

# Data Wrangling for Food Demand Forecasting

## Data Acquisition

The data was provided by the meal delivery company. The dataset, present in csv format, contains

1. Product(Meal) features such as category, sub-category, current price and discount.
2. Information for fulfillment center like center area, city information etc.
3. Historical data of demand for a product-center combination (Weeks: 1 to 145)
4. Future data of product-center combination for prediction (Weeks: 146 to 155)

## Data Analysis

Below steps were taken to analyse the dataset

| Data                                   | Analysis Steps   | Conclusion  |
|--|--|---|
| <b>Meal Information</b>                | <ul style="list-style-type: none"><li>• Missing Values</li><li>• Data Types of each column</li><li>• Distribution of category and cuisine columns</li></ul>  | There is nothing to clean.  |
| <b>Fulfilment Center Information</b>   | <ul style="list-style-type: none"><li>• Missing Values</li><li>• Data Types of each column</li><li>• Distribution of city_code, region_code and center_type columns</li><li>• Statistics of op_area column like mean, percentile, etc</li></ul>  | There is nothing to clean.  |
| <b>Training data / Historical data</b> | <ul style="list-style-type: none"><li>• Missing Values</li><li>• Data Types of each column</li><li>• Distribution of week, center_id, meal_id, emailer_for_promotion and homepage_featured columns</li><li>• Statistics of num_orders, checkout_price and base_price columns</li></ul> | <ul style="list-style-type: none"><li>• Records are missing for some week, center and meal combination</li><li>• Data contains two outliers, one with 24299 as number of orders and another as 2.97 checkout_price.</li></ul> |

|                                |   |  |
|--------------------------------|---|--|
| <b>Test data / Future data</b> | <ul style="list-style-type: none"> <li>• Missing Values</li> <li>• Data Types of each column</li> <li>• Distribution of week, center_id, meal_id, emailer_for_promotion and homepage_featured columns</li> <li>• Statistics of num_orders, checkout_price and base_price columns</li> </ul> | Records are missing for some weeks, center and meal combination. |
|--------------------------------|---|--|

## Outliers

Action on outliers will be taken during the modeling based on the performance of model with and without outliers.

## Missing Records

Missing records can be because to below two reasons:

1. There is actually no sales for that meal, center and weeks combination
2. Records were not captured due to technical error

Reason will become more clear after exploring data and then action can be taken.

## Data Merging

All three data are present in different dataframes. Hence, its required to merge them into one dataframe. Below steps were taken to merge the dataset

1. Left join on training data and meal information on meal\_id.
2. Left join on training data and fulfilment center information on center\_id.

Same steps were taken for test data.

## Derive new variables

Below are the new variables, related to number of orders, derived from existing dataset

| Variable Name        | Description  | Derived from   |
|----------------------|--|--|
| average_orders_Nweek | It is the mean of num_orders for particular meal_id and center_id in past few weeks.<br>N -> 13, 26 and 52 | <ul style="list-style-type: none"> <li>• center_id</li> <li>• meal_id</li> <li>• week</li> <li>• num_orders</li> </ul> |

|                                 |   |  |
|---------------------------------|---|--|
| average_orders_Nweek_across     | It is the mean of num_orders for particular meal_id across all centers in the past few weeks.<br>N -> 13, 26 and 52   | <ul style="list-style-type: none"> <li>• meal_id</li> <li>• week</li> <li>• num_orders</li> </ul>                      |
| average_orders_Nweek_adj        | It is the mean of num_orders for particular meal_id and center_id in past few weeks ending at 10 weeks in the past.<br>e.g:- for week 50, past weeks will be 37-40 weeks.<br>N -> 13 and 26 | <ul style="list-style-type: none"> <li>• center_id</li> <li>• meal_id</li> <li>• week</li> <li>• num_orders</li> </ul> |
| average_orders_Nweek_adj_across | It is the mean of num_orders for particular meal_id across all centers in the past few weeks ending at 10 weeks in the past.<br>N -> 13 and 26  | <ul style="list-style-type: none"> <li>• meal_id</li> <li>• week</li> <li>• num_orders</li> </ul>                      |

Below are the new variables, related to time periods, derived from week column

| Variable      | Description   |
|---------------|---|
| year          | It represents the year, group of 52 consecutive weeks, in which the record belongs.   |
| month         | It represents the month, group of 4 consecutive weeks in a year, in which the record belongs. Since, month is considered as a set of 4 weeks, there are 13 months in the dataset. |
| quarter       | It represents the quarter, group of 13 consecutive weeks in a year, in which the record belongs.  |
| week_in_month | Since, month contains set of 4 weeks, this variable represents record belongs to which of these 4 weeks.  |

Below are the new variables related to prices of meals.

| Variable        | Description   | Derived from   |
|-----------------|---|--|
| mean_base_price | It is the mean of all base_price for a particular center_id and meal_id till that week        | <ul style="list-style-type: none"><li>• center_id</li><li>• meal_id</li><li>• week (&lt;= current record)</li><li>• base_price</li></ul> |
| discount        | It is the discount (in percentage) that customers got in that week for a meal in that center. | <ul style="list-style-type: none"><li>• mean_base_price</li><li>• checkout_price</li></ul>   |