# Task 2

Team O

# Overview

- Problem Statement
- Raw Data
- Data Pre-Processing
- Outlier Handling
- Correlation Matrix

- EDA
- Performance Metrics
- Feature Importance

# Problem Statement

## Objectives

- Explore and prepare the dataset

- Train a machine learning model

- Evaluate model and present findings

# Raw Data

| | num_passengers | sales_channel | trip_type | purchase_lead | length_of_stay | flight_hour | flight_day | route | booking_origin | wants_extra_baggage | wants_preferred_seat | wants_in_flight_meals | flight_duration | booking_complete |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2 | Internet | RoundTrip | 262 | 19 | 7 | Sat | AKLDEL | New Zealand | 1 | 0 | 0 | 5.52 | 0 |
| 1 | 1 | Internet | RoundTrip | 112 | 20 | 3 | Sat | AKLDEL | New Zealand | 0 | 0 | 0 | 5.52 | 0 |
| 2 | 2 | Internet | RoundTrip | 243 | 22 | 17 | Wed | AKLDEL | India | 1 | 1 | 0 | 5.52 | 0 |
| 3 | 1 | Internet | RoundTrip | 96 | 31 | 4 | Sat | AKLDEL | New Zealand | 0 | 0 | 1 | 5.52 | 0 |
| 4 | 2 | Internet | RoundTrip | 68 | 22 | 15 | Wed | AKLDEL | India | 1 | 0 | 1 | 5.52 | 0 |

```
num_passengers          0
sales_channel           0
trip_type               0
purchase_lead           0
length_of_stay          0
flight_hour             0
flight_day              0
route                   0
booking_origin          0
wants_extra_baggage     0
wants_preferred_seat    0
wants_in_flight_meals   0
flight_duration         0
booking_complete        0
dtype: int64
```

```
(50000, 14)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 50000 entries, 0 to 49999
Data columns (total 14 columns):
 #   Column                 Non-Null Count  Dtype
---  ------                 --------------  -----
 0   num_passengers         50000 non-null  int64
 1   sales_channel          50000 non-null  object
 2   trip_type              50000 non-null  object
 3   purchase_lead          50000 non-null  int64
 4   length_of_stay         50000 non-null  int64
 5   flight_hour            50000 non-null  int64
 6   flight_day             50000 non-null  object
 7   route                  50000 non-null  object
 8   booking_origin         50000 non-null  object
 9   wants_extra_baggage    50000 non-null  int64
 10  wants_preferred_seat   50000 non-null  int64
 11  wants_in_flight_meals  50000 non-null  int64
 12  flight_duration        50000 non-null  float64
 13  booking_complete       50000 non-null  int64
dtypes: float64(1), int64(8), object(5)
memory usage: 5.3+ MB
```

# Data Pre-Processing

| | num_passengers | purchase_lead | length_of_stay | flight_hour | wants_extra_baggage | wants_preferred_seat | wants_in_flight_meals | flight_duration | booking_complete |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2 | 262 | 19 | 7 | 1 | 0 | 0 | 5.52 | 0 |
| 1 | 1 | 112 | 20 | 3 | 0 | 0 | 0 | 5.52 | 0 |
| 2 | 2 | 243 | 22 | 17 | 1 | 1 | 0 | 5.52 | 0 |
| 3 | 1 | 96 | 31 | 4 | 0 | 0 | 1 | 5.52 | 0 |
| 4 | 2 | 68 | 22 | 15 | 1 | 0 | 1 | 5.52 | 0 |

| | sales_channel | trip_type | flight_day | route | booking_origin |
|---|---|---|---|---|---|
| 0 | Internet | RoundTrip | Sat | AKLDEL | New Zealand |
| 1 | Internet | RoundTrip | Sat | AKLDEL | New Zealand |
| 2 | Internet | RoundTrip | Wed | AKLDEL | India |
| 3 | Internet | RoundTrip | Sat | AKLDEL | New Zealand |
| 4 | Internet | RoundTrip | Wed | AKLDEL | India |

```
sales_channel
Internet    43917
Mobile       5364
Name: count, dtype: int64


trip_type
RoundTrip    48779
OneWay         386
CircleTrip     116
Name: count, dtype: int64


flight_day
Mon    7988
Wed    7562
Tue    7558
Thu    7323
Fri    6685
Sun    6442
Sat    5723
Name: count, dtype: int64


route
...
Svalbard & Jan Mayen         1
Name: count, Length: 104, dtype: int64
```
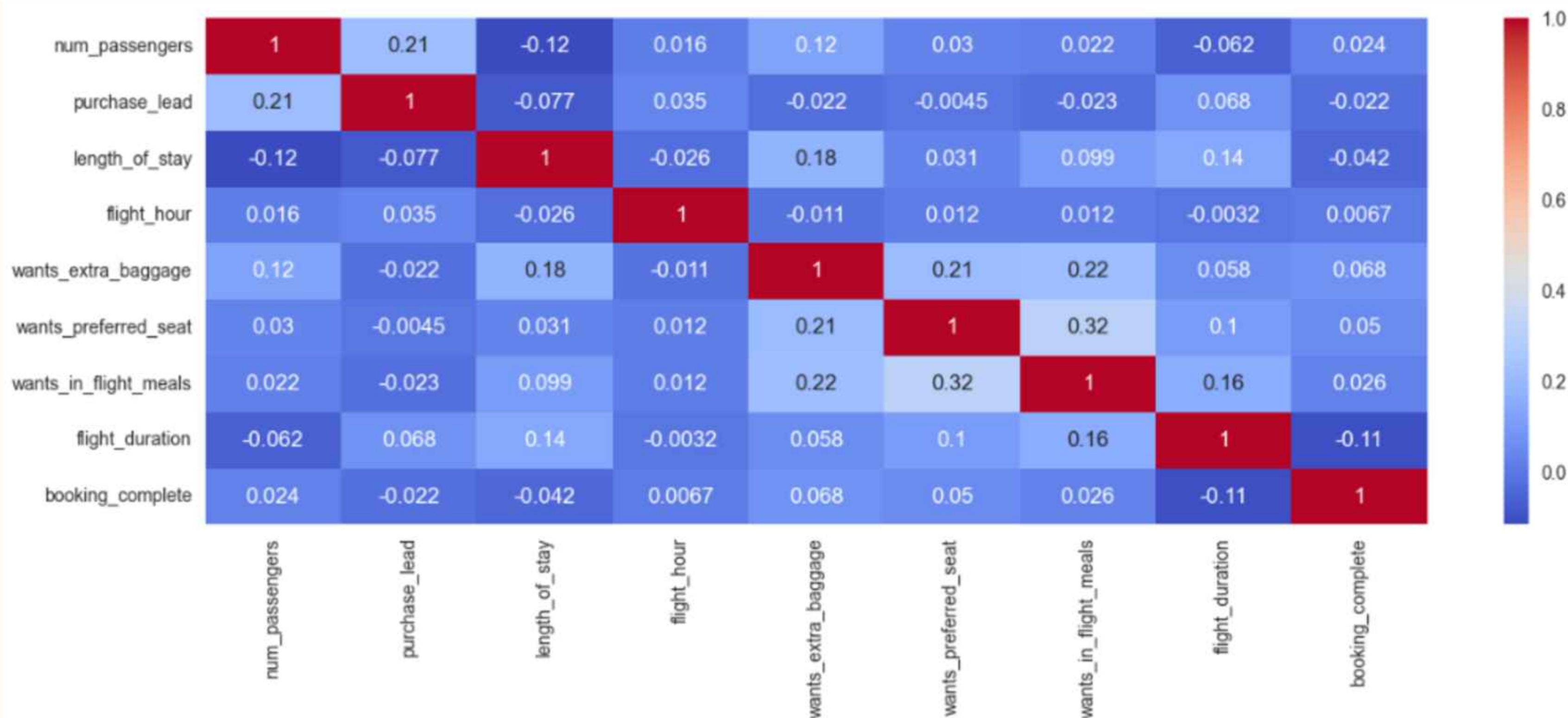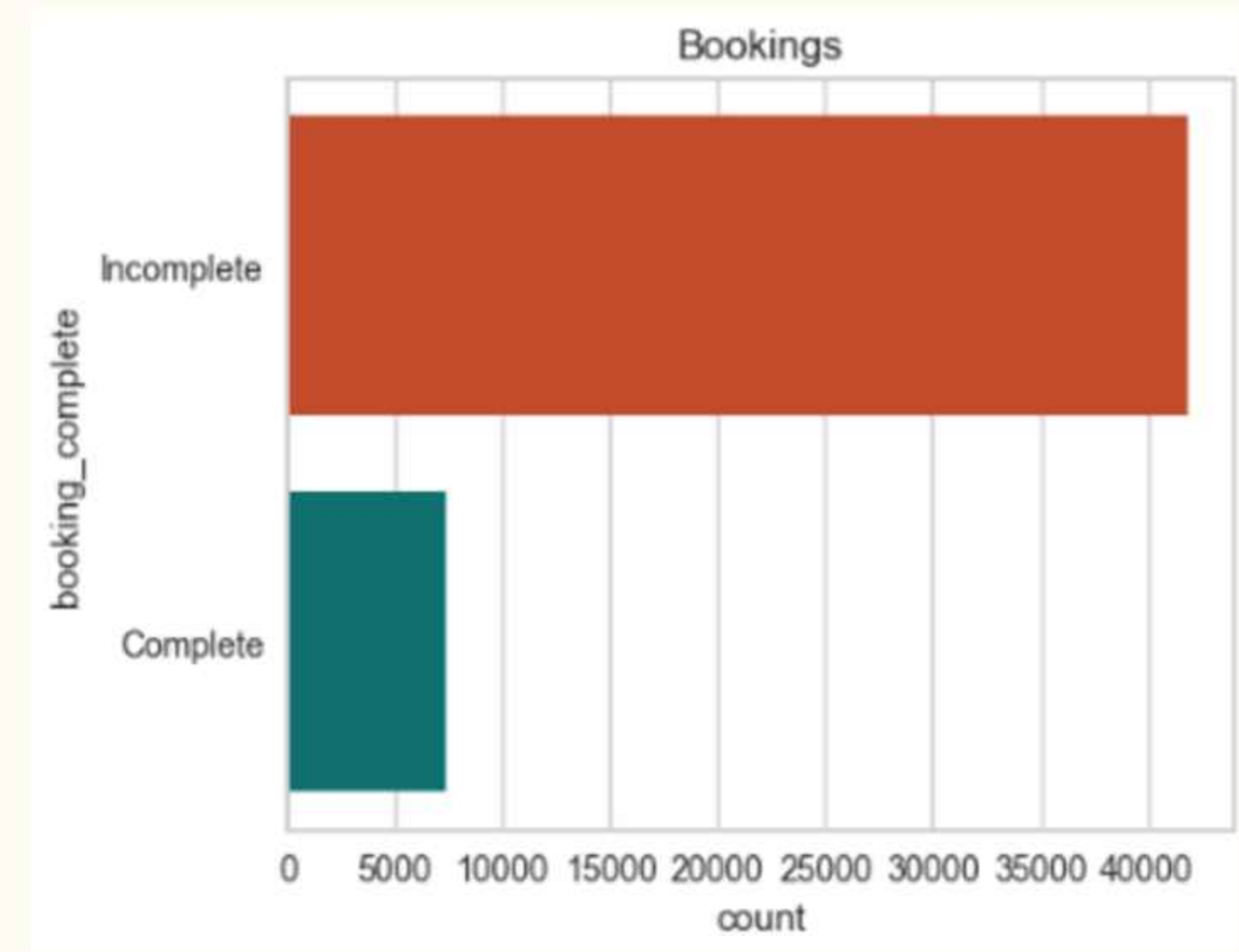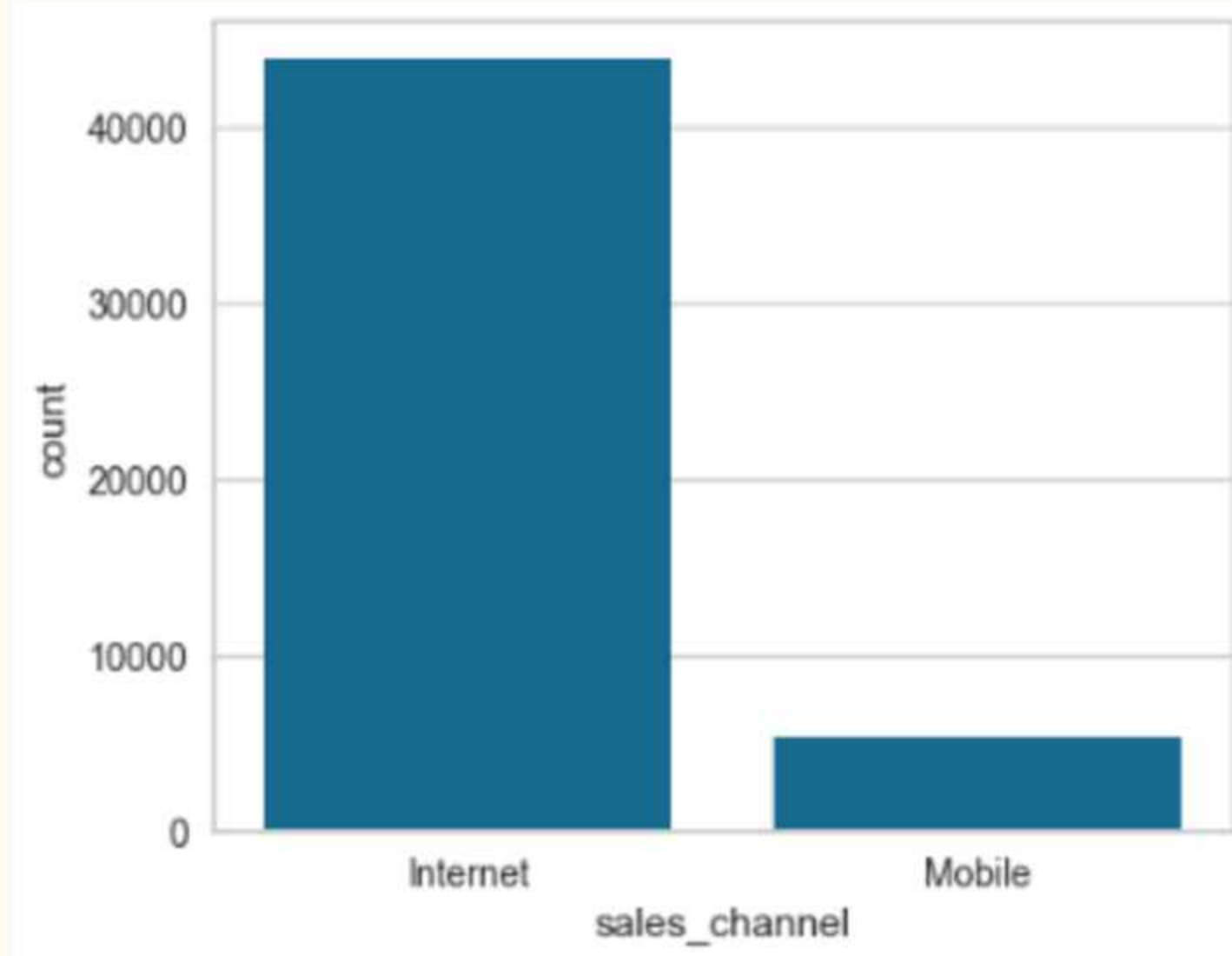
# Outlier Handling

EDA

# Performance Metrics

Accuracy: 0.9195
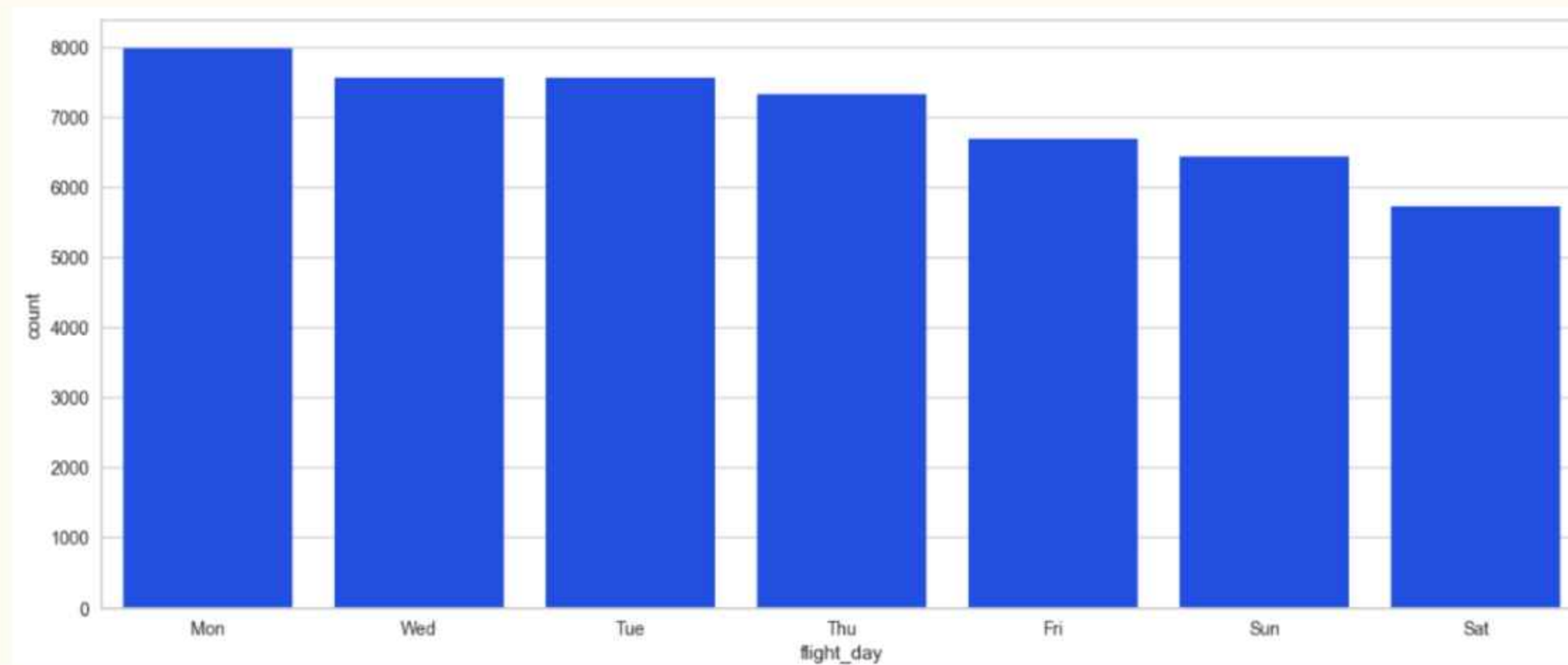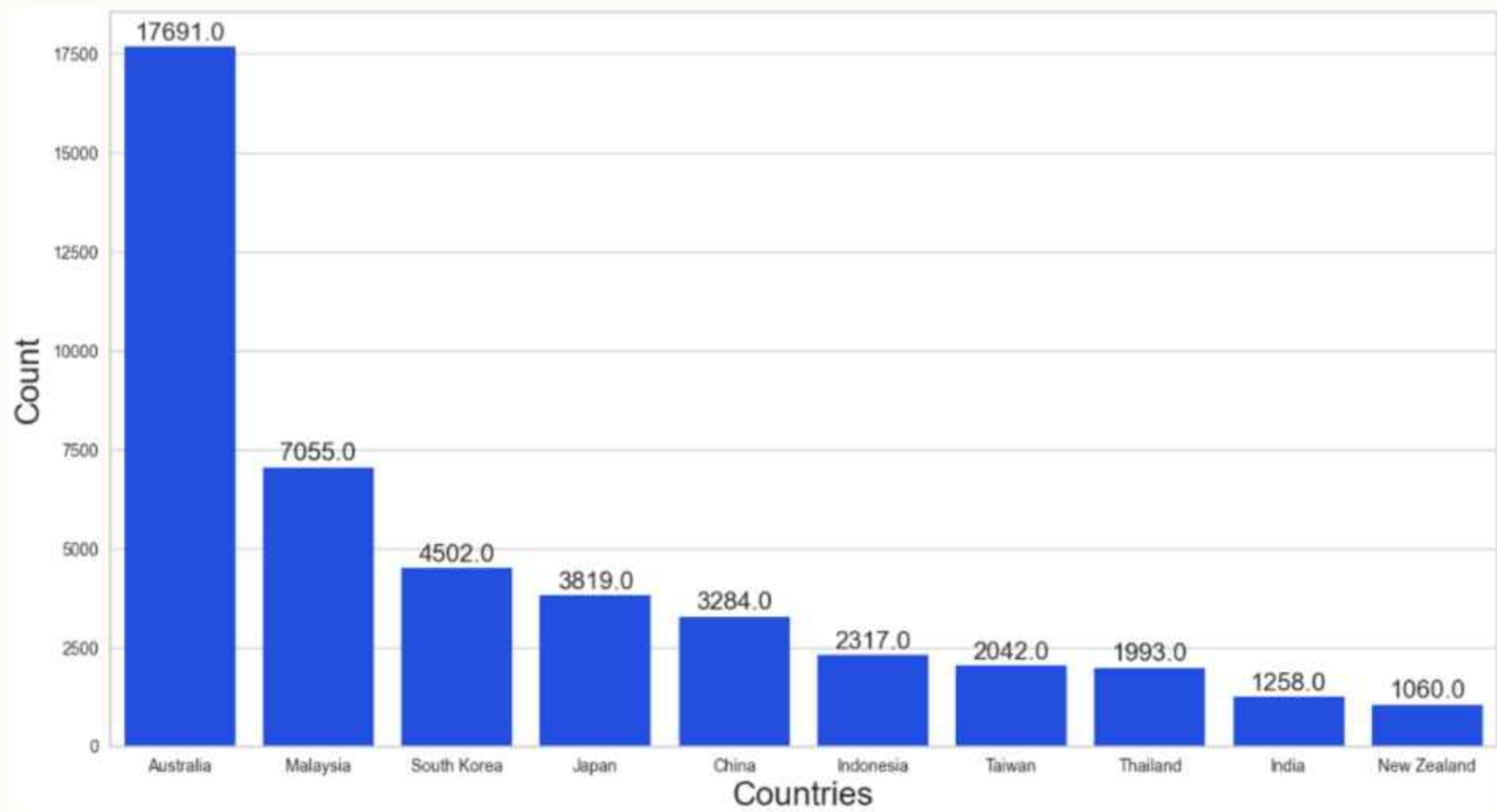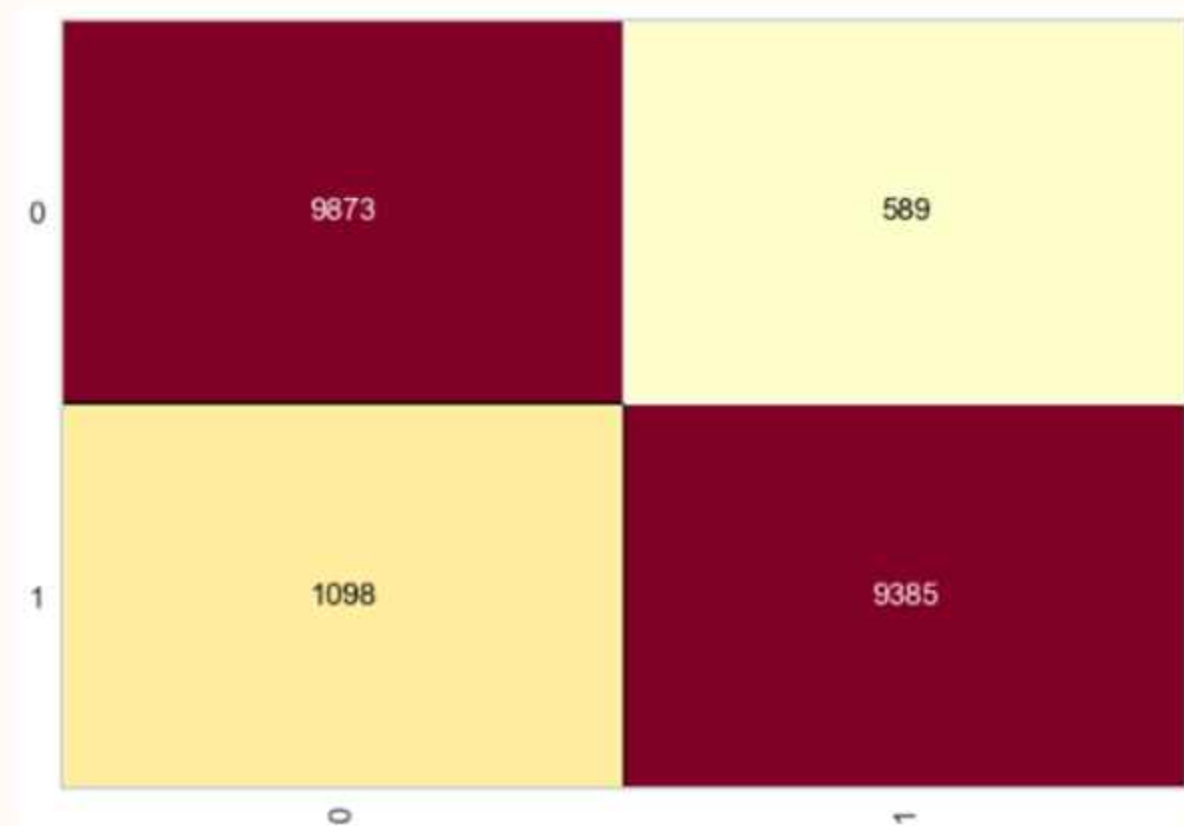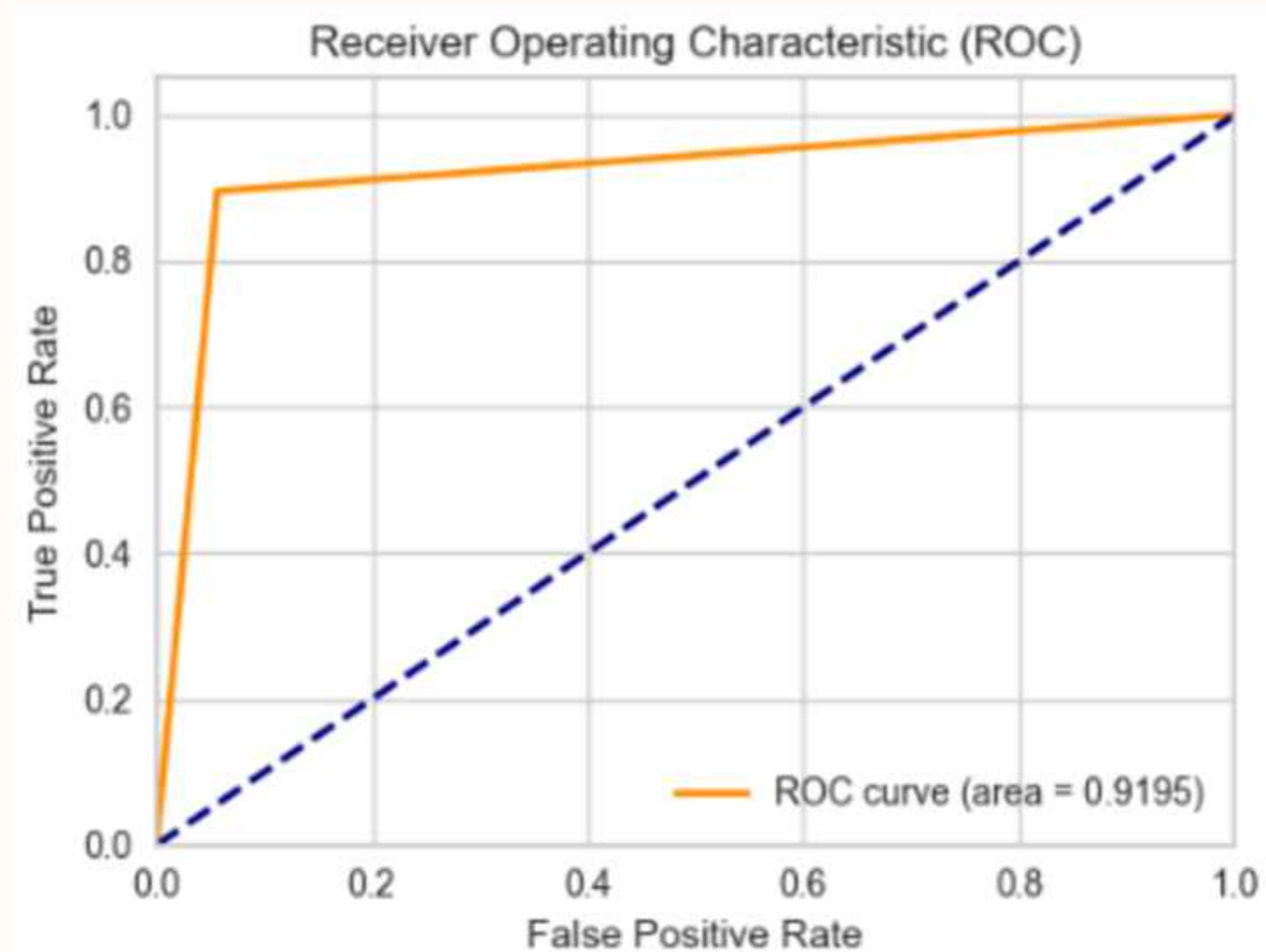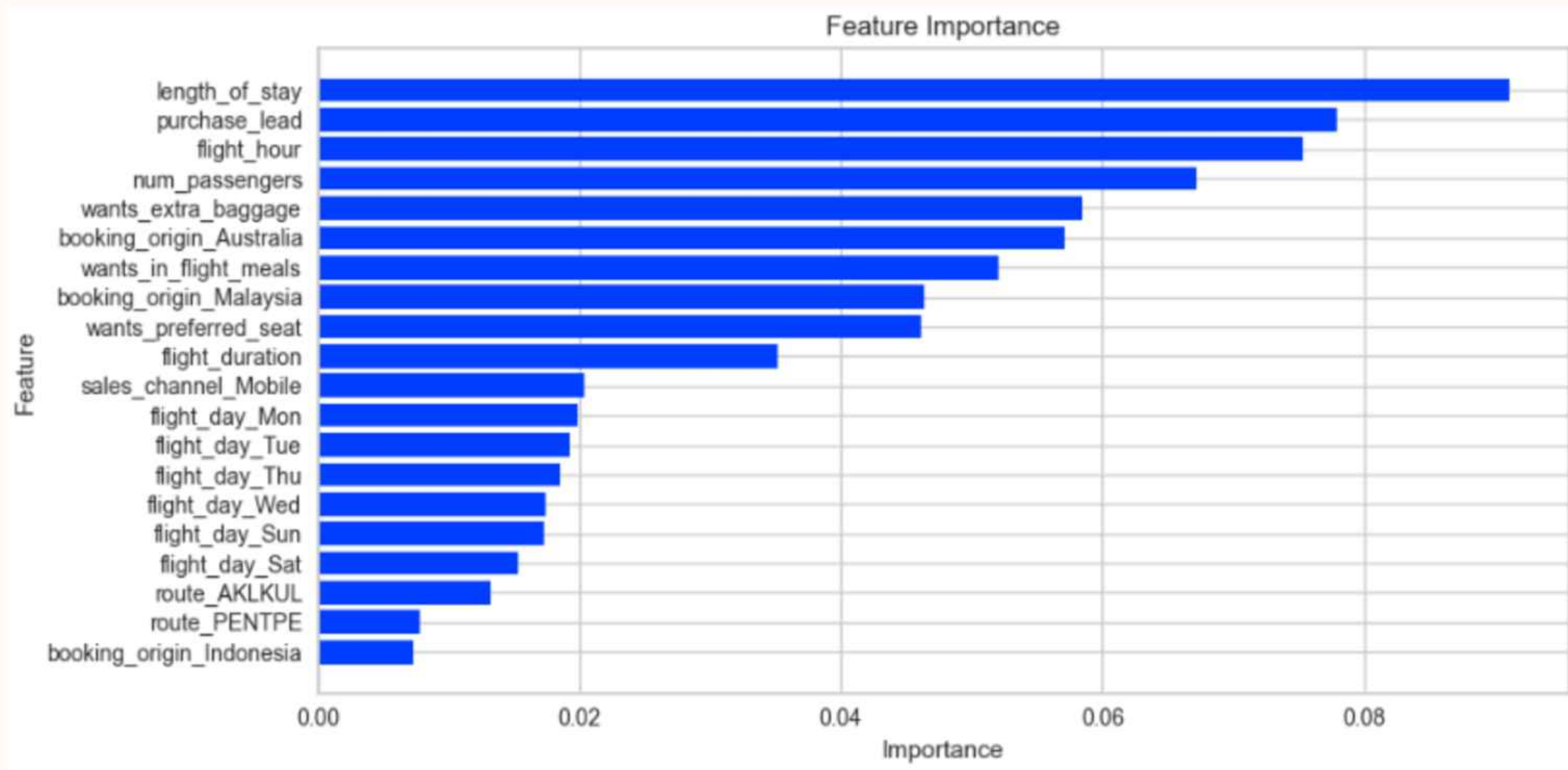AUC-ROC Score: 0.92

# Feature Importance

Feature Importance

# Thank You