

Airline Passenger Satisfaction

2440009654 - Darnell Kikoo

2440009162 - Matthew Adrianus Mulyono

2440076450 - Winson Allan Wijaya

Table of Contents

03

Problem Definition

Why did we do this?

04

Data Preparation

The chosen dataset,
and how we prepared
it.

05

EDA

...stands for
Exploratory Data
Analysis

12

Data Preprocessing

Preparing the data for
the model

16

Modeling & Evaluation

Creating the model

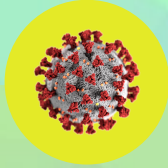
22

Results

...with the discussion



Problem Definition



COVID-19

Due to the rampant pandemic, the airline industry experienced plummeting profits.



Living Standards

However, living standards have in turn, improved as people took care of themselves more



Demand

These two factors results in airline industries desiring to improve their services.

Data Preparation

Dataset: Airline Passenger Satisfaction by TJ Klein*



Remove Outliers

We remove columns that we deem as unnecessary

Handle Missing Values

The categorical ones are replaced with its **mode**, while numericals are replaced with its **mean**.

Data Conversion

We convert categorical data to numerical and vice versa.

Exploratory Data Analysis

Airline Passenger Dataset

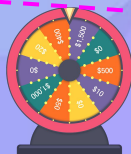


**Categorical
Data**

**Numerical
Data**



Ratings



Non-ratings

Exploratory Data Analysis

1.

Categorical Label Distribution

We find out the label distribution in both the ratings and non ratings categorical labels using bar graphs

2.

Numerical Label Distribution

This is done using *kdeplot* and *boxplot*



3.

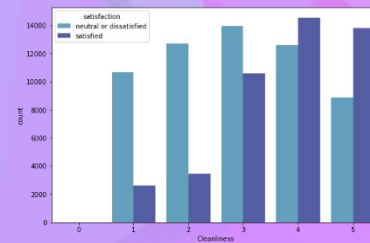
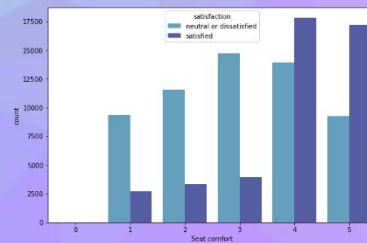
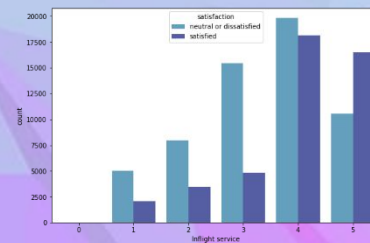
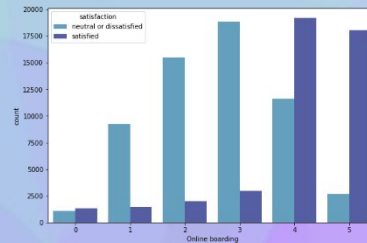
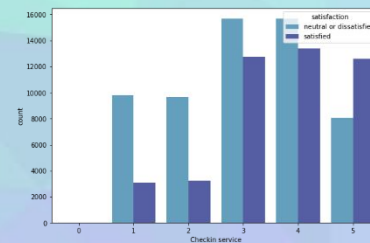
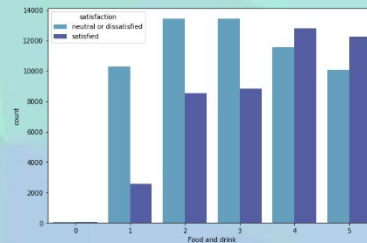
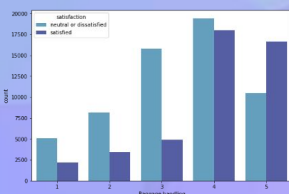
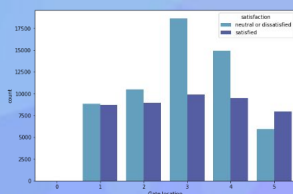
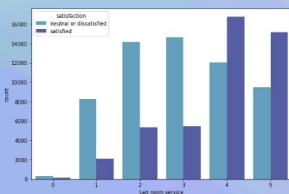
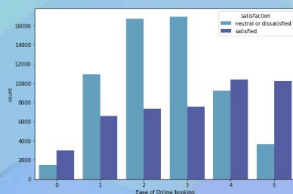
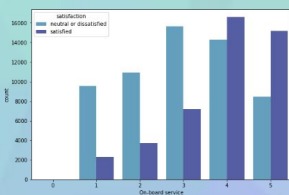
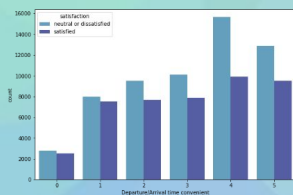
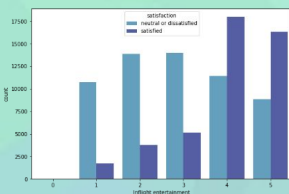
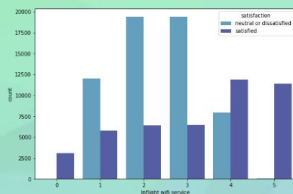
Heatmap

The heatmap is used to find correlations between the categorical rating labels

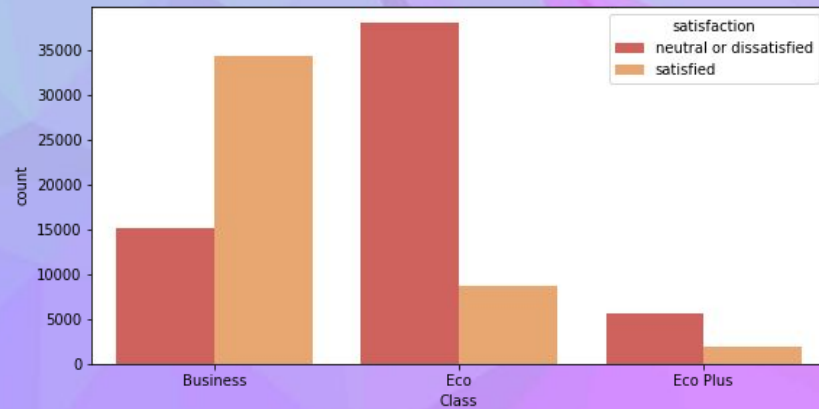
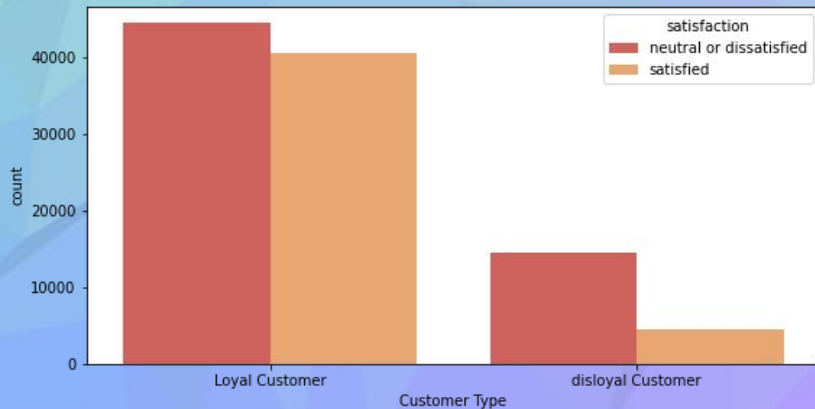
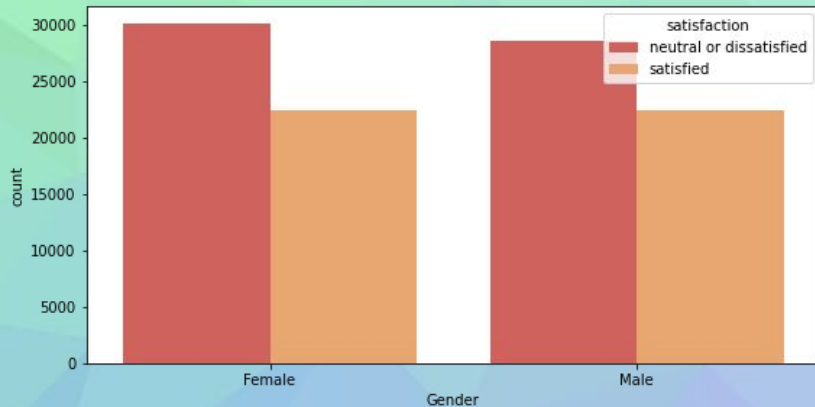
Check out the results @

https://colab.research.google.com/drive/1ozDpPvwTh44hDEafooDxac9SjwNuupHI#scrollTo=_qPqw3_BLsrS

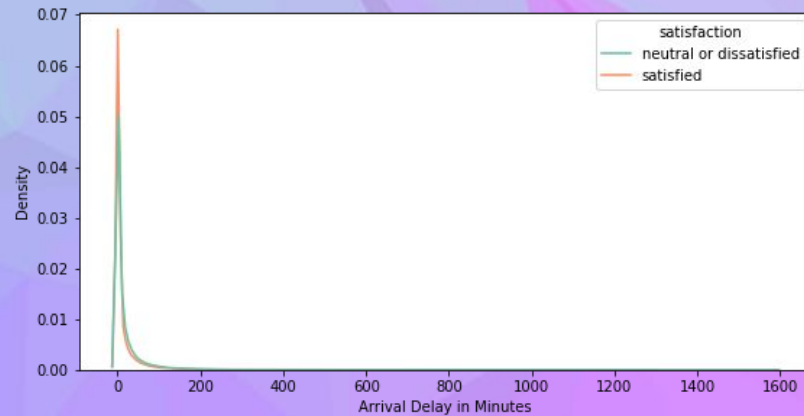
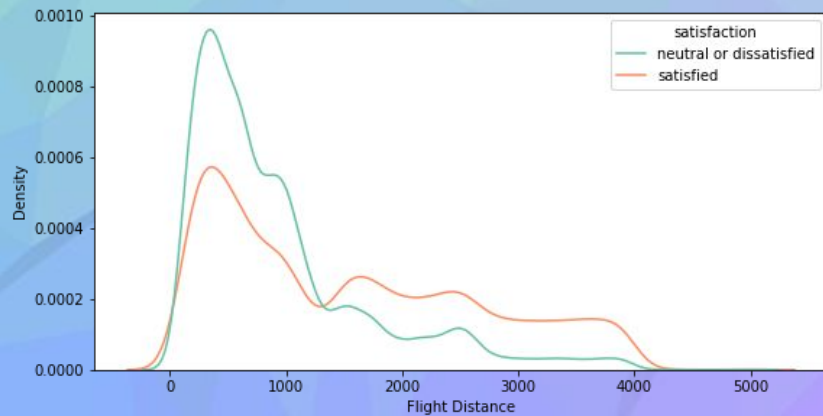
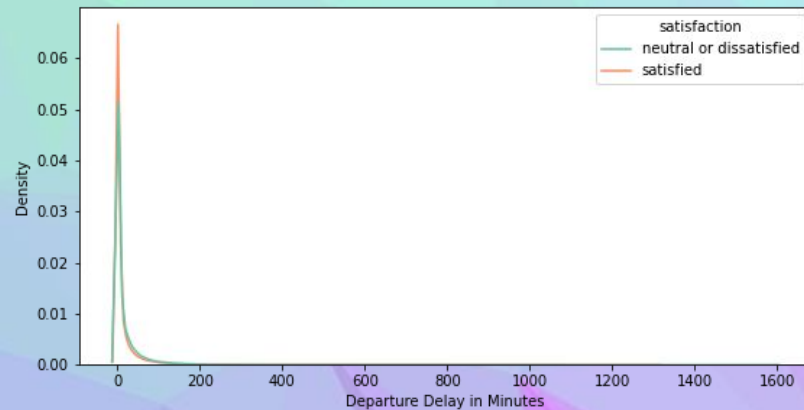
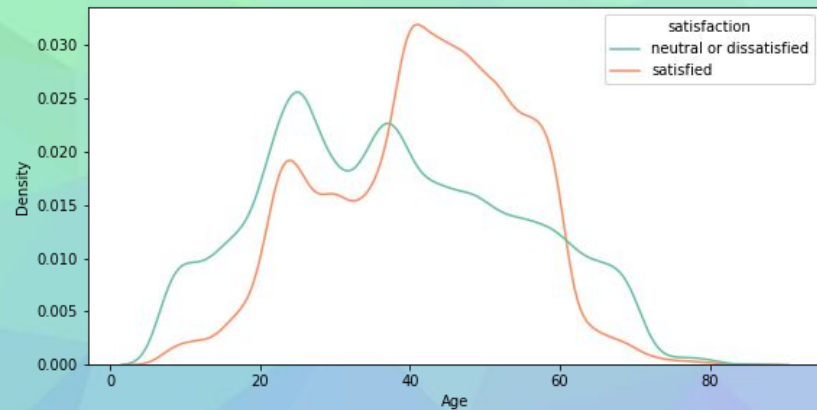
CATEGORICAL LABEL DISTRIBUTION (RATINGS)



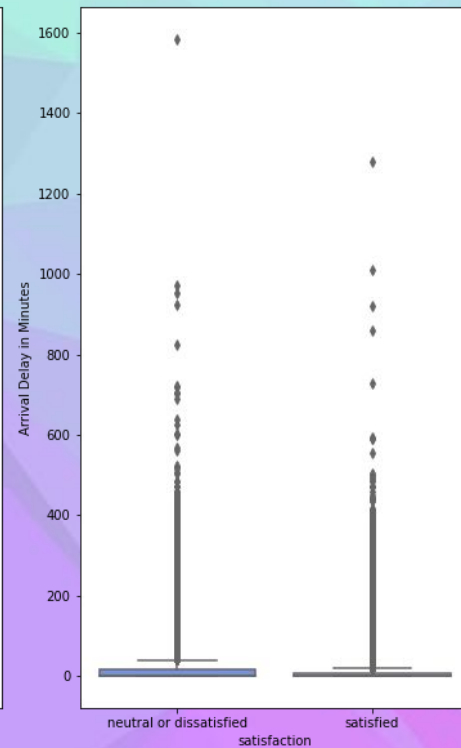
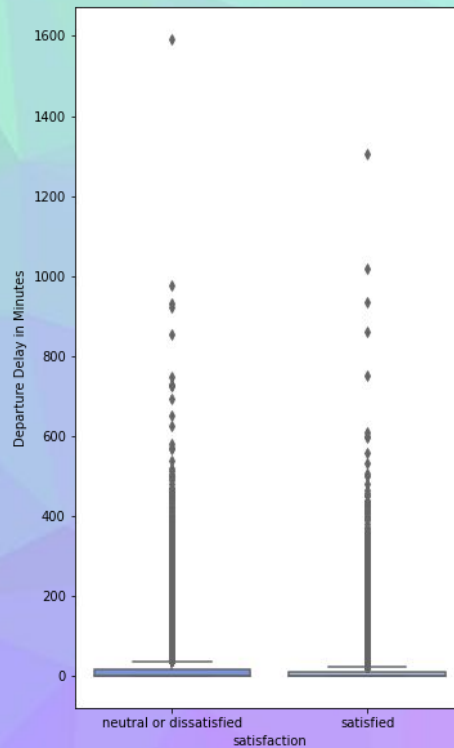
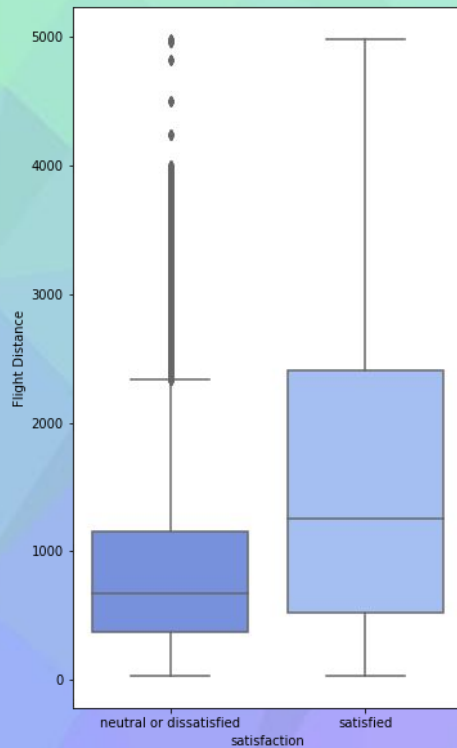
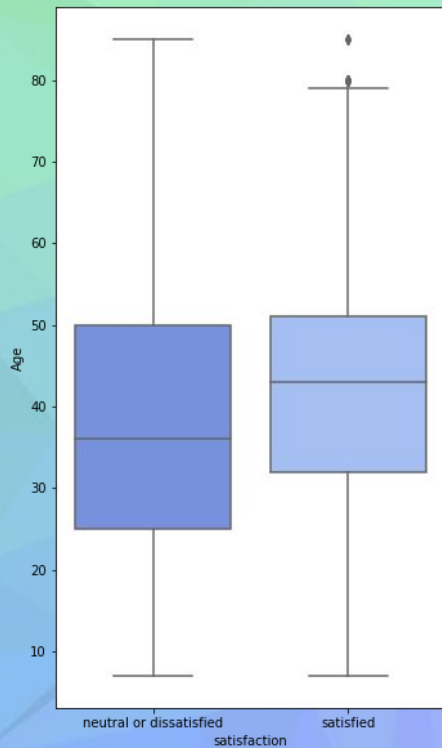
CATEGORICAL LABEL DISTRIBUTION (NON-RATINGS)



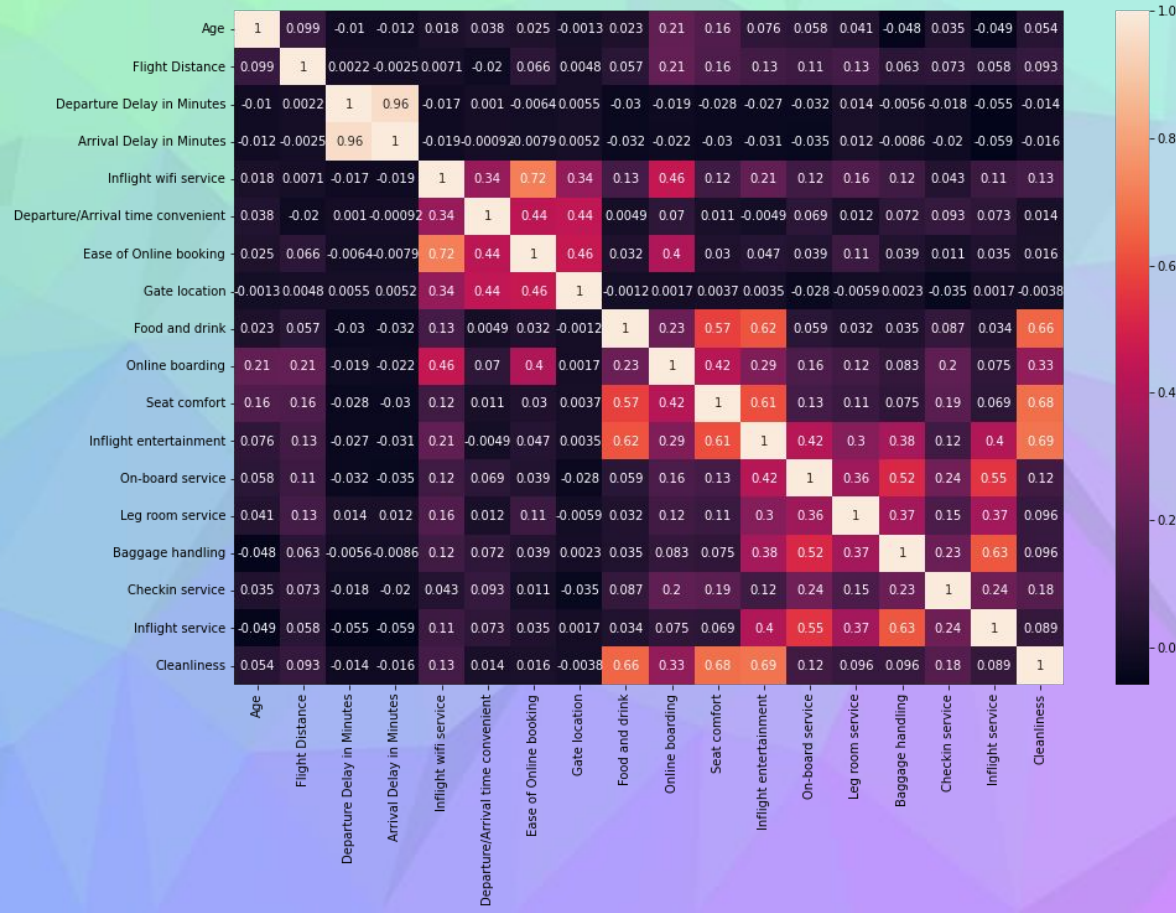
Numerical Label Distribution (Kdeplot)



Numerical Label Distribution (Boxplot)



HEATMAP



Data Preprocessing

01

Split the data

Using `train_test_split`
(train:test =4:1)

02

Remove outliers

Using Interquartile
Range (IQR) to
detect outliers

03

Scale the data

Scale numerical
data using
`MinMaxScaler`

04A

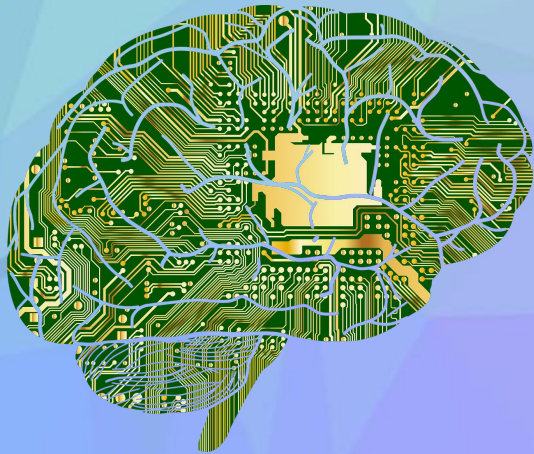
Encode Nominal Data

One-Hot encode
nominal columns using
`pandas.get_dummies`

04B

Encode Ordinal Data

Encode ordinal
columns using
`map`



Split the Dataset

```
train_df, test_df = train_test_split(df, random_state = 42, test_size = 0.2)
```

Remove Outliers

```
def removeOutlier(df, column):  
    Q1=df[column].quantile(0.25)  
    Q3=df[column].quantile(0.75)  
    IQR=Q3-Q1  
    df_final=df[~((df[column]<(Q1-1.5*IQR)) | (df[column]>(Q3+1.5*IQR)))]  
    return df_final
```


Scale Numerical Data

```
def preprocessData(scaler, df, column, type):  
    if type == 'train':  
        df[column] = scaler.fit_transform(df[[column]])  
    elif type == 'test':  
        df[column] = scaler.transform(df[column])  
    return df
```

```
minmax = MinMaxScaler()  
for i in num:  
    print(i)  
    train_df = preprocessData(minmax, train_df, i, 'train')  
    test_df = preprocessData(minmax, test_df, [i], 'test')
```

Encode Categorical Data

```
def getDummies(df2, column):  
    df2_ex = pd.DataFrame()  
    df2_ex = pd.get_dummies(df2[column])  
  
    df2 = pd.concat([df2, df2_ex], axis = 1)  
    df2 = df2.drop(column, axis = 1)  
    return df2
```



```
def ordinalEncode(df, map, column):  
    df[column] = df[column].map(map)  
    return df  
  
class_map = {'Business' : 2,  
             'Eco Plus' : 1,  
             'Eco' : 0}  
  
cType_map = {'Loyal Customer' : 1,  
             'disloyal Customer' : 0}
```


Modeling & Evaluation

Evaluation Metrics

The metrics used are accuracy, precision, recall, F1-Score, and ROC



Model redefinition

Fit the model with tuned parameters

Hyperparameter tuning

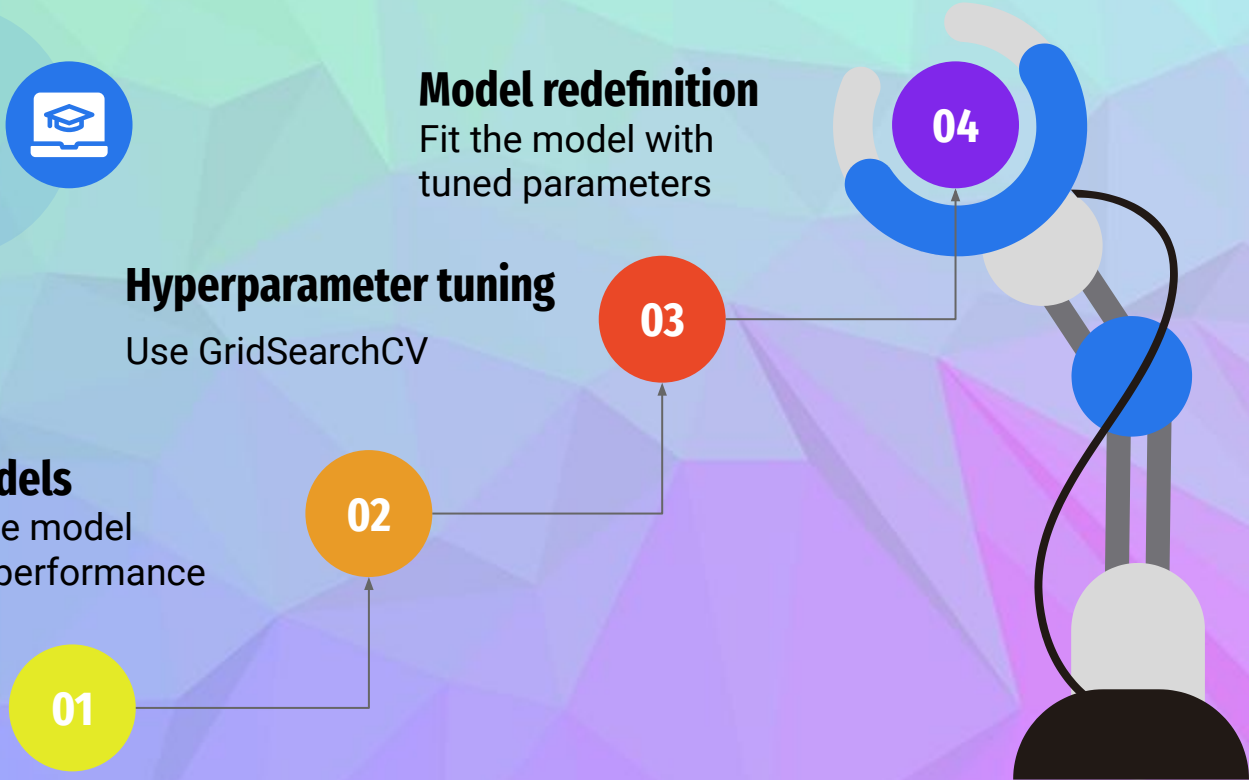
Use GridSearchCV

Filter models

Choose the model with best performance

Prepare models

Define 11 different classifier models, fit each of them and summarize it's metrics



Prepare Models

```
models = {  
    "Logistic Regression" : LogisticRegression(),  
    "Decision Tree": DecisionTreeClassifier(),  
    "LDA" : LinearDiscriminantAnalysis(),  
    "SGD" : SGDClassifier(),  
    "Gaussian": GaussianNB(),  
    "Random Forest" : RandomForestClassifier(),  
    "Gradient Boosting" : GradientBoostingClassifier(),  
    "XGBoost" : XGBClassifier(),  
    "CatBoost" : CatBoostClassifier(),  
    "LGBM" : LGBMClassifier(),  
    "KNN" : KNeighborsClassifier(n_neighbors=3)  
}
```

First of all, we initialize the models which performances we would like to compare.

```
X_train = train_df.drop("satisfaction", axis='columns')  
y_train = train_df['satisfaction']  
X_test = test_df.drop('satisfaction', axis='columns')  
y_test = test_df['satisfaction']
```

We then split the datasets into one that only contains the target column, and another that doesn't contain the target column.

Prepare Models

```
scores = []
probability = {}
for model in models:
    classifier = models[model]
    classifier.fit(X_train, y_train)
    predicts = classifier.predict(X_test)
    try:
        score = classifier.predict_proba(X_test)[:,-1]
        roc = roc_auc_score(y_test, score, average='weighted')
        probability[model] = score
    except:
        roc = 0
    scores.append([
        model,
        accuracy_score(y_test, predicts),
        f1_score(y_test, predicts, average='weighted'),
        precision_score(y_test, predicts),
        recall_score(y_test, predicts),
        roc
    ])
])
```

We then use a loop to fit the data into each model. The scores for each model are also calculated and stored in the **scores** array.

Filter Models

	Model	Accuracy	F1	Precision	Recall	ROC
0	Logistic Regression	0.866753	0.865966	0.877955	0.806793	0.916526
1	Decision Tree	0.941100	0.941135	0.927886	0.937914	0.940740
2	LDA	0.866176	0.865526	0.872320	0.812197	0.918199
3	SGD	0.869256	0.867931	0.899685	0.788266	0.000000
4	Gaussian	0.841105	0.840107	0.847517	0.775364	0.879621
5	Random Forest	0.961166	0.961084	0.971249	0.938796	0.993430
6	Gradient Boosting	0.943314	0.943213	0.947787	0.920820	0.987916
7	XGBoost	0.939656	0.939531	0.945496	0.914424	0.987658
8	CatBoost	0.962466	0.962394	0.971445	0.941663	0.994979
9	LGBM	0.962273	0.962195	0.972292	0.940340	0.994717
10	KNN	0.919542	0.919086	0.942450	0.868659	0.953444

We then compile the scores into a dataframe, then display them for comparison purposes. Here, we saw that the CatBoost model has the highest accuracy in all of the metrics (Accuracy, F1 Score, Precision, Recall, and ROC).

Hyperparameter Tuning

In order to improve the performance of the model, we tuned the hyperparameters using **GridSearchCV**.

```
grid = {'max_depth': [5, 6, 7],
        'n_estimators': [300, 400, 500],
        'learning_rate' : [0.01, 0.05, 0.1, 0.15]}

cbc = CatBoostClassifier()

gscv = GridSearchCV (estimator = cbc, param_grid = grid, scoring = 'accuracy', cv = 5)

gscv.fit(X_train, y_train)
```

We then output the best estimator, score, and parameters.

```
#returns the estimator with the best performance
print(gscv.best_estimator_, end = '\n\n')

#returns the best score
print(gscv.best_score_, end = '\n\n')

#returns the best parameters
print(gscv.best_params_)
```

The resulting output can be seen below, with the best parameters being utilized in the next stage (Model Redefinition).

```
<catboost.core.CatBoostClassifier object at 0x7f00d6544550>

0.9627804851804577

{'learning_rate': 0.1, 'max_depth': 7, 'n_estimators': 400}
```


Model Redefinition

After tuning the hyperparameters, we redefine the model with the newly acquired values.

```
final_model = CatBoostClassifier(learning_rate=0.1, max_depth = 7, n_estimators= 400)
final_model.fit(X_train, y_train)
```

We then calculate the performance of the now redefined model using the same metrics as before (Accuracy score, F1 Score, Precision Score, and Recall Score).

```
predicts = final_model.predict(X_test)
print(f"Final Model's accuracy : {accuracy_score(y_test, predicts)}")
print(f"Final Model's F1 Score : {f1_score(y_test, predicts, average = 'weighted')}")
print(f"Final Model's Precision : {precision_score(y_test, predicts)}")
print(f"Final Model's Recall : {recall_score(y_test, predicts)}")
```

Results

Catboost Model Performance

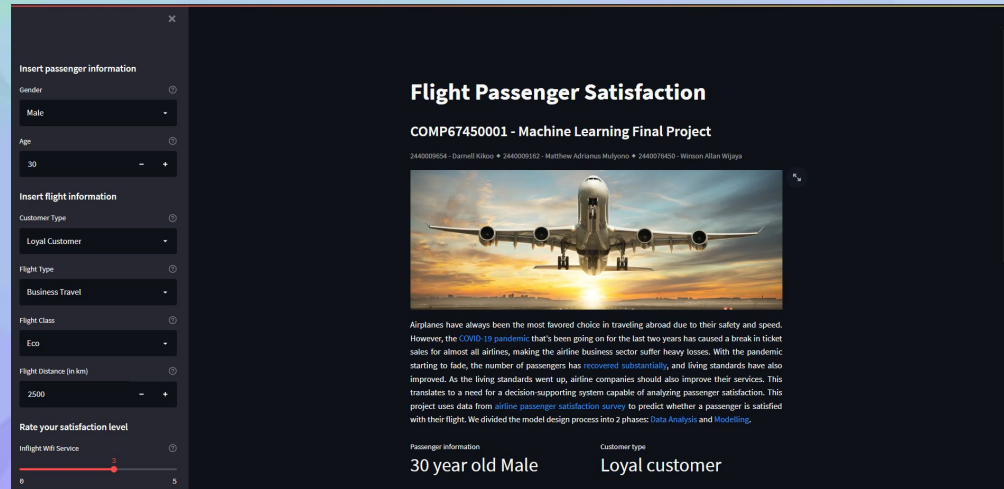
Metric	Performance
Accuracy	0.9628
Precision	0.9708
Recall	0.9431
F1 Score	0.9627

The model achieved 96.28% accuracy, supported by 97% precision, 94% recall, and 96% F1 Score.

Deployment

We deployed our model in the form of a web-application, which can be accessed through the link:

<https://flight-satisfaction.herokuapp.com/>.



The screenshot displays the 'Flight Passenger Satisfaction' web application. On the left, a dark sidebar contains input fields for passenger and flight information. The 'Insert passenger information' section includes 'Gender' (set to Male) and 'Age' (set to 30). The 'Insert flight information' section includes 'Customer Type' (set to Loyal Customer), 'Flight Type' (set to Business Travel), 'Flight Class' (set to Eco), and 'Flight Distance (in km)' (set to 2500). At the bottom of the sidebar is a 'Rate your satisfaction level' section with a slider for 'Inflight Wi-Fi Service' ranging from 0 to 5. The main content area on the right has a white background and features the title 'Flight Passenger Satisfaction', the subtitle 'COMP67450001 - Machine Learning Final Project', and a list of team members. Below this is a large image of a commercial airplane taking off. A paragraph of text discusses the impact of the COVID-19 pandemic on the airline industry and the need for a decision-supporting system. At the bottom, two summary cards show 'Passenger information' as '30 year old Male' and 'Customer type' as 'Loyal customer'.

Flight Passenger Satisfaction

COMP67450001 - Machine Learning Final Project

2460009054 - Darrell Kilian • 2460009102 - Matthew Adrianus Mulyono • 2460070450 - Winson Allan Wigaya

Airplanes have always been the most favored choice in traveling abroad due to their safety and speed. However, the COVID-19 pandemic that's been going on for the last two years has caused a break in ticket sales for almost all airlines, making the airline business sector suffer heavy losses. With the pandemic starting to fade, the number of passengers has recovered substantially, and living standards have also improved. As the living standards went up, airline companies should also improve their services. This translates to a need for a decision-supporting system capable of analyzing passenger satisfaction. This project uses data from [airline passenger satisfaction survey](#) to predict whether a passenger is satisfied with their flight. We divided the model design process into 2 phases: [Data Analysis](#) and [Modelling](#).

Passenger information
30 year old Male

Customer type
Loyal customer

Спасибо Gracias شكر Obrigado Спасибо Dank U
Grazie Euxαριστώ Danke
Merci Ngilyabonga Dank U
Dziękę Euxαριστώ
Danke Grazie תודה
Diolch 謝謝
Dank U Terima Kasih
Grazie Tack Euxαριστώ
Merci Tack Euxαριστώ