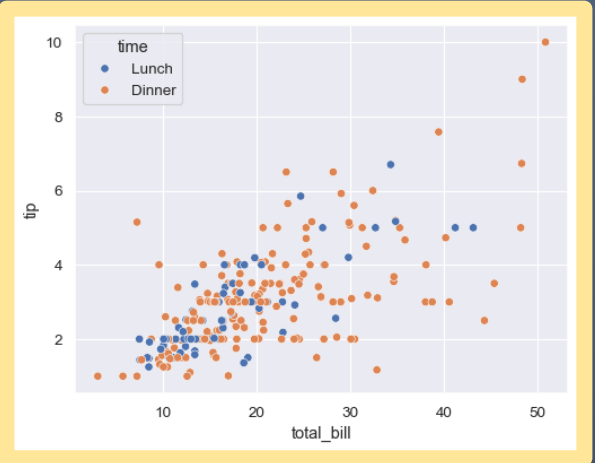
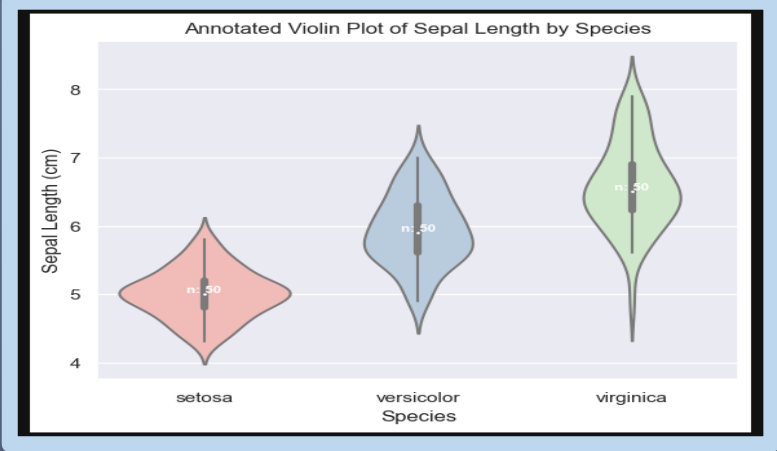
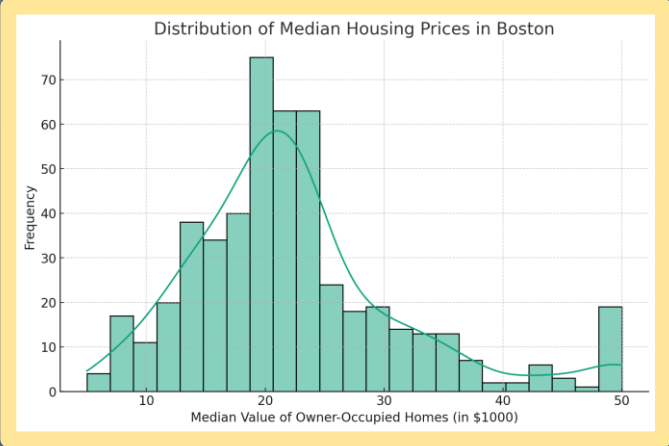
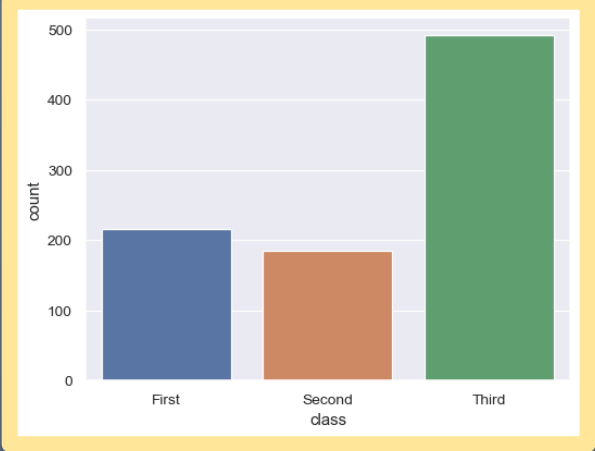
- 1) Training Pipeline
- 2) Prediction pipeline

EDA:

1) Stats-based Analysis

A) Mean, median, std etc.

Graph



## **Pre-processing:**

- 1) Handle Missing Value
- 2) Handel Duplicate Values
- 3) Outlier Handling
- 4) Handel Imbalance Data
- 5) Feature Selection
- 6) Dimension Reduction
- 7) Encoding
- 8) Scaling

### Handle Missing Value:

- 1) Random value
- 2) Forward filling/backward filling
- 3) Statistical approach (mean, mode, median)
- 4) End of the distribution
- 5) Create “your own” ML model to predict missing value

### Handle Duplicate Values:

- 1) Drop the duplicate

### Handling Imbalance Dataset

- 1) Collect more data
- 2) Under sampling
- 3) Over sampling
- 4) Cluster based oversampling

### Feature Selection

- 1) Correlation:
- 2) Variance Threshold:
- 3) Chi-Square Method:
- 4) ANOVA (Analysis of Variance):
- 5) Information Gain:

### Dimension Reduction

- 1) (PCA, LDA, TSNE)
- 2) Split/Merge/Drop/Add

### Outlier Handling:

- a) Detect the Outlier
  - 1) z-Score
  - 2) IQR
  - 3) Boxplot
  - 4) Scatter plot
  - 5) Violin plot
- b) Removing Outlier
  - 1) Drop Outlier
  - 2) Replace with other value

### Encoding

- 1) One hot encoding
- 2) Ordinal Encoding

### Scaling:

- 1) Standardization
- 2) Normalization (Min-max scaler)

# Model training

## 1) Training Pipeline

- a) Choose model based on problem
- b) Hyper parameter tuning
- c) Select the best model
- d) Save model as pkl file

## 2) Prediction pipeline

- a) Create the API to predict the new data
- b) Handover the API to user (aws, or any other cloud)
- c) API ready for prediction