# BRIEF SELF INTRO

POSTGRADUATE DIPLOMA IN EDUCATION

Major (Math Education), HKBU

ACCELERATE

Data Science and Machine Learning

Cohort 2

**2011**

**2018**

**2013**

**NOW**

CUHK BSC MATHEMATICS

Double Streams in Enrichment Math
and Applied Math

GM TEACHER

EMI secondary school
HKDSE Math, M2, Physics

# OBJECTIVES

## Four Objectives

### Insight
Looking into real data of catering industry

### Analysis
possible features that affect the number of visitors

### Prediction
Precise prediction of restaurant visitors of all kinds

### Application
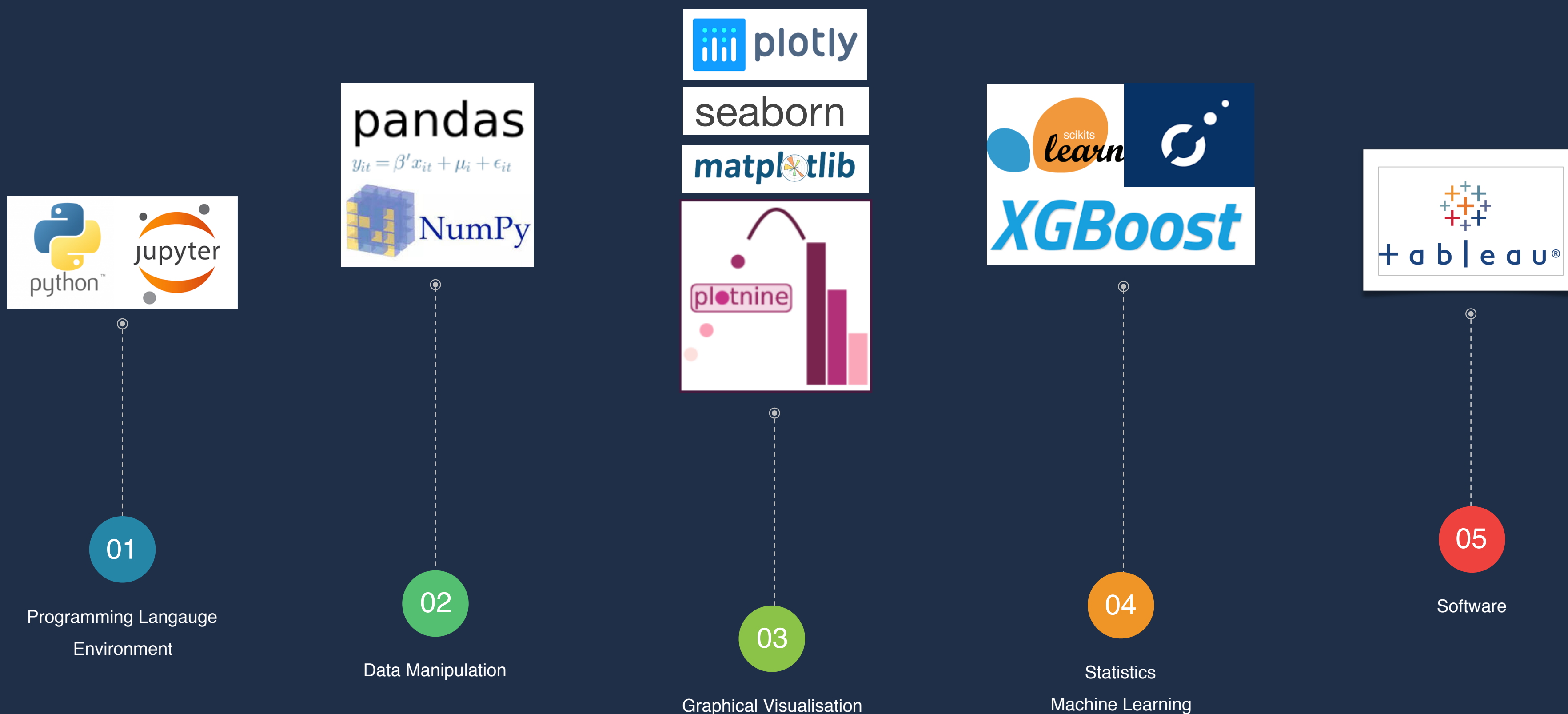Apply similar strategies in HK restaurants

# Data Source

**DATA**

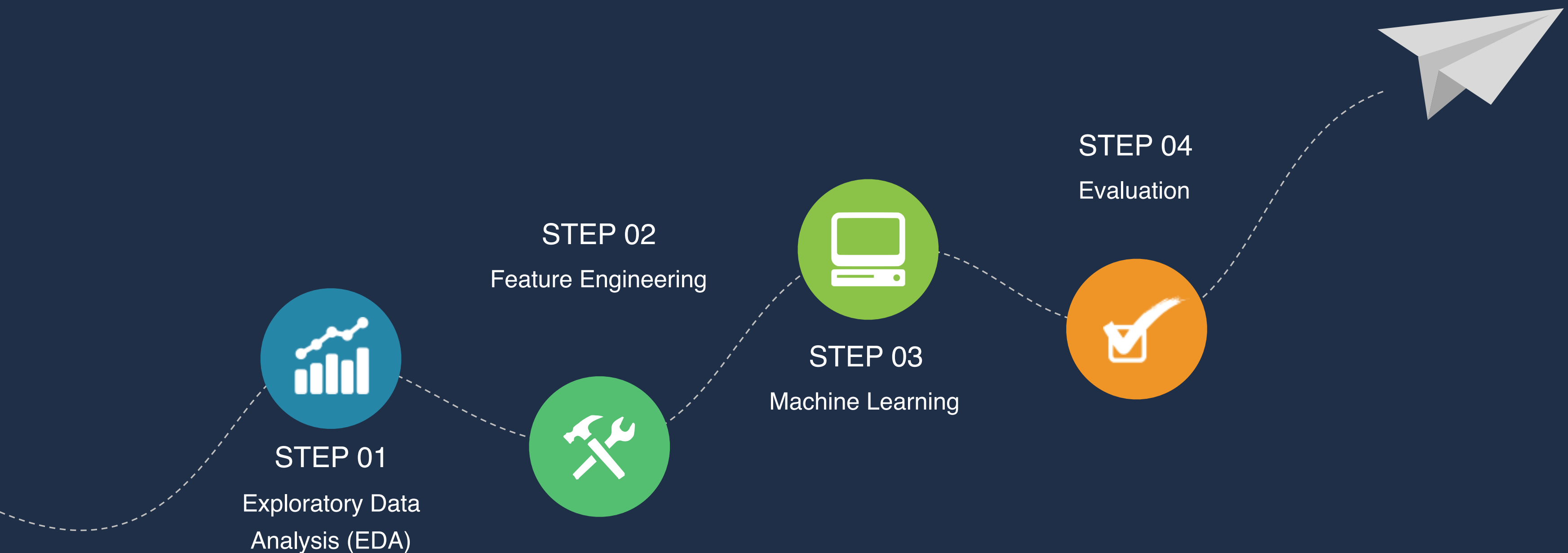AirREGI, an POS cash register app with reserve system for cafe and restaurants

2016-2017 visit and reserve data of 829 stores located in different locations of Japan

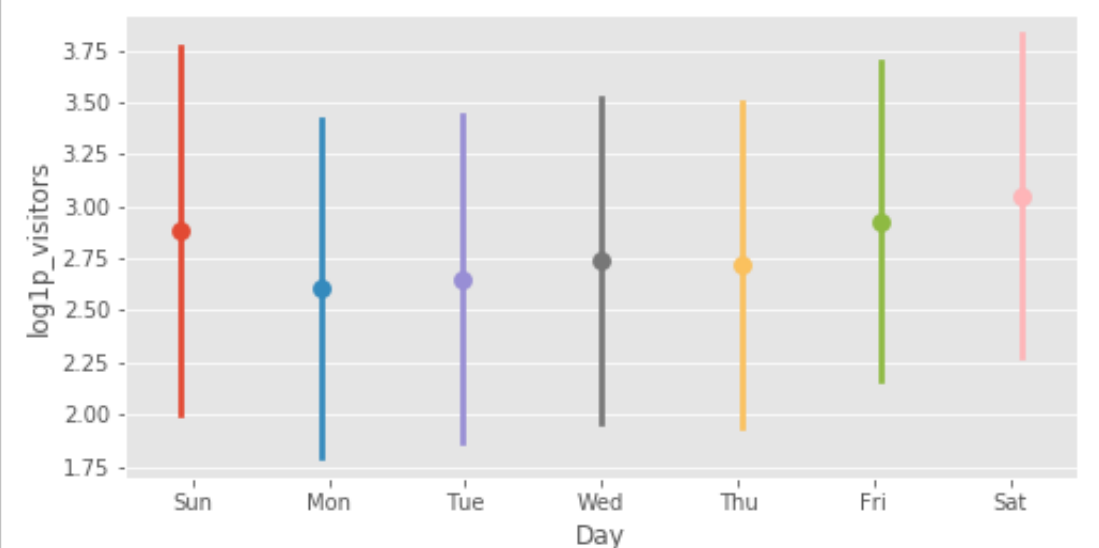Dataset can be obtained in Kaggle.

# Technical Requirement



01 — Programming Langauge Environment

02 — Data Manipulation

03 — Graphical Visualisation

04 — Statistics Machine Learning

05 — Software

# Data Science Methodology

**STEP 04**

Evaluation

**STEP 02**

Feature Engineering

**STEP 03**

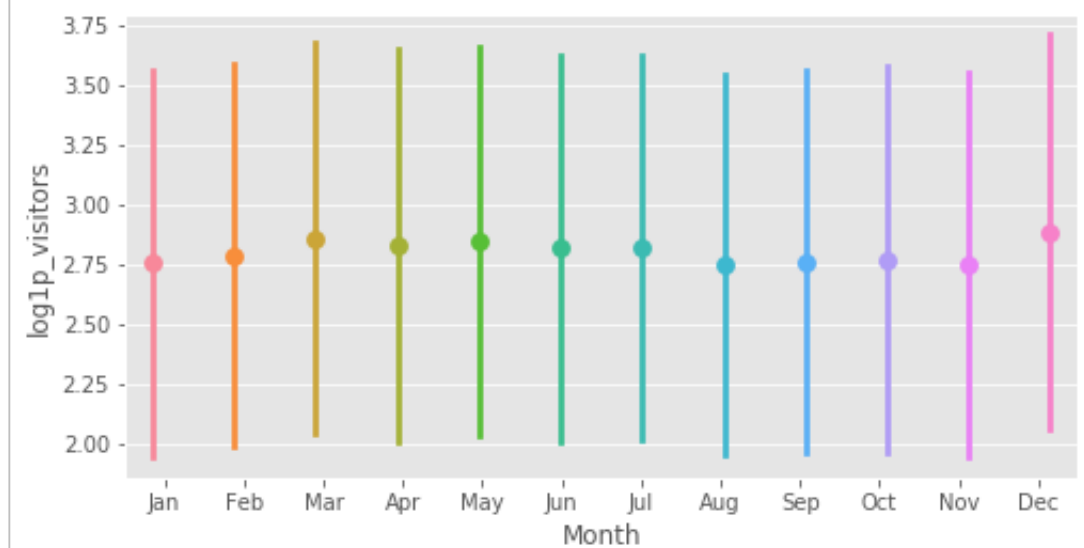Machine Learning

**STEP 01**

Exploratory Data
Analysis (EDA)

# Exploratory Data Analysis

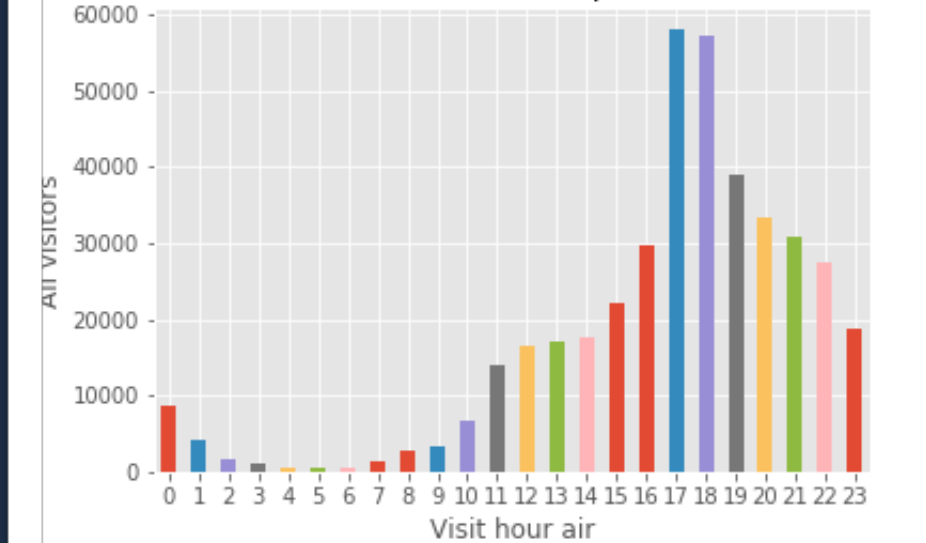An approach to summarise the main characteristics of the dataset mainly by visualisation.
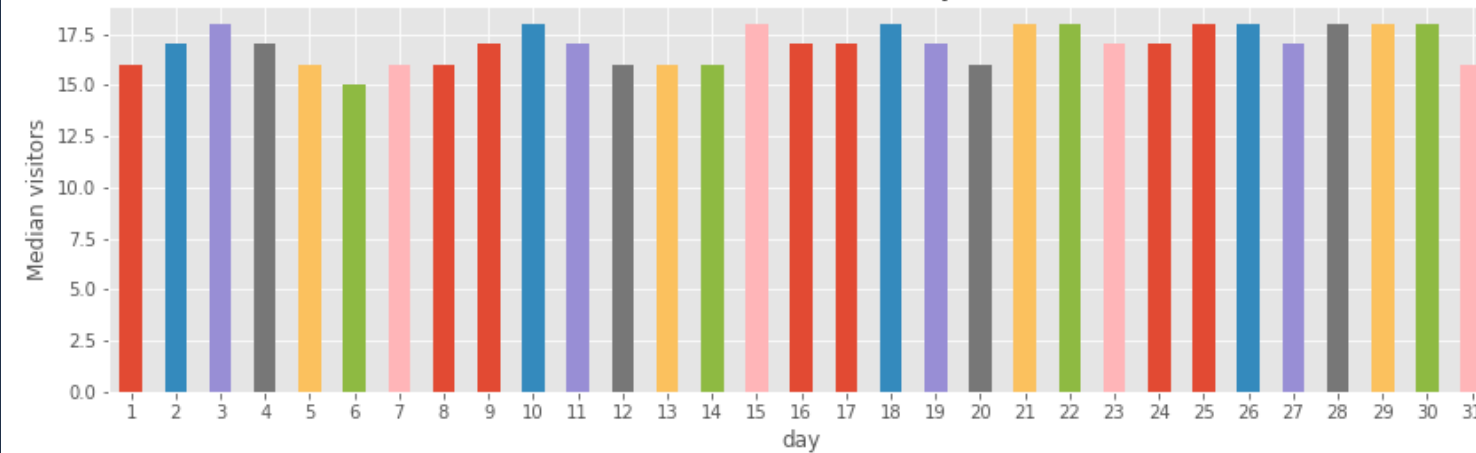
Air mean of log1p(visitors) in days of week
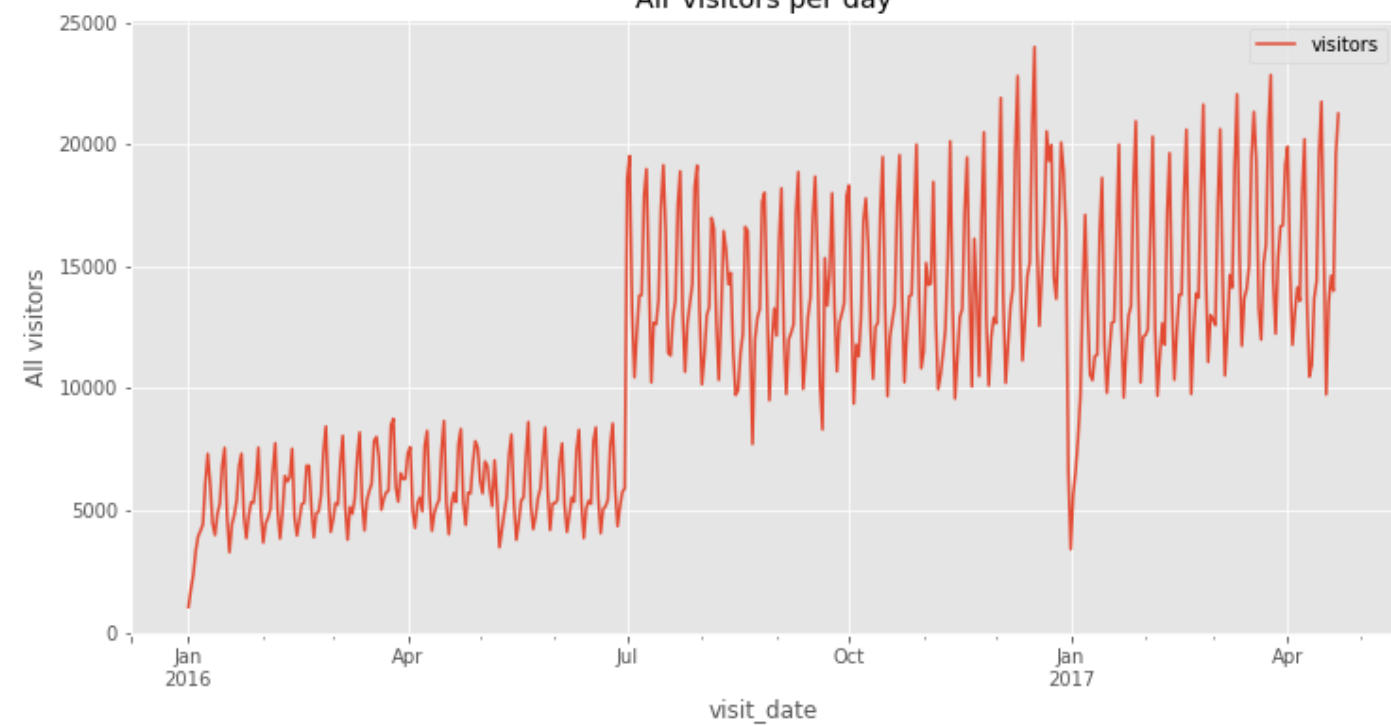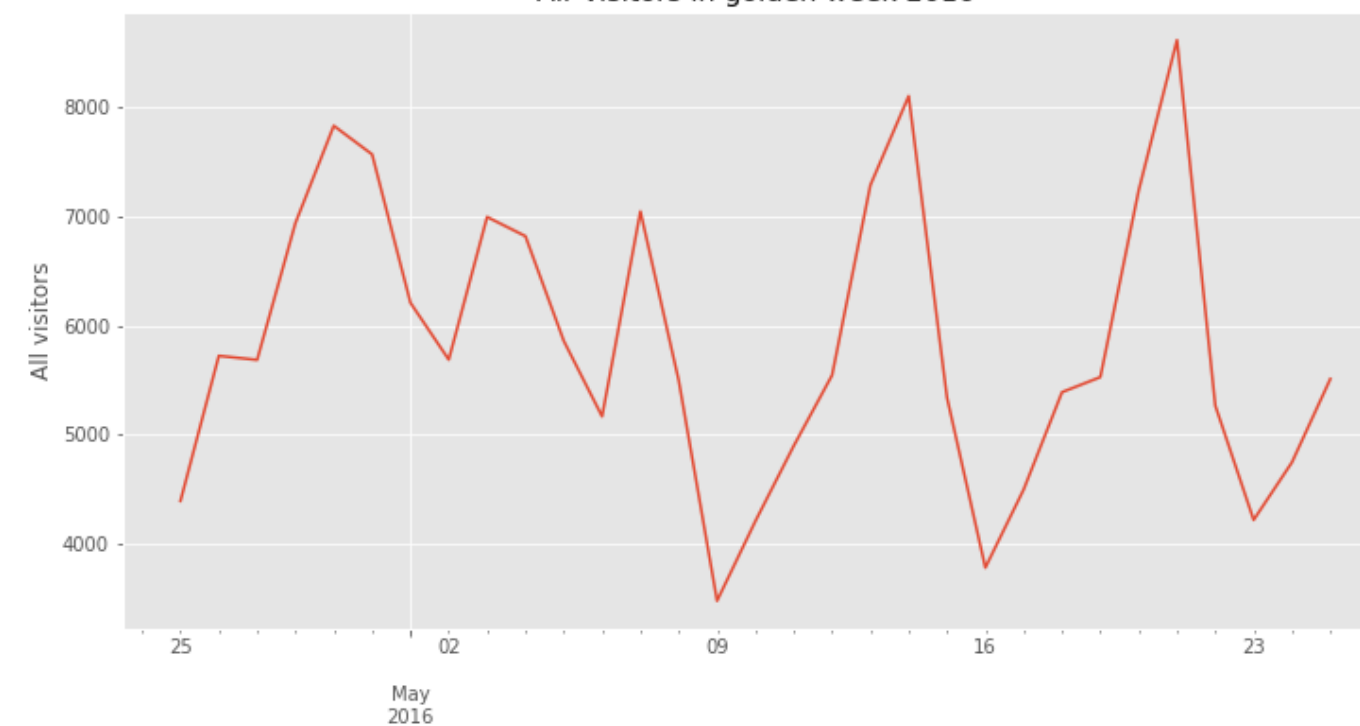
Air mean of log1p(visitors) in different months

Air no. of reserves per hour
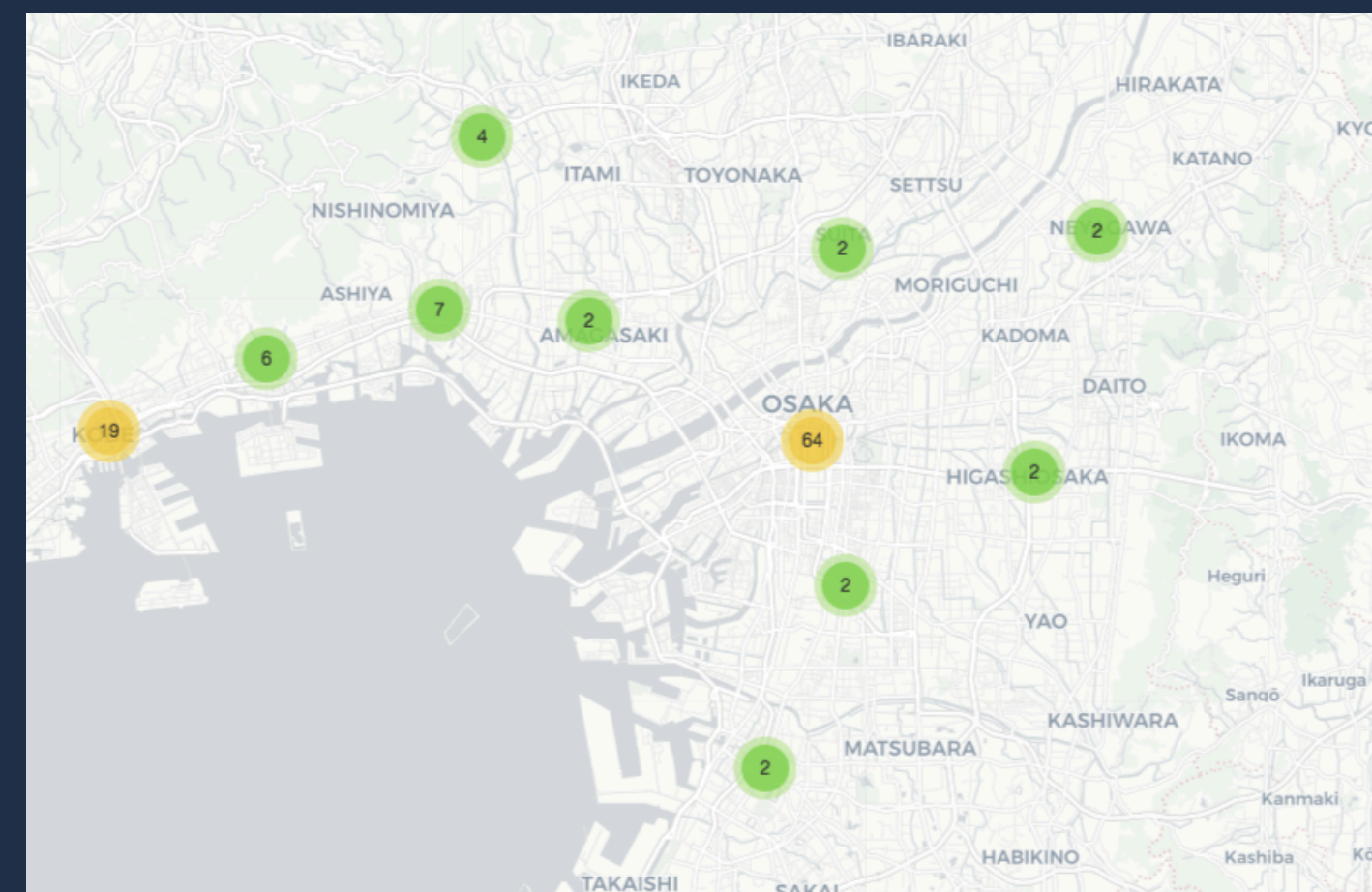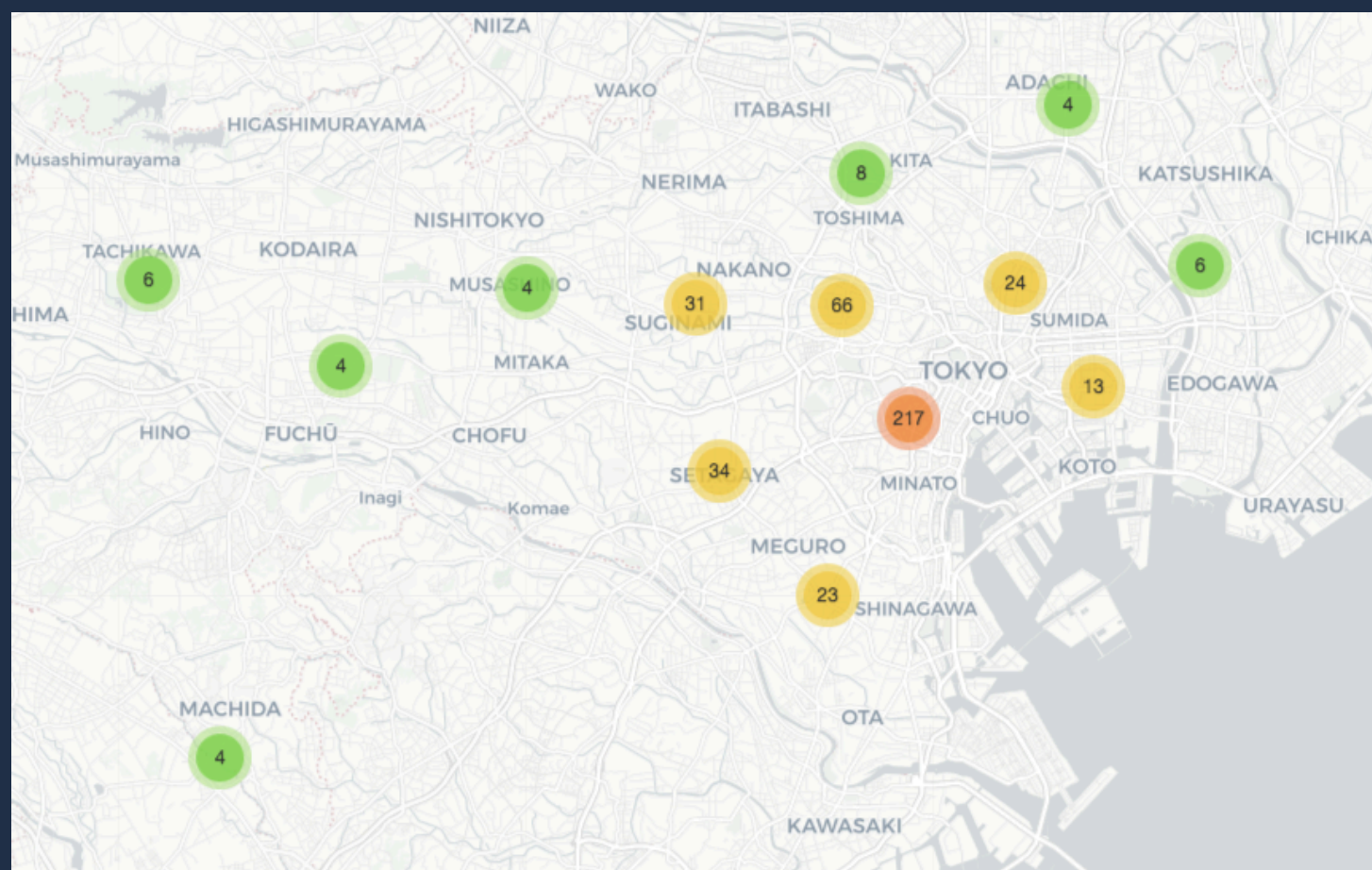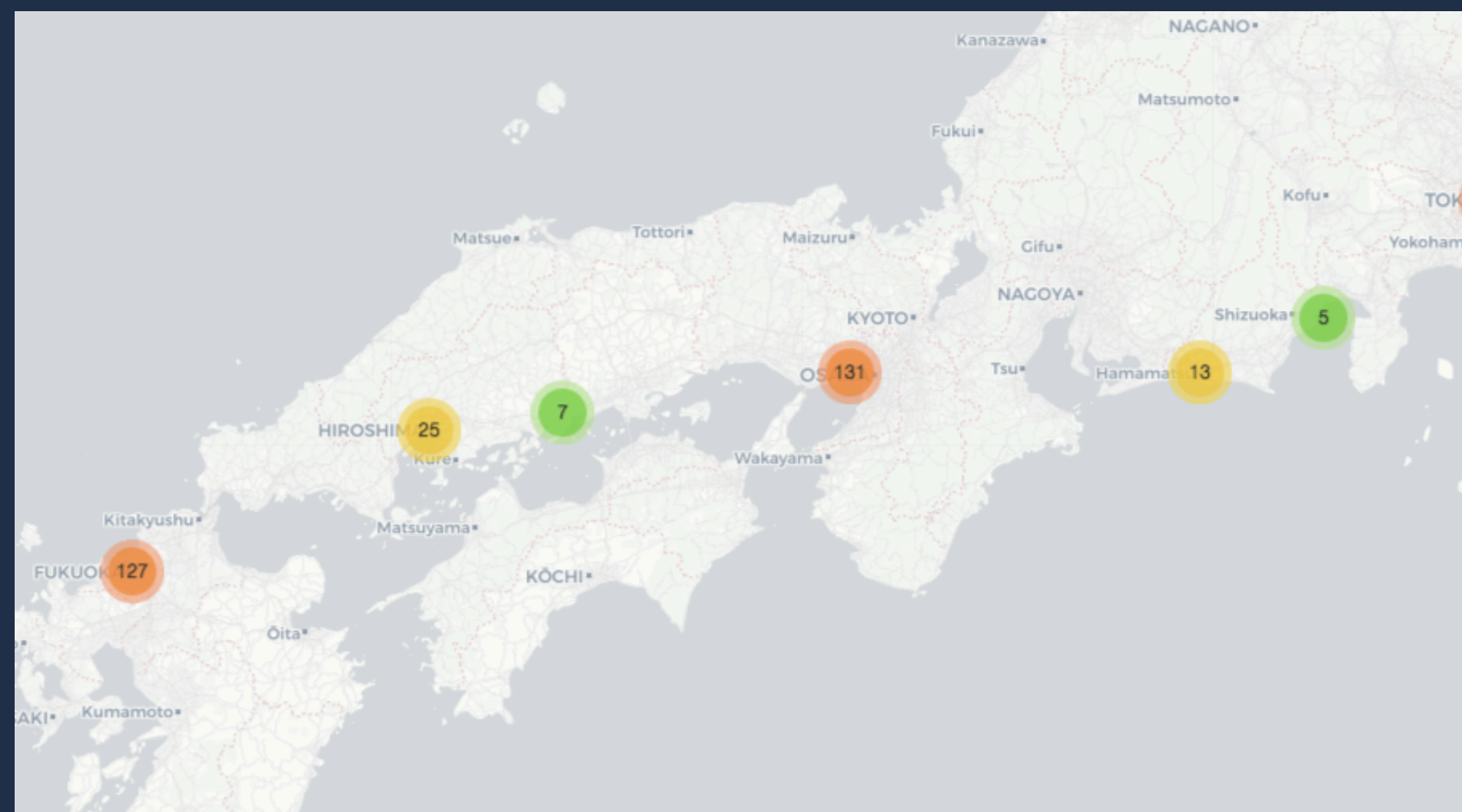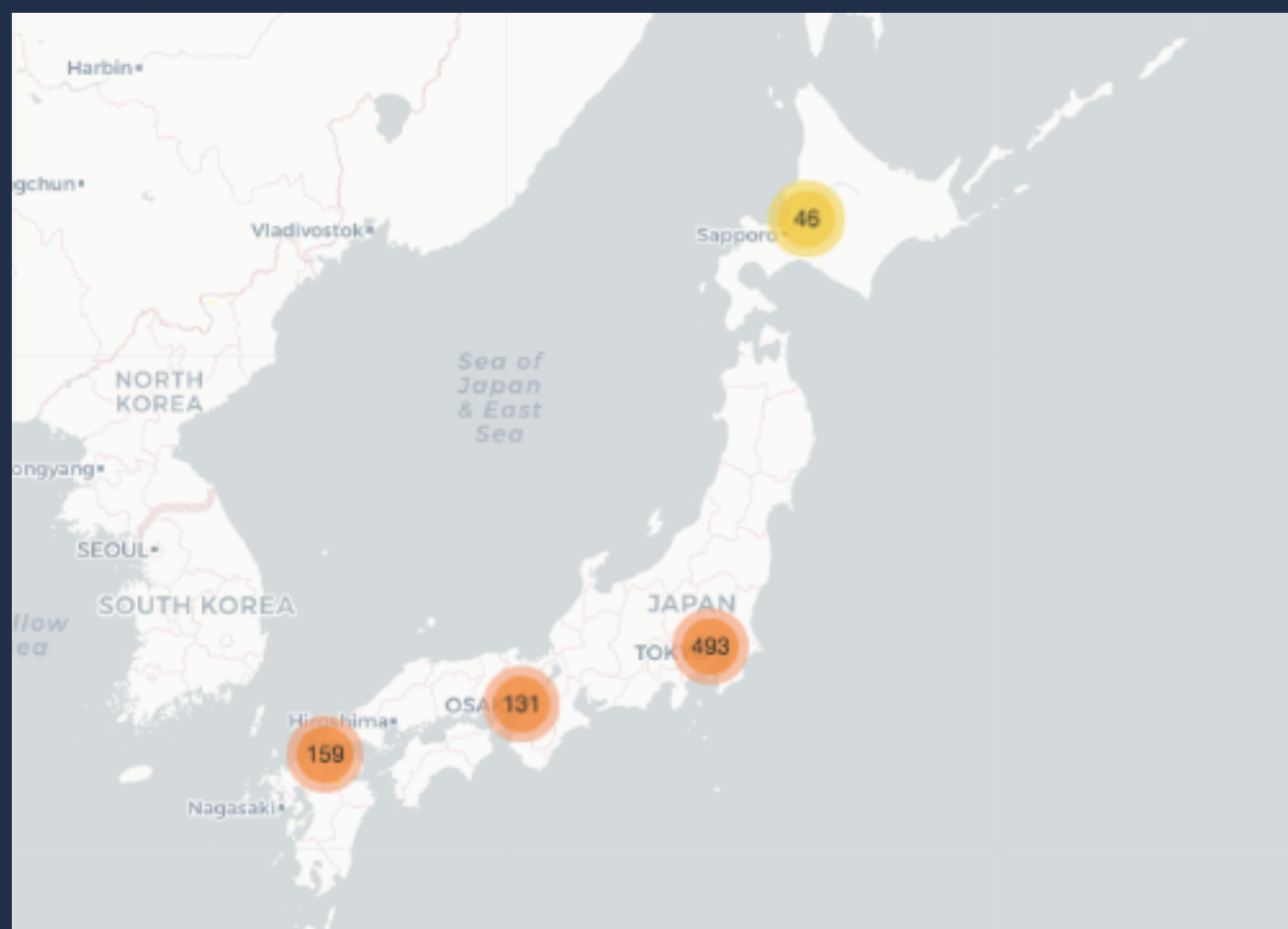
Air median visitors in different days of month
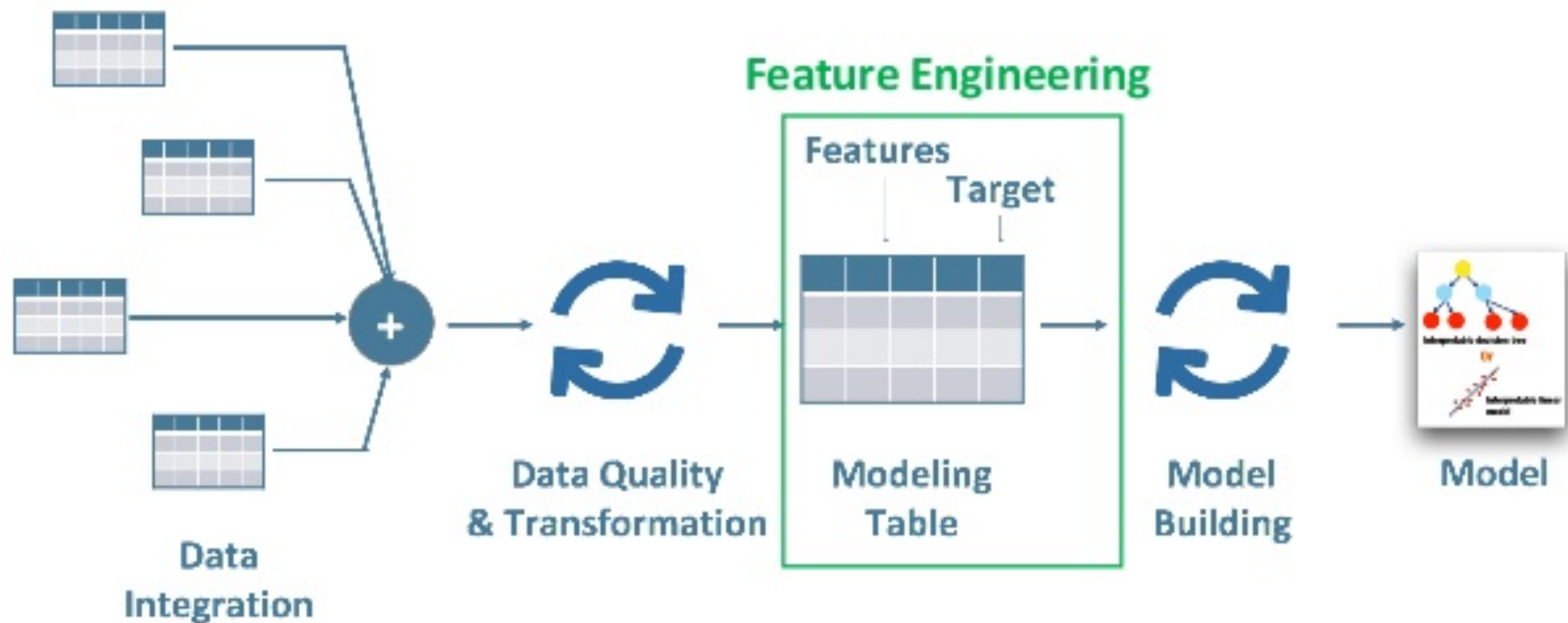
Air visitors per day

Air visitors in golden week 2016

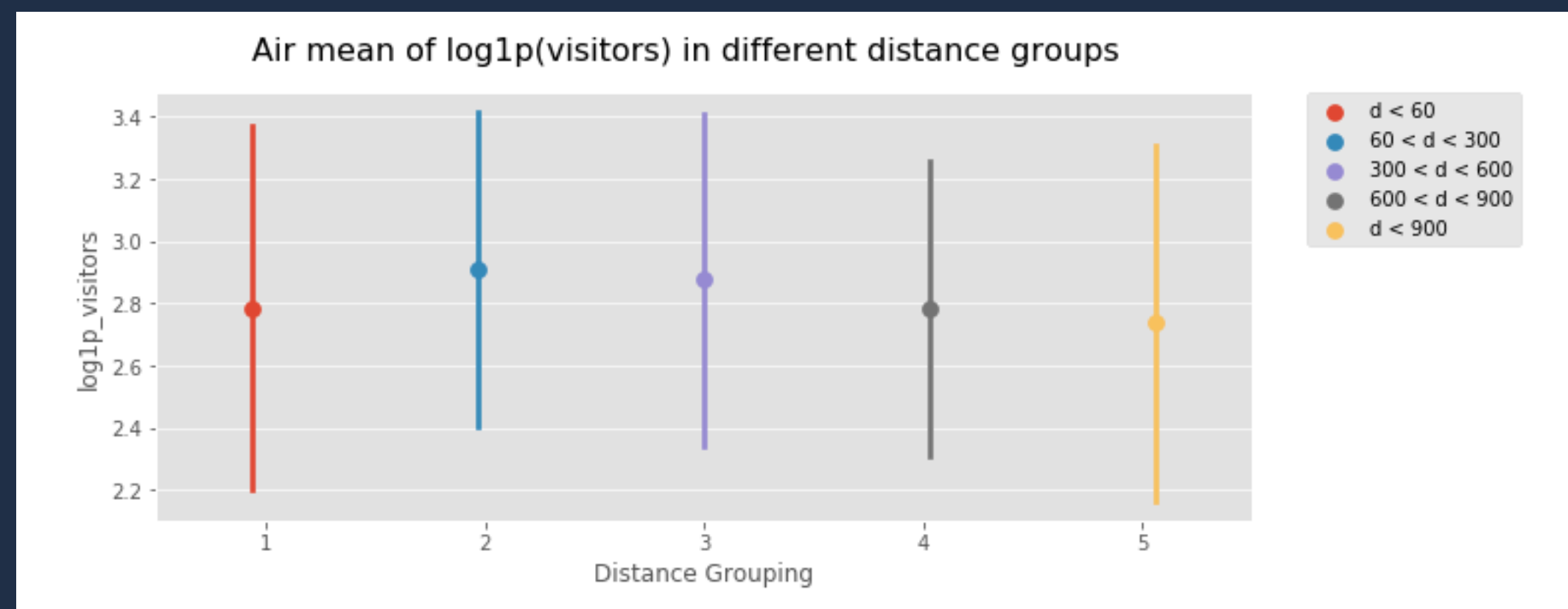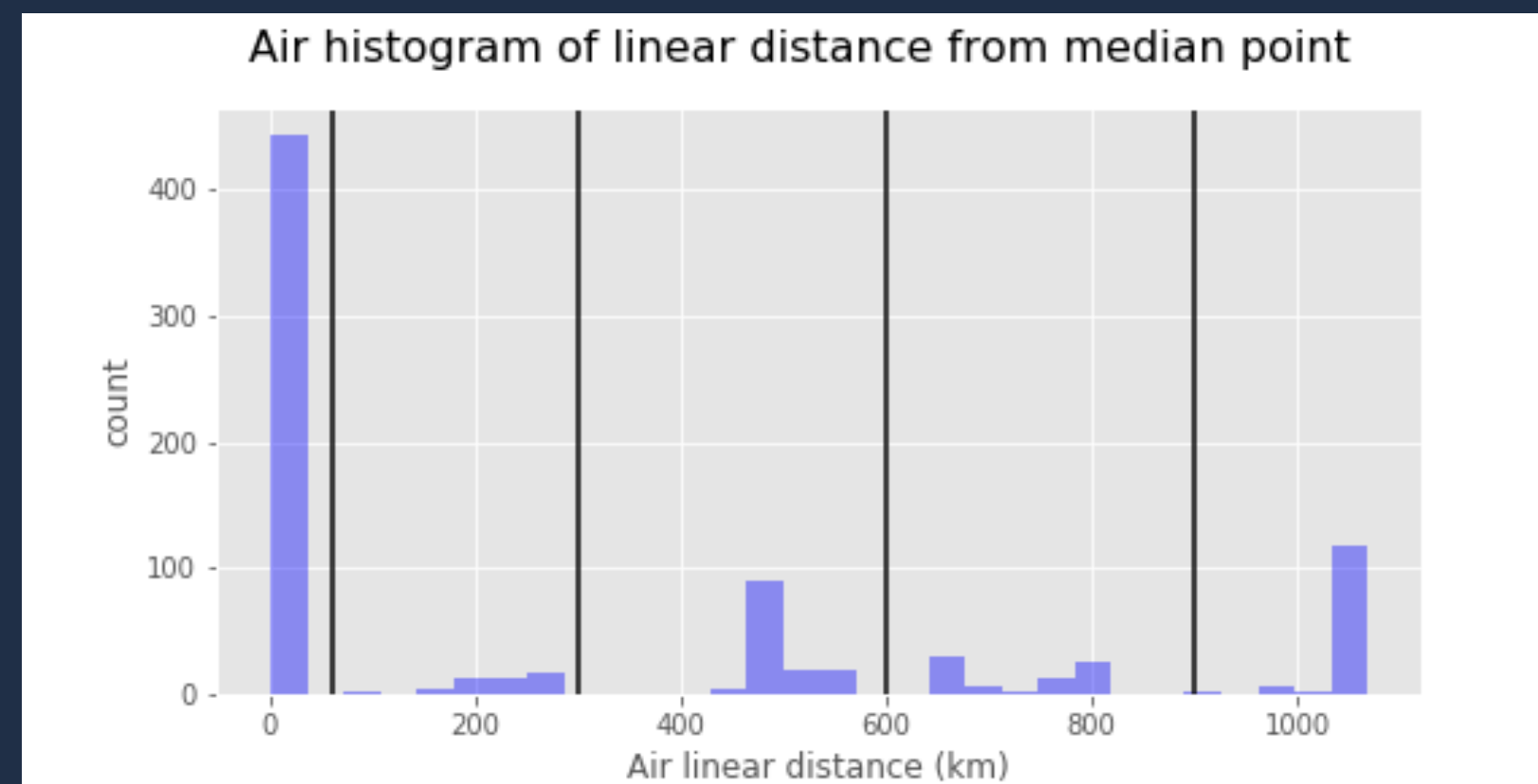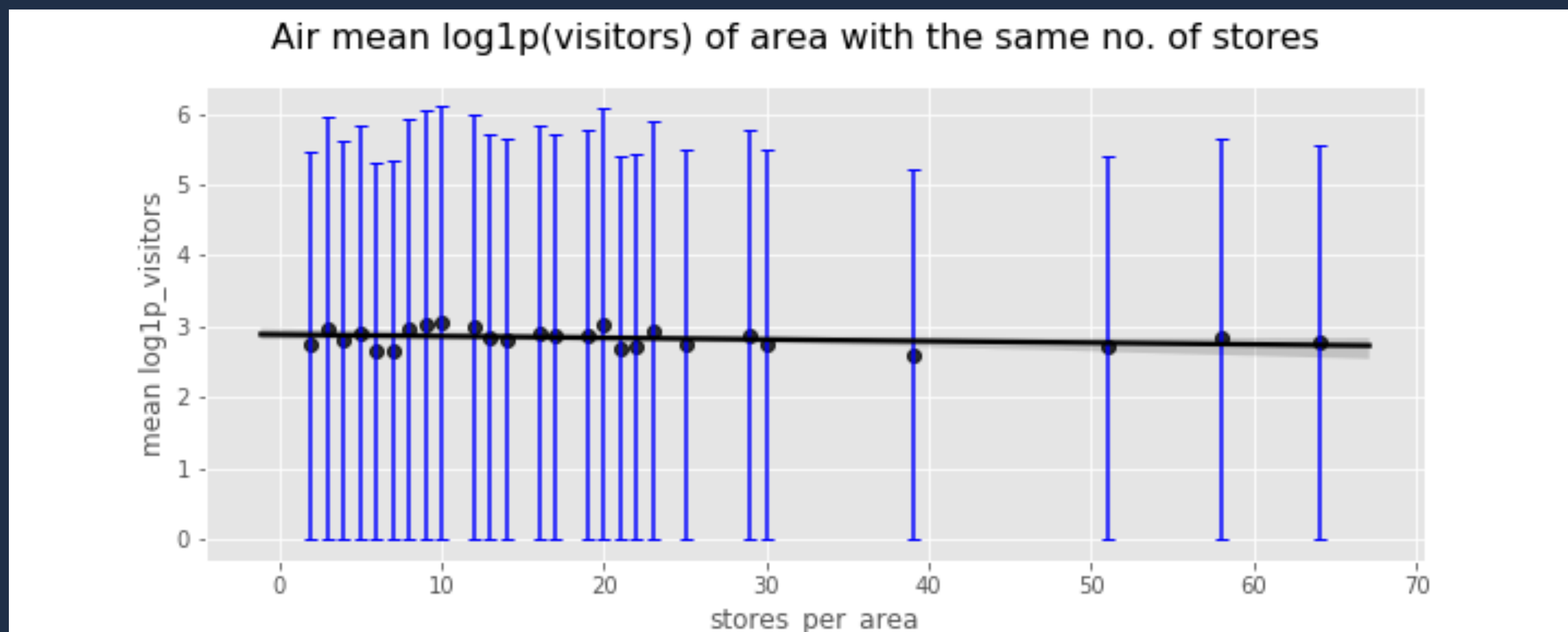# Feature Engineering

Create new features from the original ones

Air histogram of linear distance from median point



Air mean of log1p(visitors) in different distance groups

Air mean of log1p(visitors) in each prefecture



Air mean log1p(visitors) of area with the same no. of stores

# Modeling

# Facebook's Prophet ML Package

O**pen-source** forecasting model by facebook research team

Numerous applications across **Facebook**

Based on **additive model** with non-linear trends

**2018** Capture **seasonalities** of data plus **holiday** effects

**NA** Handles **missing data** and **outliers** automatically

PROPHET

# XGBoost ML Algorithm



eXtreme Gradient
Boosting

custom tree
building
algorithm

**XGBoost**
by Tianqi Chen

Used for:
• classification
• regression
• ranking
with custom loss
functions

Interfaces for
Python and R,
can be executed on
YARN

**eXtreme Gradient Boosting**

**O**pen-source gradient boosted **decision tree algorithm**

**Speedy** and high performance

Dominates **Kaggle** competitions

# RMSLE (Root Mean Squared Logarithmic Error)

$$\sqrt{\frac{1}{n} \sum_{i=1}^{n} (\log(p_i + 1) - \log(a_i + 1))^2}$$

Where:

- $n$ is the number of price quotes in the test set
- $p_i$ is your predicted price
- $a_i$ is the actual price
- $\log(x)$ is the natural logarithm

The **lower** the better

The **difference** between the predicted values and the actual values

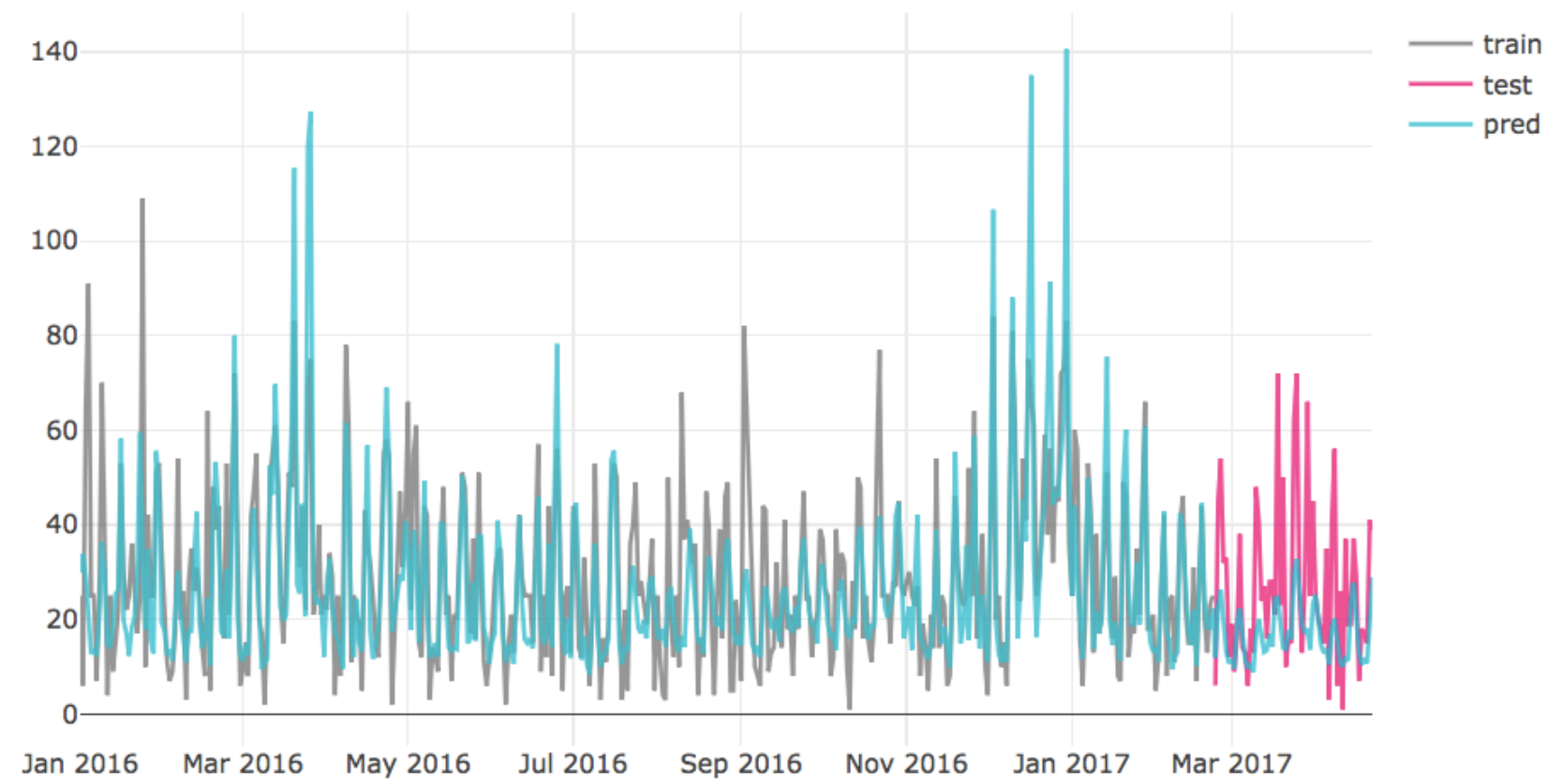Penalises an **under-predicted** estimate more

# ETL process comparison between Python and Tableau

A popular data visualising software for business intelligence

# Prophet

# XGBoost

# Prophet

# XGBoost

# thank you

To review the slide
Visit: https://https://github.com/jaycheung1096