

# Food Demand Forecast

## Business, Economic and Financial Data Project

Pierpaolo D'Odorico, Massimiliano Conte and Eddie Rossi

# Food Demand Forecasting

## The business problem:

A **meal delivery company** operates in multiple cities. They have various **fulfillment centers** in these cities for dispatching **meal orders** to their customers.

We need to **forecast** for upcoming weeks, so that these centers will **plan the stock** of raw materials accordingly.

## Task:

**Predict the demand** for the next **10 weeks!**

# Data sources

Data are collected in 3 different datasets, connected by keys.

## Datasets

- Fulfilment centers data
- Meal info data
- Sales historical data

## Fulfilment centers data

center_id	city_code	region_code	center_type	op_area
11	679	56	TYPE_A	3.7
13	590	56	TYPE_B	6.7
124	590	56	TYPE_C	4.0
66	648	34	TYPE_A	4.1
94	632	34	TYPE_C	3.6
64	553	77	TYPE_A	4.4

# Fulfilment centers data features

## Variables:

- **center\_id**: Fulfilment identifier
- **city\_code**: City id in which the center is located on
- **region\_code**: Region id in which the center is located on
- **center\_type**: Type of the center
- **op\_area**: Size of the operational area

## Unique values in dataset:

center\_id : 77, city\_code : 51, region\_code : 8, center\_type : 3,  
op\_area : 30

## Meal info data

meal_id	category	cuisine
1885	Beverages	Thai
1993	Beverages	Thai
2539	Beverages	Thai
1248	Beverages	Indian
2631	Beverages	Indian
1311	Extras	Thai

# Meal data features

## Variables:

- **meal\_id**: Meal identifier
- **category**: Category of food
- **cuisine**: Category of cuisine

## Unique values in dataset:

meal\_id : 51, category : 14, cuisine : 4

## Meal info data\*

<u>id</u>	<u>week</u>	<u>center_id</u>	<u>meal_id</u>	<u>base_price</u>	<u>num_orders</u>
1379560	1	55	1885	152.29	177
1466964	1	55	1993	135.83	270
1346989	1	55	2539	135.86	189
1338232	1	55	2139	437.53	54
1448490	1	55	2631	242.50	40
1270037	1	55	1248	252.23	28

\*binomial variables and “checkout\_price” are removed for a better dataset view.

# Meal data features

## Variables:

- **id**: Id of the single transaction
- **week**: Temporal variable, we have 145 unique weeks
- **center\_id**: Fulfilment identifier
- **meal\_id**: Meal identifier
- **checkout\_price**: Paid price for the product
- **base\_price**: Full price of the product without promotion
- **emailer\_for\_promotion**: Binomial, promotion email or not
- **homepage\_featured**: Binomial, product on web homepage
- **num\_orders**: Number of orders for the meal and center

## Meal data features

### Unique values in dataset:

id : 456548, week : 145, center\_id : 77, meal\_id : 51,  
checkout\_price : 1992, base\_price : 1907, emailer\_for\_promotion :  
2, homepage\_featured : 2, num\_orders : 1250

## Create a unique dataset

We created a unique dataset **merging by keys**.

There are **0 NA's** in the complete dataset.

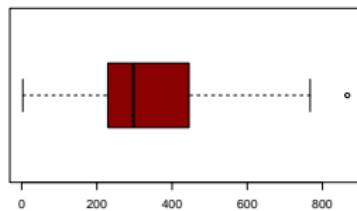
We will perform some exploratory data analysis:

- **Univariate Analysis:** Looking at single variables behaviour
- **Multivariate Analysis:** Correlation between variables

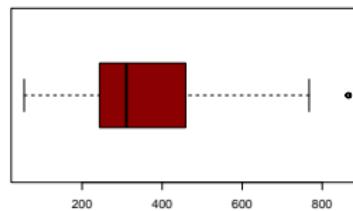
# Univariate Analysis on numerical variables

We plot boxplots of **numerical variables**:

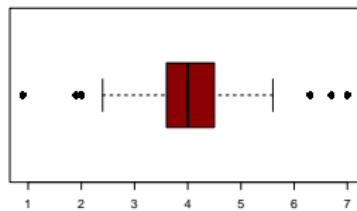
checkout\_p boxplot



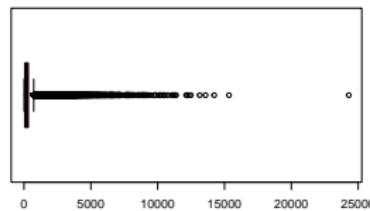
base\_p boxplot



op\_area boxplot

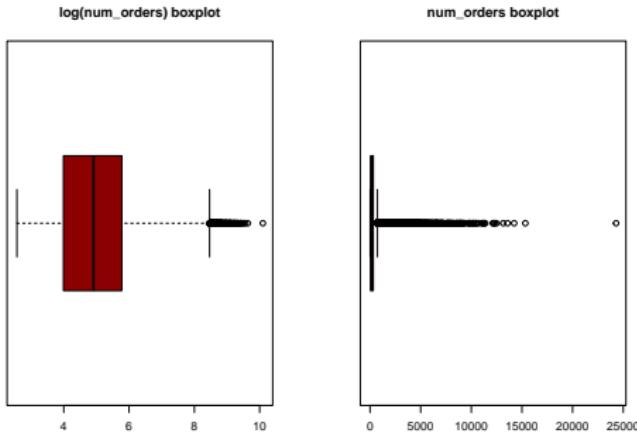


num\_orders boxplot



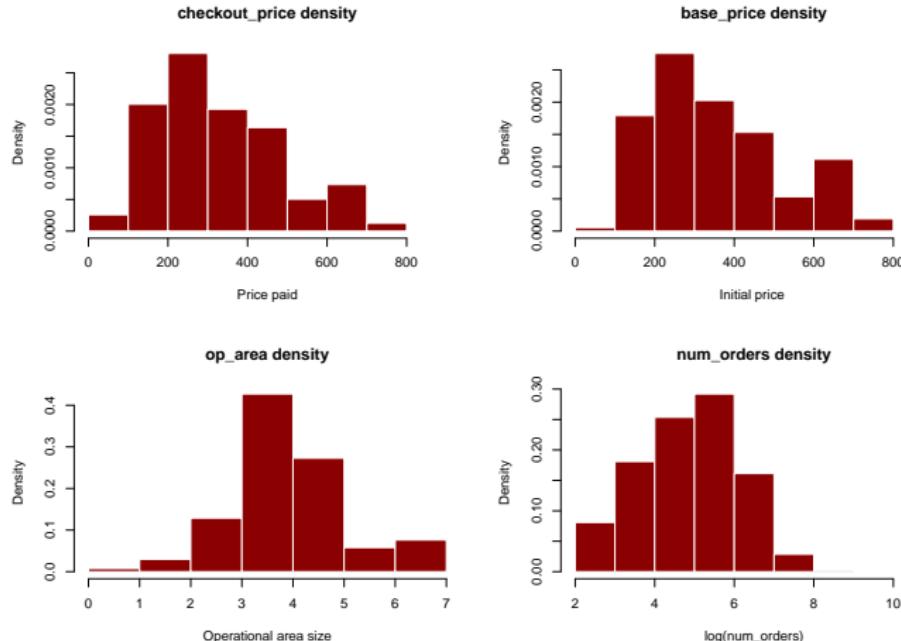
# Univariate Analysis on numerical variables

**num\_orders** boxplot vs log transformation for a better view:

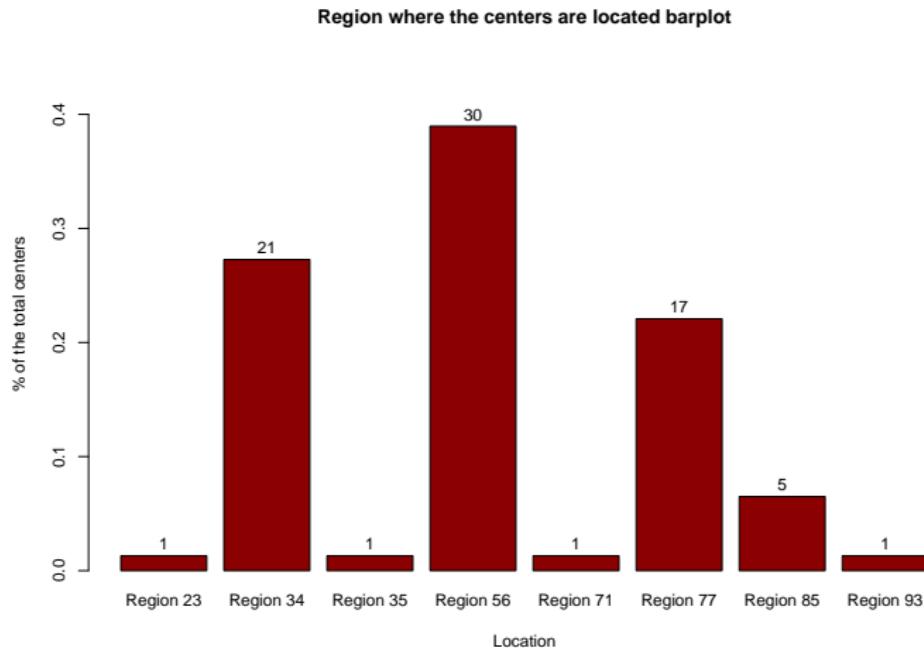


High **num\_orders** are related to some specific series from a specific center and meal with high demand. For this reason we don't consider them outliers.

# Univariate Analysis on numerical variables

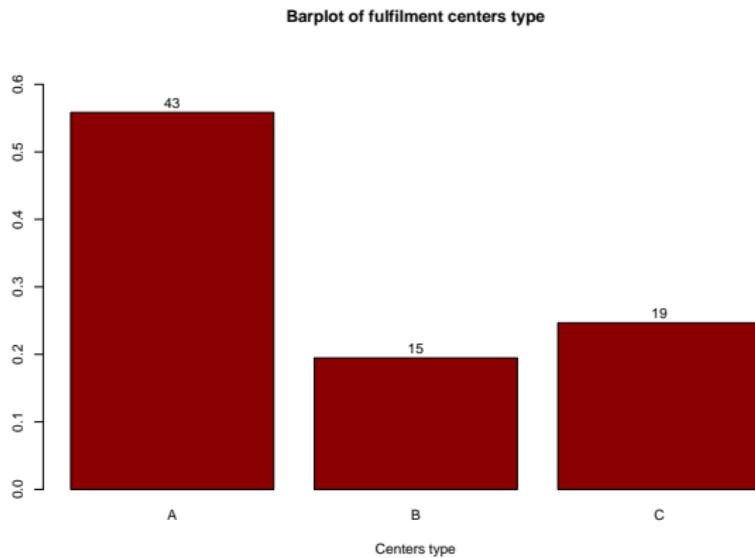


# Univariate Analysis on categorical variables

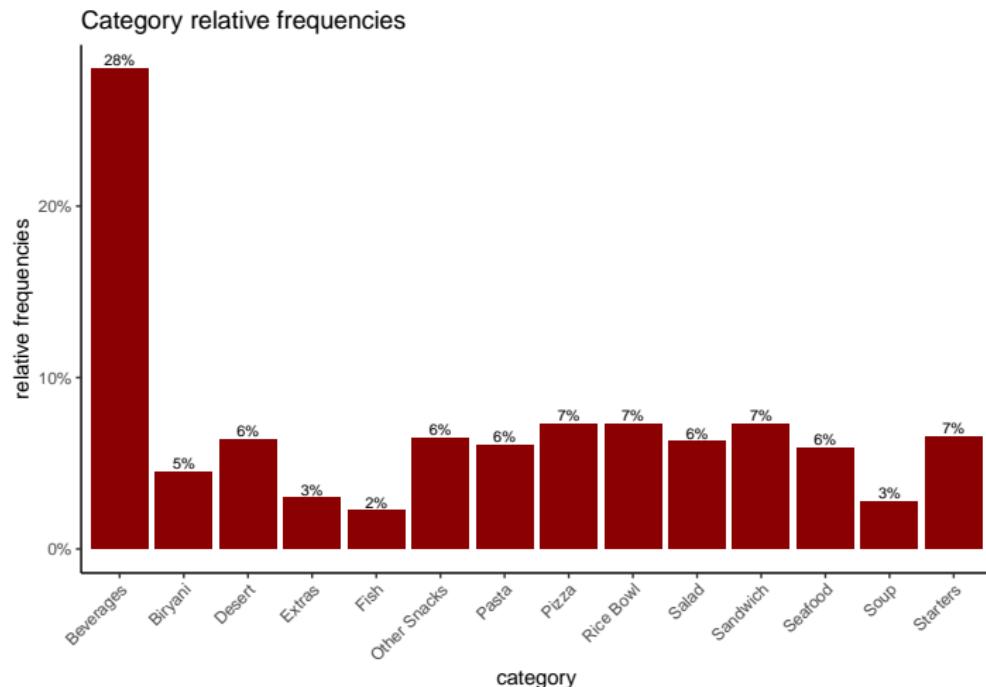


# Univariate Analysis on categorical variables

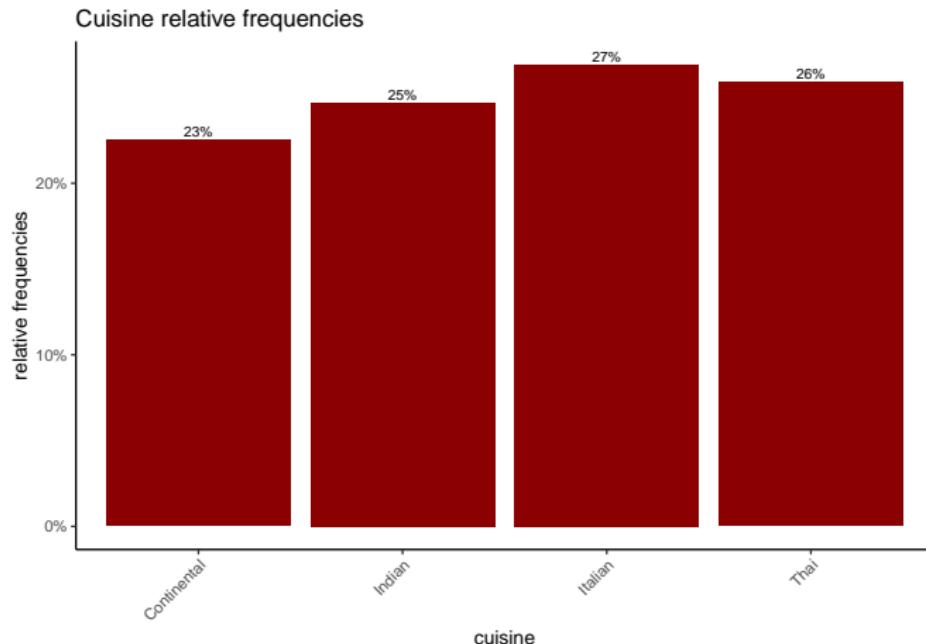
**center\_type:** There are 3 centers types: A, B and C



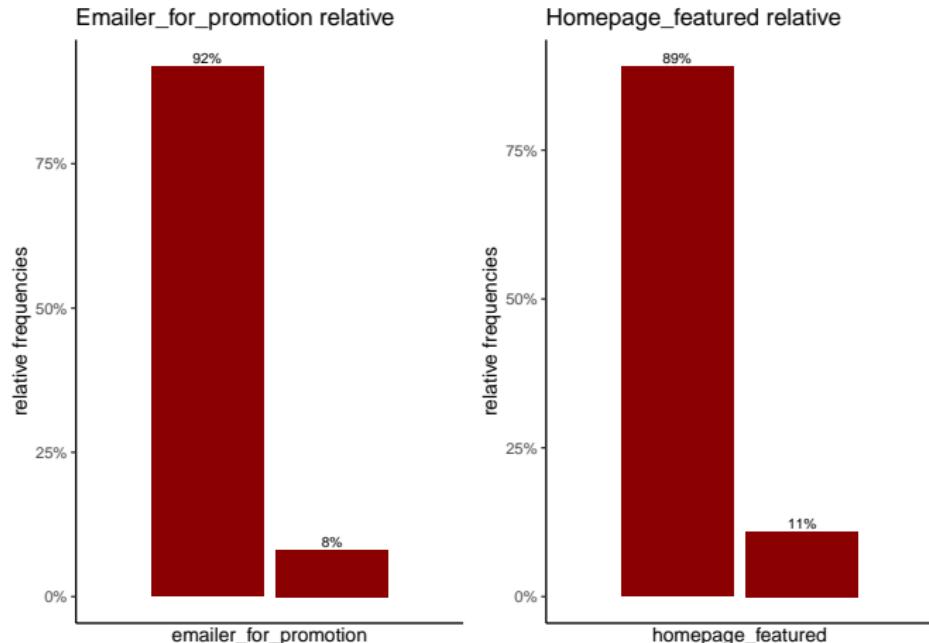
# Univariate Analysis on categorical variables



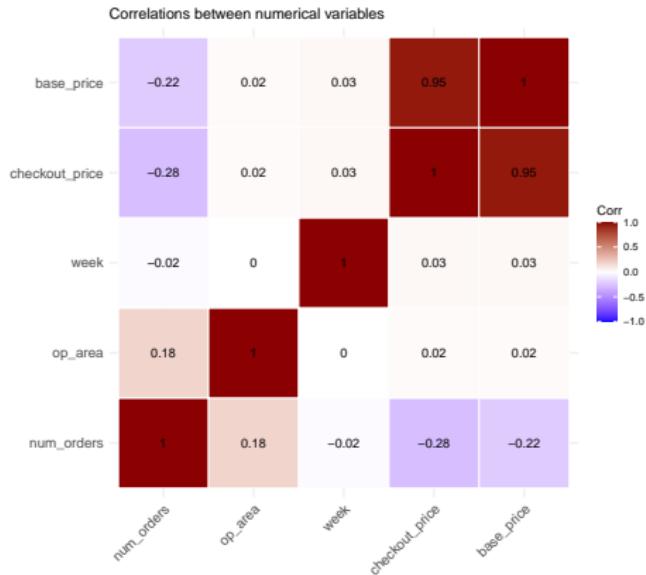
# Univariate Analysis on categorical variables



# Univariate Analysis on categorical variables

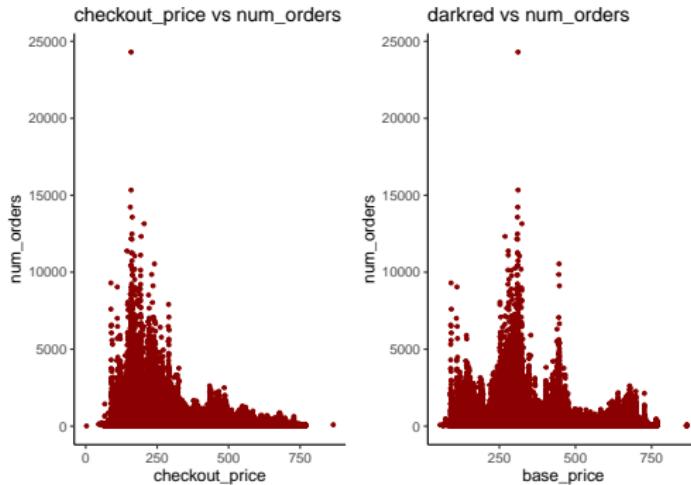


# Multivariate Analysis, correlation plot



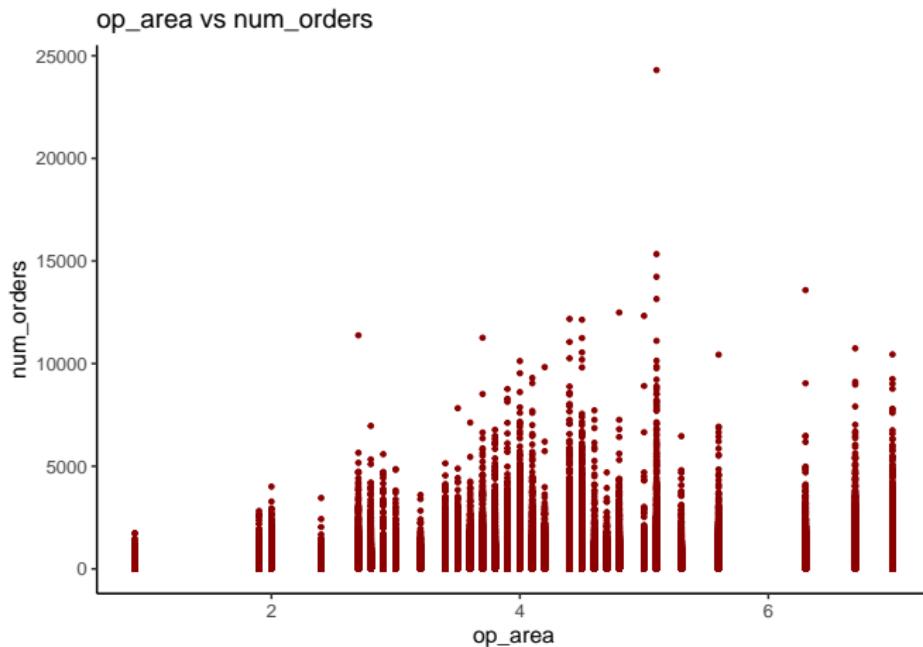
num\_orders is not highly correlated with other numerical variables.

# Multivariate Analysis, numerical variables



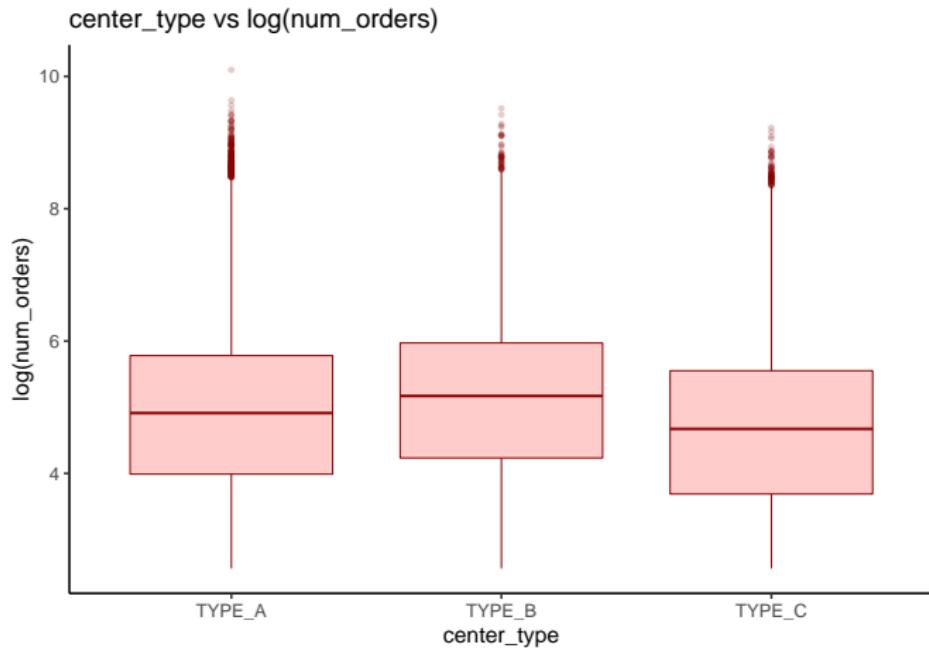
We decided to remove `checkout_price` due to the nature of the variable, in a real case scenario we can't have a checkout price because checkout means that the order is confirmed.

## Multivariate Analysis, numerical variables

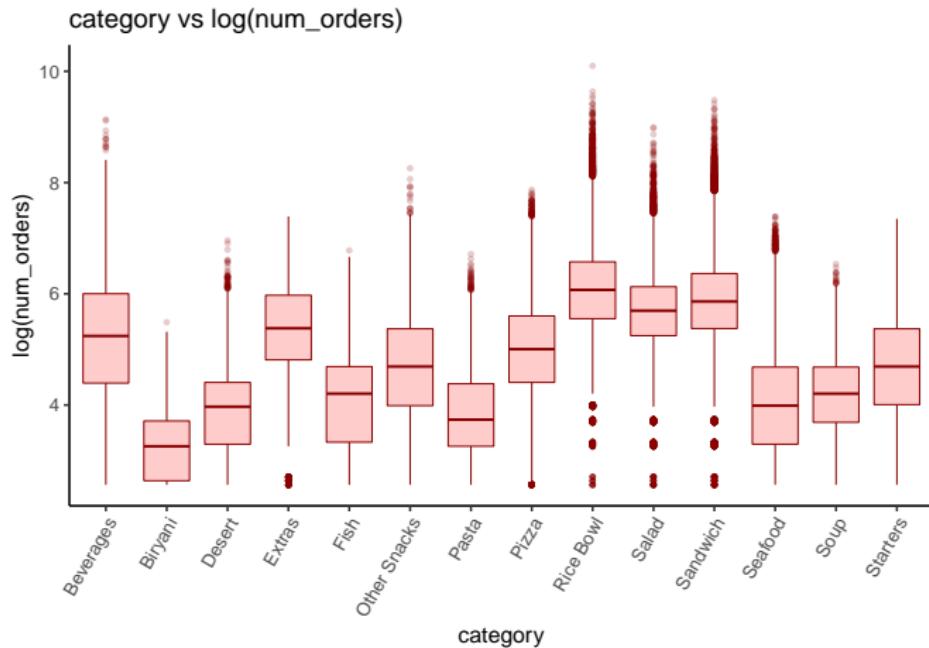


num\_orders seems to increase for bigger op\_area centers

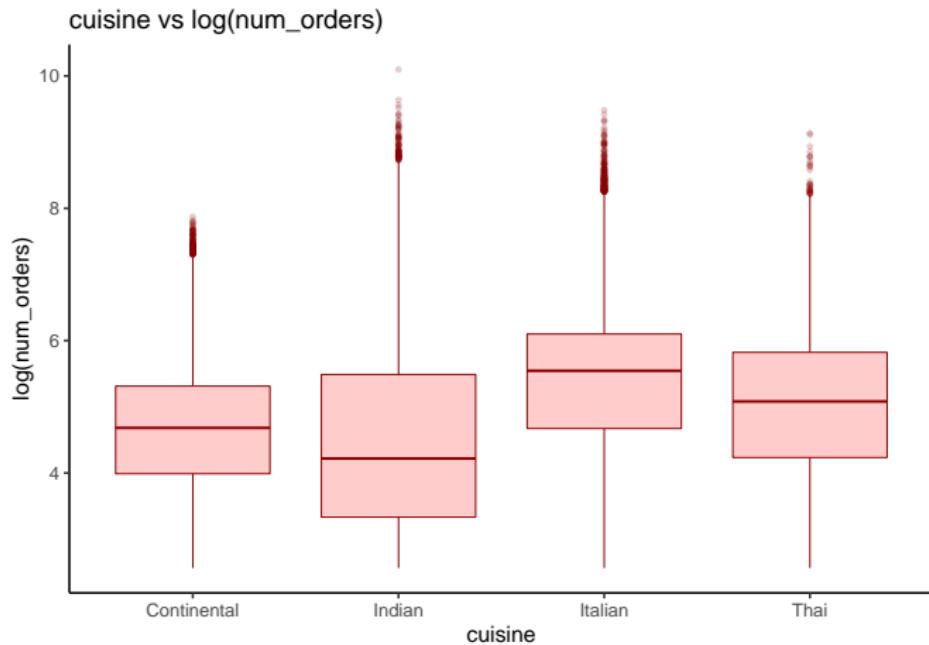
# Multivariate Analysis, categorical variables



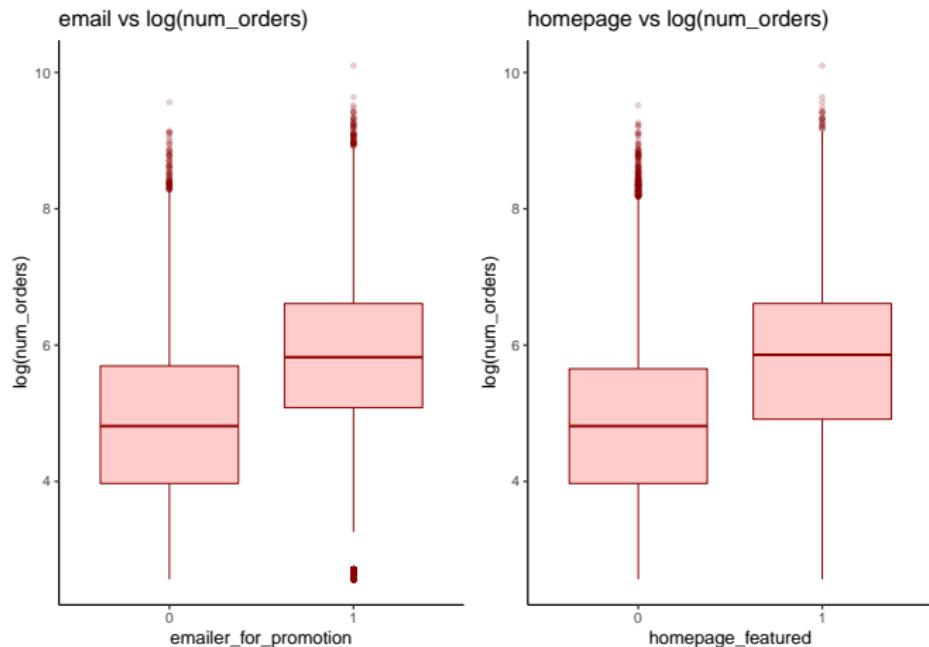
# Multivariate Analysis, categorical variables



# Multivariate Analysis, categorical variables



# Multivariate Analysis, categorical variables



Promotions ad email seem to increase the orders

## Modelling

Since we want to organize the goods for each specific fulfillment center, we need to forecast the demand for each specific center. Moreover, we also need to stratify for each unique meal, since each of them requires a different set of raw materials. We propose a two-stage approach:

- First we account for the temporal relationship using the linear model, obtaining (hopefully) i.i.d. residuals
- Then we model the obtained residuals, using some flexible method such as the gradient boosting

## Linear model

We want to fit a straight line, between demand and time, for each combination of center and meal. This mean we should fit  $N^o\text{centers} \cdot N^o\text{meals}$  ( $77 \cdot 51 = 3927$ ) linear models. But if we carefully craft some indicator variables we can specify all the simple linear models in to one single big linear model.

## Linear model

$$Y_{ij} = \beta_{0ij} + \beta_{1ij} week$$
$$\forall i = 1, \dots, 77; j = 1, \dots, 51$$

Is equivalent to:

$$Y = \beta_0 + \beta_1 week + X_{ind}\beta_{level} + X_{ind}\beta_{slope} \cdot week$$

## Linear model

where  $X_{ind}$  is a vector with  $51 \cdot 77 - 1 = 3926$  columns, and is obtained as the interaction between the dummy expansion of the categorical variables center\_id and meal\_id.

The model has  $1 + 1 + 3926 + 3926 = 7854$  scalar parameters, that in the simple formulation there are 2 parameters for each model, so  $2 \cdot 77 \cdot 51 = 7854$

# Results