# Hotel Booking Cancellation Prediction

by Ruba Alnashwan

# METHODOLOGY

- Proplem understanding
- Data collection
- Data cleaning
- Exploratory data analysis (EDA)
- Feature Engineering and selection
- Data modeling

# Proplem understanding

Overview:
In this project, we will use data from the kaggle website, which provides information hotel and the label (cancel or not) . Our goal from this project is to build classification models that predict if the customer will cancel the booking or not.


Scope:
observation represents a hotel booking between the 1st of July 2015 and 31st of August 2017, including booking that effectively arrived and booking that were canceled. The dataset contains 119390 rows and 10 columns.

# Data cleaning

check nulls

I filled in the null values in the country ,children features with mode and mean.
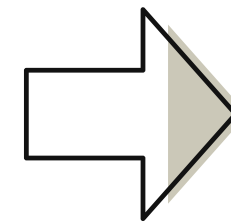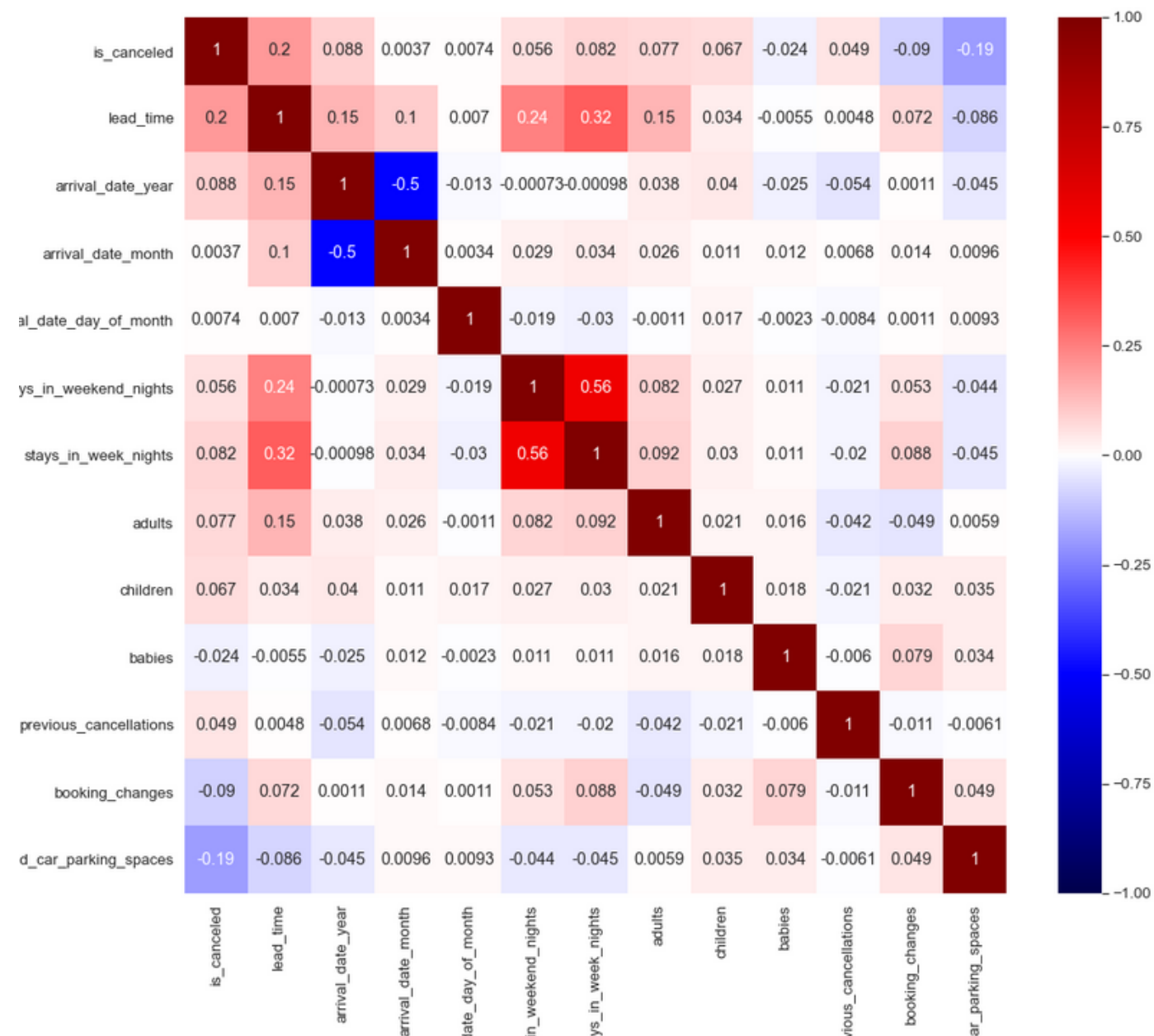
check duplicates and drop

check outliers using EDA and remove

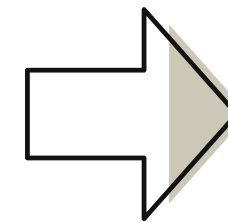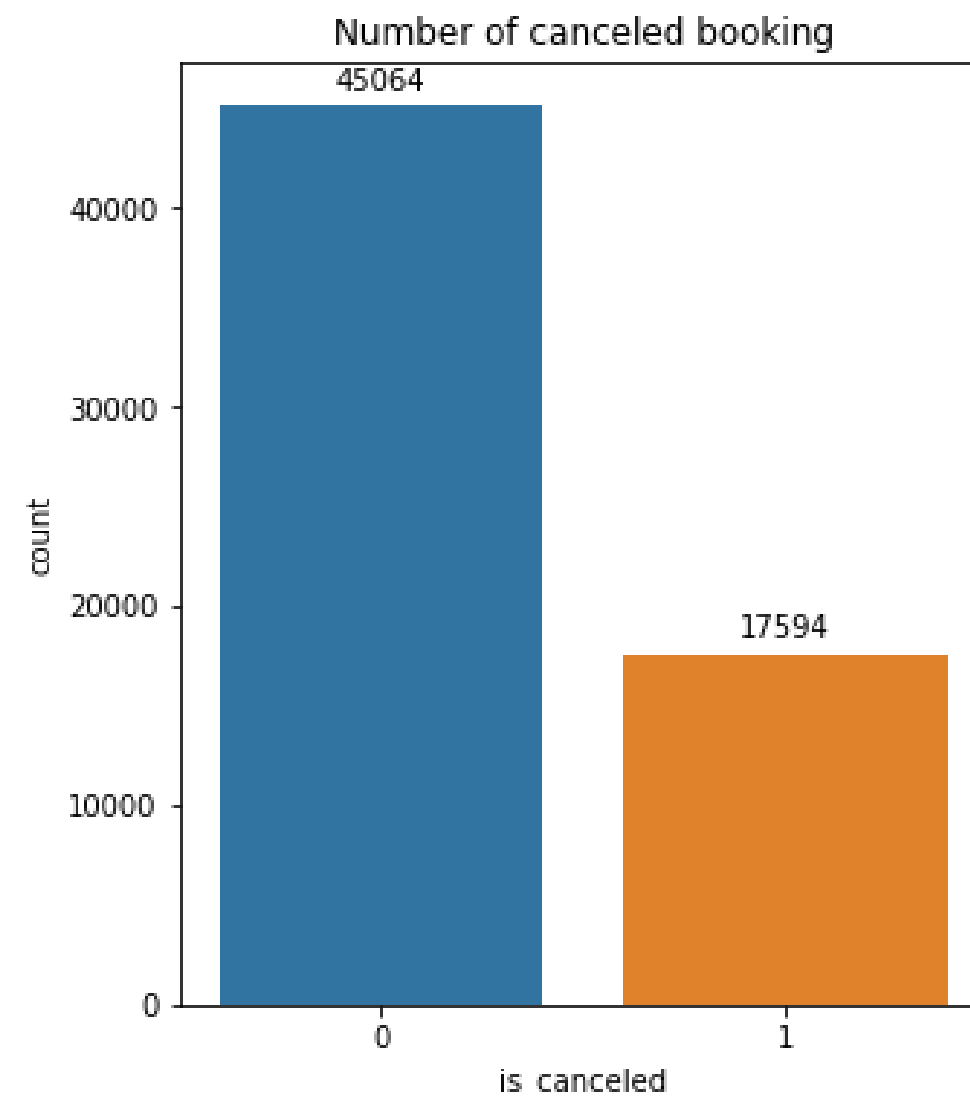replace some character to numbers like month

# Exploratory data analysis (EDA)



Observations:

This heat map shows that there is no strong relationship between the dependent variable and features .
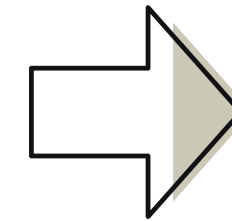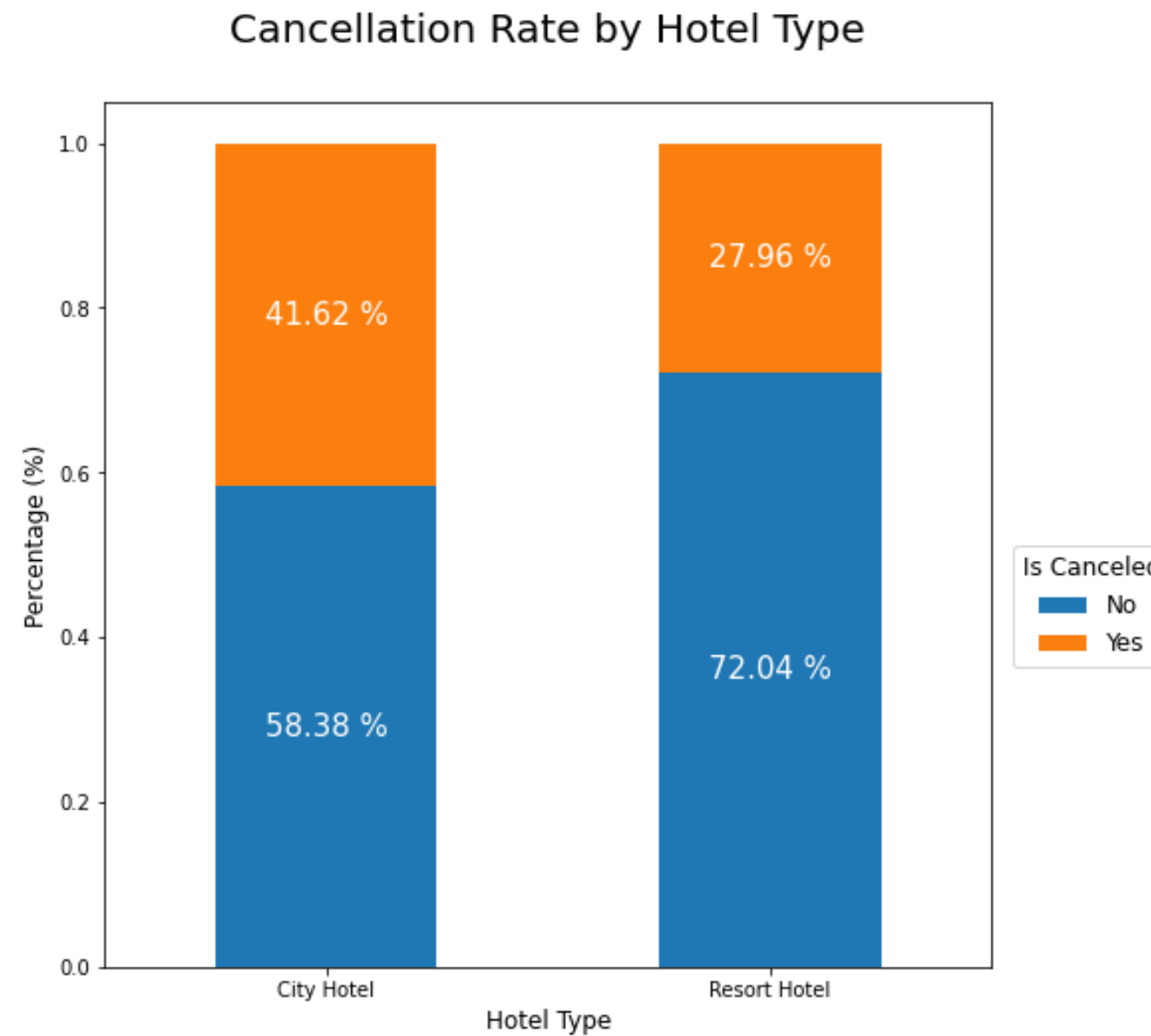
# Exploratory data analysis (EDA)



Observations:

The par plot shows the Number of canceled booking

# Exploratory data analysis (EDA)



Cancellation Rate by Hotel Type

Observations:

The par plot shows the Cancellation Rate by Hotel Type

# Feature Engineering and selection

Combine several columns into one column

Dummies varible for categorical Feature

**Scaling:** StandardScale

# Data modeling

MAIN METRIC USED FOR EVALUATION : F1

SECONDARY METRICS: PRECISION, RECALL AND ACCURACY

BASELINE MODEL USED: LOGISTIC REGRESSION

MODEL USED: RANDOM FOREST.DECISION TREE

# Data modeling

**basline**

TRAIN SCORES:
 F1 SCORE IS 0.457
 PRECISION SCORE IS 0.669
 RECALL SCORE IS 0.346

-----------------------------------------

VALIDATION SCORES:
 F1 SCORE IS 0.459
 PRECISION SCORE IS 0.665
 RECALL SCORE IS 0.350

**dummies**

TRAIN SCORES:
 F1 SCORE IS 0.573
 PRECISION SCORE IS 0.901
 RECALL SCORE IS 0.421

-----------------------------------------

VALIDATION SCORES:
 F1 SCORE IS 0.575
 PRECISION SCORE IS 0.903
 RECALL SCORE IS 0.422

**Remove outliers**

TRAIN SCORES:
 F1 SCORE IS 0.582
 PRECISION SCORE IS 0.880
 RECALL SCORE IS 0.4351

-----------------------------------------

VALIDATION SCORES:
 F1 SCORE IS 0.583
 PRECISION SCORE IS 0.859
 RECALL SCORE IS 0.441

# Data modeling

**RandomOverSampler**

TRAIN SCORES:
 F1 SCORE IS 0.674
 PRECISION SCORE IS 0.777
 RECALL SCORE IS 0.594
--------------------------------------

VALIDATION SCORES:
 F1 SCORE IS 0.625
 PRECISION SCORE IS 0.656
 RECALL SCORE IS 0.596

**TomekLinks**

TRAIN SCORES:
 F1 SCORE IS 0.597
 PRECISION SCORE IS 0.862
 RECALL SCORE IS 0.456
--------------------------------------

VALIDATION SCORES:
THE F1 SCORE IS 0.593
THE PRECISION SCORE IS 0.828
THE RECALL SCORE IS 0.4621

**SMOTE**

TRAIN SCORES:
 F1 SCORE IS 0.682
 PRECISION SCORE IS 0.763
 RECALL SCORE IS 0.616
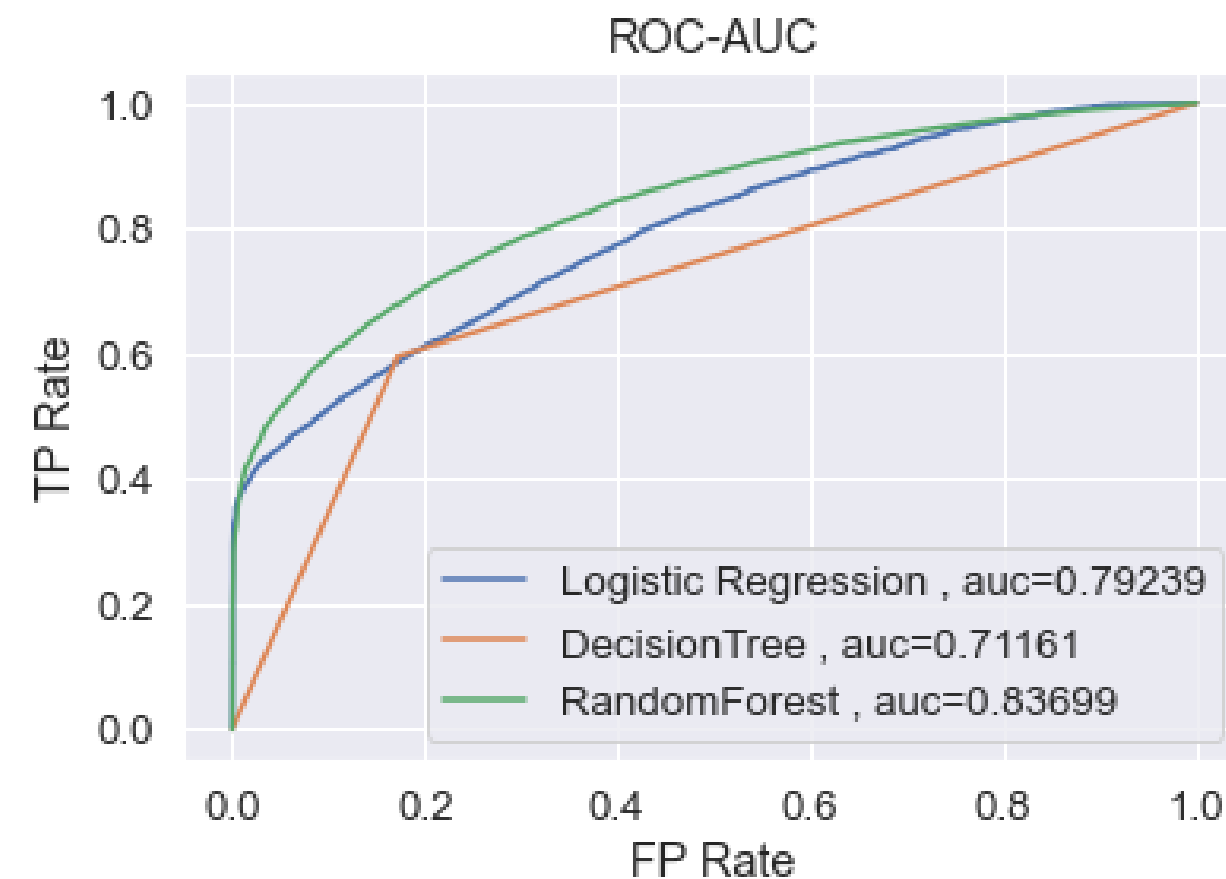--------------------------------------

VALIDATION SCORES:
 F1 SCORE IS 0.627
 PRECISION SCORE IS 0.640
 RECALL SCORE IS 0.615

# CONCLUSION

**Result the best model**



ROC-AUC

Logistic Regression , auc=0.79239
DecisionTree , auc=0.71161
RandomForest , auc=0.83699

**Difficulties:**

No strong relationship between the dependent variable and features
The data volume is large, so the training time is long

# Thank you

Any questions?