

Project

Predict Fare of Airlines Tickets

①. Data preprocessing

→ Dealing with missing values (null) { Using dropna }

→ Converting the 'object' datatype to 'Date-Time format' for the ~~columns~~ required columns like 'Arrival Time', 'Departure Time'.

* User-defined function to convert 'object' to Date-time format
OR def change_into_date_time(column)

train_df[column] = pd.to_datetime(train_df[column])

{
def change_into_date_time(col)
train_df[col] = pd.to_datetime(train_df[col])
}

(To convert multiple columns run loop with i for all column names.)

→ Splitting the timestamp into (Date) | (Month) | (Year)

→ To extract hour, minutes from Departure / Arrival Time columns.

function def extract_hour(train_df, col):
train_df[col].dt.hour
OR

train_df[col + '_hour'] = train_df[col].dt.hour

Same function of time

```
def extract_min (train_df, col):
    train_df[col].dt.
    train_df[col + '_minute'] = train_df[col].dt.minute
```

→ In 'Duration' column, we have to iterate the values of hours and minutes.

Where there is no minute add '0m' and where there is no hour, add '0h'.

eg $2h \quad 5m \rightarrow$ will make a list $\rightarrow '2h', '5m'$.
 Delimiter is ' ' (space).

```
duration = list (train_df['Duration'])
for i in range (len (duration)):
    if len (duration[i].split (' ')) == 2:
        pass
    else:
        if 'h' in duration[i]:
            duration[i] = duration[i] + '0m'
        else:
            duration[i] = '0h' + duration[i]
```

Now to extract only numeric value in 'duration' list we will use split

eg $\begin{matrix} '2h' & '5m' \\ [0][0] & [0][0] \\ [0] & [1] \end{matrix} \left\{ \begin{array}{l} \therefore \text{To extract 2 :- } \text{split}('')[0][0:-1] \\ \therefore \text{To extract 5 :- } \text{split}('')[1][0:-1] \end{array} \right.$

Put these logic in functions.

→ To differentiate between categorical data columns and numeric data columns.

Following statement is an example of a code of 'List Comprehension' in Python.

```
cat_col = [col for col in train_df.columns if
            train_df[col].dtype == 'object']

cat_col print(cat_col)
```

[Out] → All column names with categorical data in DataFrame
(This means,)

col (for col in train_df.columns) if train_df[col].dtype == 'object'

Iteration part Condition part
dtype 'o' means 'object' type

If condition == true then only column will be column will be considered in list (output)

→ Dealing with categorical data

Categorical data

Nominal Data

(Those data that have no Hierarchy/Order)



One Hot Encoding To be done

Ordinal Data

(Those data that have some hierarchy of Hierarchy)

Label encoding To be Done

✓ # Nominal Data → One Hot Encoding
 # Ordinal Data → Label Encoding

→ Visualization of 'Airlines' using boxplot.
 and
 (Total-stops)

→ Creating Dummy variables (Doing One Hot Encoding) on categorical columns → Airlines, Source, Destination

Set (drop_first = True) ← This will create $(k-1)$ dummy variables, if k is the total count.

→ This parameter helps to reduce correlations created among the dummy variables.

→ 'Route column' - Split the city codes like BLR, DEL using delimiter '→' then put them in different columns.

Perform label encoding on all the columns.

Replace the null values using 'None'

→ 'Total-stops' column i) Find unique values in column.

ii) Put them in dictionary and map it to replace with values → 0, 1, 2, 3, 4 → 0, 1, 2, 3, 4.

~~See below~~ // Do not change Order, else ERROR

② How to Deal with Outliers

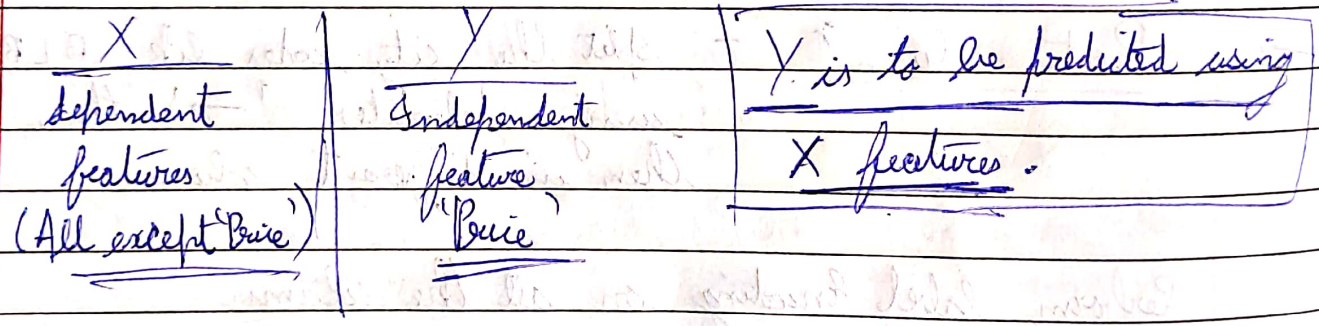
- (i) Distribution approach
- (ii) Box plot approach

Function to display the distribution plot and the box plot

```
def plot(df, col):
    fig, (ax1, ax2) = plt.subplots(2, 1)
    sns.sns.distplot(df[col], ax=ax1)
    sns.sns.boxplot(df[col], ax=ax2)
```

This function helps to get the ~~outlier~~ outliers position

③ Separating out dependent and independent feature



④ To find best features

Use mutual-info-classif module to get the importance count of the columns.

⑤ Apply Train Test Split

⑥. Use Random Forest Regressor → 81% *
Best algo

⑦. Use Linear Regression → 63%

⑧. Use KNNs → 68%

⑨. Use Decision Tree → 69%

⑩. Perform ~~hyper~~ Hyperparameter Tuning

Declaring n_estimators and max_depth

n_estimators = [int(x) for x in np.linspace(start=100, stop=1200, num=6)]
max_depth = [int(x) for x in np.linspace(start=5, stop=30, num=4)]

→ Creating dictionary to feed the parameters to the Randomized Search CV.

```
random_params = {
    'n_estimators': n_estimators
    'max_features': ['auto', 'sqrt']
    'max_depth': max_depth
    'min_sample_split': [5, 10, 15, 1000]
}
```

⑪. Use Hyper-tuned ~~Best~~ Random Forest Regressor Algo → 84%

⑫. Save model locally using pickle

⑬. Load Testing Data → Load Pre-trained Model → Preprocess Test Data → Perform Predictions → Save the new xlsx file locally.