# Rossmann Store Sales Prediction

Diandra MELO

# SUMMARY

**1**

**ROSSMANN**

**drugstore chain** with more than

# 3.000 stores

located in Europe

## BUSINESS PROBLEM

Rossmann store managers are tasked with **predicting their daily sales** for up to **six weeks in advance**

## BUSINESS UNDERSTANDING

CFO wants to **invest individually in each store**, according to their incomes

# CHALLENGES

budget definition for
**individual investment of stores**

**PROBLEM**

Inaccuracy of general sales prediction due to divergence between stores

All sales prediction is manually made by each Rossmann store

**Use of a Machine Learning method to make all stores sales predictions**

Only possible to check prediction results through a computer

**Make prediction results available through a smartphone**

# 3 METHODOLOGY

# Data Description

Sales data for **1.115** Rossmann stores

**Number of Columns: 18**

**Number of Rows: 1017209**

**DATE**

Date of sales, holidays

**STORE**

Type, assortment

**CUSTOMER**

Quantity

**PROMOTIONS**

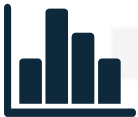Period, consecutive promo

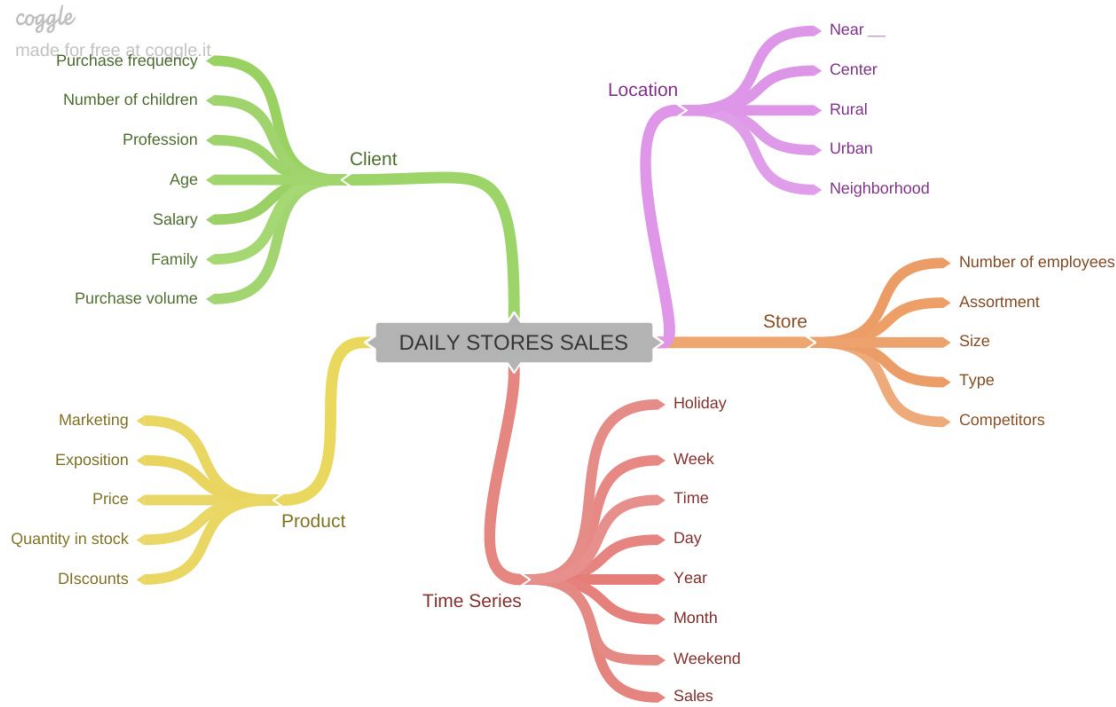**COMPETITION**

Distance, open since

# Descriptive Statistics

| | mean | median | std | min | max | range | skew | kurtosis |
|---|---|---|---|---|---|---|---|---|
| store | 558.429727 | 558.0 | 321.908651 | 1.0 | 1115.0 | 1114.0 | -0.000955 | -1.200524 |
| day_of_week | 3.998341 | 4.0 | 1.997391 | 1.0 | 7.0 | 6.0 | 0.001593 | -1.246873 |
| sales | 5773.818972 | 5744.0 | 3849.926175 | 0.0 | 41551.0 | 41551.0 | 0.641460 | 1.778375 |
| customers | 633.145946 | 609.0 | 464.411734 | 0.0 | 7388.0 | 7388.0 | 1.598650 | 7.091773 |
| open | 0.830107 | 1.0 | 0.375539 | 0.0 | 1.0 | 1.0 | -1.758045 | 1.090723 |
| promo | 0.381515 | 0.0 | 0.485759 | 0.0 | 1.0 | 1.0 | 0.487838 | -1.762018 |
| school_holiday | 0.178647 | 0.0 | 0.383056 | 0.0 | 1.0 | 1.0 | 1.677842 | 0.815154 |
| competition_distance | 5935.442677 | 2330.0 | 12547.652996 | 20.0 | 200000.0 | 199980.0 | 10.242344 | 147.789712 |
| competition_open_since_month | 6.786849 | 7.0 | 3.311087 | 1.0 | 12.0 | 11.0 | -0.042076 | -1.232607 |
| competition_open_since_year | 2010.324840 | 2012.0 | 5.515593 | 1900.0 | 2015.0 | 115.0 | -7.235657 | 124.071304 |
| promo2 | 0.500564 | 1.0 | 0.500000 | 0.0 | 1.0 | 1.0 | -0.002255 | -1.999999 |
| promo2_since_week | 23.619033 | 22.0 | 14.310064 | 1.0 | 52.0 | 51.0 | 0.178723 | -1.184046 |
| promo2_since_year | 2012.793297 | 2013.0 | 1.662658 | 2009.0 | 2015.0 | 6.0 | -0.784436 | -0.210075 |
| is_promo | 0.165966 | 0.0 | 0.372050 | 0.0 | 1.0 | 1.0 | 1.795644 | 1.224338 |

# Hypothesis Creation

Hypothesis MindMap



coggle
made for free at coggle.it

**DAILY STORES SALES**

Client
- Purchase frequency
- Number of children
- Profession
- Age
- Salary
- Family
- Purchase volume

Location
- Near __
- Center
- Rural
- Urban
- Neighborhood

Store
- Number of employees
- Assortment
- Size
- Type
- Competitors

Product
- Marketing
- Exposition
- Price
- Quantity in stock
- DIscounts

Time Series
- Holiday
- Week
- Time
- Day
- Year
- Month
- Weekend
- Sales

# Hypothesis Creation

**List of Hypothesis**

1. Stores with **more diversified products** should sell **more**
2. Stores with **closer competitors** should sell **less**
3. Stores with **longer active promotions** should sell **more**
4. Stores with **more consecutive promotions** should sell **more**
5. Stores, during **school holidays**, should sell **more**
6. Stores, **in Christmas period**, should sell **more** than during other holidays
7. Stores, during **weekends**, should sell **less**
8. At the **first half of the month**, stores should sell **more**
9. In the **first semester of the year**, stores should sell **less**
10. **Throughout the years**, stores should sell **more**
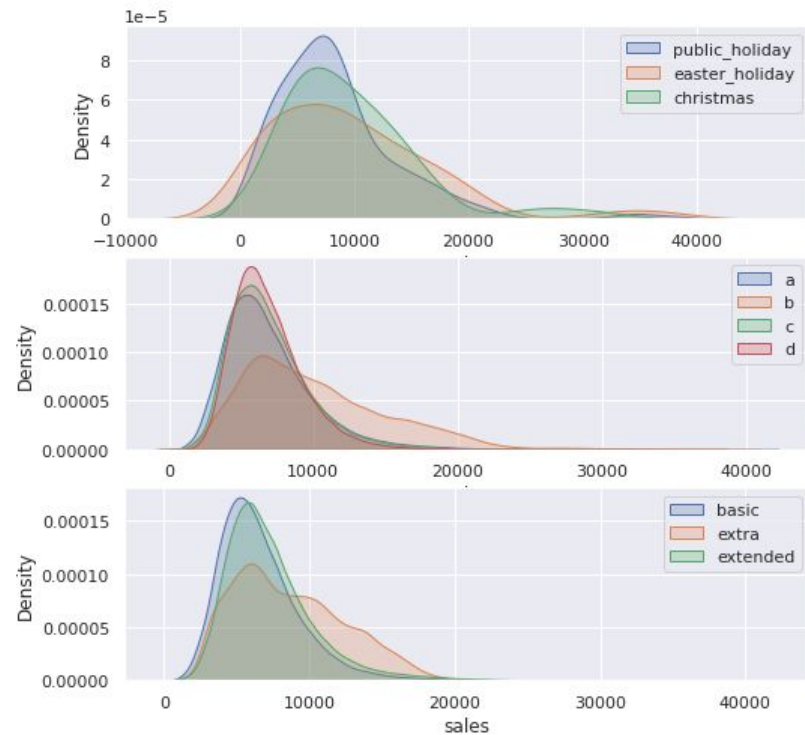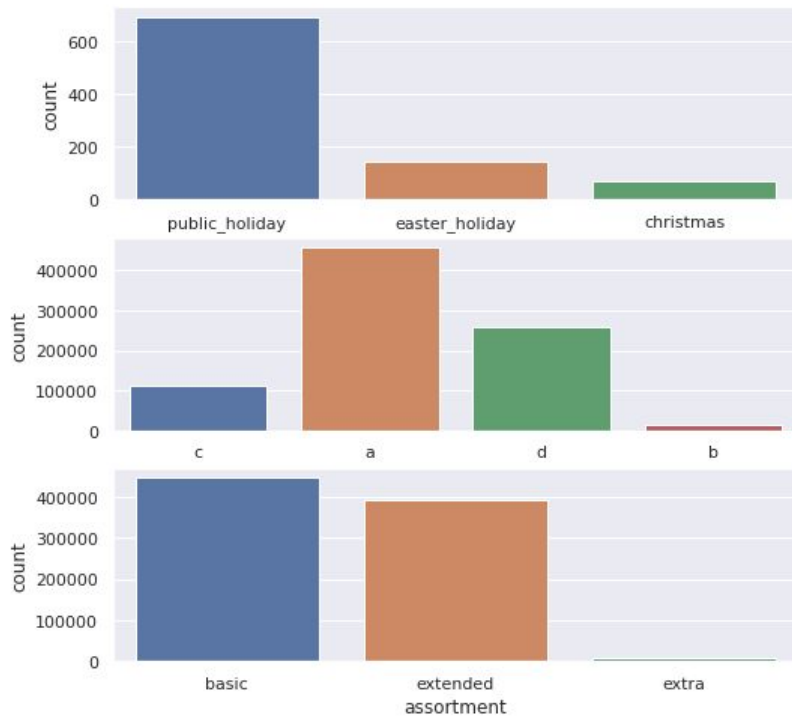
# Exploratory Data Analysis

Univariate Analysis

Rossmann Store Sales Prediction – Diandra MELO

# Exploratory Data Analysis

Univariate Analysis

# Exploratory Data Analysis
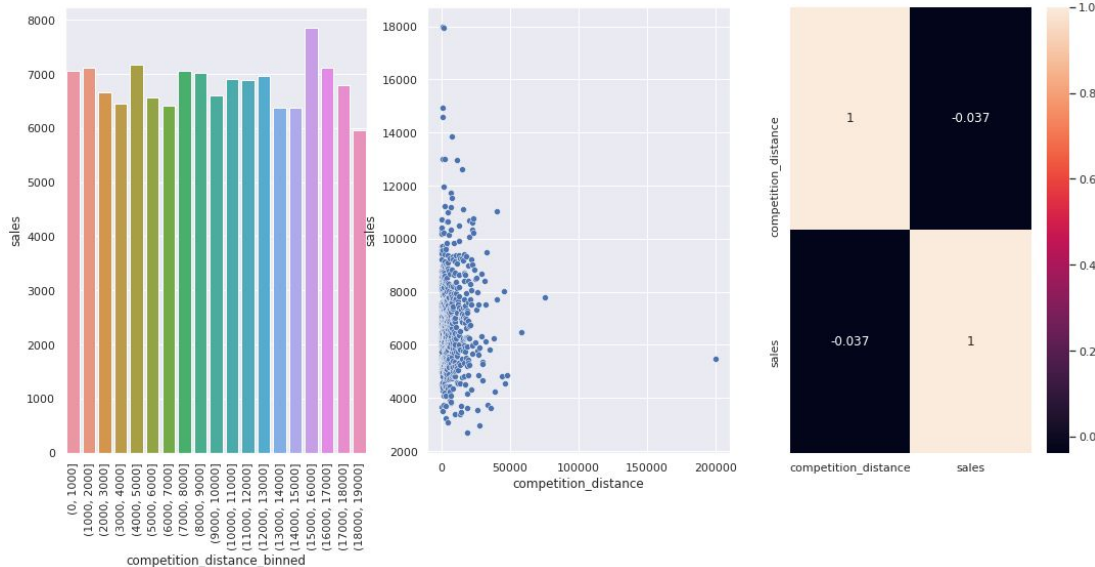
Univariate Analysis

# Exploratory Data Analysis

## Hypothesis 02

### Stores with <u>closer competitors</u> should sell <u>less</u>



This hypothesis turned to be **FALSE**.

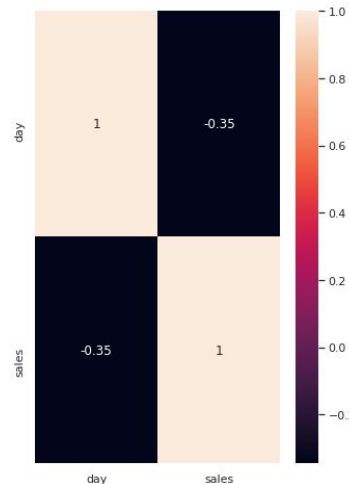Mean of sales does not vary much: in fact, the highest mean value is found for higher distances.
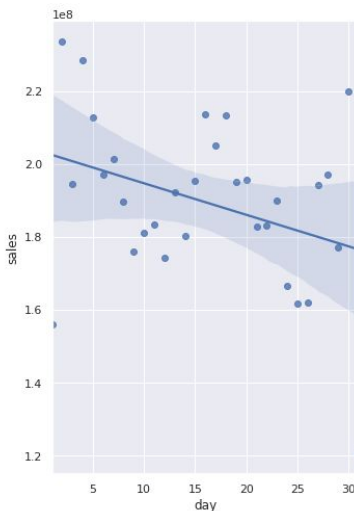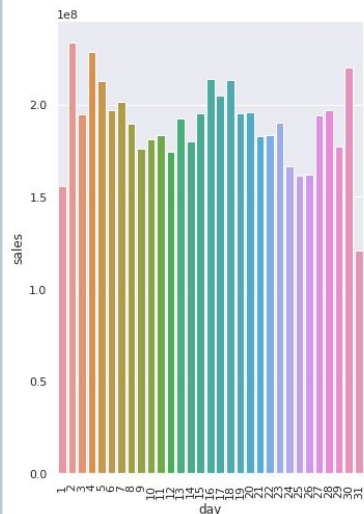
Not a strong correlation (heat map).

# Exploratory Data Analysis

**Hypothesis 08**

## At the first half of the month, stores should sell more
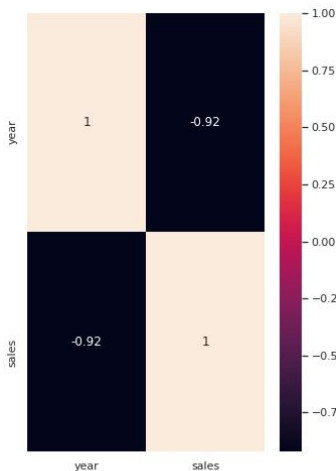


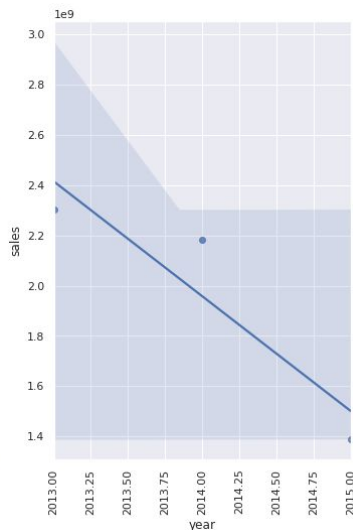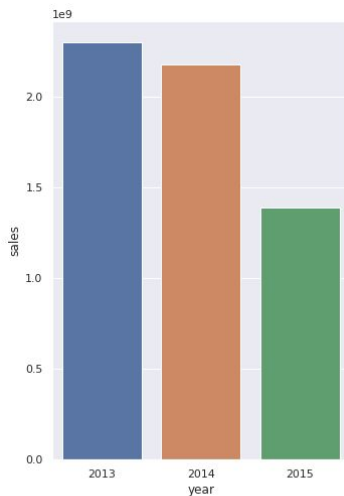This hypothesis turned to be **TRUE**.

◇ The linear regression slope exhibits well this tendency.

Rossmann Store Sales Prediction – Diandra MELO

# Exploratory Data Analysis

**Throughout the years, stores should sell more**

## Hypothesis 10



This hypothesis turned to be **FALSE**.

◇ Even though data of 2015 is incomplete (until July), it was considered to observe the tendency of sales reduction throughout years.

# Exploratory Data Analysis
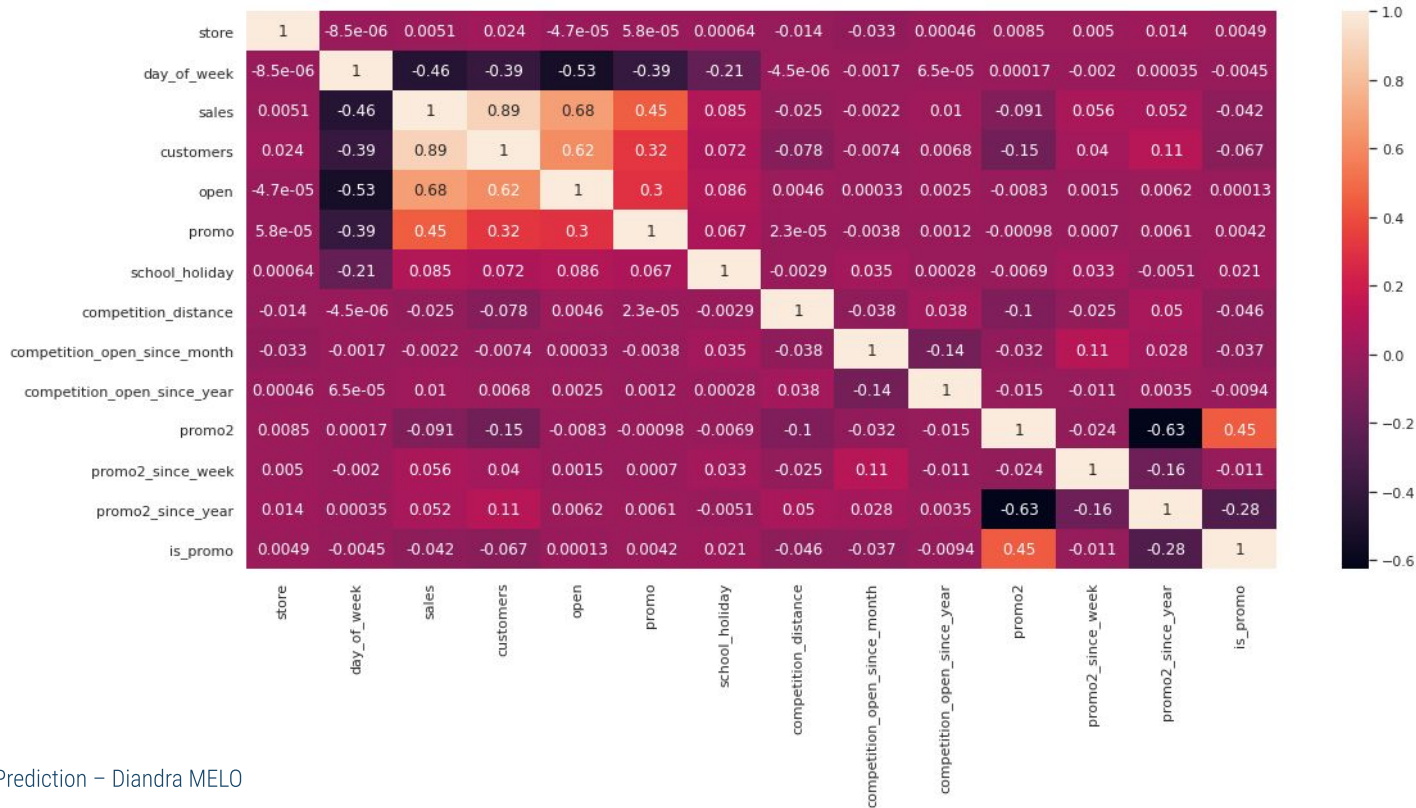
**3.4**

**Hypothesis Final Table**

| Hypothesis | Conclusion | Relevance |
|---|---|---|
| H1 – Stores with more diversified products should sell more | TRUE | Medium |
| H2 – Stores with closer competitors should sell less | FALSE | Low |
| H3 – Stores with longer active promotions should sell more | TRUE | Medium |
| H4 – Stores with more consecutive promotions should sell more | FALSE | Low |
| H5 – Stores, during school holidays, should sell more | TRUE | Low |
| H6 – Stores, in Christmas period, should sell more than during other holidays | FALSE | Medium |
| H7 – Stores, during weekends, should sell less | TRUE | High |
| H8 – At the first half of the month, stores should sell more | TRUE | High |
| H9 – In the first semester of the year, stores should sell less | TRUE | High |
| H10 – Throughout the years, stores should sell more | FALSE | High |

Rossmann Store Sales Prediction – Diandra MELO

# Exploratory Data Analysis

## Multivariate Analysis

### Numerical Attributes



Rossmann Store Sales Prediction – Diandra MELO

# Exploratory Data Analysis

**Multivariate Analysis**

**Categorical Attributes**



Rossmann Store Sales Prediction – Diandra MELO

# Machine Learning Modelling

## Used Models

- Linear Regression
- Random Forest
- Lasso
- XGBoost

## Model's Performance

| Model Name | MAE CV | MAPE CV | RMSE CV |
|---|---|---|---|
| Linear Regression | 2081.73 +/- 295.63 | 0.3 +/- 0.02 | 2952.52 +/- 468.37 |
| Lasso | 2116.38 +/- 341.5 | 0.29 +/- 0.01 | 3057.75 +/- 504.26 |
| RF Regressor | 837.68 +/- 219.1 | 0.12 +/- 0.02 | 1256.08 +/- 320.36 |
| XGBoost Regressor | 1030.28 +/- 167.19 | 0.14 +/- 0.02 | 1478.26 +/- 229.79 |

### MODEL'S CHOICE

**XGBoost Regressor**

◇ Lighter than RF Regressor, with not much difference in error

# 4 CONCLUSION

# Business Performance

## 285,860,497.77 $

Total Stores Sales Prediction for the next 6 weeks

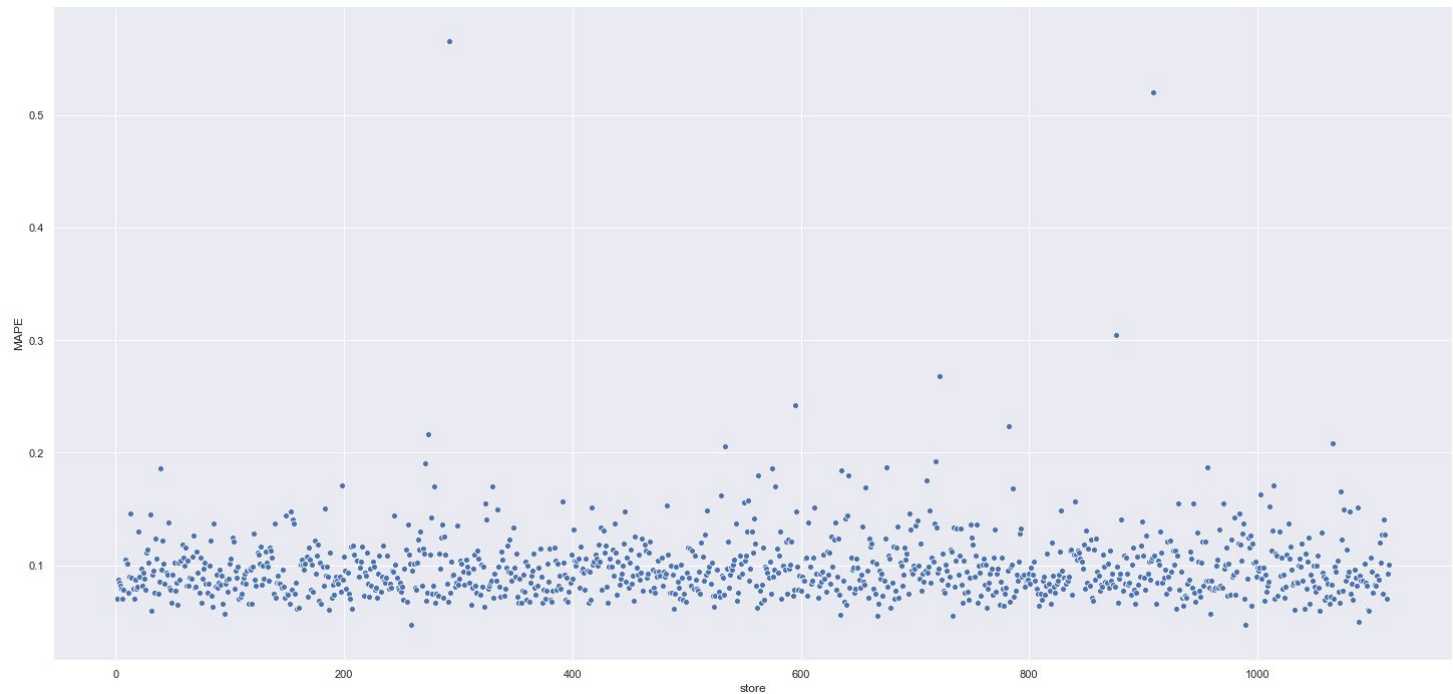## ± 745,482.07 $

Margin of error (Worst/Best Scenario)

## < 10 %
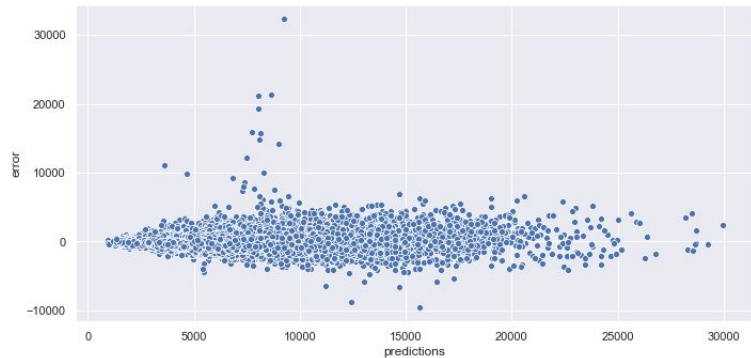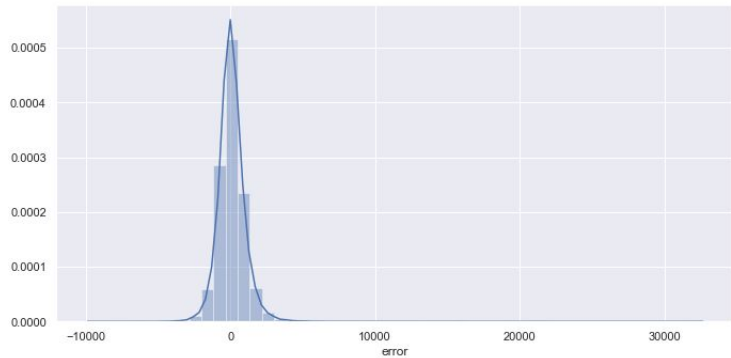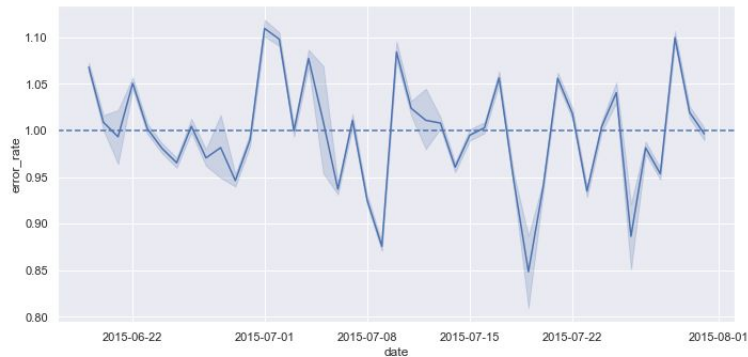
Average percentage error of the stores sales

# Business Performance

# Machine Learning Performance



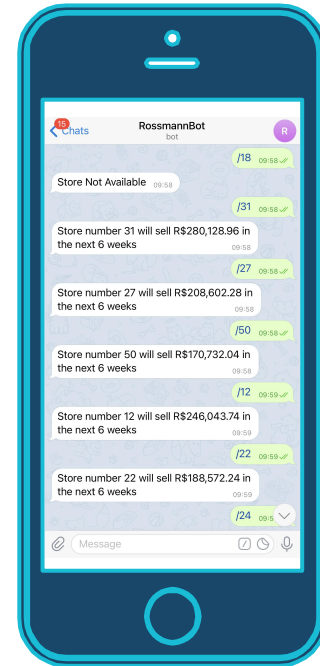Rossmann Store Sales Prediction – Diandra MELO

# DEMONSTRATION

## Model Deployment - Heroku

◈ Telegram Bot

# 6. Next Steps

◇   Improve Model's performance;

◇   Test usability of deployment;

◇   Check for other ML models that could fit well into this study case.

Rossmann Store Sales Prediction – Diandra MELO

# THANKS!

## Any questions?

You can find me at:

◇   LinkedIn: Diandra Melo
◇   GitHub: @diandramelo
◇   Gmail: dcxsmelo@gmail.com

# Credits

I would like to give an special thanks to **Meigarom Lopes** for providing the orientation needed for me to achieve these results, as well as for improving my Data Science knowledge throughout the course **Data Science in Production**.

◇ Forecast sales competition *Rossmann Store Sales* from kaggle.com

◇ Presentation template by SlidesCarnival