

```
In [12]: import pandas as pd

In [13]: df = pd.read_csv("C:\USRR DATA\Data Sciences\machine_learning\Projects\Forage\customer_booking.csv", encoding='latin1')
df.head()

Out[13]:
```

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

Data Understanding

```
In [14]: df.info()

<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

In [15]: df['flight_day'].unique()

Out[15]: array(['Sat', 'Wed', 'Thu', 'Mon', 'Sun', 'Tue', 'Fri'], dtype=object)

In [16]: mapping = {
    "Mon": 1,
    "Tue": 2,
    "Wed": 3,
    "Thu": 4,
    "Fri": 5,
    "Sat": 6,
    "Sun": 7,
}

df['flight_day'] = df['flight_day'].map(mapping)

In [17]: df['flight_day'].unique()

Out[17]: array([6, 3, 4, 1, 7, 2, 5], dtype=int64)

In [18]: df.describe()

Out[18]:
```

	num_passengers	purchase_lead	length_of_stay	flight_hour	flight_day	wants_extra_baggage	wants_preferred_seat	wants_in_flight_meals	flight_duration
count	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000	50000.000000
mean	1.591240	84.940480	23.044656	9.06634	3.814420	0.668780	0.296960	0.427140	7.277561
std	1.020165	90.451378	33.88767	9.41266	1.992792	0.470657	0.456923	0.494668	1.496863
min	1.000000	0.000000	0.00000	0.00000	1.000000	0.000000	0.000000	0.000000	4.670000
25%	1.000000	21.000000	5.00000	5.00000	2.000000	0.000000	0.000000	0.000000	5.620000
50%	1.000000	51.000000	17.00000	9.00000	4.000000	1.000000	0.000000	0.000000	7.570000
75%	2.000000	115.000000	28.00000	13.00000	5.000000	1.000000	1.000000	1.000000	8.830000
max	9.000000	867.000000	778.00000	23.00000	7.000000	1.000000	1.000000	1.000000	9.500000

```
In [19]: df[['sales_channel', 'trip_type', 'booking_origin', 'route']].describe()

Out[19]:
```

	sales_channel	trip_type	booking_origin	route
count	50000	50000	50000	50000
unique	2	3	104	799
top	Internet	RoundTrip	Australia	AKLJUL
freq	44382	49497	17672	2680

```
In [120]: df_categories=['wants_extra_baggage', 'wants_preferred_seat', 'flight_day', 'wants_in_flight_meals',
                    'booking_complete']
for category in df_categories:
    df['category_'+category] = df[category].astype('category')
df.describe(include='category')

Out[120]:
```

	flight_day	wants_extra_baggage	wants_preferred_seat	wants_in_flight_meals	booking_complete
count	50000	50000	50000	50000	50000
unique	7	2	2	2	2
top	1	1	0	0	0
freq	8102	33439	35152	28643	42522

Feature Engineering

```
In [111]: df['Origin_Country']=df['route'].str[:3]
df['Destination_Country']=df['route'].str[3:]

In [121]: Timing=[0,12,24]
labels=['AM','PM']
df['Day Timing']=pd.out(df['flight_hour'],bins=Timing,labels=labels,right=False)
df.head()

Out[121]:
```

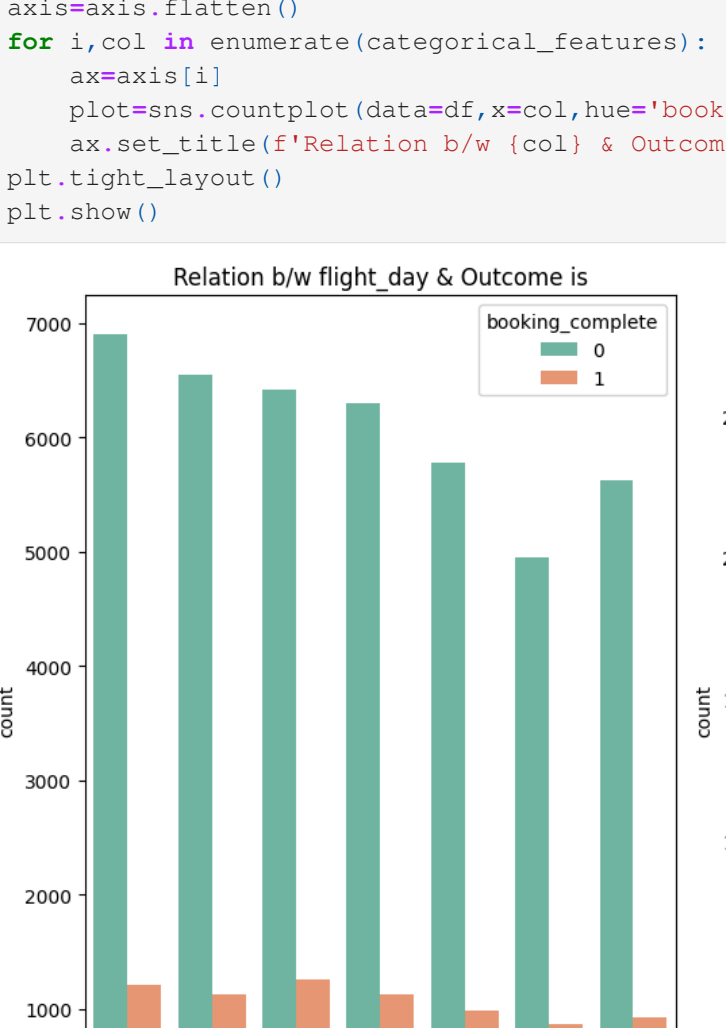
	num_passengers	sales_channel	trip_type	purchase_lead	length_of_stay	flight_hour	flight_day	route	booking_origin	wants_extra_baggage	wants_preferred_seat
0	2	Internet	RoundTrip	262	19	7	Sat	AKLDEL	New Zealand		1
1	1	Internet	RoundTrip	112	20	3	6	AKLDEL	New Zealand		0
2	2	Internet	RoundTrip	243	22	17	3	AKLDEL	India		1
3	1	Internet	RoundTrip	96	31	4	6	AKLDEL	New Zealand		0
4	2	Internet	RoundTrip	68	22	15	3	AKLDEL	India		1

```
In [131]: import matplotlib.pyplot as plt
import seaborn as sns

In [131]: plt.figure(figsize=(4,5))
sns.countplot(data=df,x='booking_complete',palette='Set2')
plt.title('Distribution of target variable(Booking complete)')
for p in plt.patches:
    count=int(p.get_height())
    plot.annotate(str(count), (p.get_x()+p.get_width()/2,p.get_height()+1),ha='center',va='bottom')
plt.tight_layout()

C:\Users\Qadri Laptops\AppData\Local\Temp\ipykernel_2488\1326082283.py:2: FutureWarning:
Passing 'palette' without assigning 'hue' is deprecated and will be removed in v0.14.0. Assign the 'x' variable to 'hue' and set 'legend=False' for the same effect.
plot=sns.countplot(data=df,x='booking_complete',palette='Set2')
```

Distribution of target variable(Booking complete)



```
In [141]: categorical_features=['flight_day', 'wants_extra_baggage', 'wants_preferred_seat', 'wants_in_flight_meals',
                             'Day Timing']
fig,ax=plt.subplots(2,2,figsize=(16,14))
ax[0].axis.flatten()
for i,col in enumerate(categorical_features):
    ax=ax+i
    plot=sns.countplot(data=df,x=col,hue='booking_complete',palette='Set2',ax=ax)
    ax.set_title(f'Relation b/w {col} & Outcome is')
    plt.tight_layout()
plt.show()
```



```
In [151]: categorical_features=['flight_day', 'wants_extra_baggage', 'wants_preferred_seat', 'wants_in_flight_meals',
                             'Day Timing']
for feature in categorical_features:
    output=df.groupby('booking_complete')[feature].value_counts(normalize=True)
    print(f'Value Count in feature {feature} is:\n {round(output*100,ndigits=2)}\n')
```

Value Count in feature flight-day is:

booking_complete flight-day

0 16.22 1

2 15.39

3 15.10

4 14.85

5 13.59

6 13.22

7 11.51

1 16.74

2 15.40

3 15.10

4 15.00

5 12.25

6 12.40

7 11.51

Name: proportion, dtype: float64

Value Count in feature wants_extra_baggage is:

booking_complete wants_extra_baggage

0 65.53

1 34.47

0 74.53

1 25.47

Name: proportion, dtype: float64

Value Count in feature wants_preferred_seat is:

booking_complete wants_preferred_seat

0 71.26

1 28.74

0 64.84

1 35.16

Name: proportion, dtype: float64

Value Count in feature wants_in_flight_meals is:

booking_complete wants_in_flight_meals

0 57.84

1 42.16

0 54.36

1 45.64

Name: proportion, dtype: float64

Value Count in feature Day Timing is:

booking_complete Day Timing

0 AM 66.91

PM 33.09

1 AM 64.66

PM 35.34

Name: proportion, dtype: float64

C:\Users\Qadri Laptops\AppData\Local\Temp\ipykernel_15380\31151519346.py:4: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=True to retain current behavior or observed=True to adopt the future default and silence this warning.

output=df.groupby('booking_complete')[feature].value_counts(normalize=True)

C:\Users\Qadri Laptops\AppData\Local\Temp\ipykernel_15380\31151519346.py:4: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=True to retain current behavior or observed=True to adopt the future default and silence this warning.

output=df.groupby('booking_complete')[feature].value_counts(normalize=True)

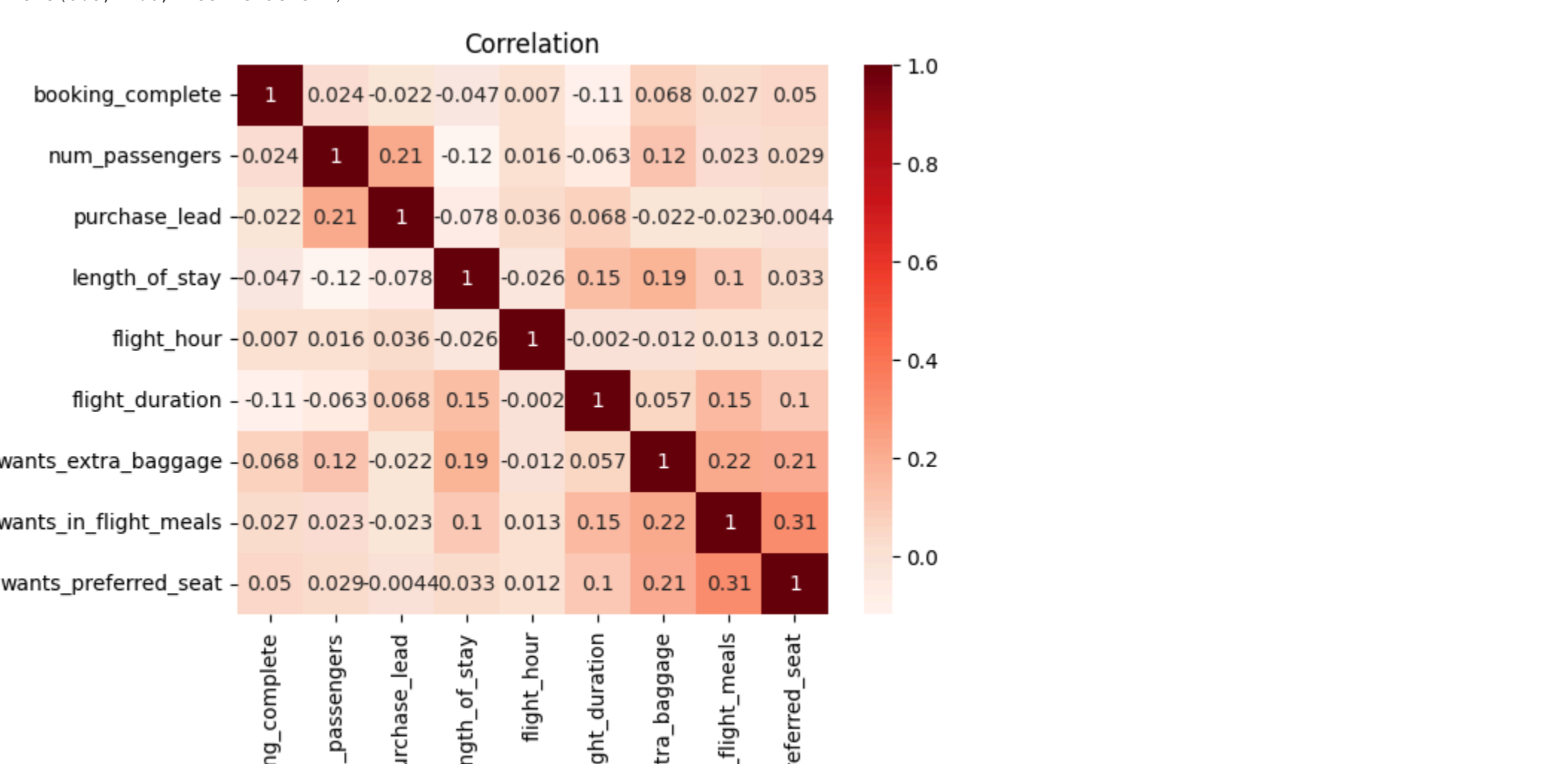
C:\Users\Qadri Laptops\AppData\Local\Temp\ipykernel_15380\31151519346.py:4: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=True to retain current behavior or observed=True to adopt the future default and silence this warning.

output=df.groupby('booking_complete')[feature].value_counts(normalize=True)

C:\Users\Qadri Laptops\AppData\Local\Temp\ipykernel_15380\31151519346.py:4: FutureWarning: The default of observed=False is deprecated and will be changed to True in a future version of pandas. Pass observed=True to retain current behavior or observed=True to adopt the future default and silence this warning.

output=df.groupby('booking_complete')[feature].value_counts(normalize=True)

```
In [161]: numerical_features=['purchase_lead', 'length_of_stay', 'flight_duration', 'num_passengers']
fig,ax=plt.subplots(2,2,figsize=(16,14))
ax[0].axis.flatten()
for i,col in enumerate(numerical_features):
    ax=ax+i
    plot=sns.boxplot(data=df,x=col,hue='booking_complete',palette='Set3',ax=ax)
    ax.set_title(f'Relation b/w {col} & Outcome is')
    plt.tight_layout()
plt.show()
```



```
In [162]: df=df[df['length_of_stay']<365]
df.shape

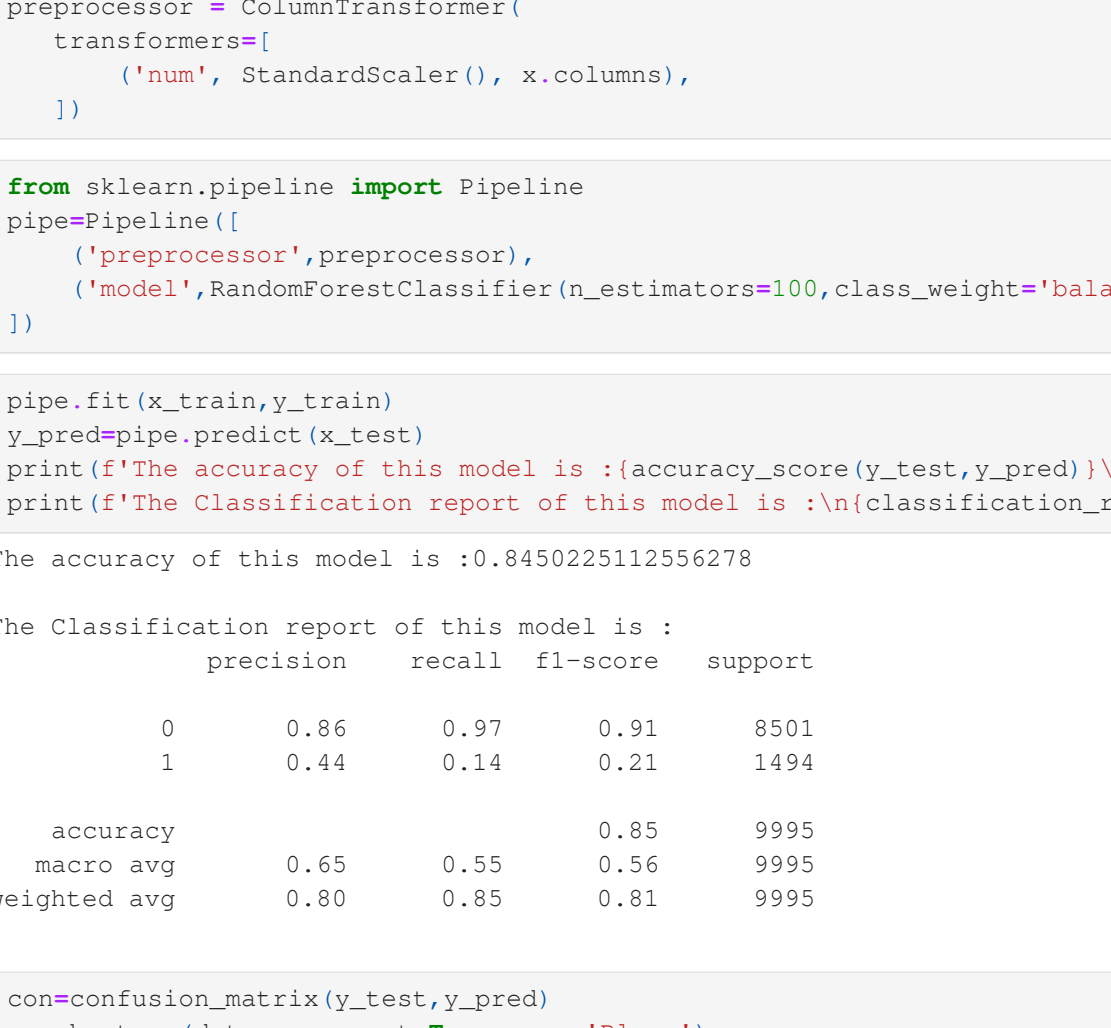
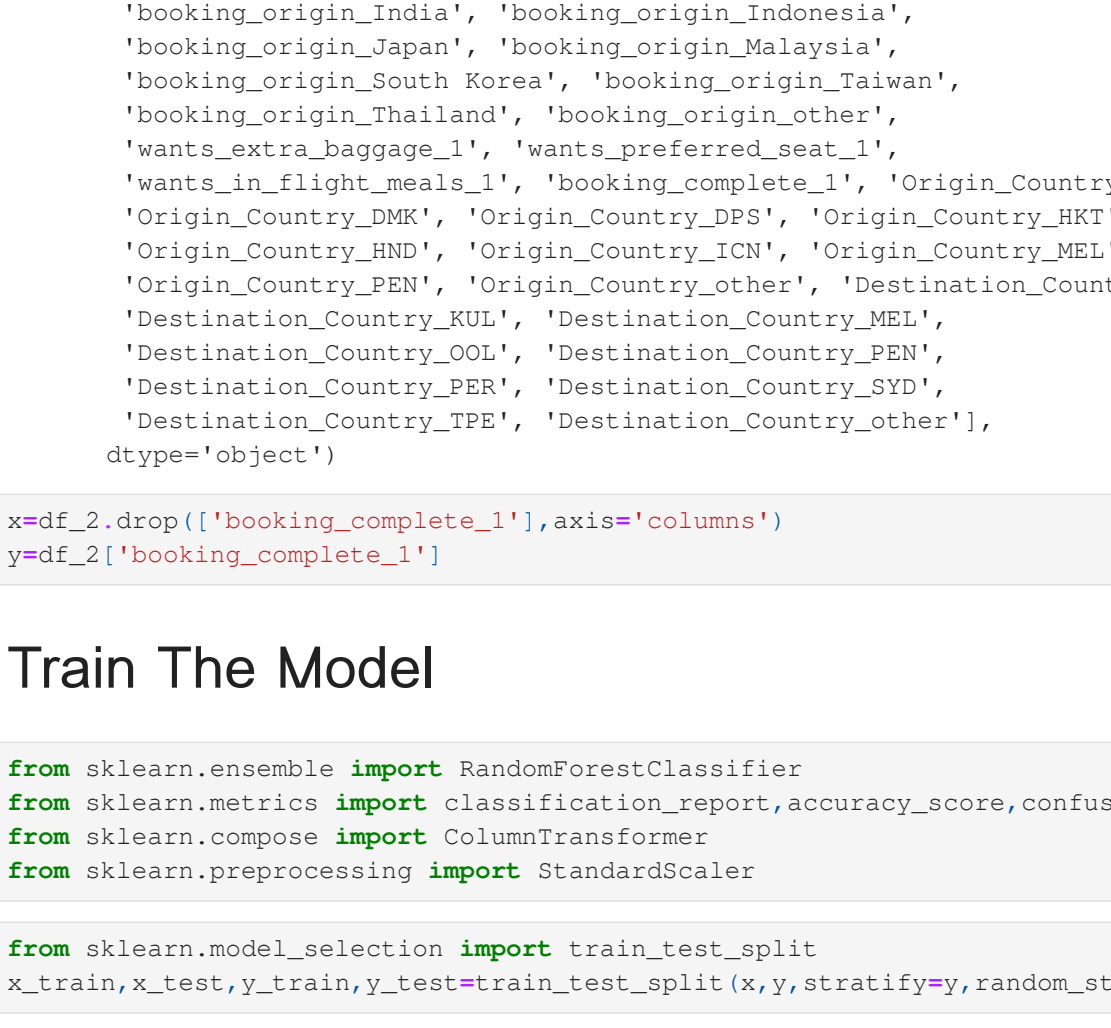
Out[162]: (49971, 17)

In [163]: len(df['booking_origin'].unique())

Out[163]: 104
```

```
In [191]: def categorical_feature(df, feature_list):
    for feature in feature_list:
        df[feature].value_counts().head(10).plot(kind='bar')
        plt.title(f'Top 10 {feature}')
        plt.xlabel(feature)
        plt.ylabel('Count')
        plt.xticks(rotation=45)
        plt.show()
```

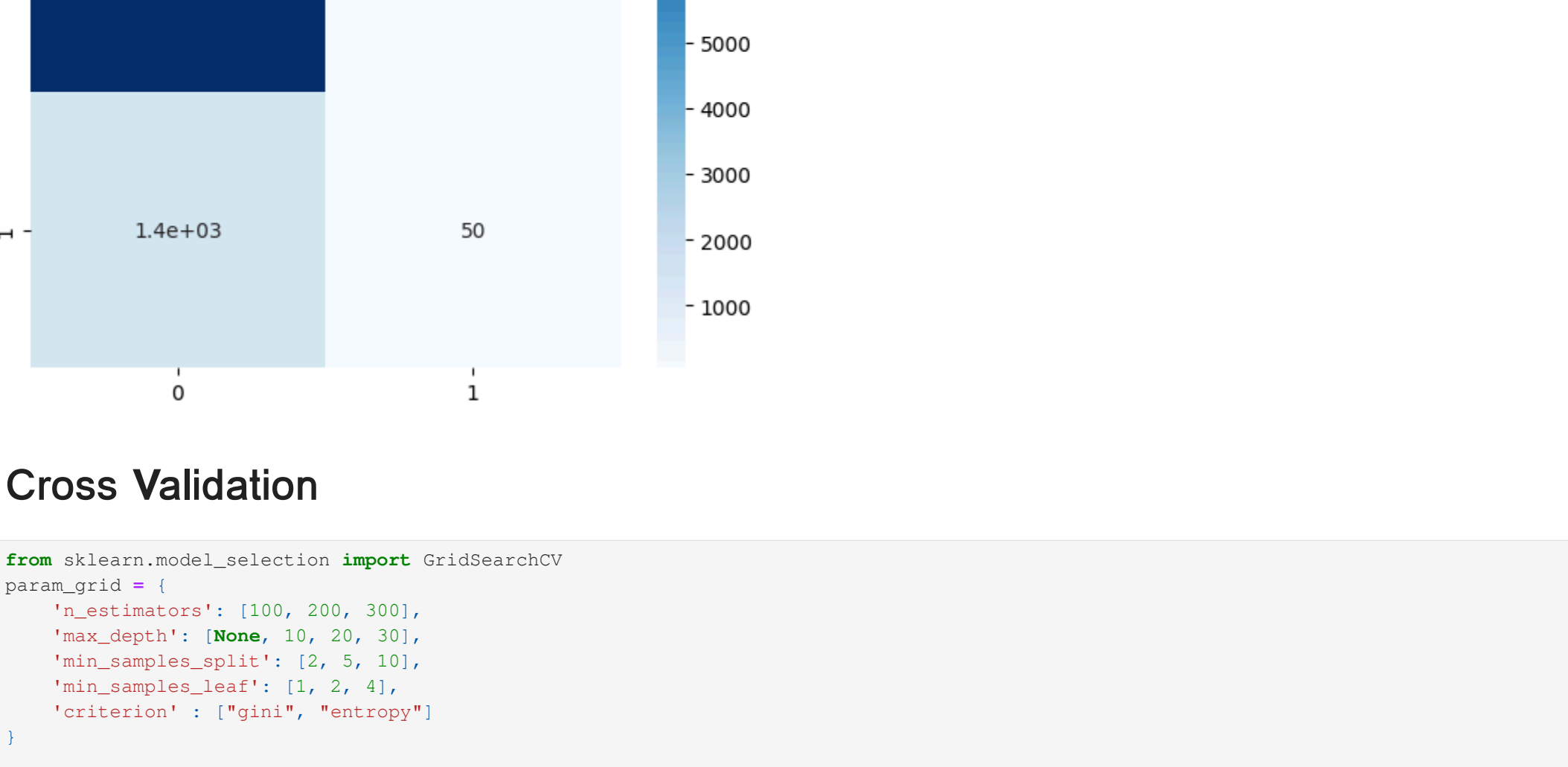
```
In [194]: categorical_feature(df, ['booking_origin', 'Origin_Country', 'Destination_Country'])
```



```
In [107]: df_correlation=df[['booking_complete', 'num_passengers', 'purchase_lead', 'length_of_stay',
                          'flight_hour', 'flight_duration', 'wants_extra_baggage', 'wants_in_flight_meals',
                          'wants_preferred_seat']]

In [111]: sns.heatmap(df_correlation.corr(),annot=True,cmap='Reds')
plt.title('Correlation')

Out[111]: Text(0.5, 1.0, 'Correlation')
```



```
In [113]: df_idx=df[['trip_type', 'flight_hour', 'sales_channel', 'Day Timing', 'route'],axis='columns']
df_idx.head()

Out[113]:
```

	num_passengers	purchase_lead	length_of_stay	flight_day	booking_origin	wants_extra_baggage	wants_preferred_seat	wants_in_flight_meals	flight_duration
0	2	262	19	6	New Zealand	0	0	0	5.52
1	1	112	20	3	AKLDEL	1	0	0	5.52
2	2	243	22	3	India	1	1	0	5.52
3	1	96	31	6	New Zealand	0	0	0	5.52
4	2	68	22	3	India	1	0	1	5.52

```
In [117]: def conversion(df, feature_list):
    for feature in feature_list:
        top_features=df[feature].value_counts().head(10).index
        df[feature]=df[feature].apply(lambda x: x if x in top_features else 'other')
    return df

In [123]: df_2=conversion(df_1, ['booking_origin', 'Origin_Country', 'Destination_Country'])

In [125]: df_2=pd.get_dummies(data=df_2,drop_first=True)

In [130]: df_2.columns

Out[130]: Index(['num_passengers', 'purchase_lead', 'length_of_stay', 'flight_duration',
        'flight_day_2', 'flight_day_3', 'flight_day_4', 'flight_day_5',
        'flight_day_6', 'flight_day_7', 'flight_day_8', 'flight_day_9',
        'booking_origin_India', 'booking_origin_Indonesia',
        'booking_origin_Japan', 'booking_origin_Malaysia',
        'booking_origin_South Korea', 'booking_origin_Taiwan',
        'booking_origin_Thailand', 'booking_origin_other',
        'wants_extra_baggage_1', 'wants_preferred_seat_1',
        'wants_in_flight_meals_1', 'booking_complete_1', 'Origin_Country_DMK',
        'Origin_Country_HKT', 'Origin_Country_JPN', 'Origin_Country_MEL',
        'Origin_Country_PEN', 'Origin_Country_PUL', 'Origin_Country_RMX',
        'Destination_Country_DMK', 'Destination_Country_HKT',
        'Destination_Country_JPN', 'Destination_Country_MEL',
        'Destination_Country_PEN', 'Destination_Country_PUL',
        'Destination_Country_RMX', 'Destination_Country_SIN',
        'Destination_Country_OTHER'],
        dtype='object')
```

Train The Model

```
In [184]: from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report, accuracy_score, confusion_matrix
from sklearn.compose import ColumnTransformer
from sklearn.preprocessing import StandardScaler

In [185]: from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test=train_test_split(x,y,stratify=y,random_state=0,test_size=0.2)

In [186]: preprocessor = ColumnTransformer(
    transformers=[
        ('num', StandardScaler(), x.columns),
    ])

In [187]: from sklearn.pipeline import Pipeline
pipe=Pipeline([
    ('preprocessor', preprocessor),
    ('model', RandomForestClassifier(n_estimators=100, class_weight='balanced'))
])

In [188]: pipe.fit(x_train,y_train)
y_pred=pipe.predict(x_test)
print(f'The accuracy of this model is : {accuracy_score(y_test,y_pred)}\n')
print(f'The Classification report of this model is : {classification_report(y_test,y_pred)}\n')
```

The accuracy of this model is 10.8450225112556278

The Classification report of this model is :

precision recall f1-score support

0 0.86 0.97 0.91 8501

1 0.44 0.14 0.21 1494

accuracy 0.85 0.85 0.85 9995

macro avg 0.65 0.51 0.49 9995

weighted avg 0.79 0.85 0.79 9995

```
In [204]: from sklearn.metrics import confusion_matrix
sns.heatmap(confusion_matrix(y_test,y_pred),cmap='Blues')
plt.title('Confusion Matrix')

Out[204]: Text(0.5, 1.0, 'Confusion Matrix')
```

Confusion Matrix



Cross Validation

```
In [197]: from sklearn.model_selection import GridSearchCV
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth': [None, 10, 20, 30],
    'min_samples_split': [2, 5, 10],
    'min_samples_leaf': [1, 2, 4],
    'criterion': ['gini', 'entropy']}

R = RandomForestClassifier(random_state=42)
grid_search = GridSearchCV(estimator=R, param_grid=param_grid, cv=5, n_jobs=-1, verbose=0)
grid_search.fit(x_train,y_train)
best_params = grid_search.best_params_
print(f'Best Hyperparameters: (best_params)')
```

C:\Users\Qadri Laptops\PycharmProjects\numpy\venv\lib\site-packages\numpy\ma\core.py:2820: RuntimeWarning: Invalid value encountered in c

est_data = np.ma.array(data, dtype=ndarray.dtype, copy=False)

Best Hyperparameters: {'criterion': 'gini', 'max_depth': 20, 'min_samples_leaf': 4, 'min_samples_split': 2, 'n_estimators': 300}

```
In [221]: y_pred=grid_search.predict(x_test)
print(f'Best Hyperparameter accuracy is : {grid_search.best_score_}\n')
print(f'Classification report of hyperparameter is : {classification_report(y_test,y_pred)}\n')
print(f'Confusion matrix is : {confusion_matrix(x_test,y_pred)}\n')
```

Best Hyperparameter accuracy is 10.852396430639624

The Classification report of hyperparameter is :

precision recall f1-score support

0 0.85 0.99 0.92 8501

1 0.45 0.03 0.06 1494

accuracy 0.85 0.85 0.85 9995

macro avg 0.65 0.51 0.49 9995

weighted avg 0.79 0.85 0.79 9995

```
In [204]: Confusion matrix is :
[[8440 61]
 [1446 50]]
```

Top Features Using Training

```
In [202]: feature_names = pipe.named_steps['model'].get_feature_names_out()
importance = pipe.named_steps['model'].feature_importances_
importances_df = pd.DataFrame({
    'feature': feature_names,
    'importance': importance
}).sort_values(by='importance', ascending=False)
importance_df.head(10)
```

Feature Importance

1 num_purchase_lead 0.244491

2 num_length_of_stay 0.174109

0 num_num_passengers 0.055143

14 num_booking_origin_Malaysia 0.063784

3 num_flight_duration 0.045665

21 num_wants_in_flight_meals_1 0.030880

2 num_wants_preferred_seat_1 0.027069

19 num_wants_extra