# Recruit Restaurant Visitor Data Analysis and Forecasting Project



Running a thriving local restaurant isn't always as charming as first impressions appear. There are often all sorts of unexpected troubles popping up that could hurt business.

One common problem is that restaurants need to know how many customers to expect each day to effectively purchase ingredients and schedule staff members. This forecast isn't easy to make because many unpredictable factors affect restaurant attendance, like weather and local competition. It's even harder for newer restaurants with little historical data.

Recruit Holdings has unique access to key datasets that could make automated future customer prediction possible. The datasets comes from a restaurant review service (Hot Pepper Gourmet) and restaurant point of sales service (AirREGI, and reservation log management (Restaurant Board).

Recruit challenges people to use reservation and visitation data to predict the total number of visitors to a restaurant for future dates. This information will help restaurants be much more efficient and allow them to focus on creating an enjoyable experience for their customers.

In this paper, a key findings from a data analysis of the available datasets are presented, and a prediction algorithm is constructed with ARIMA for different types of restaurants.

The results indicate that there is a lot of interesting properties of the datasets, which could be used to improve processes. In addition, the prediction is shown to perform remarkably well in predicting the amount of visitors, in general, to different types of restaurants.

The original source of the problem and datasets can be found at:
https://www.kaggle.com/c/recruit-restaurant-visitor-forecasting

- For a start, all the datasets available will be presented, for an idea of what characteristics we can analyze

## Description of Datasets
Want to present the available datasets and typical values

### AirREGI Reservations
- Contains information of reservations done in the Air system

| | |
|---|---|
| **Store ID:** Identification of the restaurant | air_6b15edd1b4fbb96a |
| **Visit Date:** The date for the reservation | 2016-01-02 17:00:00 |
| **Reservation Date:** The date the reservation was done | 2016-01-01 22:00:00 |
| **Visitors:** The amount of spots reserved | 3 |

### Hot Pepper Gourmet Reservations
- Contains information of reservations done in the HPG system

| | |
|---|---|
| **Store ID:** Identification of the restaurant | hpg_33ec1499d6b13141 |
| **Visit Date:** The date for the reservation | 2016-01-01 17:00:00 |
| **Reservation Date:** The date the reservation was done | 2016-01-01 15:00:00 |
| **Visitors:** The amount of spots reserved | 2 |

### Date Information
- Contains information of dates

| | |
|---|---|
| **Calendar Date:** A date | 2016-01-01 |
| **Day of the Week:** Which day of the week | Friday |
| **Holiday:** Whether the date is a holiday | 1 |

### AirREGI Store Information
- Contains information of the stores in the air system

| | |
|---|---|
| **Store ID:** Identification of the restaurant | air_0f0cdeee6c9bf3d7 |
| **Store Genre:** The type of restaurant | Italian/French |
| **Area:** The area the restaurant is located at | Hyogo-ken Kobe-shi Kumoidori |
| **Latitude:** The latitude of the restaurant's location | 34.6951242 |
| **Longitude:** The longitude of the restaurant's location | 135.19785249999998 |

### AirREGI Visitors Information
- Contains information of the visitors at restaurants for different dates

| | |
|---|---|
| **Store ID:** Identification of the restaurant | air_ba937bf13d40fb24 |
| **Visit Date:** A date the restaurant is open | 2016-01-13 |
| **Visitors:** The amount of visitors for a date | 25 |

- To give a rough idea of how many people people, on average, visit and book reservations for restaurants, we may consider the distribution of visitor and reservation values



Histogram of reservations, visitors, and unplanned visitors for restaurants, for a given day
Outliers removed

· Amount of reservations done for a restaurant, for a given day, is often smaller than 10 but there exists
 cases where reservations reach larger values like 10, 20, 30 - perhaps they corresponds to holiday/weekend days
· Amount of visitors, generally, exceeds the number of reservations, implying that restaurants often receive
more customers than the amount of reservations - indicates consumers seldom have to worry about a restaurant running out of places
· Further illustrating the relation between reservations and visitors, the amount of visitors with no underlying
reservations are mostly non-negative, implying again that restaurants tends to receive more customers than the amount of reservations they receive
· Also note that from the amount of reservations and visitors, we can see that there exists a lot of cases where
a restaurant may receive visitors without a single reservation, for a particular day - which might corresponds to weekdays
or other low-traffic inducing days
· Days with high amount of visitors most likely correspond to weekends or holidays, which might explain
why these cases occur much less than days with low amount of visitors (which probably, in turn, corresponds to weekdays)

# Recruit Restaurants - Forecasting - Data Analysis - Type of restaurants in areas, amount of visitors to every type of restaurant over the week

- As all areas are different, there is an interest in considering what type, and how many, restaurants exists in each area - perhaps there is some area with a lot of Dining Bars



**Amount of Restaurants of Each Genre in Each Area**
Each box corresponds to one genre and one area, with the color denoting the amount of restaurants

· There is one area with a lot of Cafe/Sweets - might correspond to a place where spending moderate amount of money on food/drinks is really popular, like a tourist place, shopping center, or close to perhaps schools
· Areas that have a large amount of restaurants of one type usually have a large/moderate amount of restaurants of other types - might indicate that the underlying area is a popular place to eat/drink, for example inner cities, or areas with a large nightlife
· As expected from the analysis of the amount of restaurants of each type, Cafe/Sweets, Izakaya, and Italian/French restaurants are spread, more or less, over all different areas. The large spread of Cafe/Sweets and Izakaya over different areas indicates that informal restaurants are attractive among Japanese people

- Intuitively, the traffic to restaurants should change over the week, with perhaps an uptick of visitors during the weekend. This can be considered by the average amount of visitors for each type of restaurant over the week



**Average Amount of Visitors for Each Genre of Restaurants for Each Day of the Week**

| | Monday | Tuesday | Wednesday | Thursday | Friday | Saturday | Sunday |
|---|---|---|---|---|---|---|---|
| Asian | 34 | 33 | 35 | 33 | 39 | 45 | 48 |
| Bar/Cocktail | 9.8 | 11 | 11 | 11 | 14 | 18 | 16 |
| Cafe/Sweets | 20 | 19 | 20 | 20 | 20 | 29 | 30 |
| Creative cuisine | 19 | 19 | 21 | 21 | 25 | 30 | 29 |
| Dining bar | 15 | 15 | 17 | 17 | 22 | 25 | 20 |
| International cuisine | 16 | 19 | 19 | 20 | 25 | 41 | 37 |
| Italian/French | 19 | 20 | 21 | 21 | 26 | 27 | 25 |
| Izakaya | 19 | 19 | 22 | 20 | 28 | 29 | 23 |
| Japanese food | 17 | 17 | 19 | 18 | 22 | 23 | 21 |
| Karaoke/Party | 15 | 18 | 21 | 21 | 22 | 48 | 57 |
| Okonomiyaki/Monja/Teppanyaki | 18 | 20 | 20 | 21 | 24 | 29 | 29 |
| Other | 17 | 17 | 18 | 19 | 20 | 25 | 24 |
| Western food | 17 | 20 | 21 | 21 | 24 | 27 | 25 |
| Yakiniku/Korean food | 18 | 18 | 21 | 19 | 26 | 24 | 21 |

· Asian has the highest amount of visitors in general all weekdays, by a wide margin. Implies that Asian-oriented restaurants attract a lot of Japanese people. Might be because Japan is in Asia, and a generic Asian restaurant might offer a diverse set of meals for most people's preference, as opposed to e.g. Italian/French which attract people who like European food.
· International cuisine has a big jump in visitors during the weekend. A possible reason is that people want to eat something different toward the weekend, and since Japan is an Asian country, international food might be an attractive choice. For the same reason that Westerner might want to eat Asian food during the weekend
· Karaoke/Party, likewise, has a big jump in visitors during the weekend. Most likely, the weekend is an appropriate day to go out with people to either a a karaoke or party
· The weekday with the most visitors for all genres is friday (except weekend). No surprise as this is the last working day before the weekend, and people feel like eating something special to celebrate the weekend

# Recruit Restaurants - Forecasting - Data Analysis - Type of restaurants in Japan, amount of restaurants in areas

- To get an idea of what type of restaurants exists in Japan, and in what proportions, we consider an analysis of the air_genre_name attribute in the datasets



**Amount of Restaurants in Different Genres**

| Genre | Percentage | Avg Visitors |
|---|---|---|
| Izakaya | 23.8 % | 23 |
| Cafe/Sweets | 21.8 % | 22 |
| Dining bar | 13.0 % | 18 |
| Italian/French | 12.3 % | 23 |
| Bar/Cocktail | 9.5 % | 13 |
| Japanese food | 7.6 % | 19 |
| Other | 3.3 % | 20 |
| Yakiniku/Korean food | 2.8 % | 21 |
| Western food | 1.9 % | 22 |
| Okonomiyaki/Monja/Teppanyaki | 1.7 % | 29 |
| Creative cuisine | 1.6 % | 24 |
| Karaoke/Party | 0.2 % | 31 |
| International cuisine | 0.2 % | 23 |
| Asian | 0.2 % | 39 |

Legend:
- Percentage of Restaurants in Genre
- Average Amount of Visitors for Restaurants in Genre (daily)

· Izakaya: Informal Japanese Pub
· Okonomiyaki: Japanese Savory Pancakes
· Monja: Japanese Pan-fried Batter
· Teppanyaki: Japanese Cuisine that uses an iron graddle to cook food
· Yakiniku: Japanese Grilled Meat Cuisine

· Izakaya & Cafe/Sweets are among the most frequent type of restaurants - most likely because they are both informal, and they are most likely more affordable 'restaurants', as compared to restaurants offering lunch and/or dinner
· There is a large amount of Italian/French restaurants, which might imply that Italian and French food is popular in Japan - note that Italian restaurants might include Pizzerias
· Western food and international cuisine seems to be less frequent - which might imply that western food, aside Italian and French, are not too popular in Japan

· The high amount of average visitors for karaoke party might indicate that large groups might visit karaoke parties together (e.g. friends, workplace friends, family)
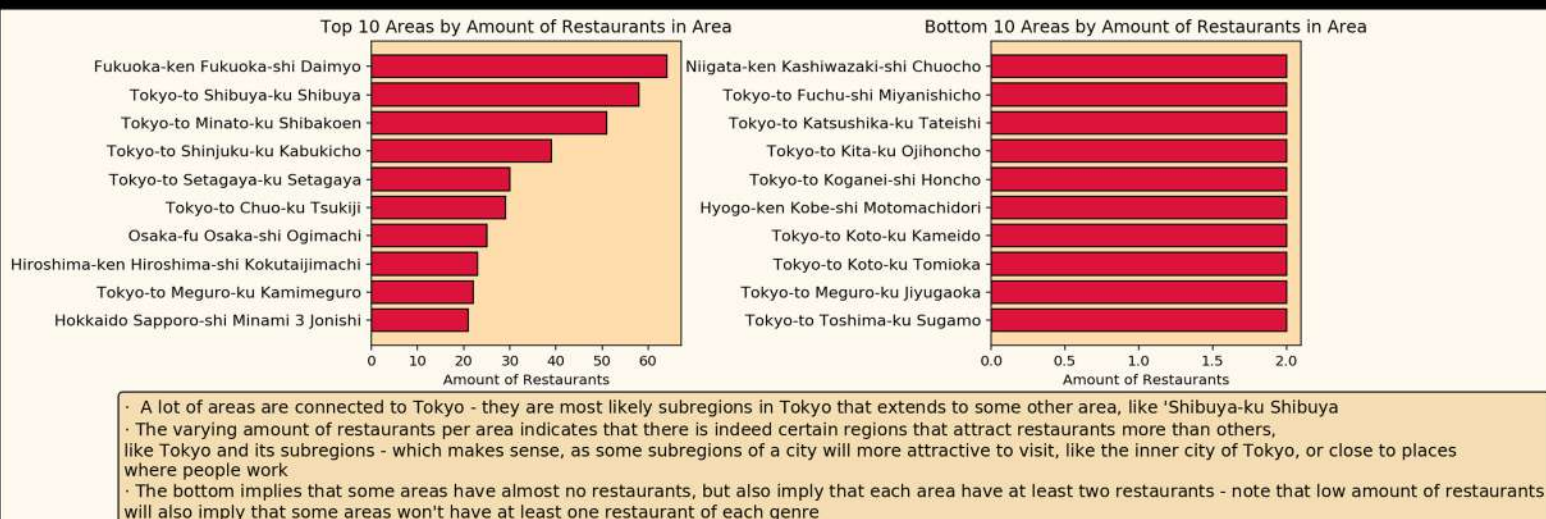· The low average visitors for Bar/Cocktail, relative Izakaya, might indicate that informal pubs are more popular in Japan
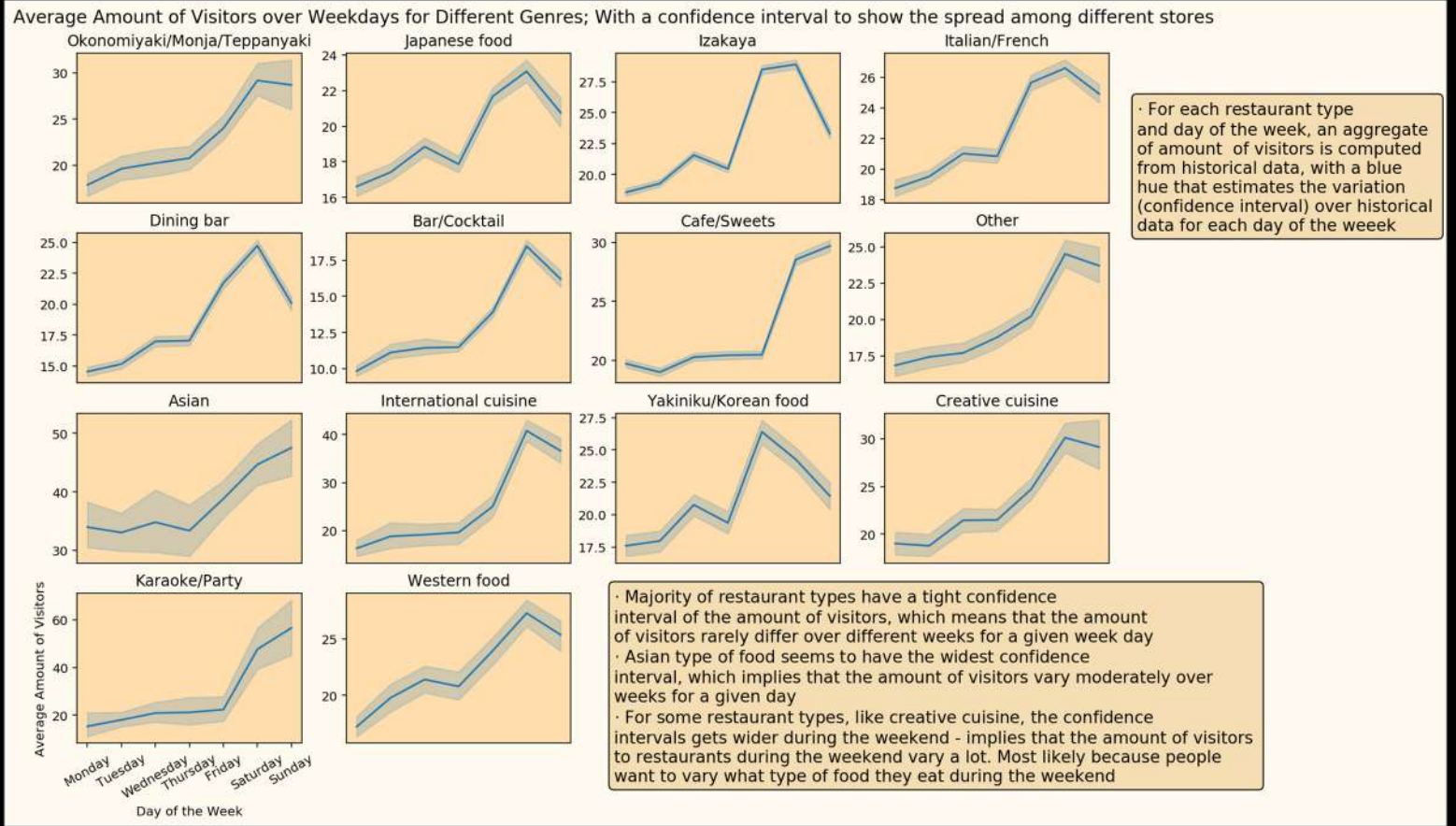
- Another interesting property is what areas of Japan exists, and how many restaurants exists in each area. Because the amount of areas exceed 100, only the top ten and top bottom areas by count are considered



**Top 10 Areas by Amount of Restaurants in Area**
- Fukuoka-ken Fukuoka-shi Daimyo
- Tokyo-to Shibuya-ku Shibuya
- Tokyo-to Minato-ku Shibakoen
- Tokyo-to Shinjuku-ku Kabukicho
- Tokyo-to Setagaya-ku Setagaya
- Tokyo-to Chuo-ku Tsukiji
- Osaka-fu Osaka-shi Ogimachi
- Hiroshima-ken Hiroshima-shi Kokutaijimachi
- Tokyo-to Meguro-ku Kamimeguro
- Hokkaido Sapporo-shi Minami 3 Jonishi

**Bottom 10 Areas by Amount of Restaurants in Area**
- Niigata-ken Kashiwazaki-shi Chuocho
- Tokyo-to Fuchu-shi Miyanishicho
- Tokyo-to Katsushika-ku Tateishi
- Tokyo-to Kita-ku Ojihoncho
- Tokyo-to Koganei-shi Honcho
- Hyogo-ken Kobe-shi Motomachidori
- Tokyo-to Koto-ku Kameido
- Tokyo-to Koto-ku Tomioka
- Tokyo-to Meguro-ku Jiyugaoka
- Tokyo-to Toshima-ku Sugamo

· A lot of areas are connected to Tokyo - they are most likely subregions in Tokyo that extends to some other area, like 'Shibuya-ku Shibuya'
· The varying amount of restaurants per area indicates that there is indeed certain regions that attract restaurants more than others, like Tokyo and its subregions - which makes sense, as some subregions of a city will more attractive to visit, like the inner city of Tokyo, or close to places where people work
· The bottom implies that some areas have almost no restaurants, but also imply that each area have at least two restaurants - note that low amount of restaurants will also imply that some areas won't have at least one restaurant of each genre
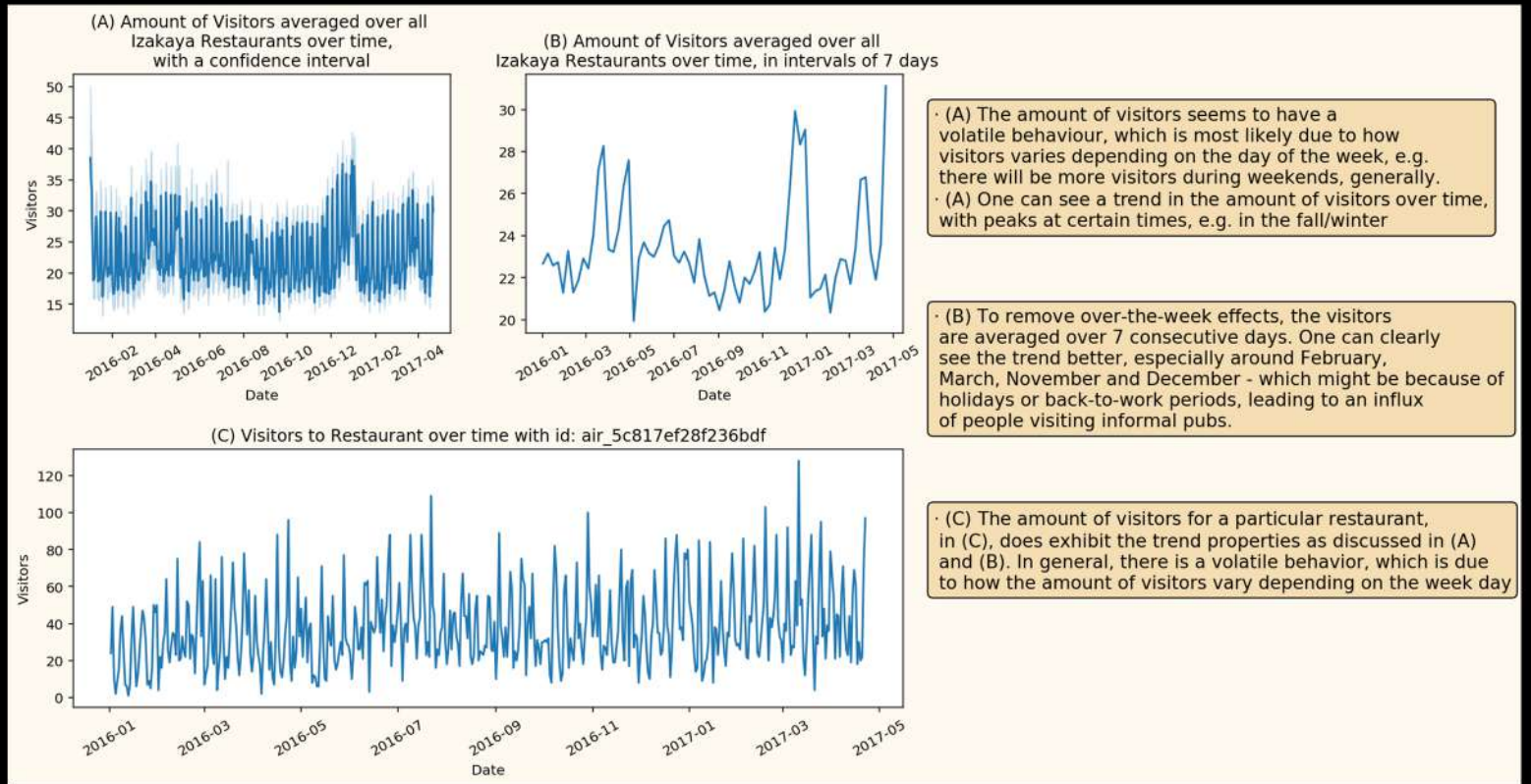
# Recruit Restaurants - Forecasting - Data Analysis - Average Amount of Visitors over the week, time-series for a particular restaurant

- To analyze how the amount of visitors varies over each day of the week, the average of amount of visitors is computed for each genre over all days of the week



Average Amount of Visitors over Weekdays for Different Genres; With a confidence interval to show the spread among different stores

· For each restaurant type and day of the week, an aggregate of amount of visitors is computed from historical data, with a blue hue that estimates the variation (confidence interval) over historical data for each day of the weeek

· Majority of restaurant types have a tight confidence interval of the amount of visitors, which means that the amount of visitors rarely differ over different weeks for a given week day
· Asian type of food seems to have the widest confidence interval, which implies that the amount of visitors vary moderately over weeks for a given day
· For some restaurant types, like creative cuisine, the confidence intervals gets wider during the weekend - implies that the amount of visitors to restaurants during the weekend vary a lot. Most likely because people want to vary what type of food they eat during the weekend
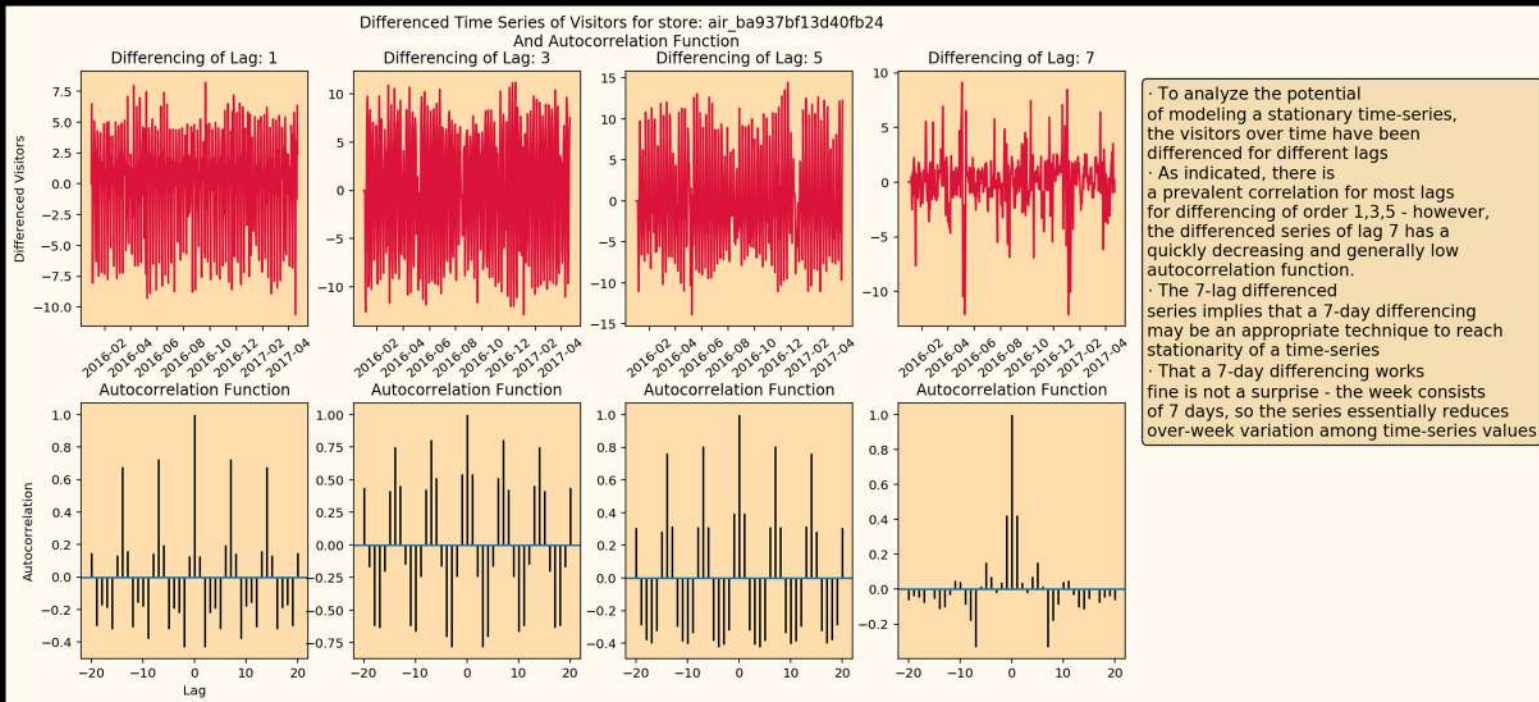
- For visitors in a particular genre, an idea is to analyze restaurants in the Izakaya genre. In particular, it is of interest to analyze the amount of visitors over time for a particular store in the Izakaya genre



· (A) The amount of visitors seems to have a volatile behaviour, which is most likely due to how visitors varies depending on the day of the week, e.g. there will be more visitors during weekends, generally.
· (A) One can see a trend in the amount of visitors over time, with peaks at certain times, e.g. in the fall/winter

· (B) To remove over-the-week effects, the visitors are averaged over 7 consecutive days. One can clearly see the trend better, especially around February, March, November and December - which might be because of holidays or back-to-work periods, leading to an influx of people visiting informal pubs.

· (C) The amount of visitors for a particular restaurant, in (C), does exhibit the trend properties as discussed in (A) and (B). In general, there is a volatile behavior, which is due to how the amount of visitors vary depending on the week day
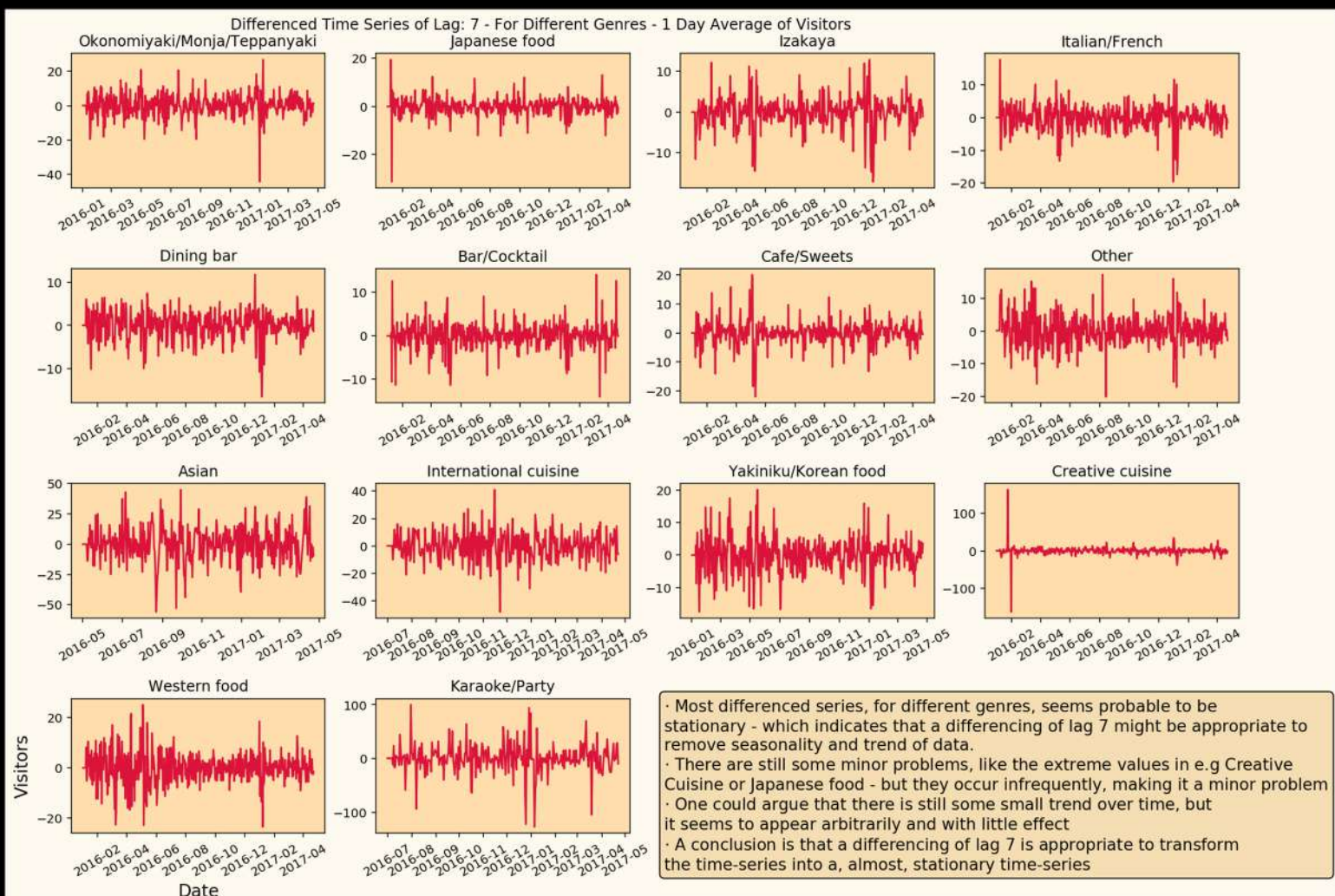
# Recruit Restaurants - Forecasting - Prediction Analysis - Time series analysis

- Before applying mathematical models to our time-series data, a key idea is to analyze whether the time-series can be transformed to a stationary, or almost, time-series. One typical approach is to difference the time-series, and check for weak stationarity (constant mean over time, and time-independent autocorrelation function)



Differenced Time Series of Visitors for store: air_ba937bf13d40fb24 And Autocorrelation Function

· To analyze the potential of modeling a stationary time-series, the visitors over time have been differenced for different lags
· As indicated, there is a prevalent correlation for most lags for differencing of order 1,3,5 - however, the differenced series of lag 7 has a quickly decreasing and generally low autocorrelation function.
· The 7-lag differenced series implies that a 7-day differencing may be an appropriate technique to reach stationarity of a time-series
· That a 7-day differencing works fine is not a surprise - the week consists of 7 days, so the series essentially reduces over-week variation among time-series values

- To further illustrate the effect of differencing time-series with a lag of 7, we consider this type of transformation to the average amount of visitors (daily) for each type of restaurant, to see if the assumption of stationarity is plausible



Differenced Time Series of Lag: 7 - For Different Genres - 1 Day Average of Visitors

· Most differenced series, for different genres, seems probable to be stationary - which indicates that a differencing of lag 7 might be appropriate to remove seasonality and trend of data.
· There are still some minor problems, like the extreme values in e.g Creative Cuisine or Japanese food - but they occur infrequently, making it a minor problem
· One could argue that there is still some small trend over time, but it seems to appear arbitrarily and with little effect
· A conclusion is that a differencing of lag 7 is appropriate to transform the time-series into a, almost, stationary time-series

# Recruit Restaurants - Forecasting - Prediction Analysis - Optimal ARIMA models, and prediction on a week of days with optimal ARMA models

- With our 7-day differenced time-series of the daily average visitors of each genre, the next step is to fit it to optimal ARIMA models. Based on considering different AR and MA orders, and comparing the prediction performance with two metrics on a test week

| Optimal Parameters for Different Genres - MAE | | | |
|---|---|---|---|
| Genre Name | p | q | MAE |
| Okonomiyaki/... | 5 | 3 | 2.66 |
| Japanese food | 6 | 3 | 1.37 |
| Izakaya | 5 | 3 | 1.65 |
| Italian/French | 7 | 3 | 2.58 |
| Dining bar | 6 | 1 | 0.64 |
| Bar/Cocktail | 7 | 1 | 0.96 |
| Cafe/Sweets | 4 | 2 | 1.79 |
| Other | 6 | 1 | 1.88 |
| Asian | 5 | 1 | 12.4 |
| International cuisine | 4 | 2 | 5.66 |
| Yakiniku/Korean food | 5 | 2 | 1.25 |
| Creative cuisine | 6 | 3 | 3.8 |
| Western food | 2 | 3 | 2.65 |
| Karaoke/Party | 6 | 3 | 9.86 |

| Optimal Parameters for Different Genres - RMSE | | | |
|---|---|---|---|
| Genre Name | p | q | RMSE |
| Okonomiyaki/... | 7 | 2 | 3.46 |
| Japanese food | 7 | 3 | 1.91 |
| Izakaya | 6 | 2 | 2.0 |
| Italian/French | 7 | 3 | 3.0 |
| Dining bar | 6 | 2 | 0.86 |
| Bar/Cocktail | 7 | 3 | 1.2 |
| Cafe/Sweets | 7 | 2 | 2.65 |
| Other | 6 | 1 | 2.09 |
| Asian | 4 | 2 | 14.61 |
| International cuisine | 4 | 2 | 6.5 |
| Yakiniku/Korean food | 5 | 2 | 1.87 |
| Creative cuisine | 6 | 3 | 4.59 |
| Western food | 2 | 3 | 3.22 |
| Karaoke/Party | 6 | 3 | 11.45 |

**ARIMA(p,d,q) Model - with statsmodels package (Python)**
Based on utilizing d=7, i.e. a differencing of lag 7
To find optimal p and q, two metrics, MAE and RMSE, are evaluated on on a test set consisting of a week, given an ARIMA model trained on a training set.
For the search of optimal p, q, all ARIMA models for p=1,2,...,7 and q=1,2,...,7 have been considered

In above tables, the optimal choices of p and q for each genre is presented.
MAE is the mean of the residuals, with respect to true and predicted time series values.
RMSE is the square root of the mean of the residuals squared

- With our optimal models, based on MAE, the last step is to evaluate the models on a week of days, the validation set, to really measure the performance of our models



Optimal ARIMA Models, based on MAE, Predictions on a Single Week (Validation Dates); For Different Genres, 1 Day Average of Visitors

Most predictions follow the trend of the true values, and is in a lot of cases not too far apart from the true values.
The predicted values are pretty bad for a single date in Bar/Cocktail, because of some reason.
The predictions are, in particular, good for Western Food and Italian/French.