

Airline Passenger Satisfaction

CS513 Project Presentation

Section B Group 6



Shivang Medhekar
CWID: 10478510



Mingfang Liang
CWID: 10473184



Yutong Wei
CWID: 2005208



Jiani Yu
CWID: 10474729

Problem Statement

- Airline Market
 - 5000 airlines
 - \$785.6 bn
- Airline Challenges
 - Many Features
 - What to improve ?
 - How to differentiate ?



Work Flow



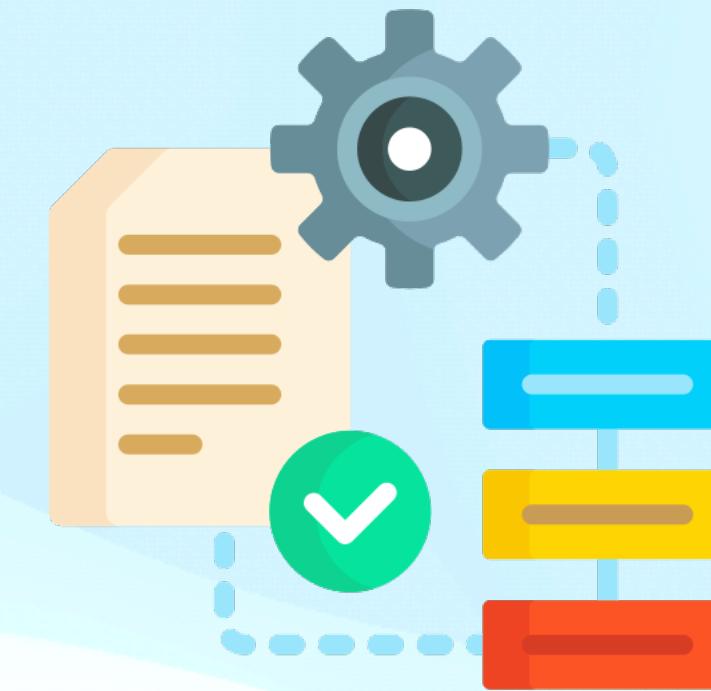
Data Exploration



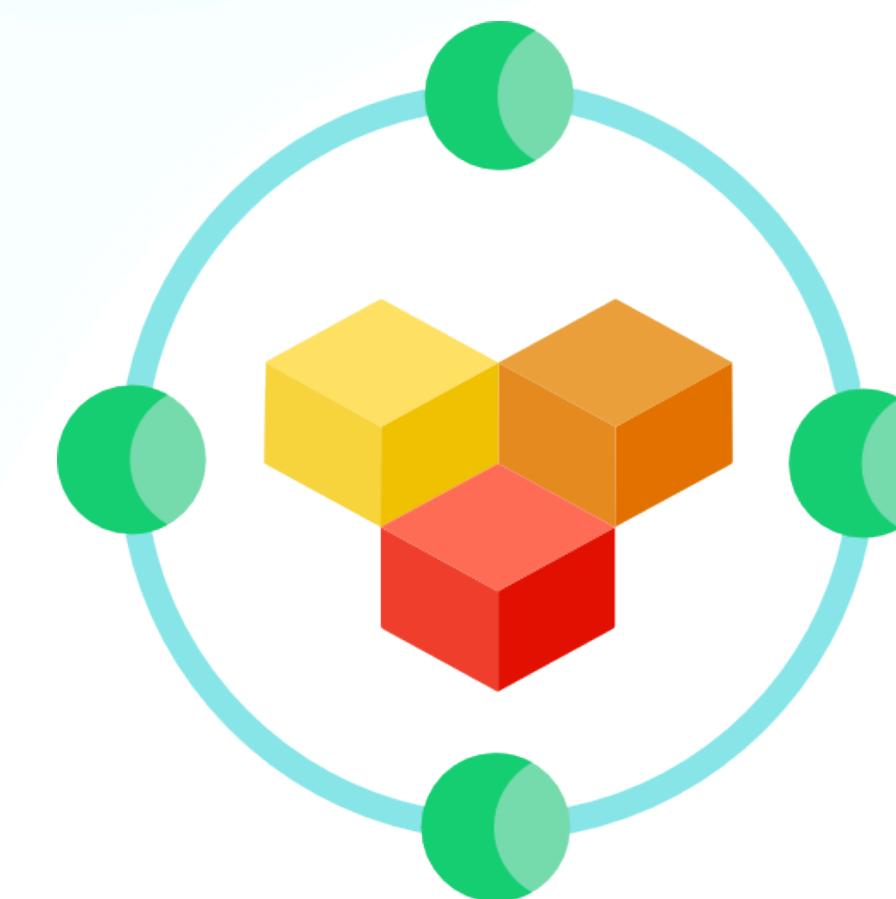
Data Analysis



Data Cleaning



Data Preprocessing



Modeling

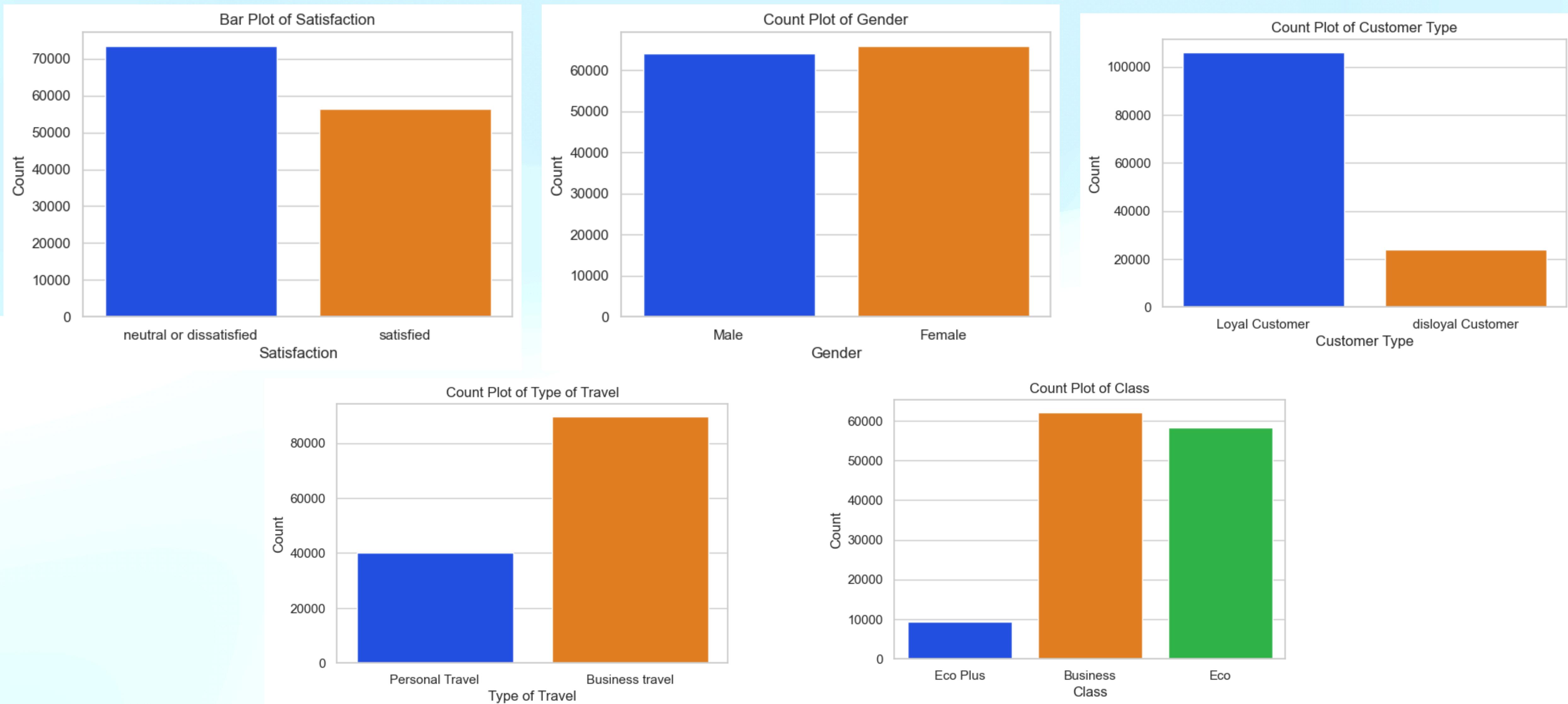


Evaluation

Dataset Exploration

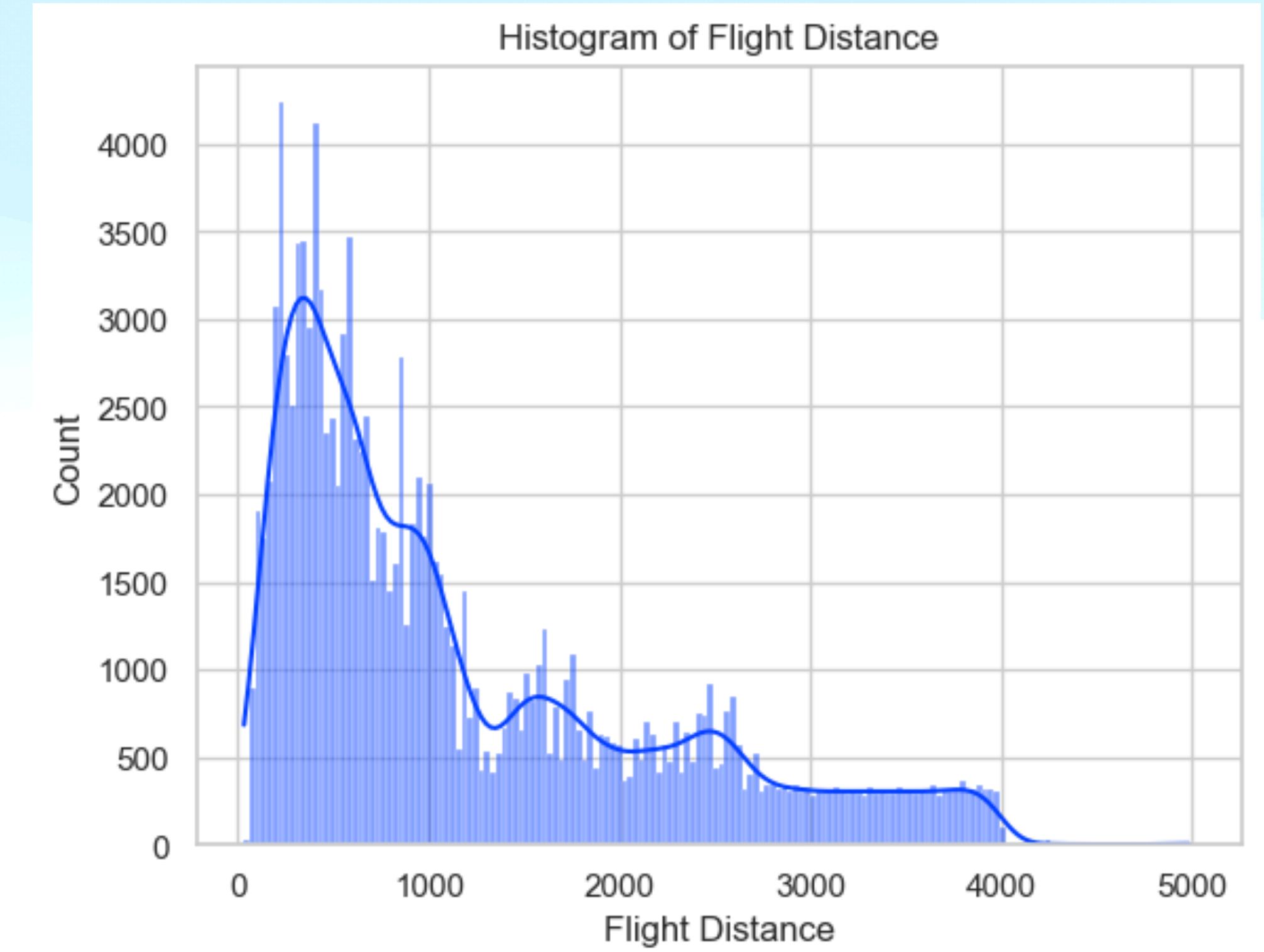
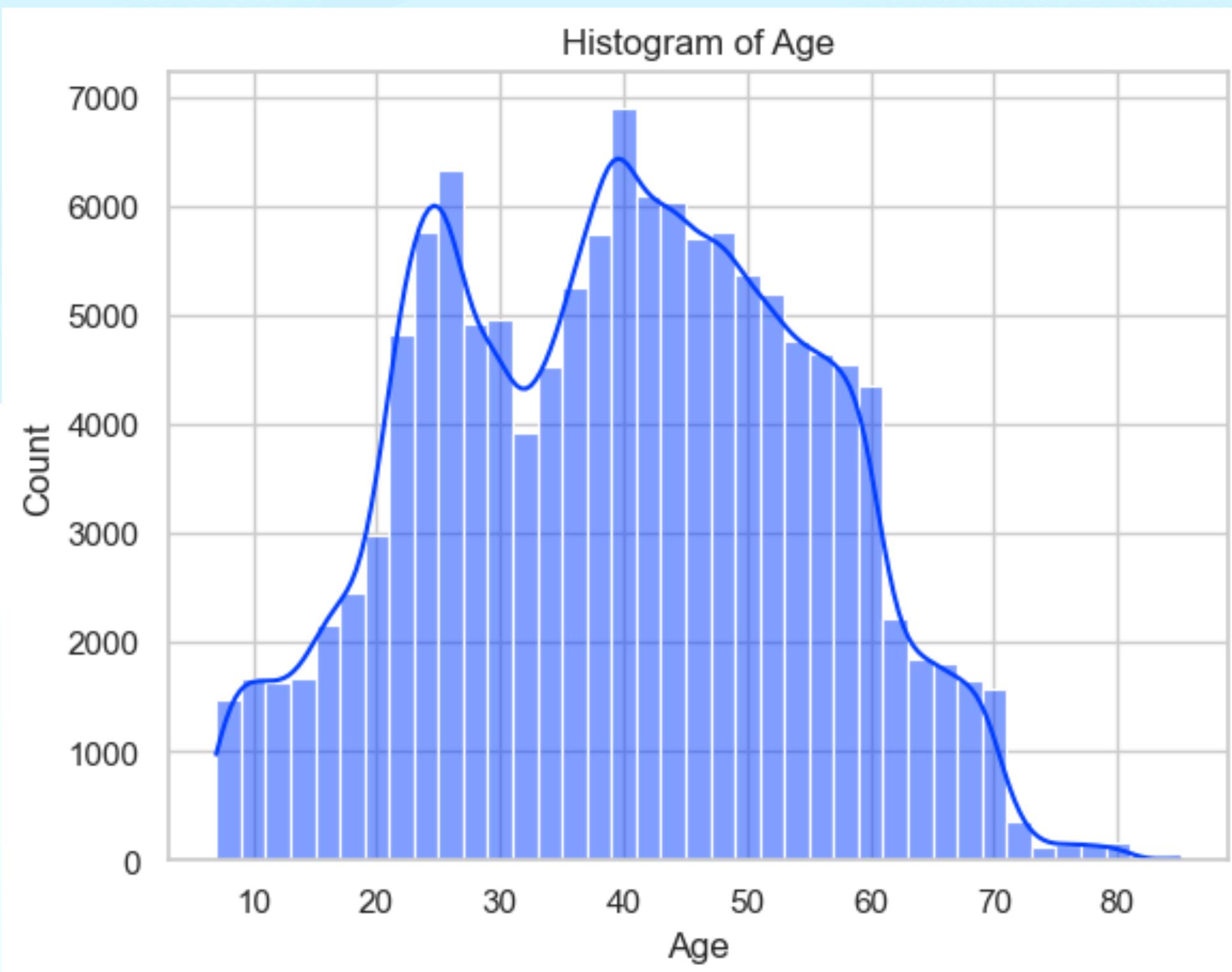
- **Gender:** Gender of the passengers (Female, Male)
 - **Customer Type:** The customer type (Loyal customer, disloyal customer)
 - **Age:** The actual age of the passengers
 - **Type of Travel:** Purpose of the flight of the passengers (Personal Travel, Business Travel)
 - **Class:** Travel class in the plane of the passengers (Business, Eco, Eco Plus)
 - **Flight distance:** The flight distance of this journey
 - **Inflight wifi service:** Satisfaction level of the inflight wifi service (0:Not Applicable;1-5)
 - **Departure/Arrival time convenient:** Satisfaction level of Departure/Arrival time convenient
 - **Ease of Online booking:** Satisfaction level of online booking
 - **Gate location:** Satisfaction level of Gate location
 - **Food and drink:** Satisfaction level of Food and drink
 - **Online boarding:** Satisfaction level of online boarding
 - **Seat comfort:** Satisfaction level of Seat comfort
 - **Inflight entertainment:** Satisfaction level of inflight entertainment
 - **On-board service:** Satisfaction level of On-board service
 - **Leg room service:** Satisfaction level of Leg room service
 - **Baggage handling:** Satisfaction level of baggage handling
 - **Check-in service:** Satisfaction level of Check-in service
 - **Inflight service:** Satisfaction level of inflight service
 - **Cleanliness:** Satisfaction level of Cleanliness
 - **Departure Delay in Minutes:** Minutes delayed when departure
 - **Arrival Delay in Minutes:** Minutes delayed when Arrival
 - **Satisfaction:** Airline satisfaction level(Satisfaction, neutral or dissatisfaction)
- Dataset: airline passenger satisfaction survey.
 - 23 Variables and 129,880 Observations

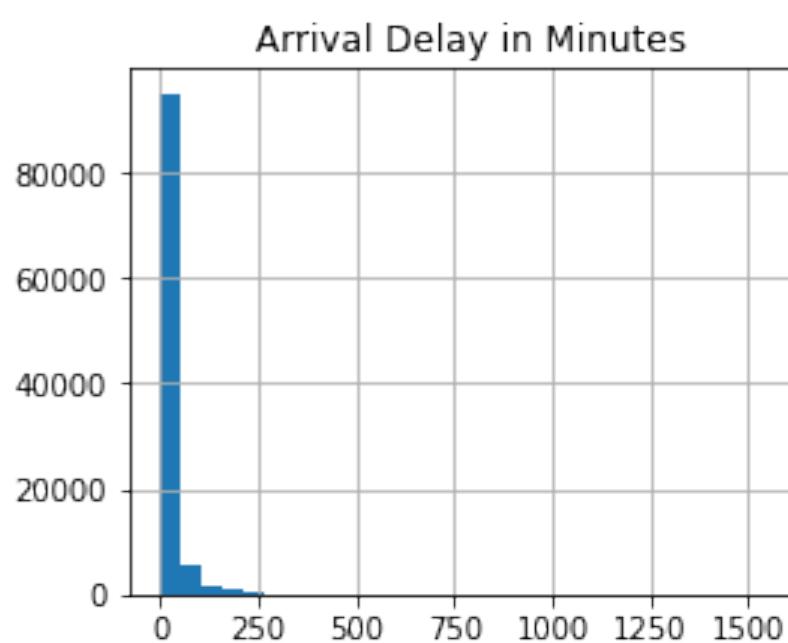
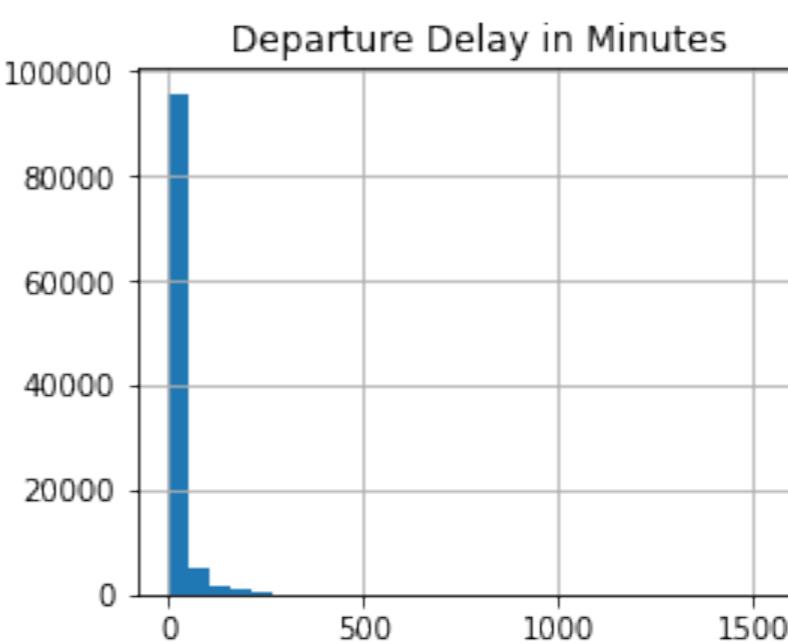
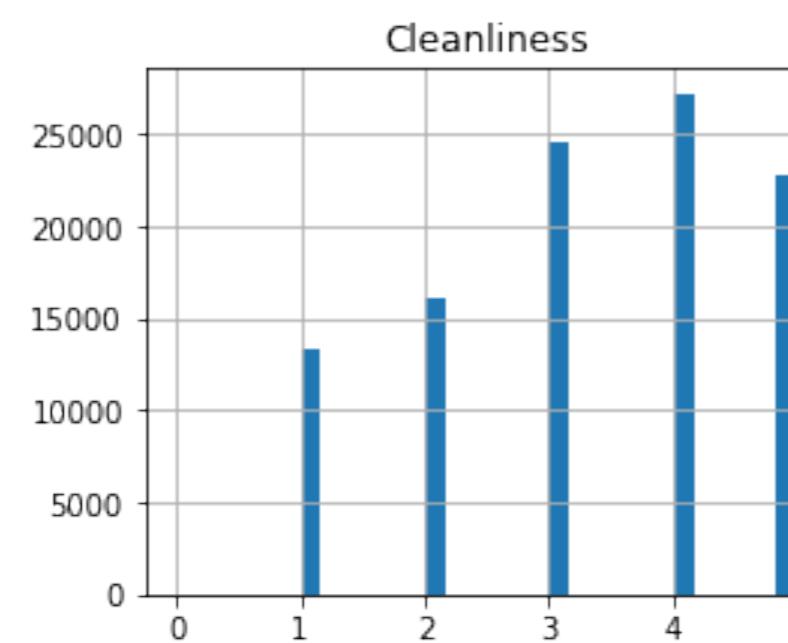
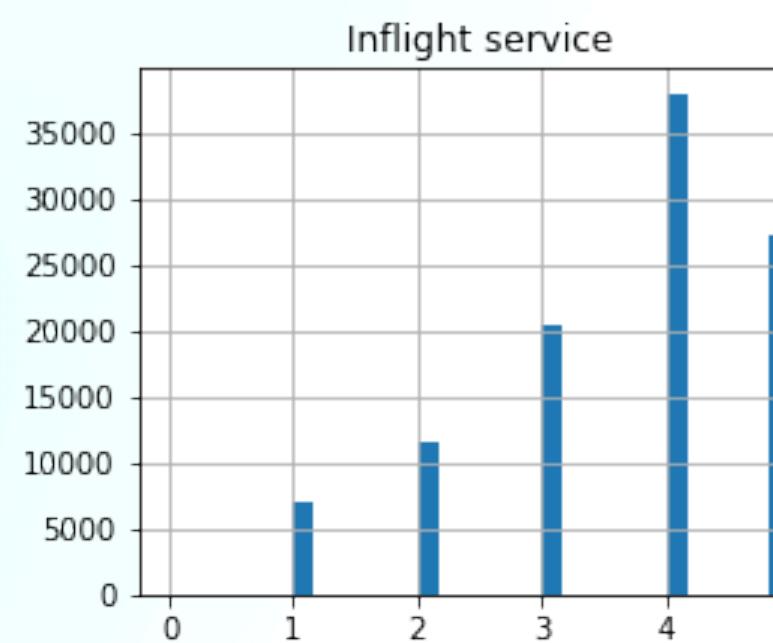
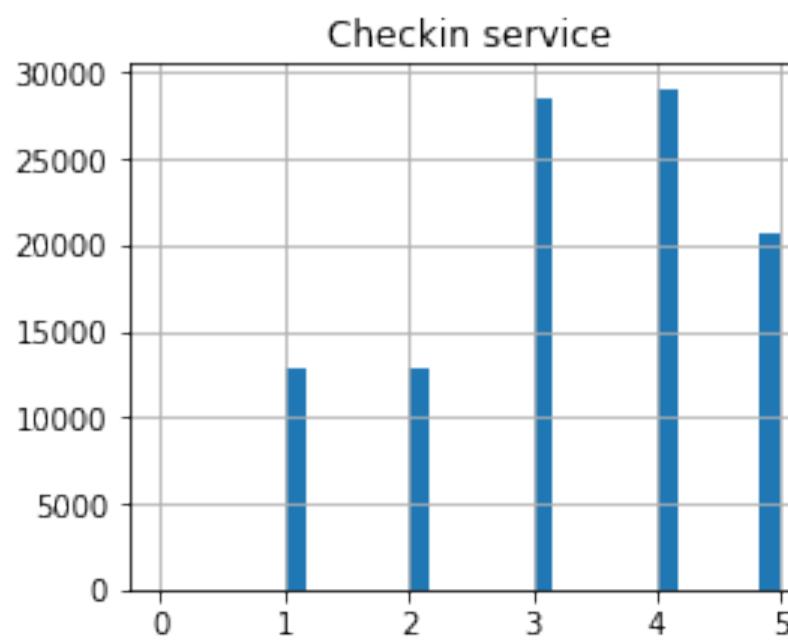
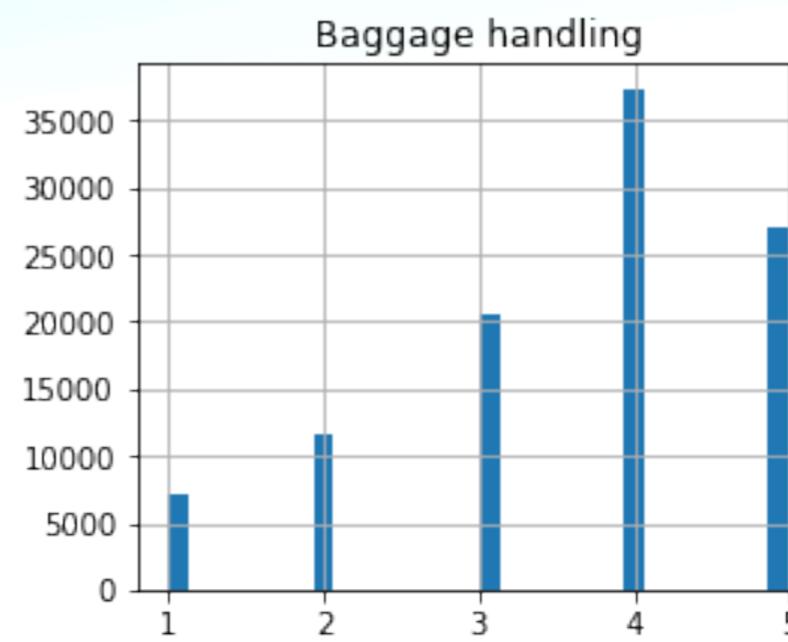
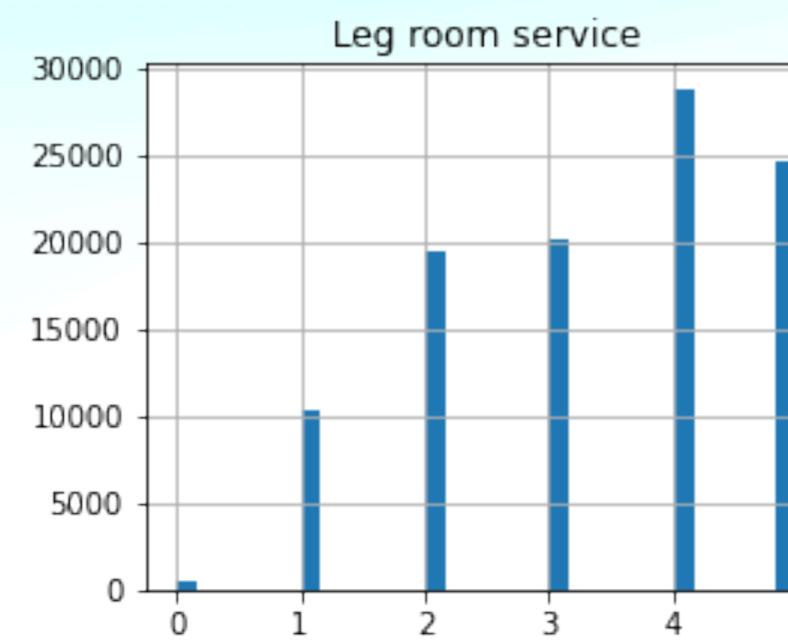
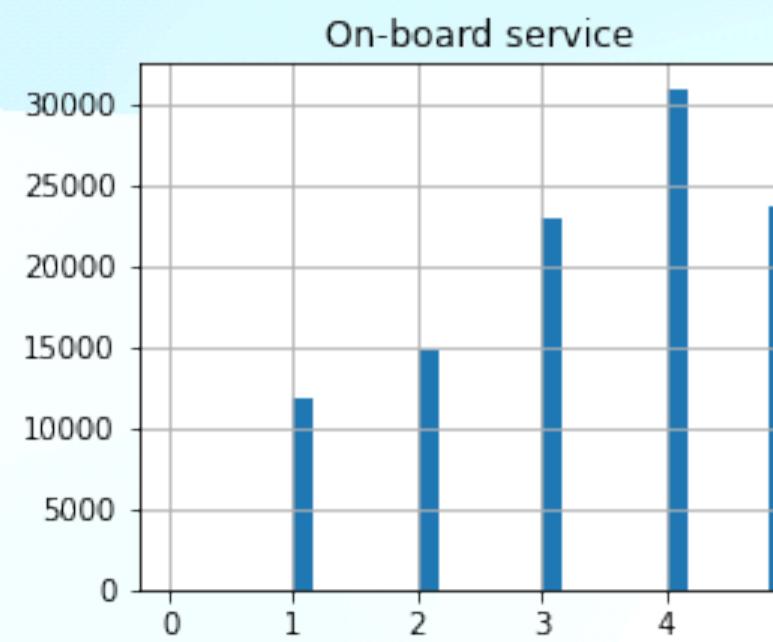
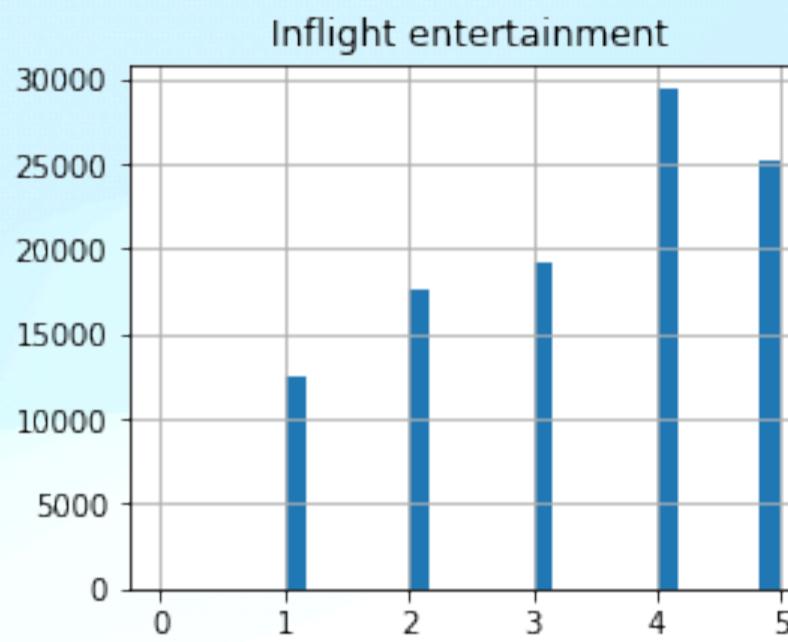
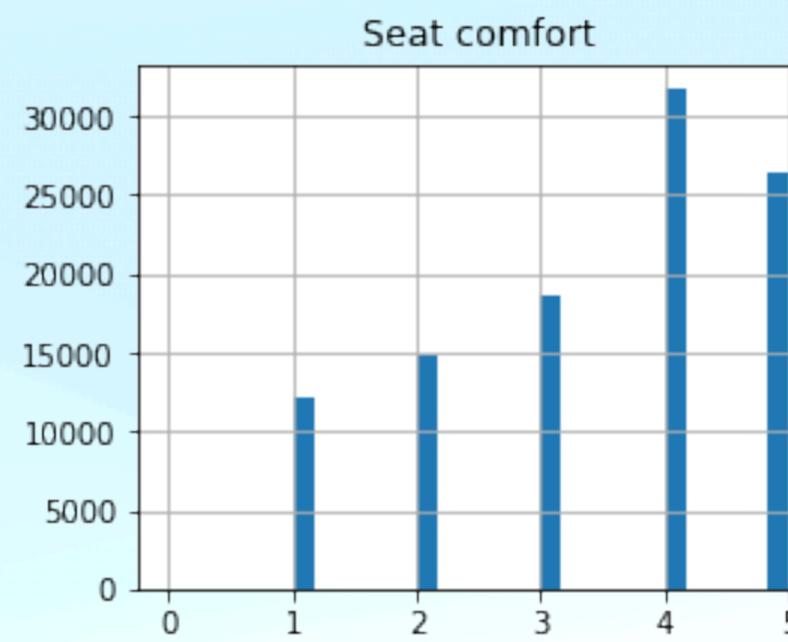
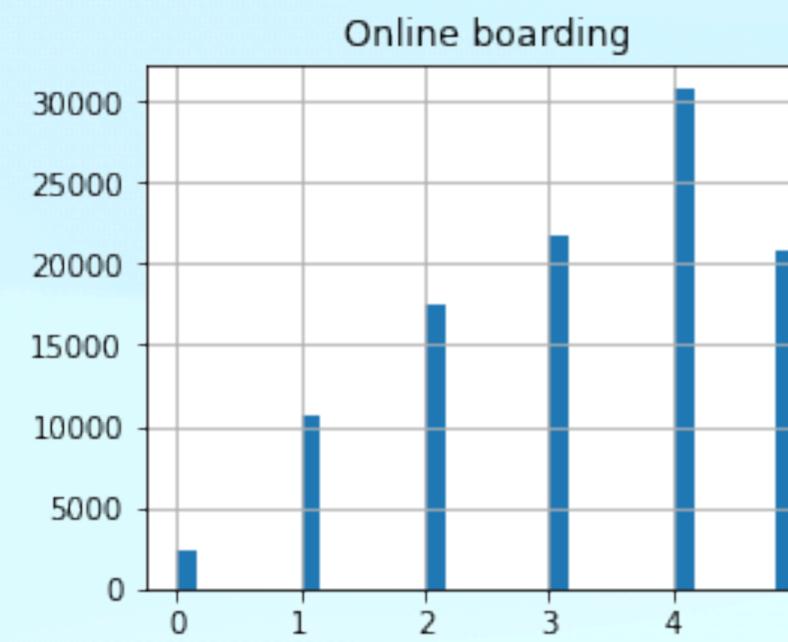
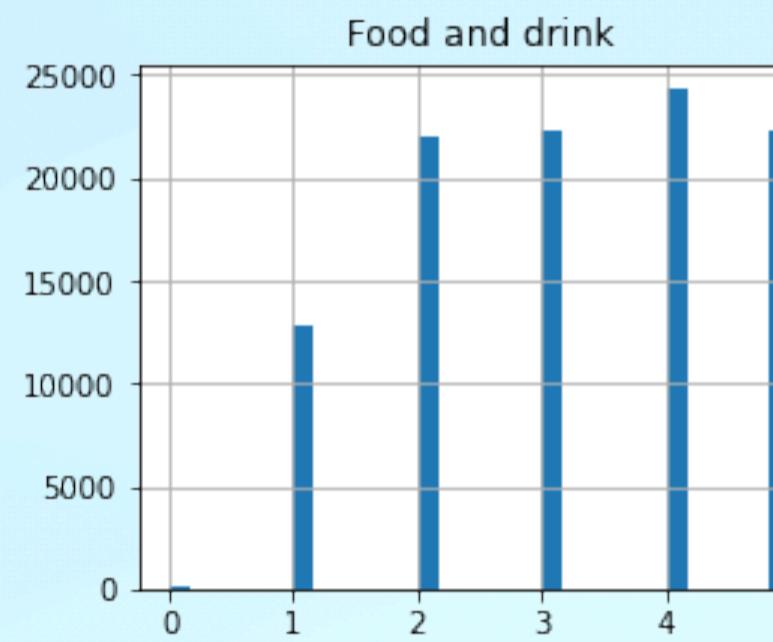
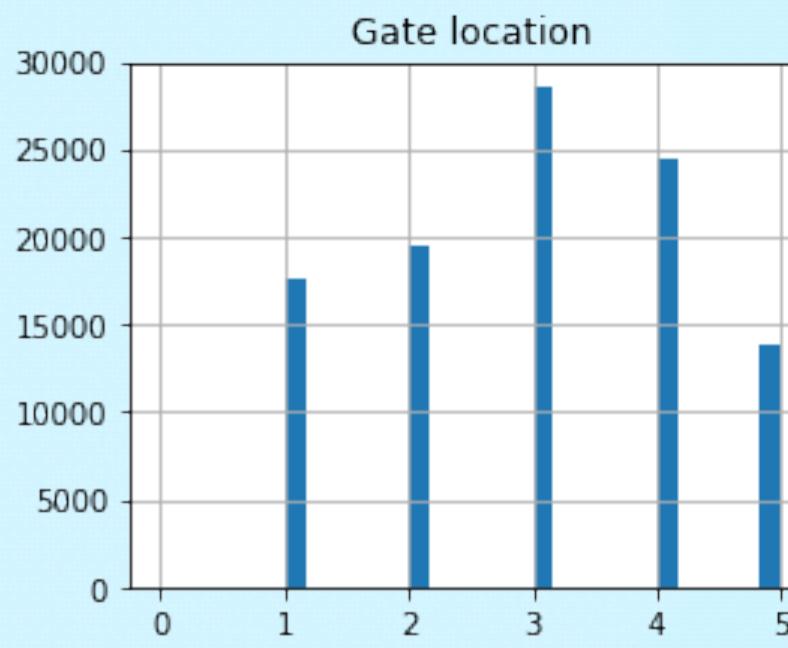
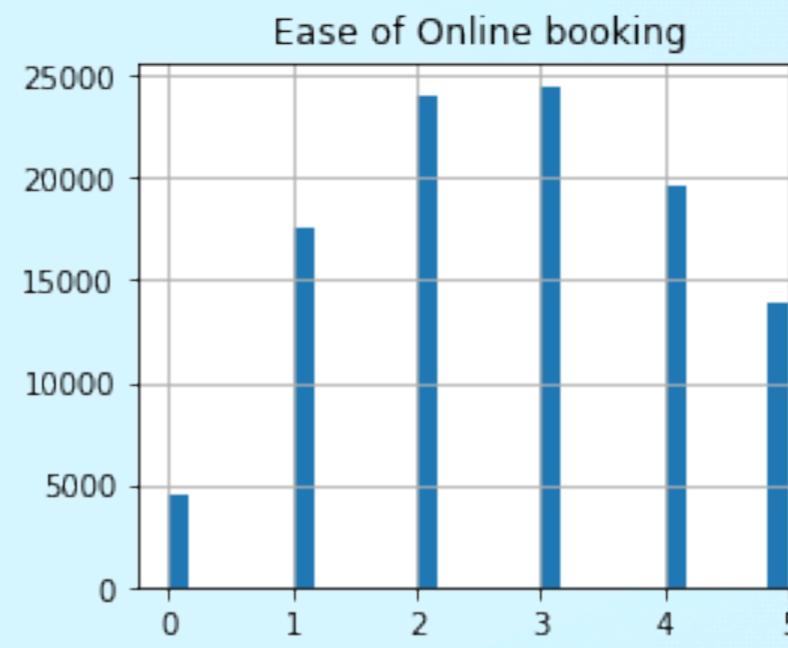
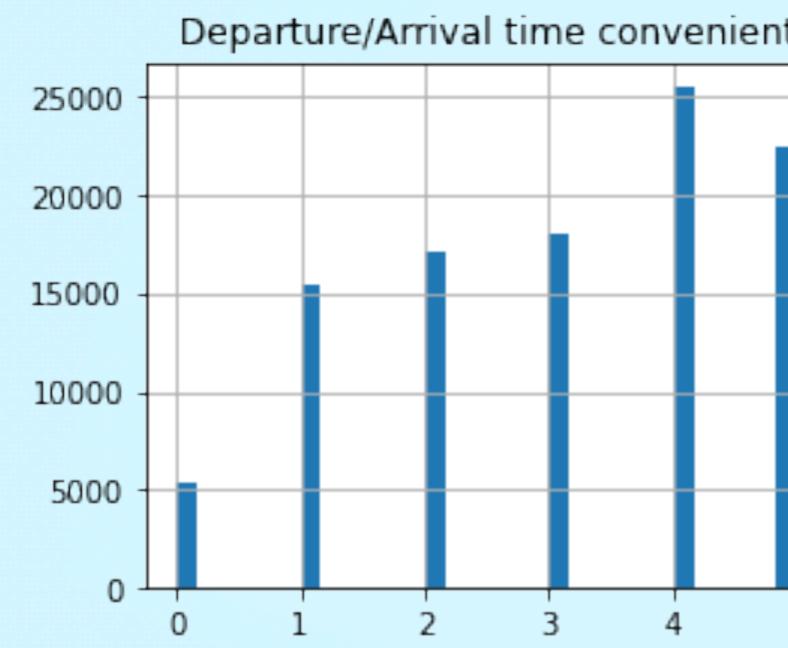
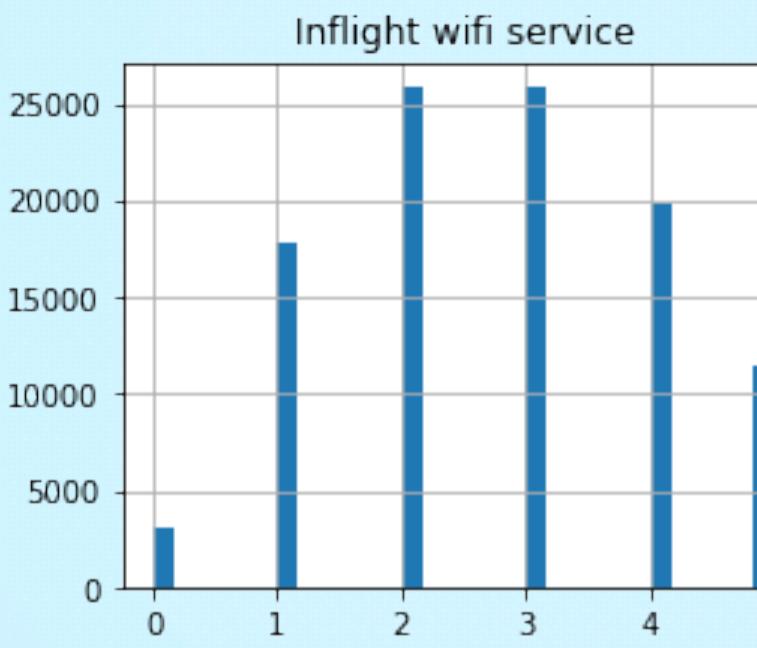
Data Analysis: Categorical Data



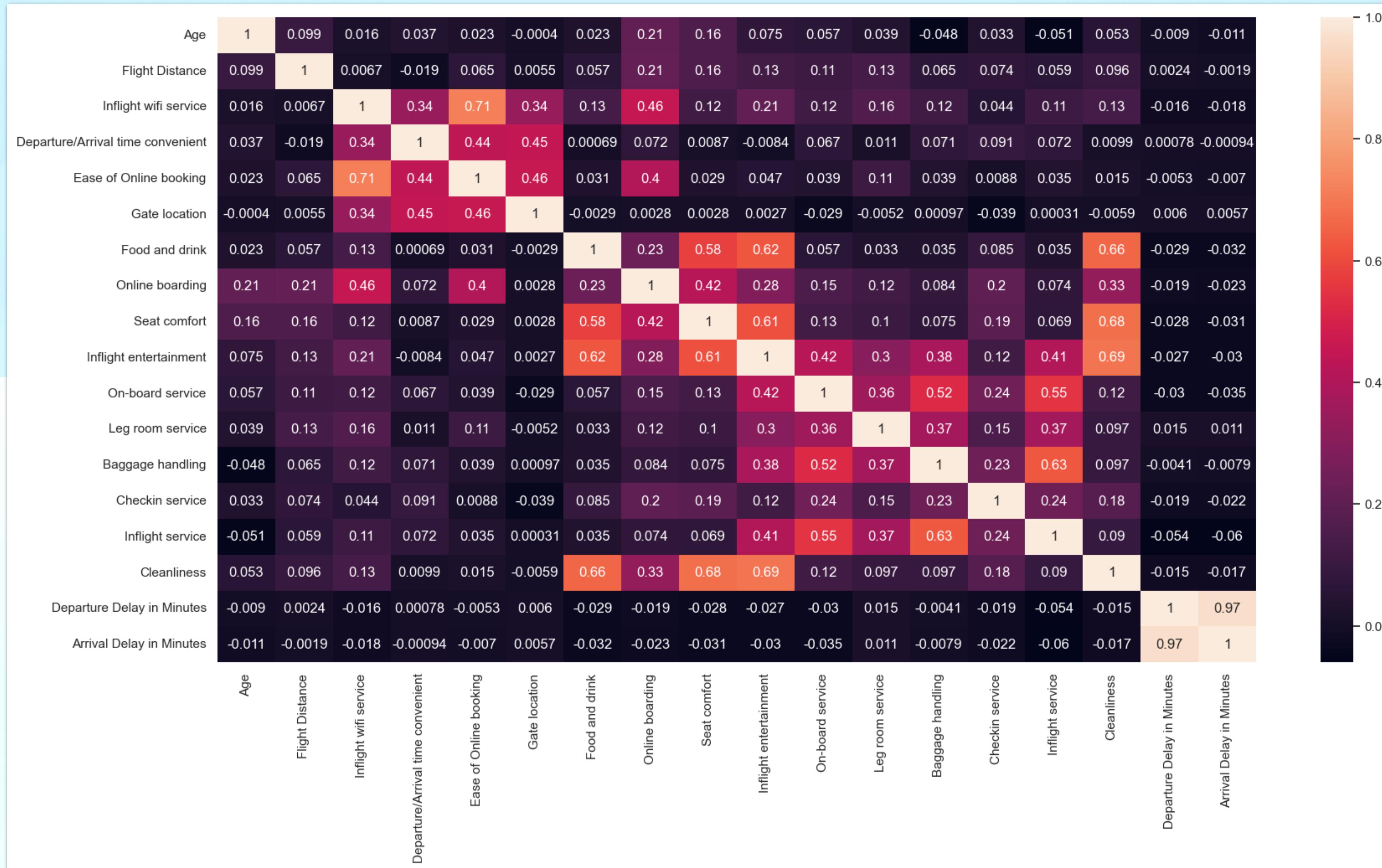
Data Analysis

Numerical Data





Heat Map



Data Cleaning

NA Values

dataset.isnull().any()	
	data
Unnamed: 0	False
id	False
Gender	False
Customer Type	False
Age	False
Type of Travel	False
Class	False
Flight Distance	False
Inflight wifi service	False
Departure/Arrival time convenient	False
Ease of Online booking	False
Gate location	False
Food and drink	False
Online boarding	False
Seat comfort	False
Inflight entertainment	False
On-board service	False
Leg room service	False
Baggage handling	False
Checkin service	False
Inflight service	False
Cleanliness	False
Departure Delay in Minutes	False
Arrival Delay in Minutes	True
satisfaction	False

dataset.isnull().sum()	
	data
Gender	0
Customer Type	0
Age	0
Type of Travel	0
Class	0
Flight Distance	0
Inflight wifi service	0
Departure/Arrival time convenient	0
Ease of Online booking	0
Gate location	0
Food and drink	0
Online boarding	0
Seat comfort	0
Inflight entertainment	0
On-board service	0
Leg room service	0
Baggage handling	0
Checkin service	0
Inflight service	0
Cleanliness	0
Departure Delay in Minutes	0
Arrival Delay in Minutes	393
satisfaction	0

Replaced NA with
Mean of the
Column
Arrival Data in
Minutes

Data Preprocessing

Binary Encoding

- Customer Type - ['Loyal Customer', 'disloyal Customer']
- Gender - ['Male', 'Female']
- Type of Travel - ['Personal Travel', 'Business travel']
- Satisfaction - ['neutral or dissatisfied', 'satisfied']

One Hot Encoding

- Class - ['Eco Plus', 'Business', 'Eco']

Data Preprocessing

Data Split and Scale

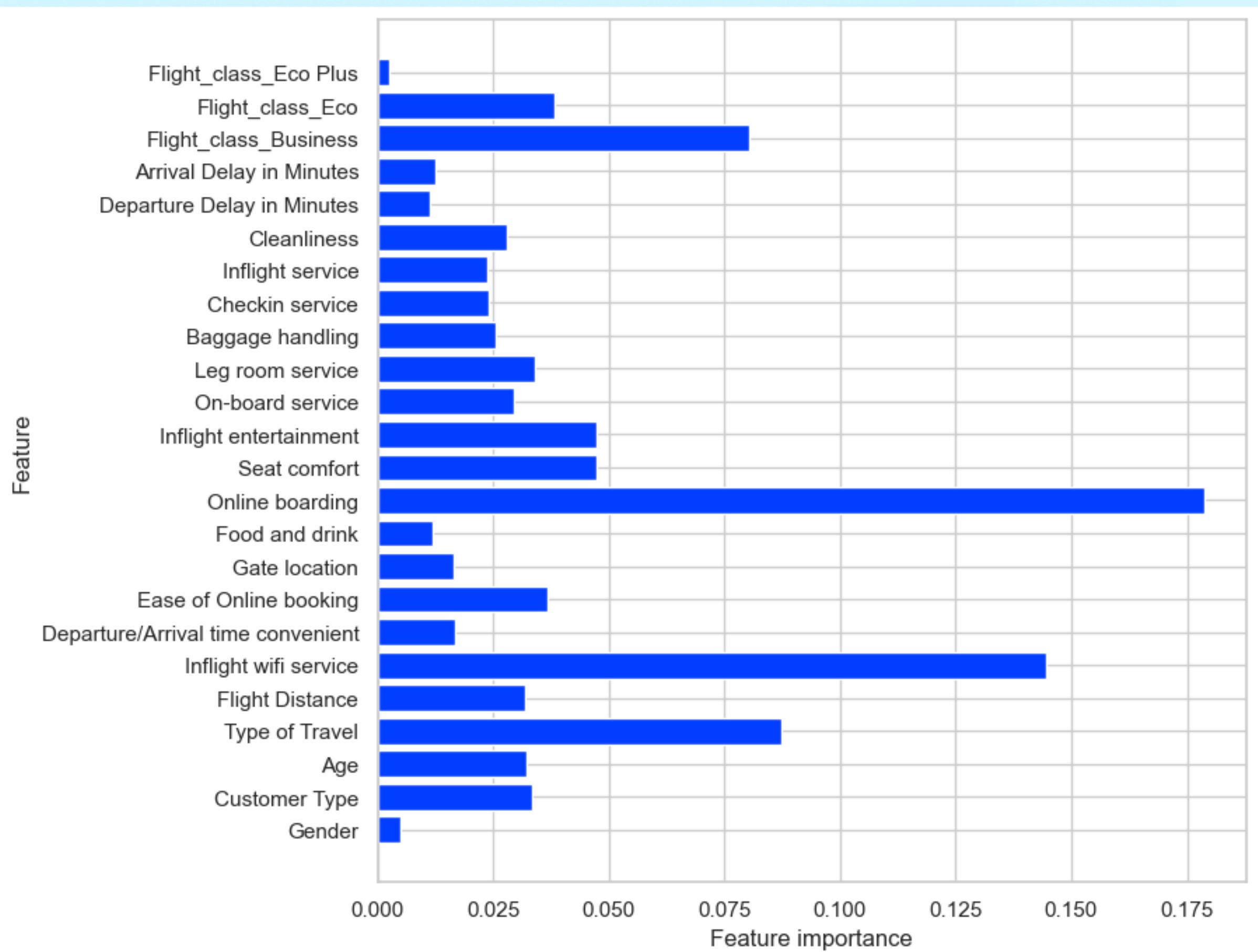
Split Data in 70 % Test and 30 % Train
And Scale between -1 and 1 using
MinMaxScaler

Data Preprocessing

Feature Selection

Recursive feature elimination with cross-validation RFECV with RandomForestClassifier

1. Online boarding
2. Inflight wifi service
3. Type of Travel
4. Flight_class_Business
5. Seat comfort
6. Inflight entertainment
7. Flight_class_Eco
8. Ease of Online booking
9. Leg room service
10. Customer Type
11. Age
12. Flight Distance
13. On-board service
14. Cleanliness
15. Baggage handling
16. Checkin service
17. Inflight service
18. Departure/Arrival time convenient
19. Gate location
20. Arrival Delay in Minutes
21. Food and drink
22. Departure Delay in Minutes
23. Gender
24. Flight_class_Eco Plus



Logistic Regression with Grid SearchCV

Logistic Regression

- A Machine Learning classification algorithm
- Used to predict the probability of certain classes based on some dependent variables.

GridSearchCV

- A function that is in sklearn's model_selection package.
- specify the different values for each hyperparameter
- training and testing using cross validation of your dataset – hence the acronym “CV” in GridSearchCV.

LogisticRegression with GridSearchCV

```
In [39]: def get_logr(mode = 'load'):

    path = DIRPATH + "/Models/"
    filename = 'log_r.sav'

    if mode == 'load':
        log_r = pickle.load(open(path + filename, 'rb'))
    else:

        log_r_param_grid = {
            'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000],
            'penalty': ['l2'],
        }

        log_r = GridSearchCV(
            estimator = LogisticRegression(max_iter = 1000),
            param_grid = log_r_param_grid,
            cv = 5,
            verbose = 1,
            scoring = 'accuracy',
            n_jobs = -1
        )
        log_r.fit(X_train, y_train)

        pickle.dump(log_r, open(path + filename, 'wb'))

    return log_r
```

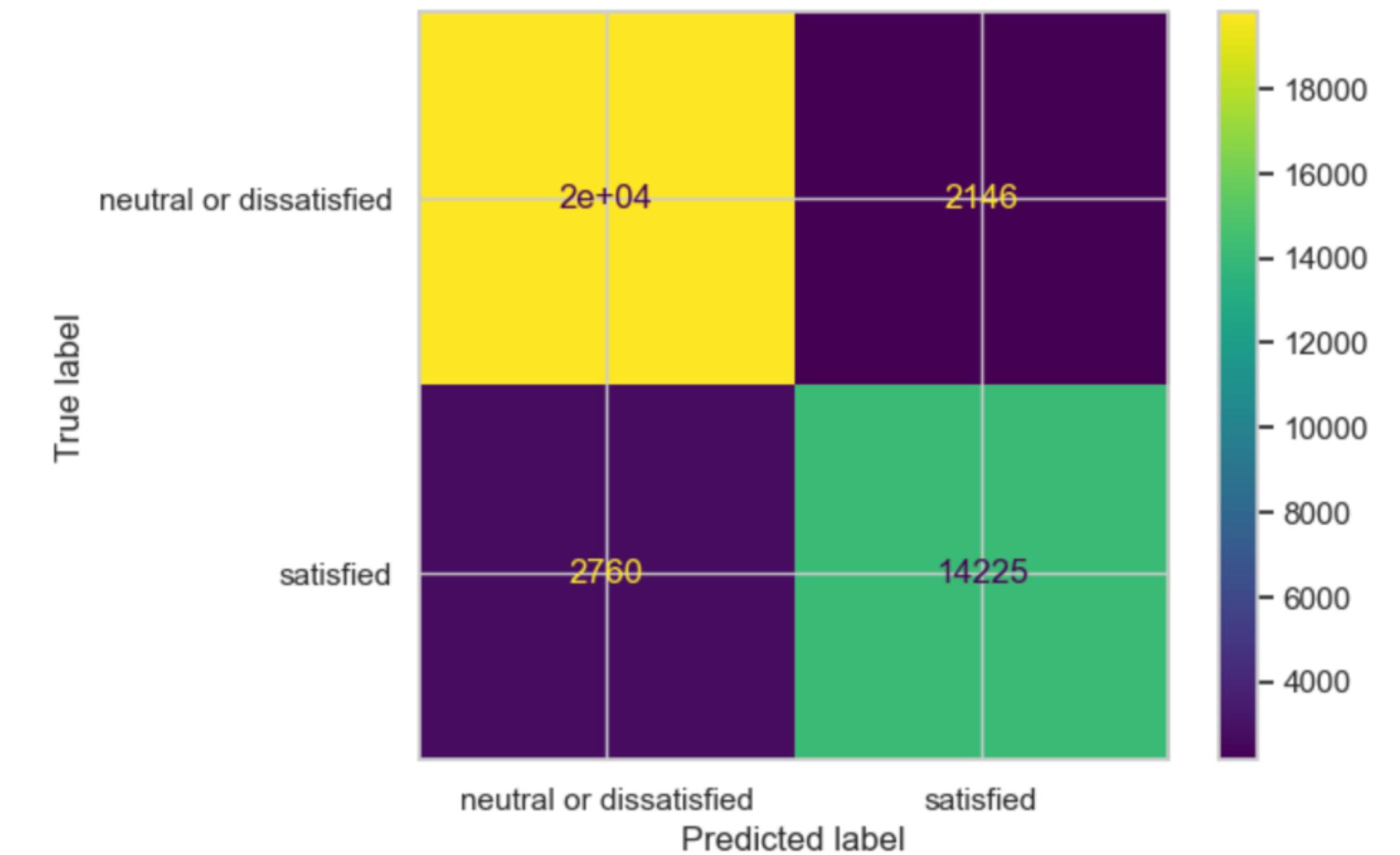
```
In [40]: log_r = get_logr()
```

```
In [41]: predictions_lr = log_r.predict(X_test)
print_metrics(y_true=y_test, y_pred=predictions_lr)
```

AUC-ROC: 0.8699325123187744

Classification:

	precision	recall	f1-score	support
neutral or dissatisfied	0.88	0.90	0.89	21979
satisfied	0.87	0.84	0.85	16985
accuracy			0.87	38964
macro avg	0.87	0.87	0.87	38964
weighted avg	0.87	0.87	0.87	38964



Random Forest with Randomized Search CV

- *Random Forest*
 - *Supervised Machine Learning Algorithm*
 - *used widely in Classification and Regression problems*
- **Randomized Search CV**
 - implements a “**fit**” and a “**score**” method
 - a fixed number of parameter settings is sampled from the specified distributions

```
Random Forest with Randomized Search CV

In [45]: def get_rf_rs(mode = "load"):

    path = DIRPATH + "/Models/"
    filename = 'rf_rs.sav'

    if mode == "load":
        rf_rs = pickle.load(open(path + filename, 'rb'))

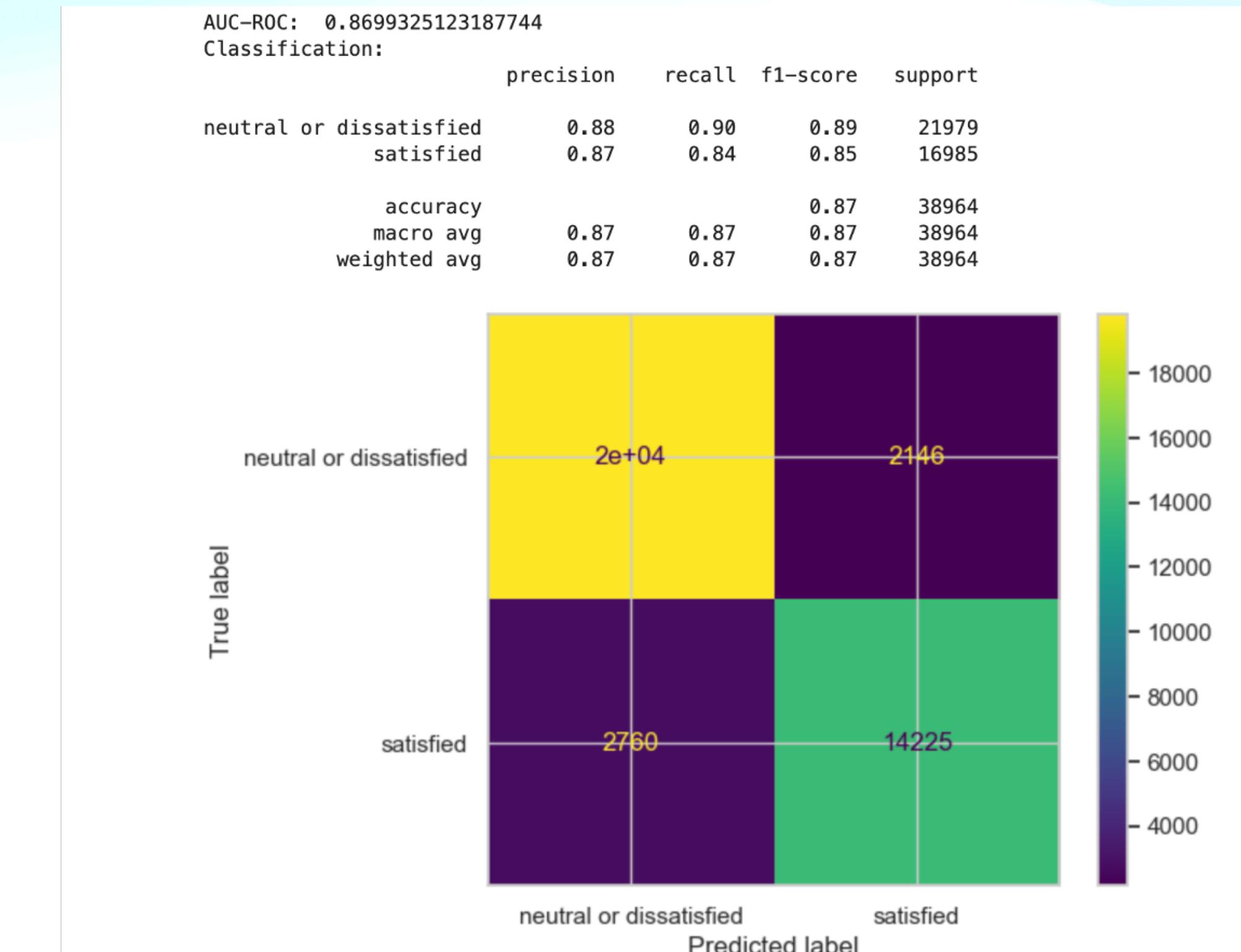
    else:
        rf_parameters_grid = {
            'max_depth': [5, 15],
            'min_samples_leaf': [2, 8],
            'n_estimators': [50, 100],
            'max_features': [5, 10]
        }

        # define grid search
        rf_rs = RandomizedSearchCV(
            estimator = RandomForestClassifier(),
            param_distributions = rf_parameters_grid,
            cv = 10,
            verbose = 1,
            n_jobs = -1
        )
        rf_rs.fit(X_train, y_train)
        pickle.dump(rf_rs, open(path + filename, 'wb'))

    return rf_rs

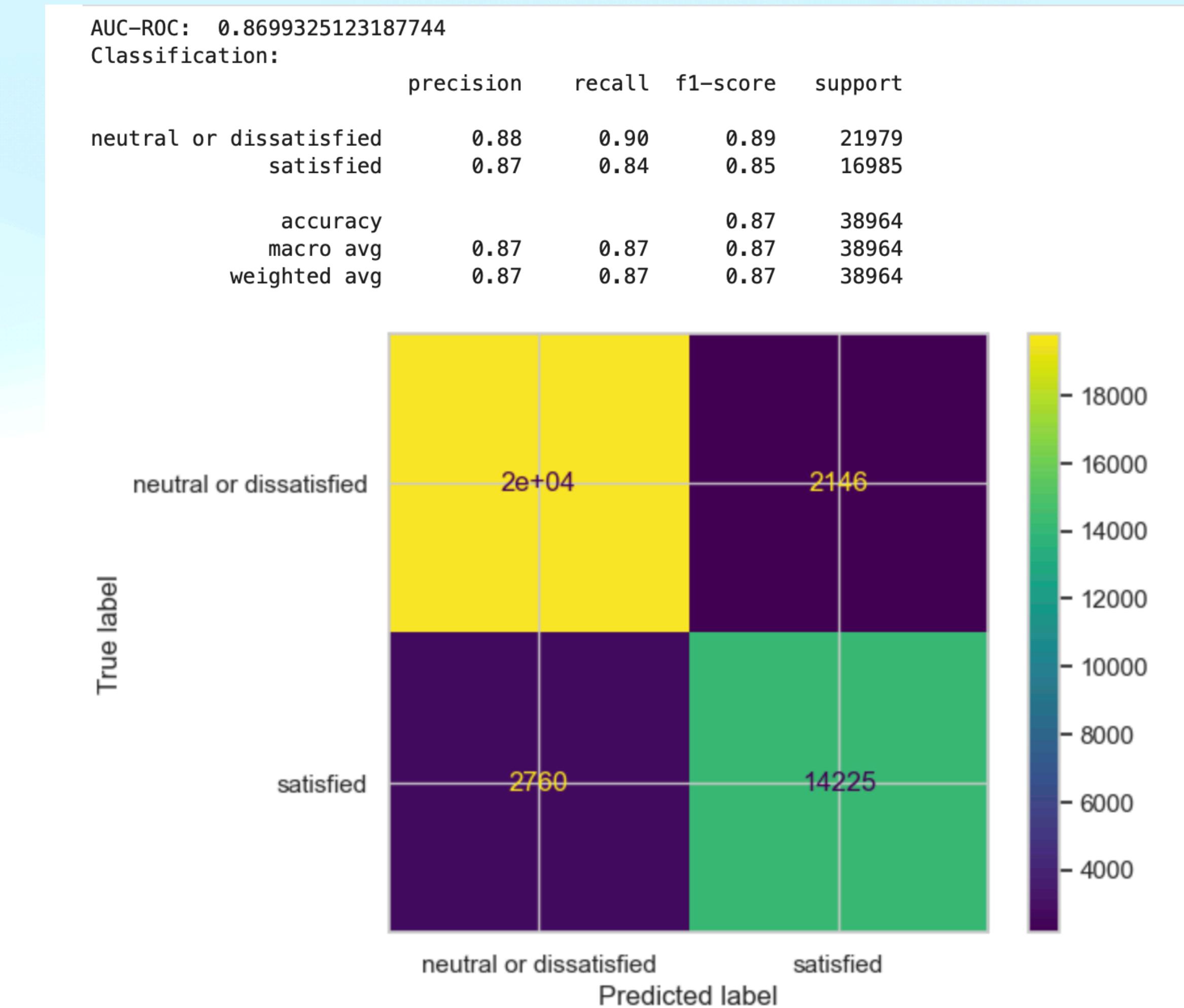
In [46]: rf_rs = get_rf_rs()

In [47]: predictions_rf_rs = rf_rs.predict(X_test)
print_metrics(y_true = y_test, y_pred = predictions_rf_rs)
```



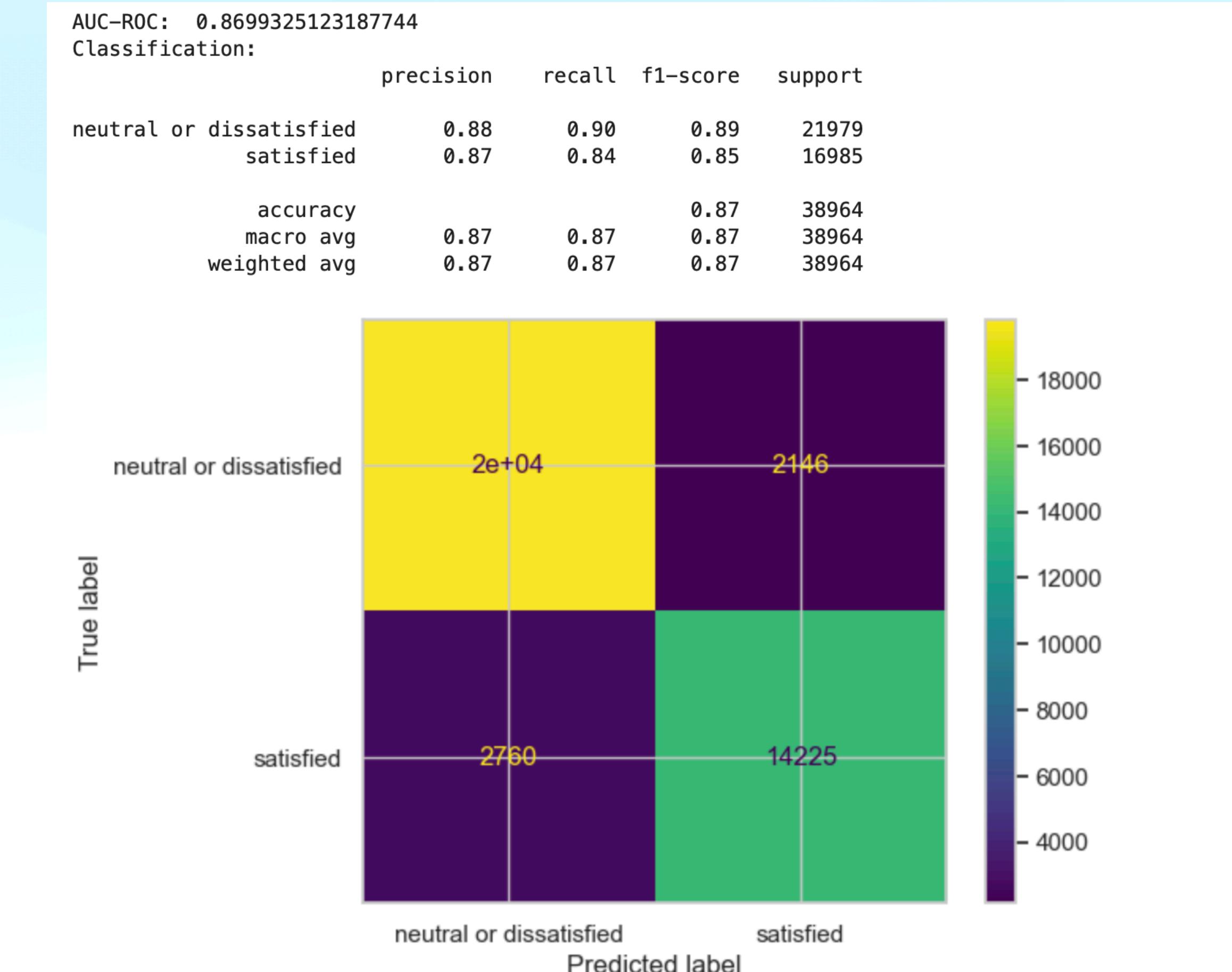
K-nearest neighbor with Grid Search CV

- Neighbors-based classification is a type of instance-based learning or non-generalizing learning.
- Key parameters search using grid search
 - `n_neighbors: int, default=5`
 - `Weights: {'uniform', 'distance'}` or `callable, default='uniform'`



Decision Tree Classifier with Grid Search CV

- Decision Trees (DTs) are a non-parametric supervised learning method used for classification and regression.
- Key parameters search using grid search
 - `min_samples_split`: *int or float, default=2*
 - `max_depth`: *int, default=None*



Support Vector Machine with Grid Search CV

Support vector machines (SVMs) are a set of supervised learning methods used for classification, regression and outliers detection.

The advantages of support vector machines are:

- Effective in high dimensional spaces.
- Still effective in cases where number of dimensions is greater than the number of samples.
- Uses a subset of training points in the decision function (called support vectors), so it is also memory efficient.
- Versatile: different Kernel functions can be specified for the decision function. Common kernels are provided, but it is also possible to specify custom kernels.

```

def get_svm(mode = 'load'):

    path = DIRPATH + "/Models/"
    filename = 'svc.sav'

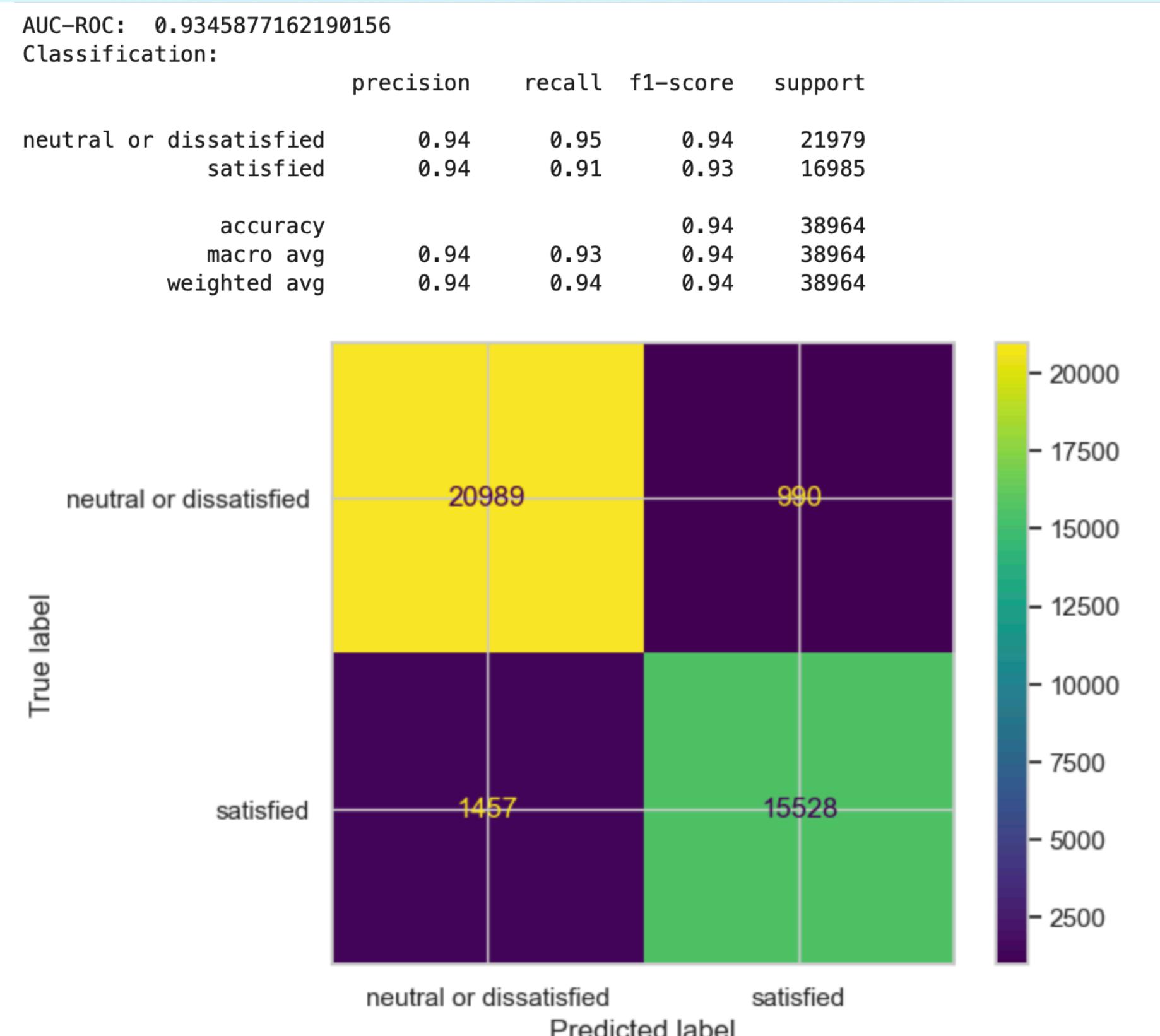
    if mode == "load":
        svc = pickle.load(open(path + filename, 'rb'))
    else:
        # defining parameter range
        param_grid = {'C': [0.1, 1, 10],
                      'kernel': ['rbf'],
                      'gamma': ['scale', 'auto']}

        svc = GridSearchCV(SVC(random_state=False), param_grid, verbose = 3, cv=5, n_jobs=-1)

        svc.fit(X_new_train, y_train)
        pickle.dump(svc, open(path + filename, 'wb'))

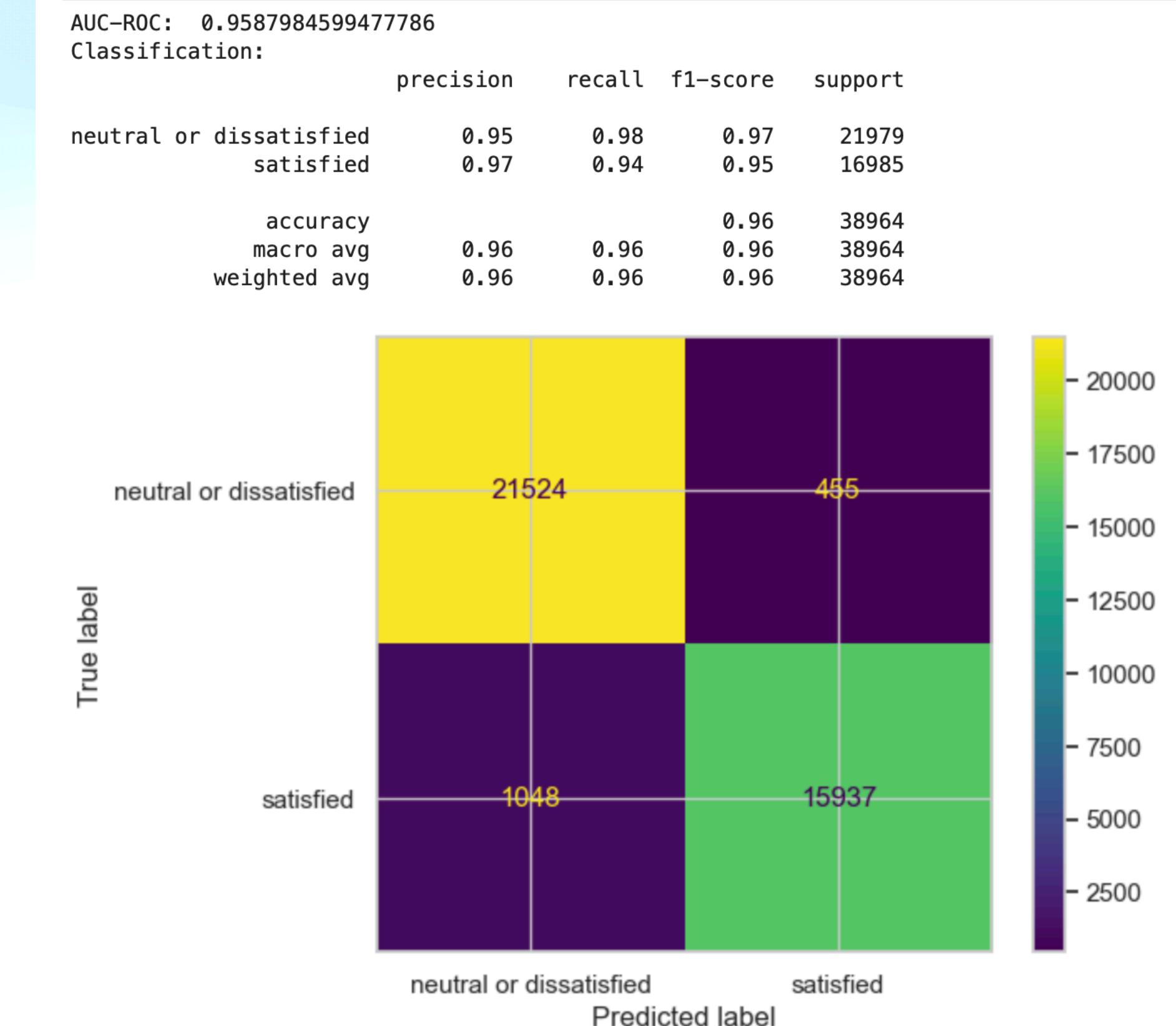
    return svc

```



Bagging along with Random Forest

Random forests or random decision forests is an ensemble learning method for classification, regression and other tasks that operates by constructing a multitude of decision trees at training time.



Multilayer Classifier

Metal device set to: Apple M2
Model: "sequential_9"

Layer (type)	Output Shape	Param #
<hr/>		
dense_20 (Dense)	(None, 22)	550
dropout_5 (Dropout)	(None, 22)	0
dense_21 (Dense)	(None, 11)	253
dense_22 (Dense)	(None, 11)	132
dense_23 (Dense)	(None, 1)	12
<hr/>		

Total params: 947

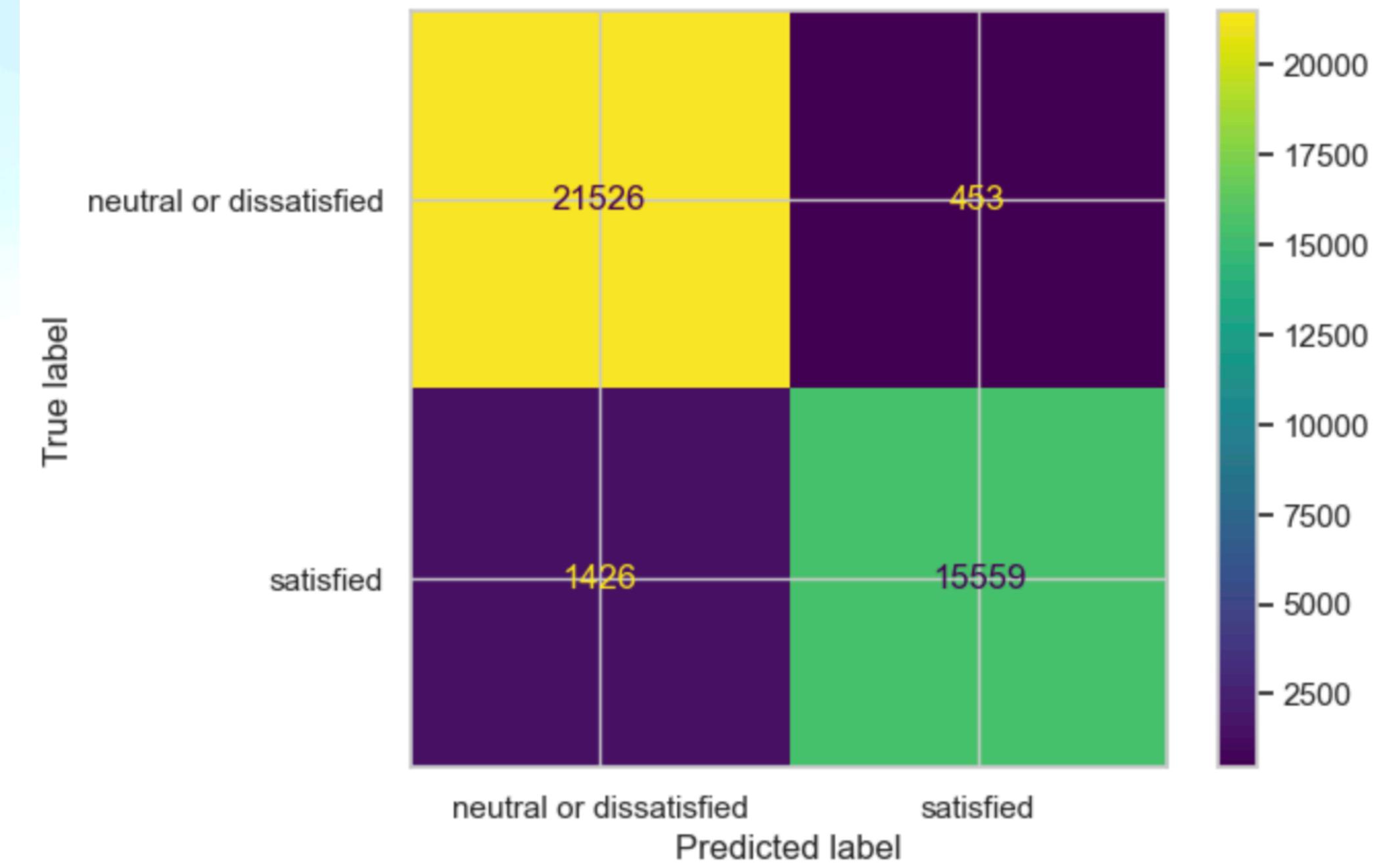
Trainable params: 947

Non-trainable params: 0

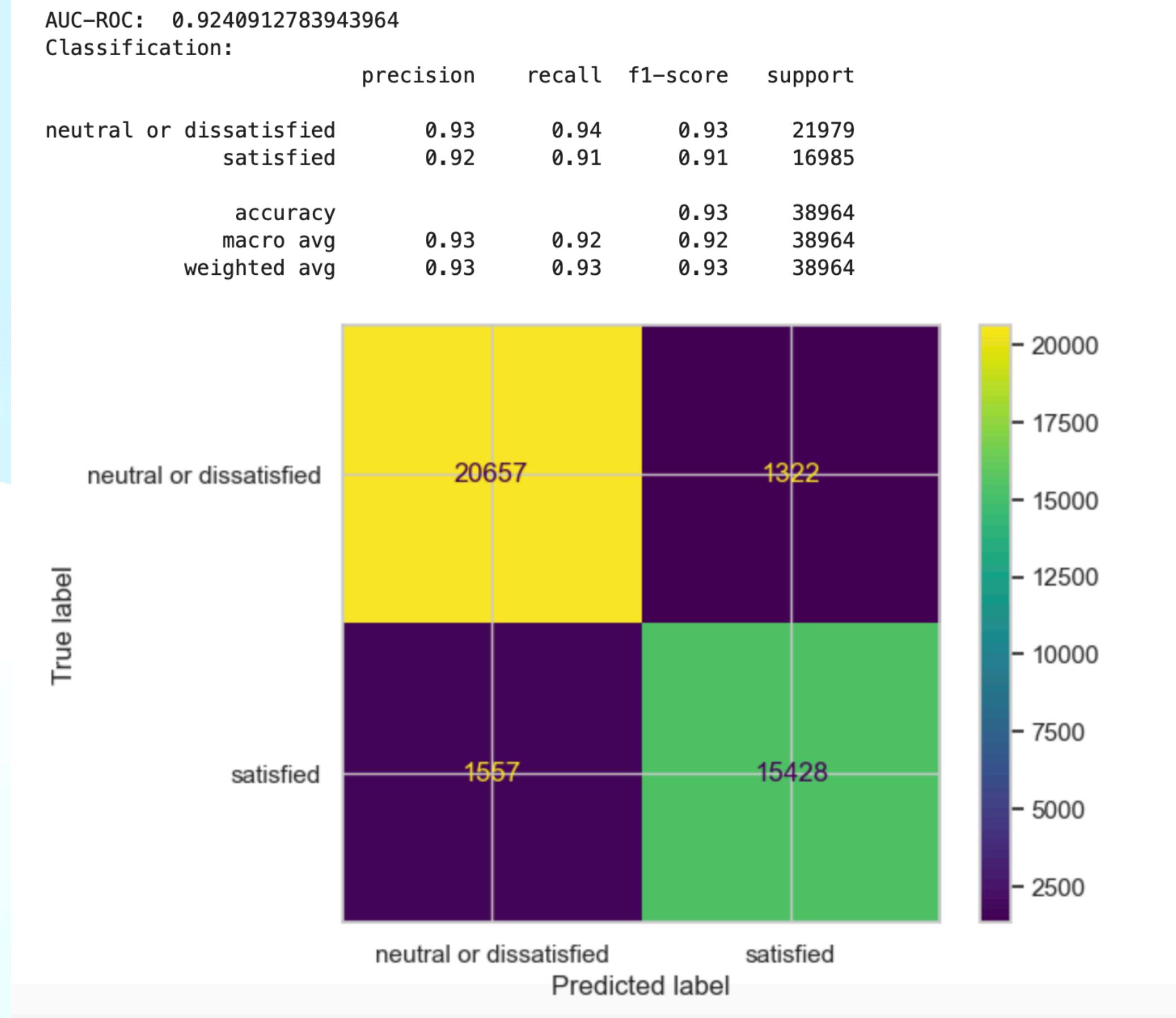
AUC-ROC: 0.9477164925124624

Classification:

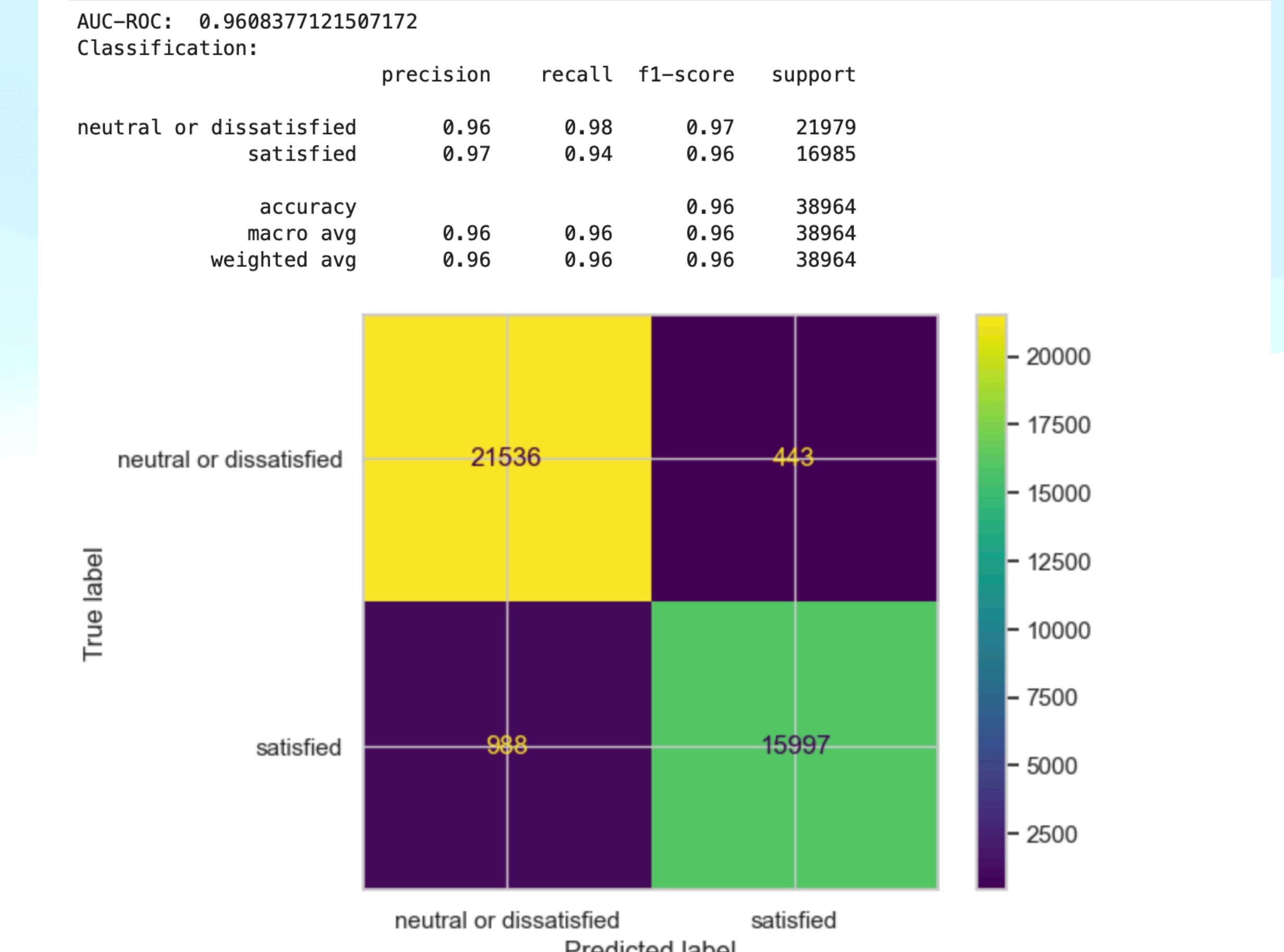
	precision	recall	f1-score	support
neutral or dissatisfied	0.94	0.98	0.96	21979
satisfied	0.97	0.92	0.94	16985
accuracy			0.95	38964
macro avg	0.95	0.95	0.95	38964
weighted avg	0.95	0.95	0.95	38964



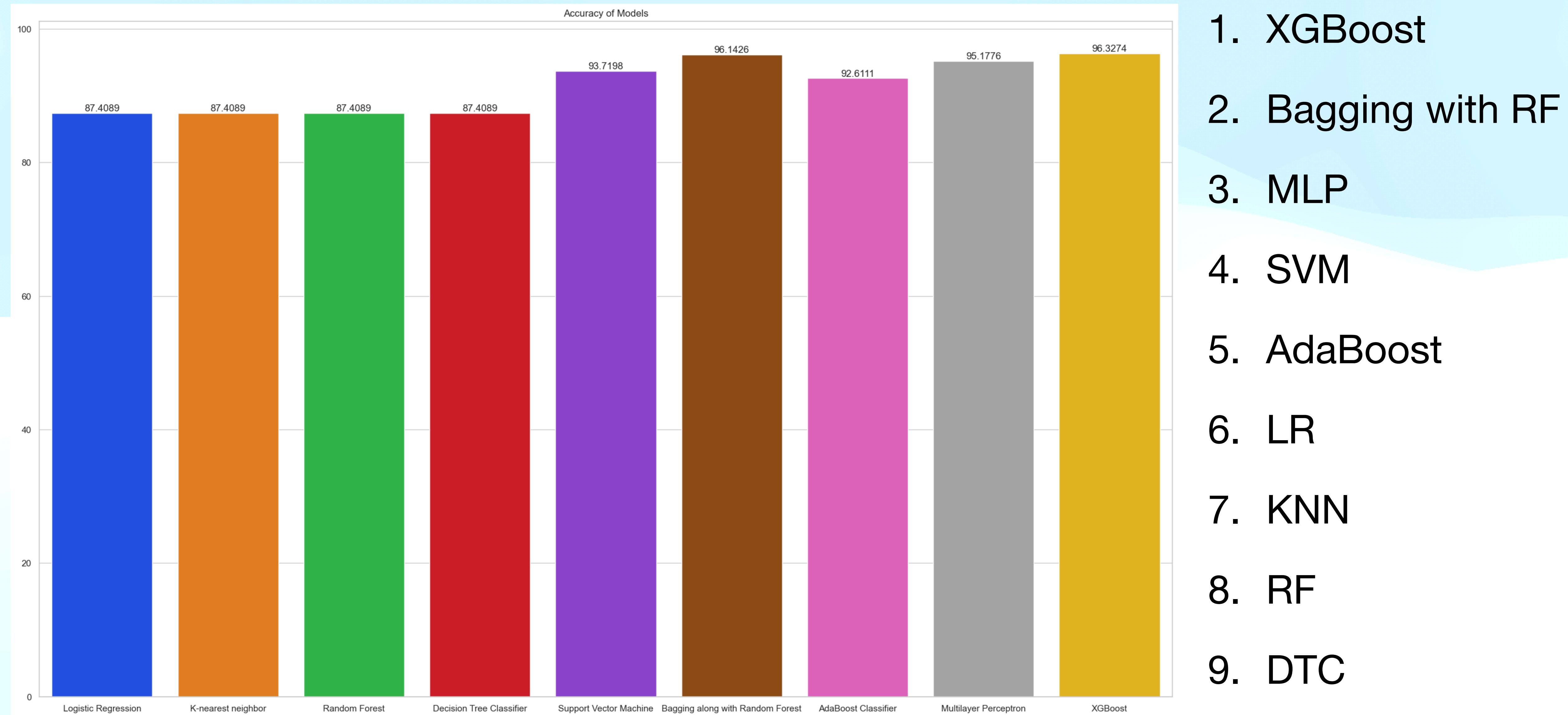
AdaBoost Classifier



Gradient Boosting XGBoost

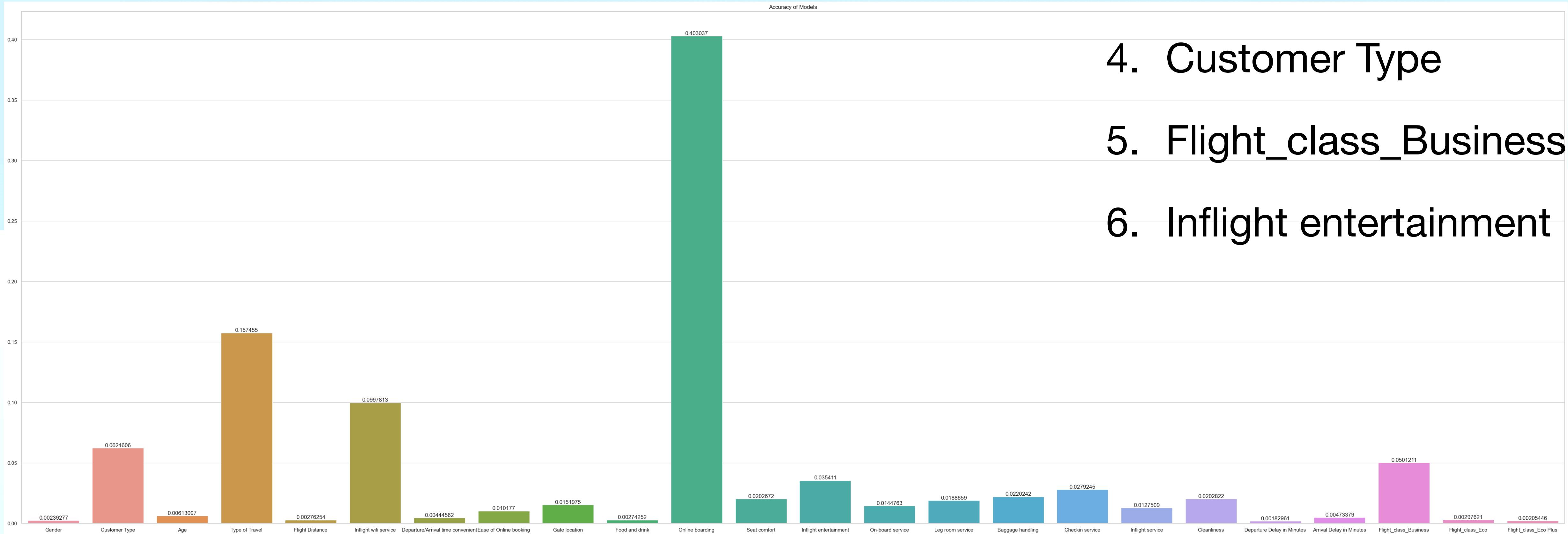


Accuracy of Models



Important features XGB

1. Online boarding
2. Type of Travel
3. Inflight wifi service
4. Customer Type
5. Flight_class_Business
6. Inflight entertainment





Thank You