

RESTAURANT'S VISITOR FORECASTING

IBM ADVANCED DATA SCIENCE CAPSTONE PROJECT



Introduction

ETL

EDA

Model

Results

VISITOR FORECASTING

- What is the problem?
- Why is it an important problem?
- So, what is the solution?
- What is needed for solution?
- How is the solution achieved?



Introduction

ETL

EDA

Model

Results

EXTRACT - TRANSFORM - LOAD



Introduction

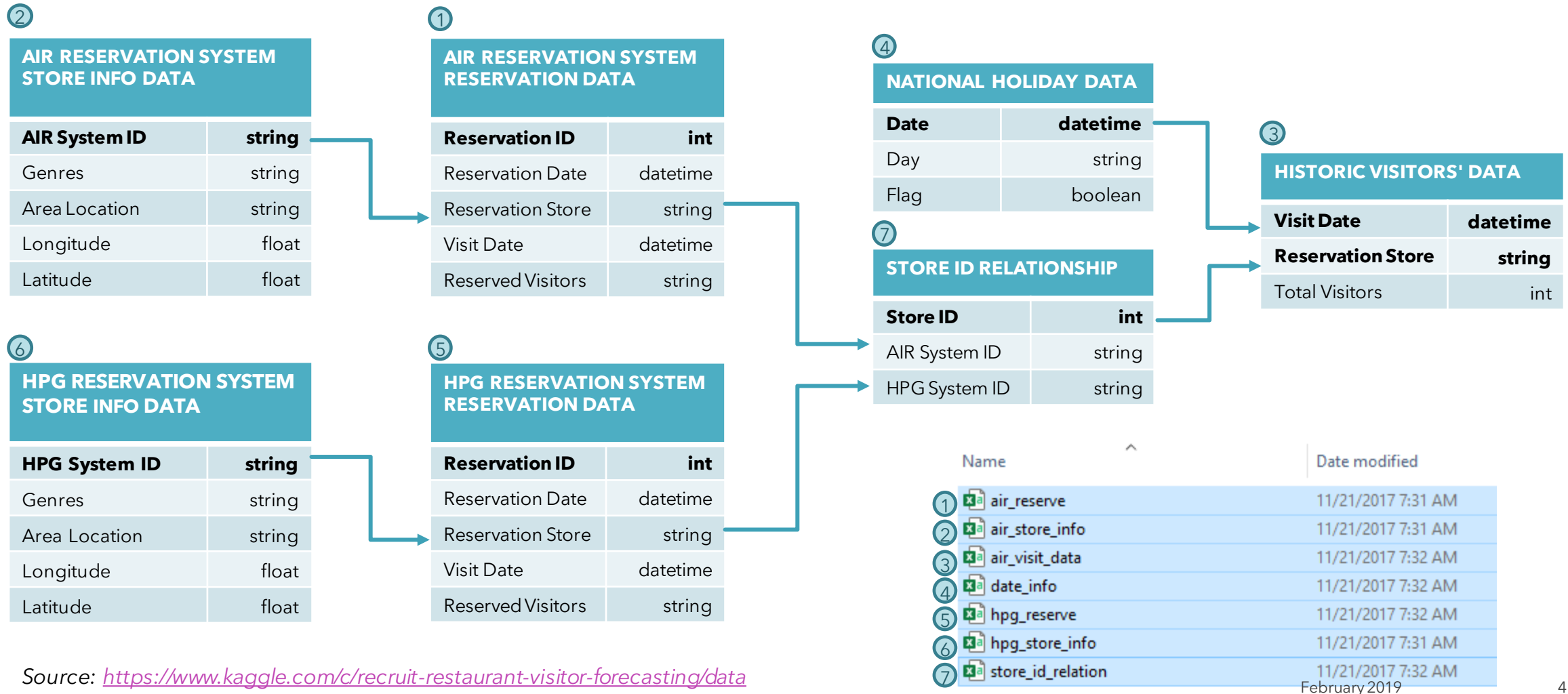
ETL

EDA

Model

Results

ETL - DATA SOURCE



Source: <https://www.kaggle.com/c/recruit-restaurant-visitor-forecasting/data>

Introduction

ETL

EDA

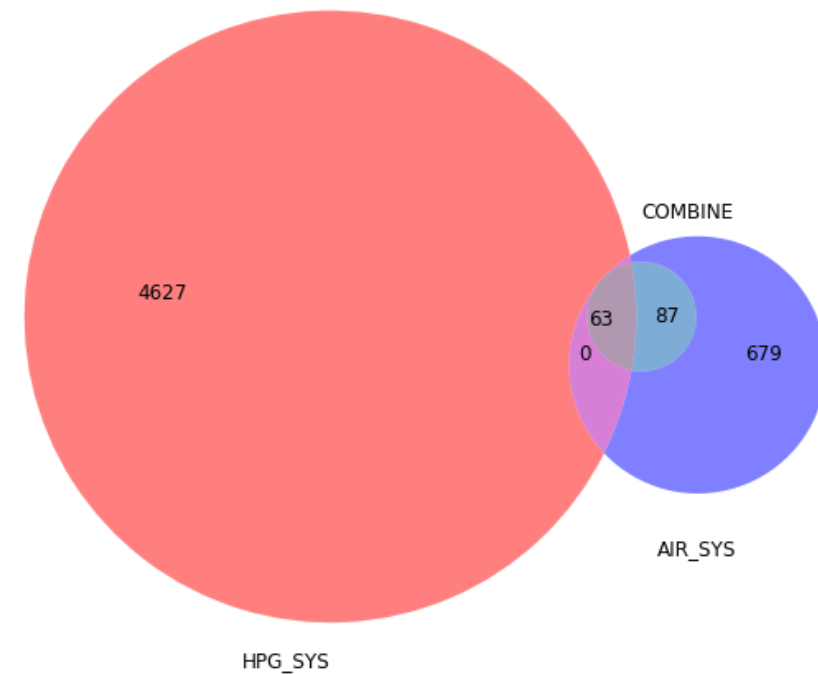
Model

Results

RESTAURANTS BY SYSTEMS

-	AIR_SYS	HPG_SYS
Unique	679	4627
Combine	150	63
Explicit	0	87

System-wise distribution of hotels



Introduction

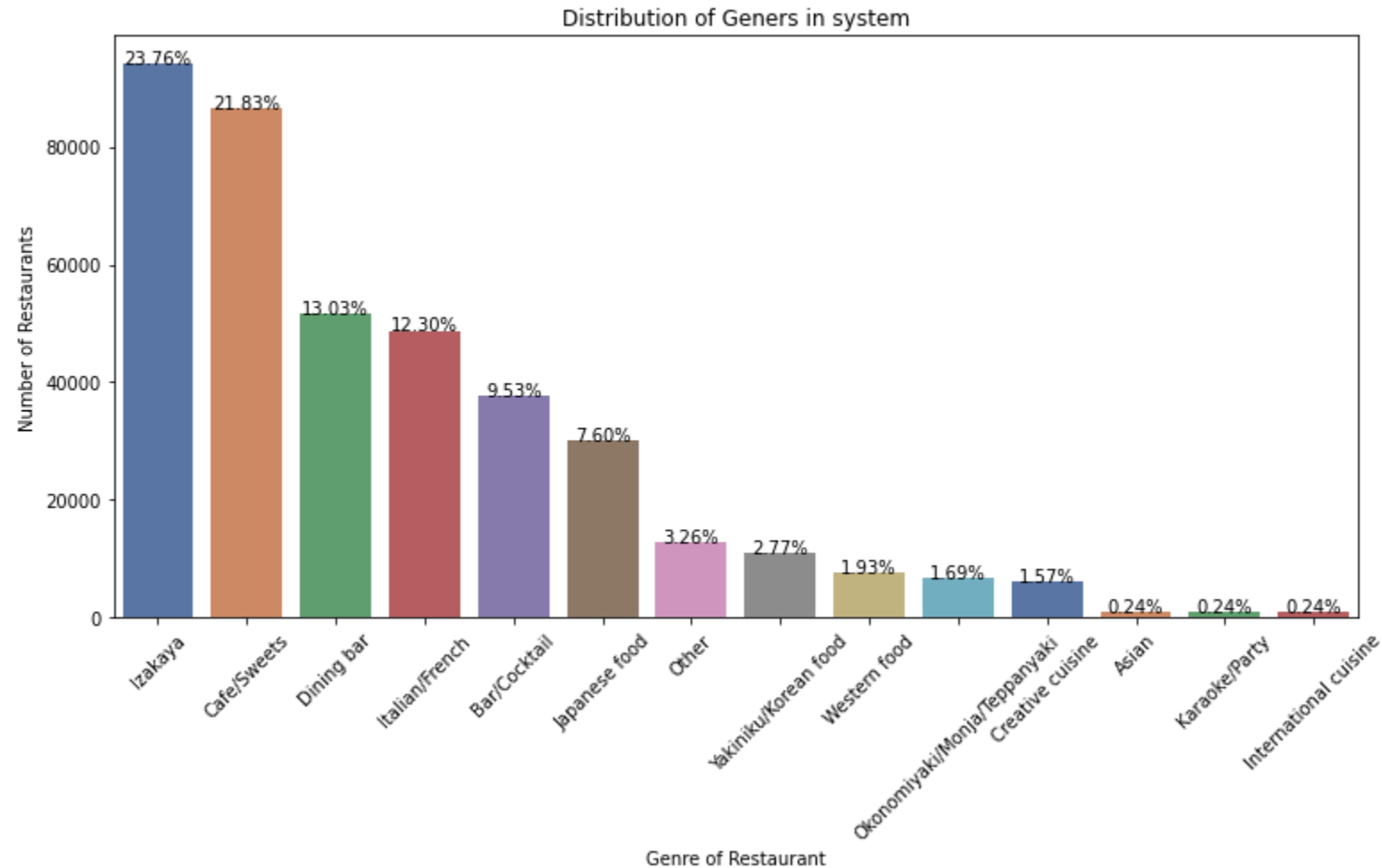
ETL

EDA

Model

Results

RESTAURANTS BY GENRES



Introduction

ETL

EDA

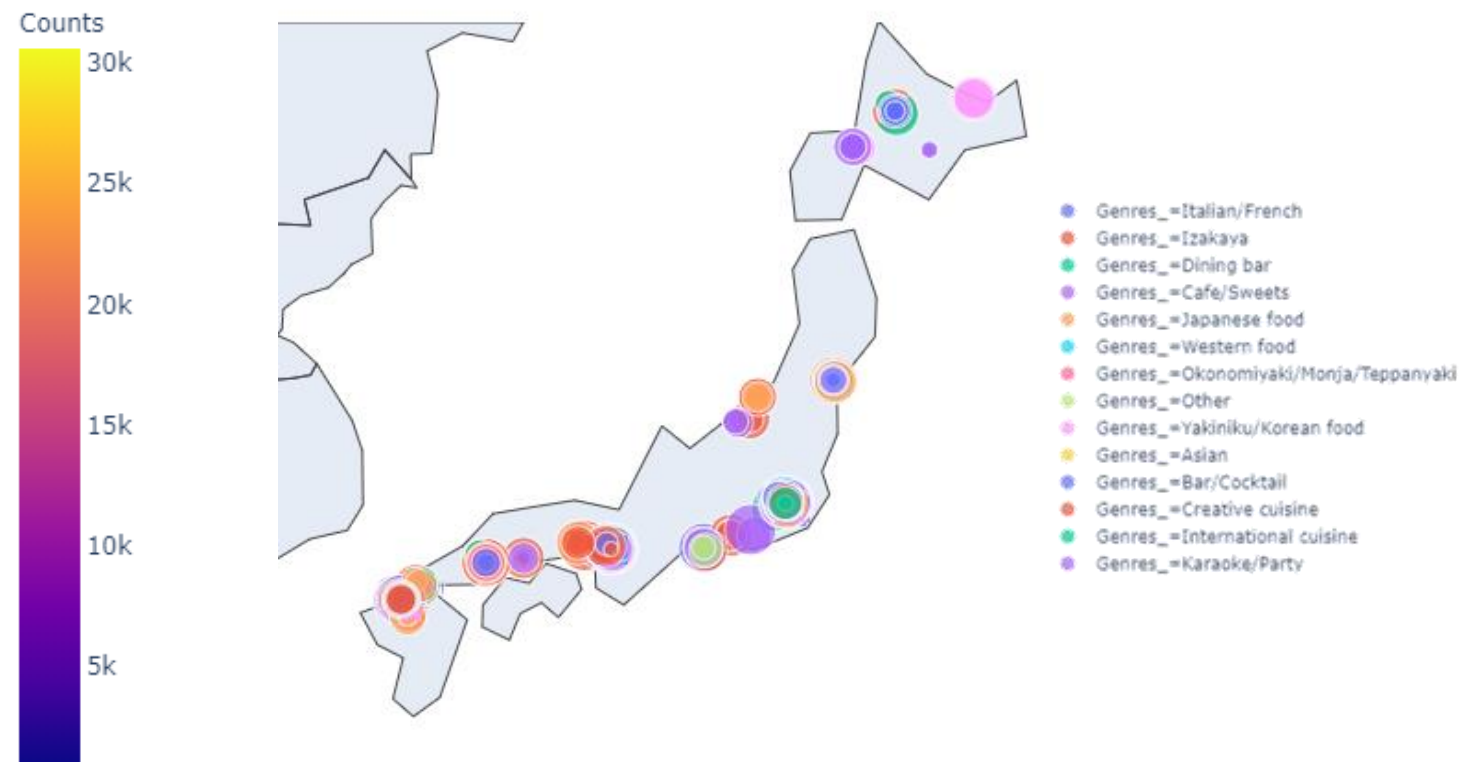
Model

Results

RESTAURANTS AND VISITORS BY LOCATION



Restaurants by Area



Visitors by Area

Introduction

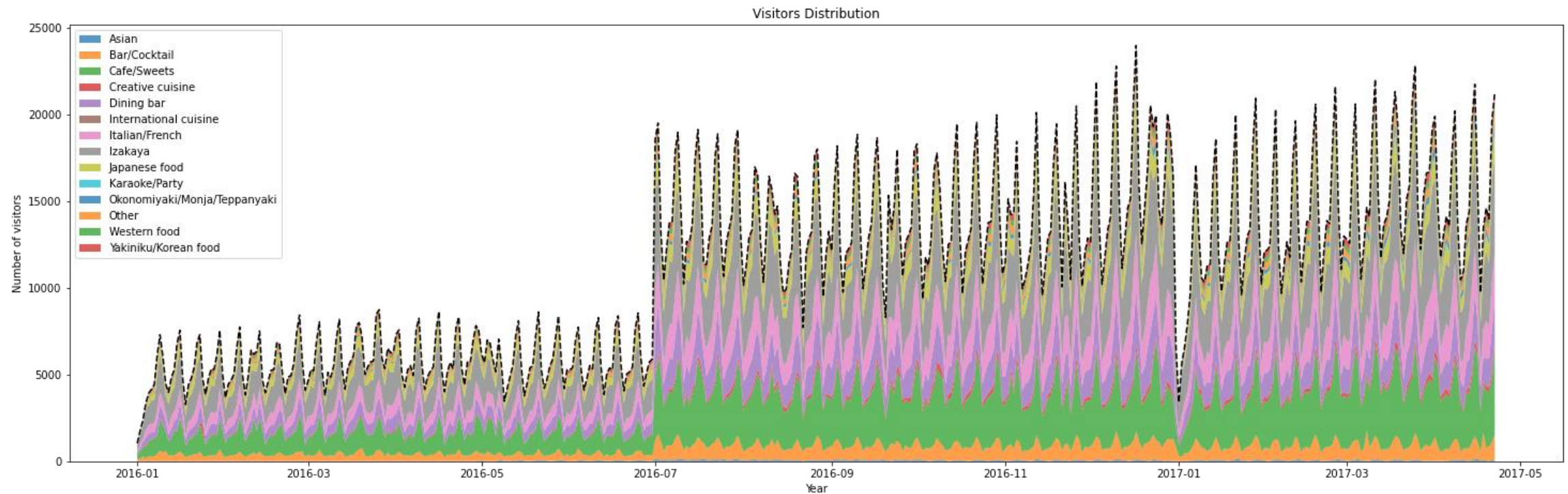
ETL

EDA

Model

Results

VISITORS TIMESERIES PLOT



Introduction

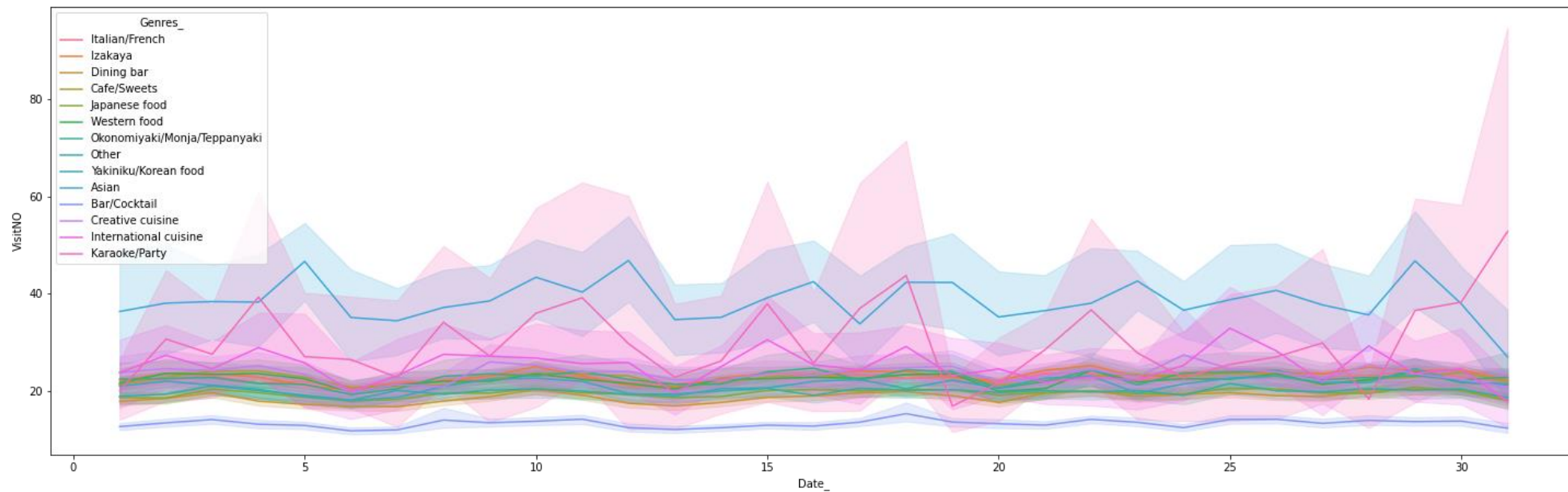
ETL

EDA

Model

Results

TRENDS IN VISITORS - OVER MONTH



Introduction

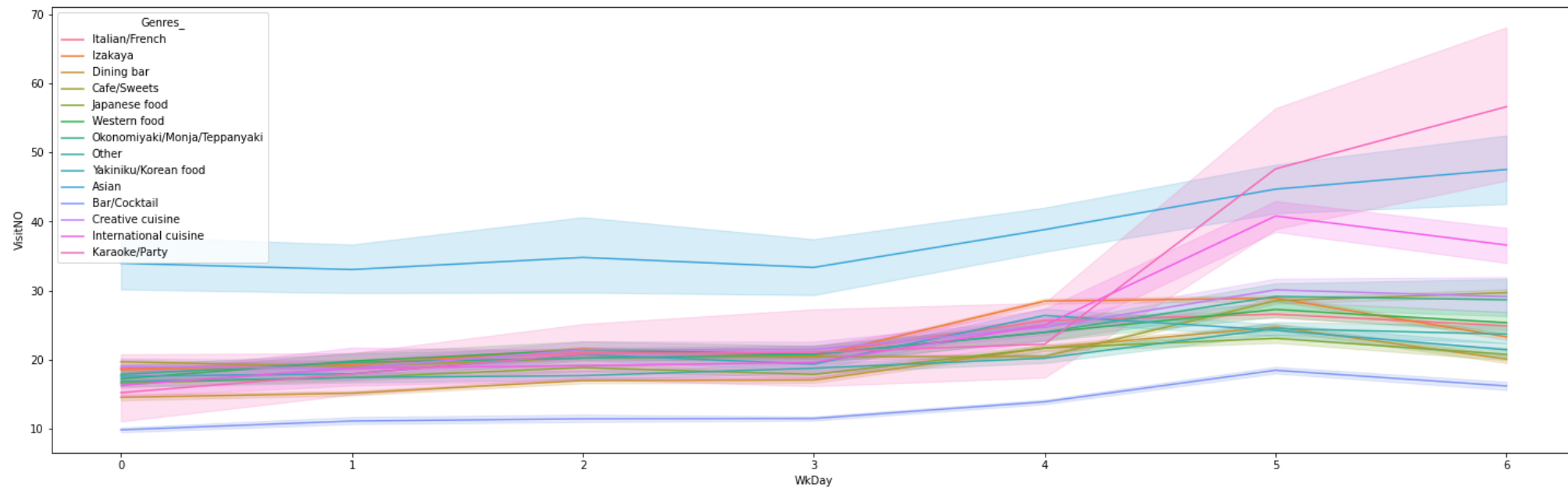
ETL

EDA

Model

Results

TRENDS IN VISITORS - OVER WEEK



Introduction

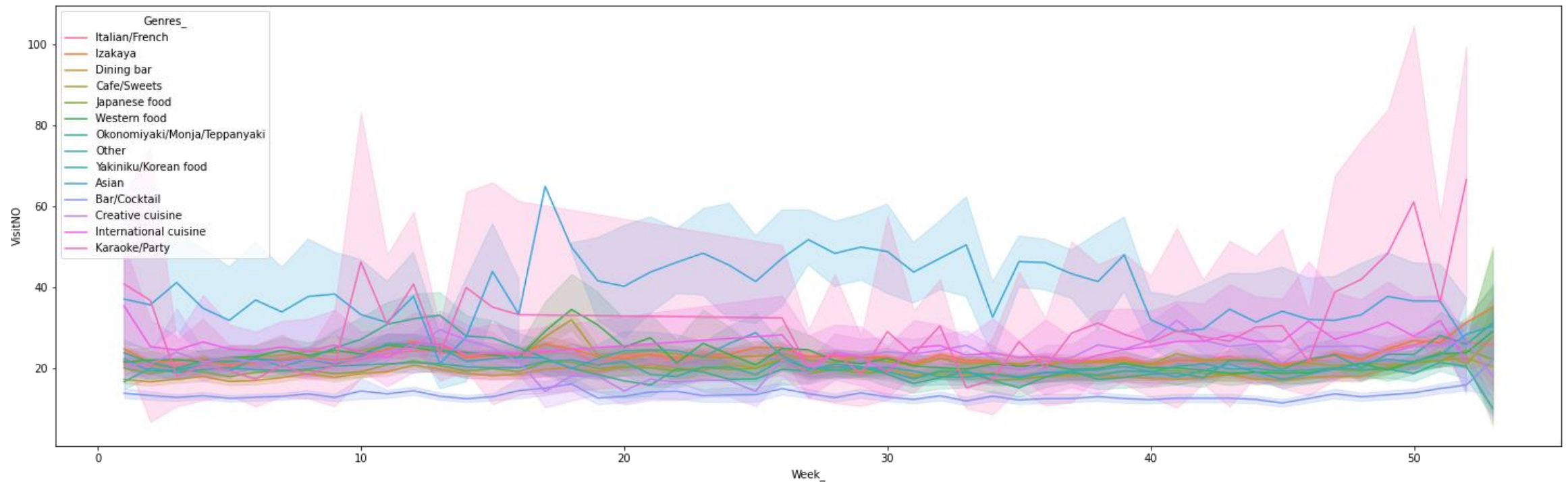
ETL

EDA

Model

Results

TRENDS IN VISITORS - OVER YEAR



Introduction

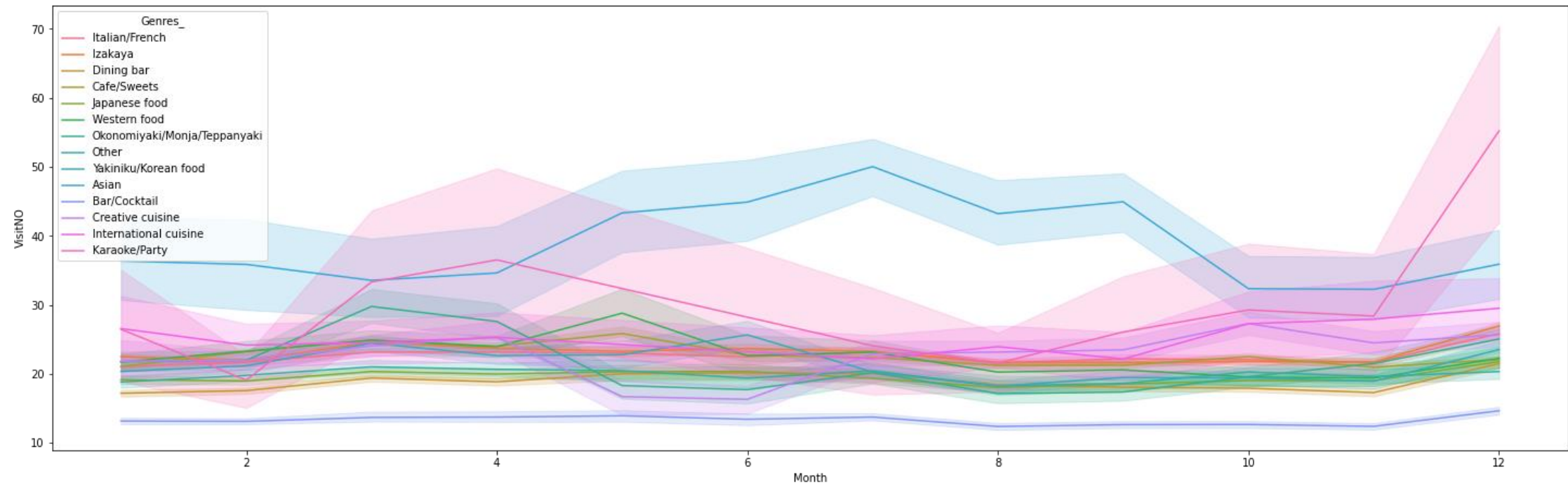
ETL

EDA

Model

Results

TRENDS IN VISITORS - OVER YEAR



Introduction

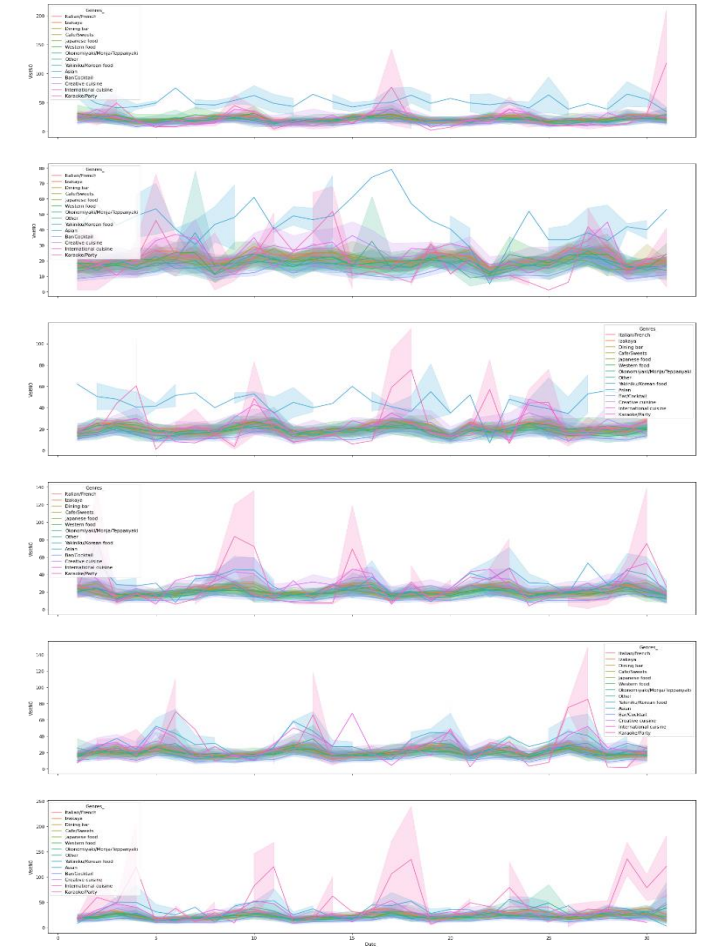
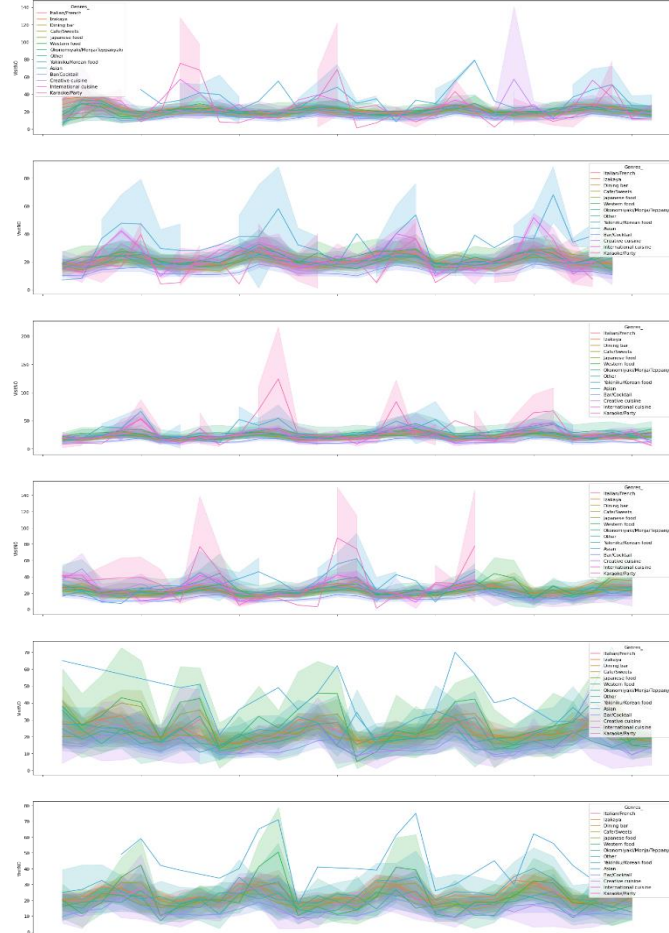
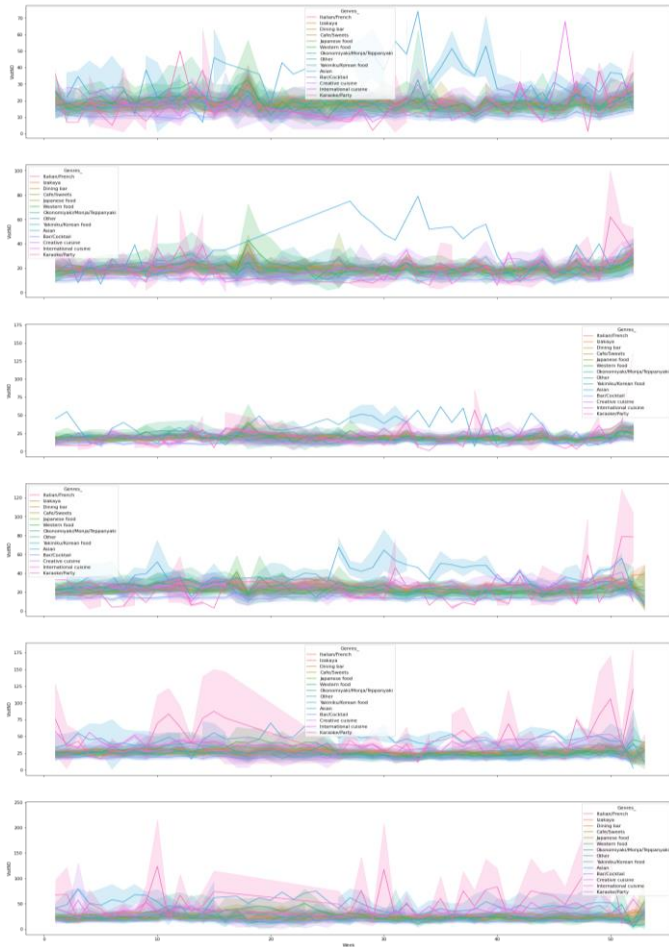
ETL

EDA

Model

Results

VISITORS VISITING PATTERN ACROSS YEAR



February 2019

13

Introduction

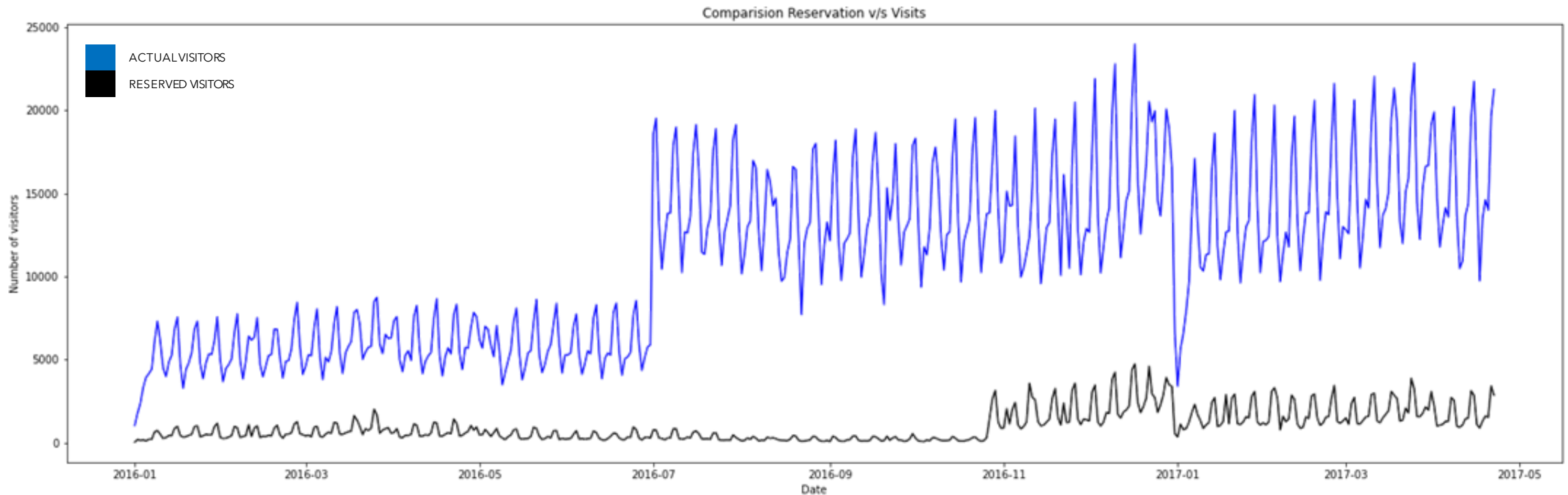
ETL

EDA

Model

Results

VISITORS - ACTUAL V/S RESERVED



Introduction

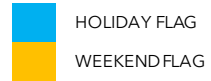
ETL

EDA

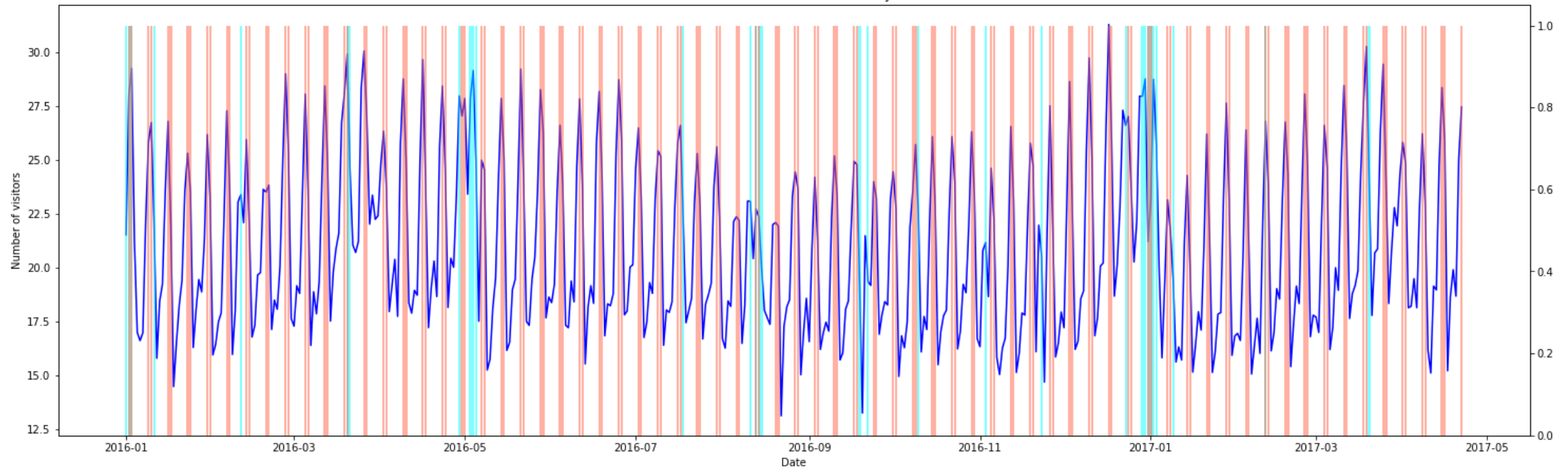
Model

Results

VISITORS - ON HOLIDAYS & WEEKENDS



Visitors on Weekend and Holidays



Introduction

ETL

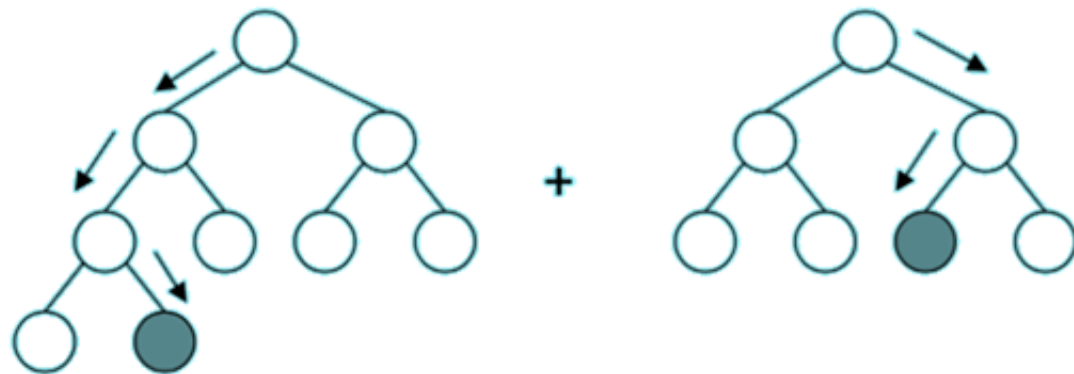
EDA

Model

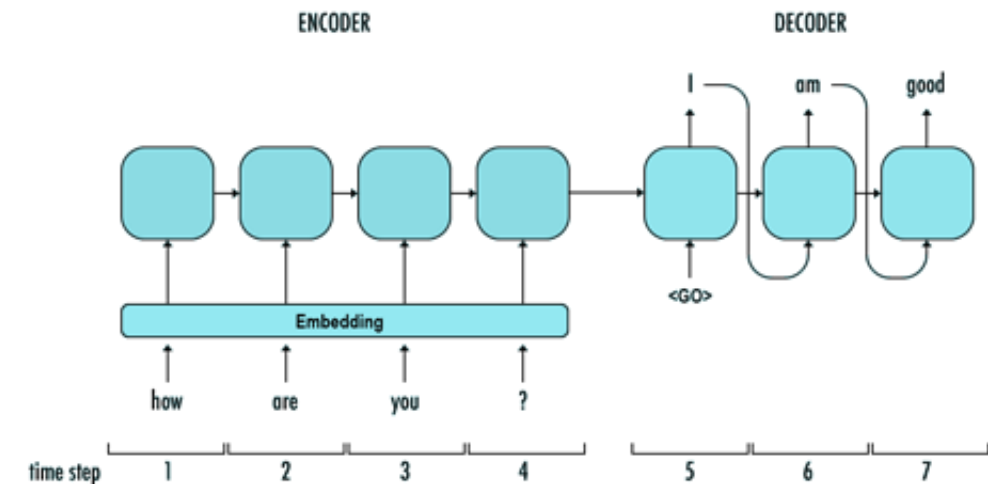
Results

MODEL DEFINITION

- CLASSICAL ML MODEL
- Gradient Booster Regressor



- DEEP LEARNING MODEL
- Seq2Seq LSTM Encode-Decoder



Introduction

ETL

EDA

Model

Results

MODEL EVALUATION

- Data split 80:10:10 train-test-validation split
- Root mean squared logarithmic error as metric

$$RMSLE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\log(1 + p) - \log(1 + a))^2}$$

n is the total number of observations | p_i is your prediction of target | a_i is the actual target for i .

- Robustness to the effect of the outliers
- Measurement of relative error
- Biased penalty for overestimation

Introduction

ETL

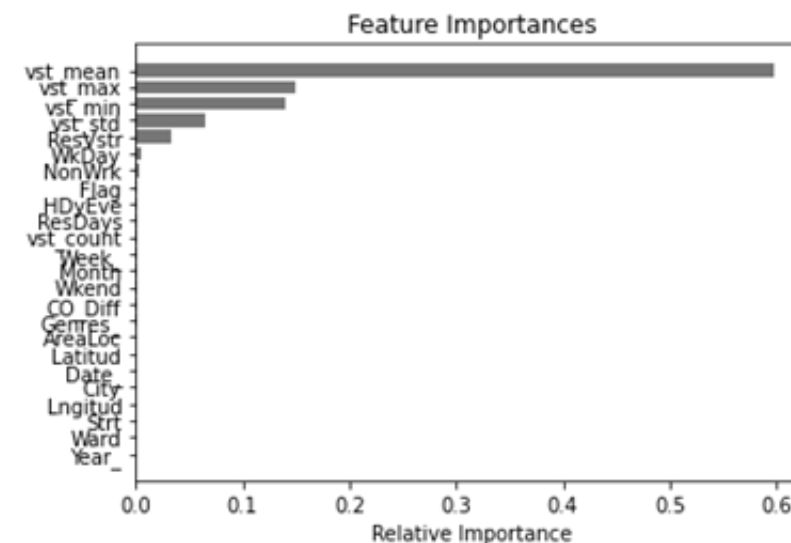
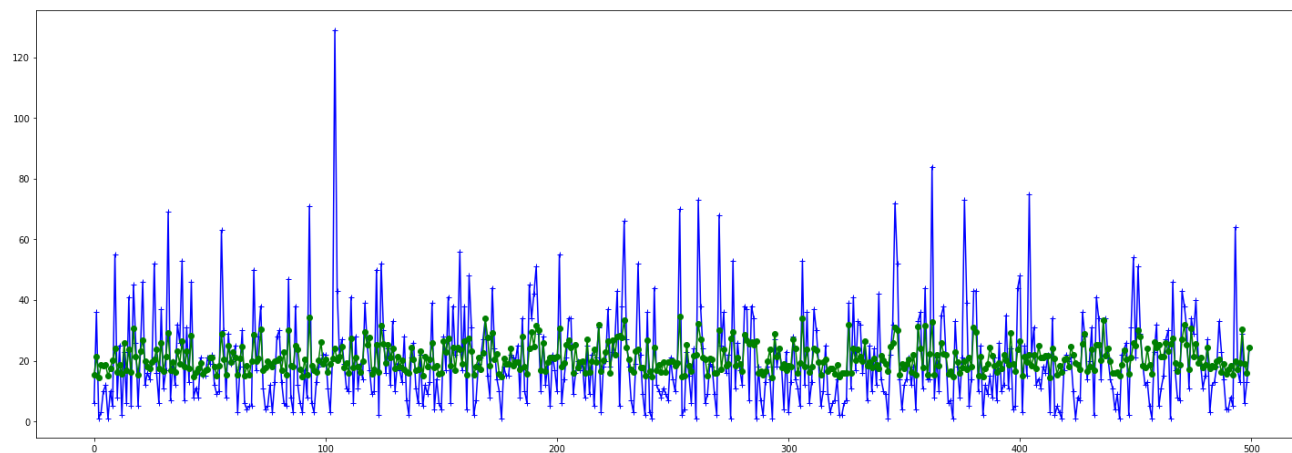
EDA

Model

Results

GRADIENT BOOST MODEL

- K-fold cross-validation and training
- Prediction is averaged over 5 folds
- Hyper Parameter tuning performed



Performance of base model:
0.71

Performance of fine model:
0.72

Introduction

ETL

EDA

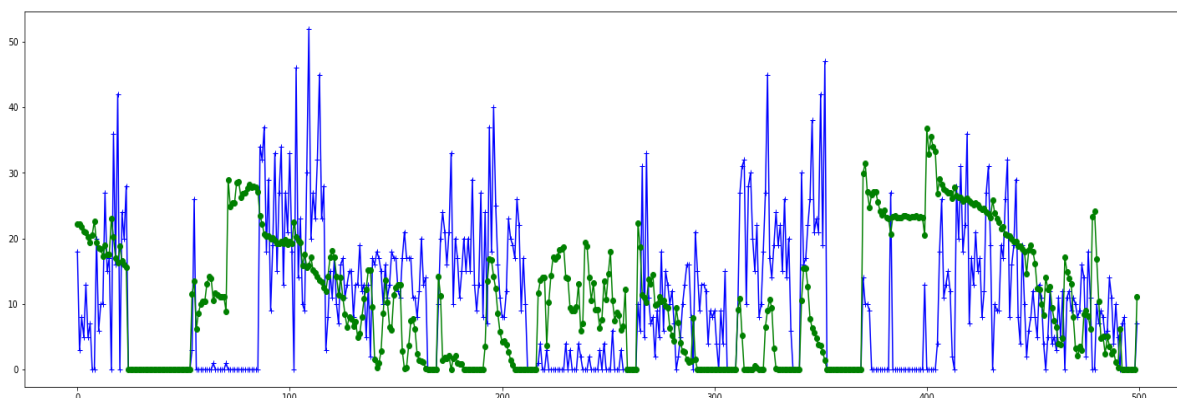
Model

Results

LSTM ENCODE-DECODER

- One layer of encoder
- Two layers of decoder units
- Two iterations with different features

Layer (type)	Output Shape	Param #	Connected to
input_5 (InputLayer)	[(None, None, 32)]	0	[]
input_6 (InputLayer)	[(None, None, 32)]	0	[]
lstm_3 (LSTM)	[(None, 64), (None, 64), (None, 64)]	24832	['input_5[0][0]']
lstm_4 (LSTM)	(None, None, 64)	24832	['input_6[0][0]', 'lstm_3[0][1]', 'lstm_3[0][2]']
lstm_5 (LSTM)	[(None, None, 64), (None, 64), (None, 64)]	33024	['lstm_4[0][0]']
time_distributed_1 (TimeDistributed)	(None, None, 1)	65	['lstm_5[0][0]']
Total params: 82,753 Trainable params: 82,753 Non-trainable params: 0			



Performance of first model:
1.61

Performance of second model:
2.09

Introduction

ETL

EDA

Model

Results

SUMMARY

- GBM works better
- Further tasks: Tuning LSTM for better performance
 - Activation Function
 - Number of layers
 - Number of hidden units in each layer
 - Optimizer
- Links below:
 - Architectural decision document :
[Recruite Restaurants Visitors Forecasting ADD Document.pdf](#)
 - Entity relationship diagram:
[Database Documentation.pdf](#)
 - Jupyter Notebook:
[IBM_Capstone.ipynb](#)

Algorithm	Variation	RSMLE	Visual
Gradient Boost	Before tuning	0.7174	--
	After tuning	0.7204	OK
Encoder-Decoder	With 3 prev days	1.6358	--
	With 7 prev days	2.0922	--