

Sales Forecasting for Promotional Flyers





Agenda

- 1 Introduction
- 2 Data Preprocessing
- 3 Feature Engineering
- 4 Model exploration
- 5 Final approach
- 6 Discussions



1

Introduction

Our goal & datasets



Predict the sales units for each product featured in promotional flyers from November to December 2023.



Anagrafica_volantini,
Risultati_prodotti_volantini,
Storico_quantità
Gfk_caratteristiche,
KPI_prodotti volantini.
Icecat_caratteristiche



2

Preprocessing

Data manipulation



Prezzo_listino



Recovered the missing prices by averaging the prices of products with the same art_code, for the remaining NaNs, we used historical revenue divided by historical quantity.

Prezzo_promo



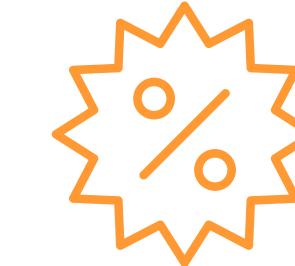
For the missing values, it was calculated using the percentage discount and the list price.

Stock_pz



All negative values have been replaced with 0.

Sconto_perc



Measured the absolute value for the negative ones, while multiplying by 100 all those less than one. For the missing values, we calculated the discount using the promotional price and the list price.

Focus on gfk_caratteristiche



Initially, we did not know how to use this dataset because it contained a large number of categorical variables.



We noticed that for the product group, the type of characteristic was similar within the same column.



We then split it by product group and performed ordinal encoding, avoiding the creation of a large number of variables.

Final DataSet



Anagrafica volantini, Risultati prodotti volantini, Storico quanità, Gfk caratteristiche and Kpi prodotti volantini got merged together on “art_cod” and “codice volantino” and then split for product group

Resulting in 5 different datasets, one for each product group:

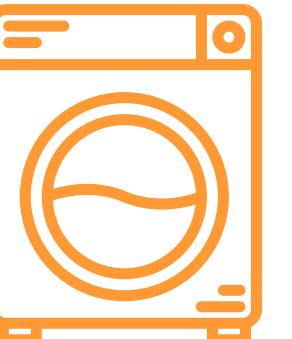
Smartphone



PC



Wash



TV



Core_wear



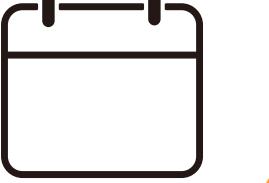


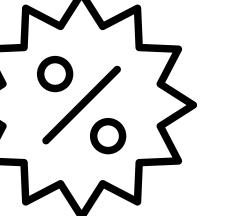
3

Feature engineering

New variables



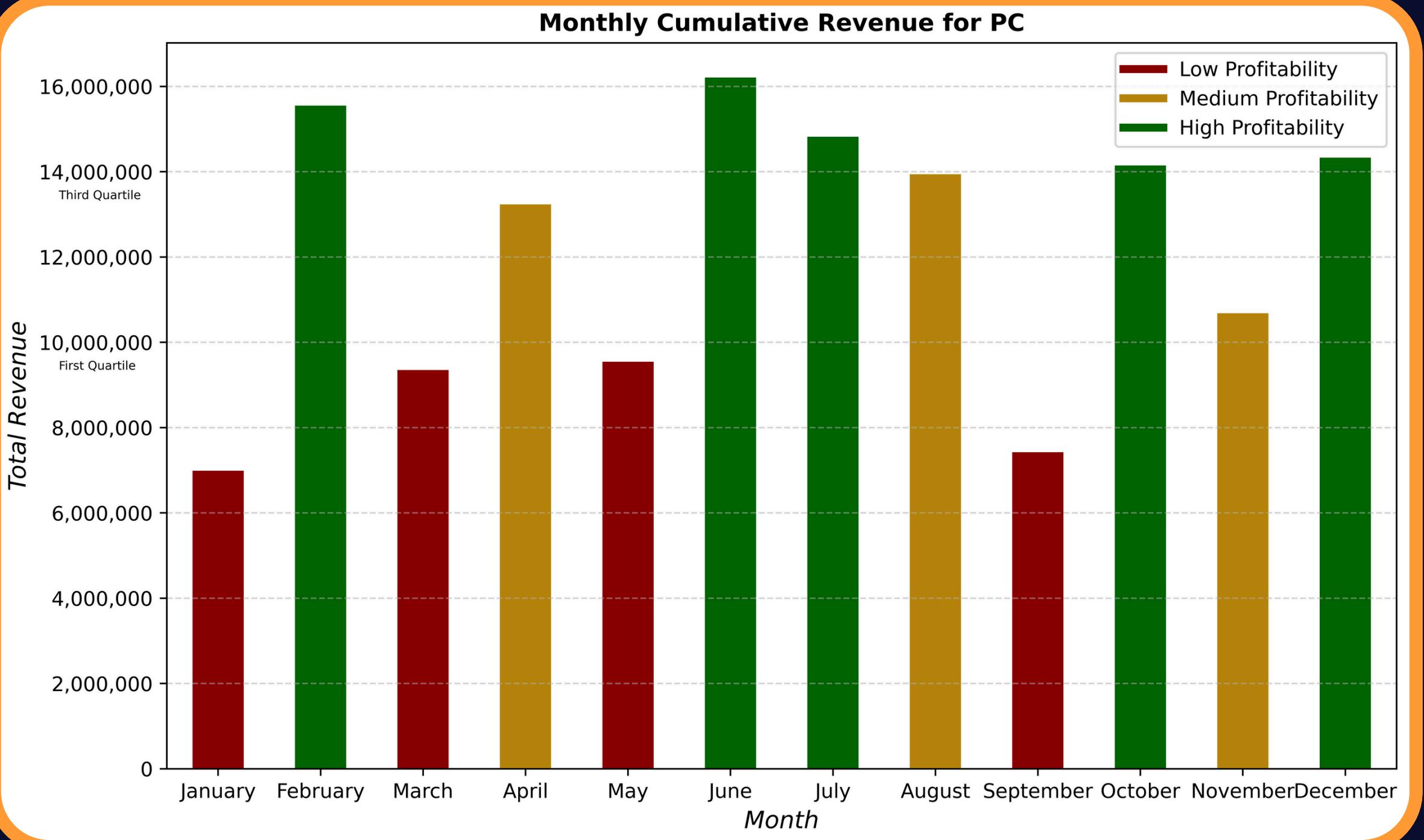
 Flyer duration
in days

 Average
discount by
campaign
name

 Average
Monthly
Revenue

 Average
discount per
flyer

Average Monthly Revenue



- 1 Computed as the **monthly cumulative revenue** of each product group
- 2 Then grouped the months into **profitability classes**: low, medium, and high.

Average Discount by Campaign Name



- 1 We grouped the **campaign names** into different macro areas (e.g., all campaigns regarding ***Black Friday*, *Scontissimi*, *Tasso 0***, etc.)
- 2 Computed the **average discount** related to each group of campaigns
- 3 Finally, the column “**Campaign Name**” was dropped



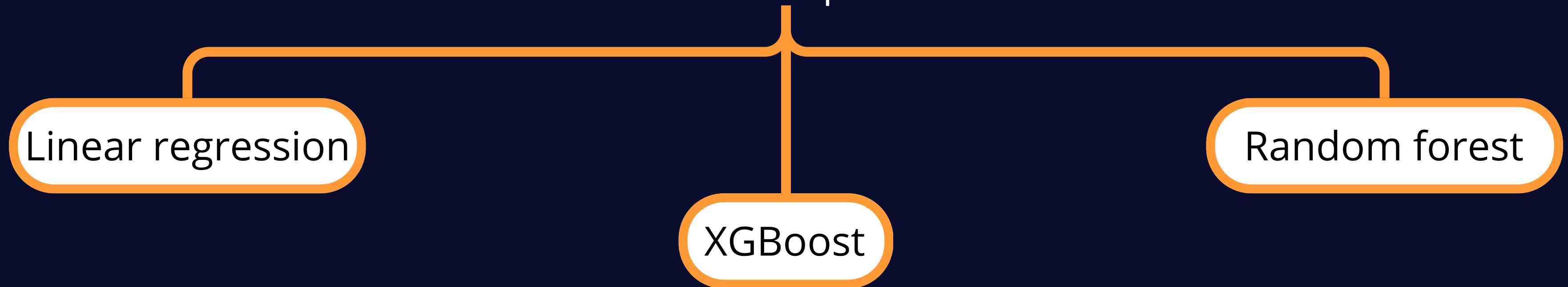
4

Model exploration

Our approach



We decided to implement 3 different models and a benchmark
to evaluate their performance



Our benchmark consists of multiplying the average daily sales per
product by the days the flyer is online

Work phases



The work can be divided into 3 different phases. These phases differ from one another based on the train and test split.

1

5-folds cross validation using **the entire dataset**

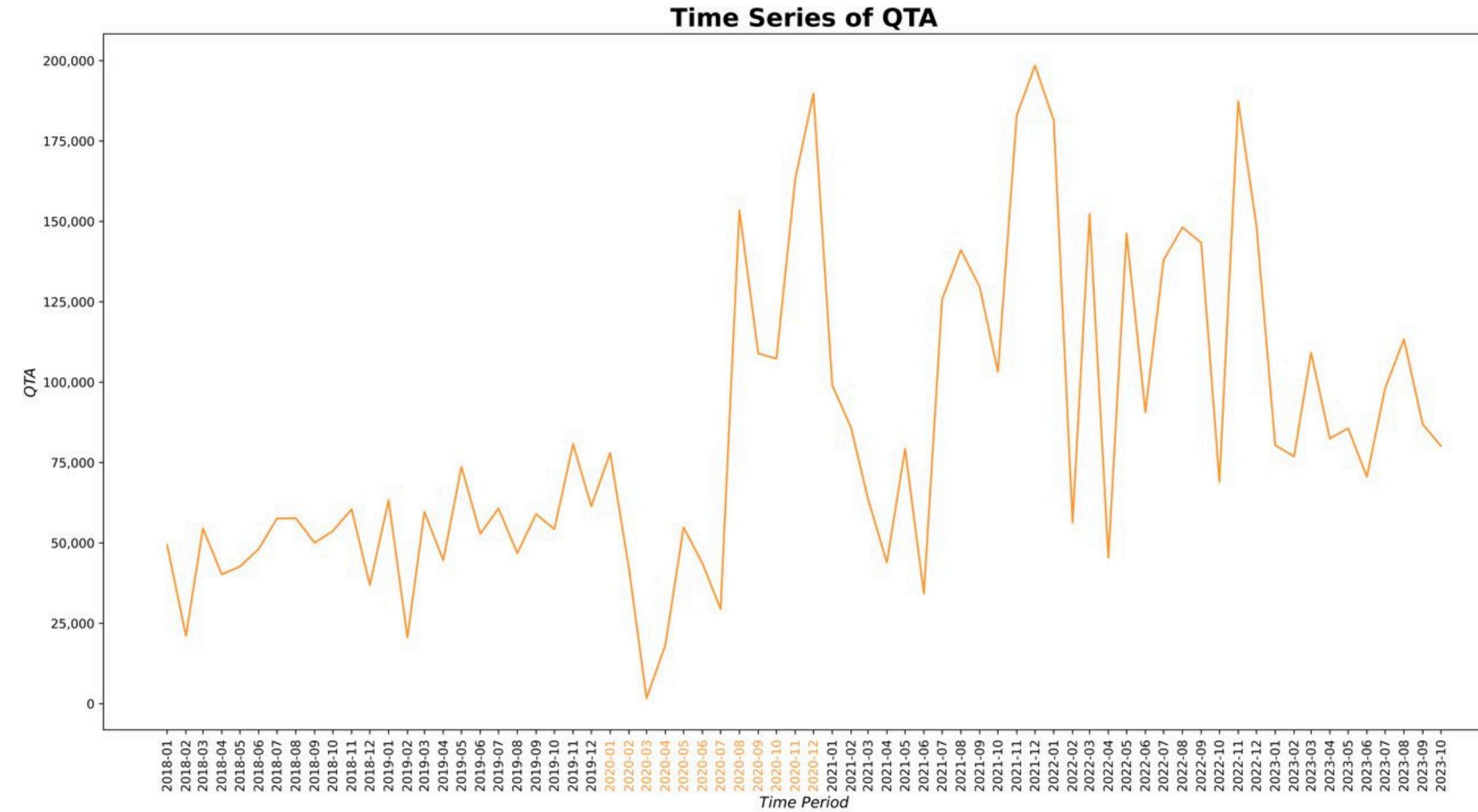
2

5-folds cross validation **excluding 2020**

3

The **entire dataset** was included in the training set, **except data from 2020**, and used as test set **August, September and October 2023**

What about 2020?



The year **2020** was **removed** from our analysis under the assumption that since it was the period of the **COVID-19** pandemic, sales would not follow a normal trend, thus **impacting the performance of our models**.

This particular trend can also be seen from the graph above.



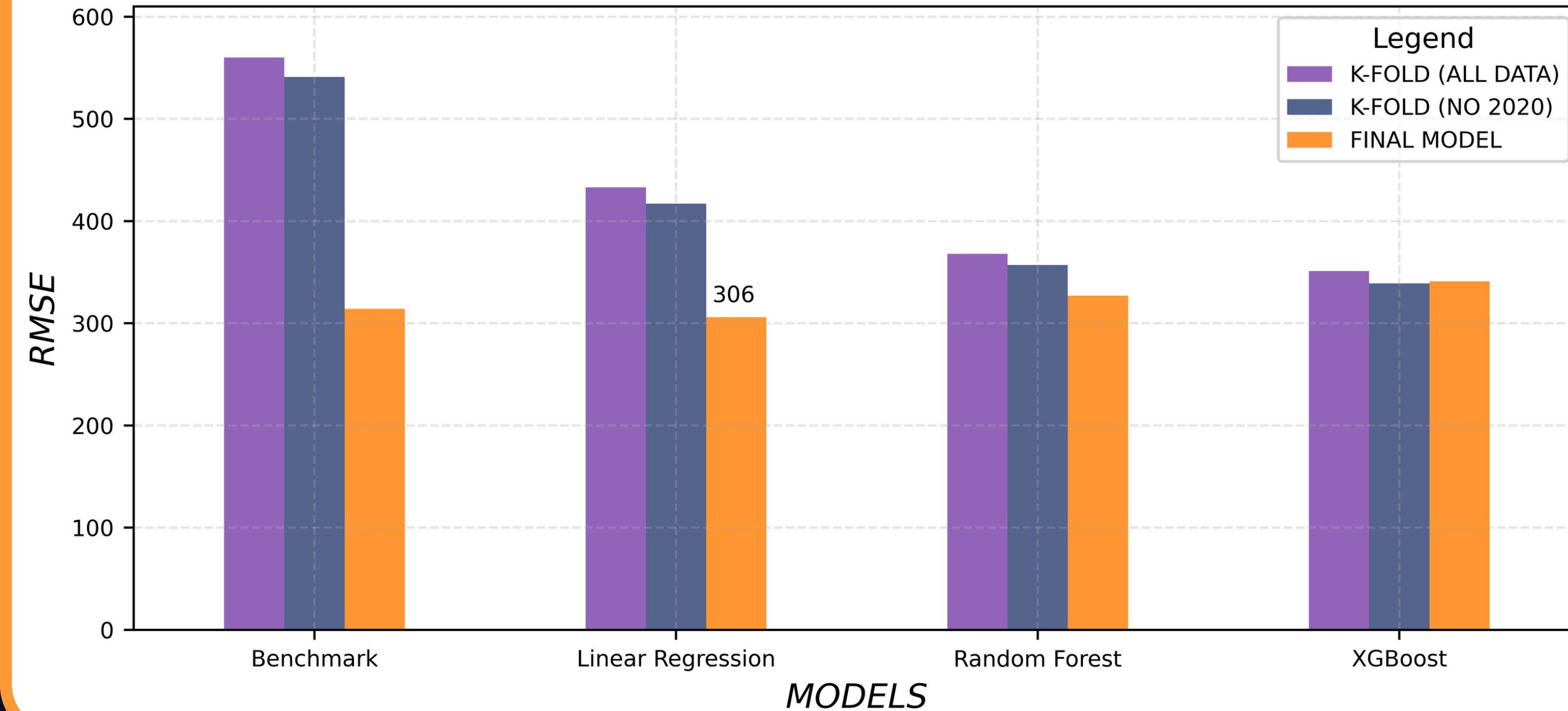
5

Final models

Smartphone



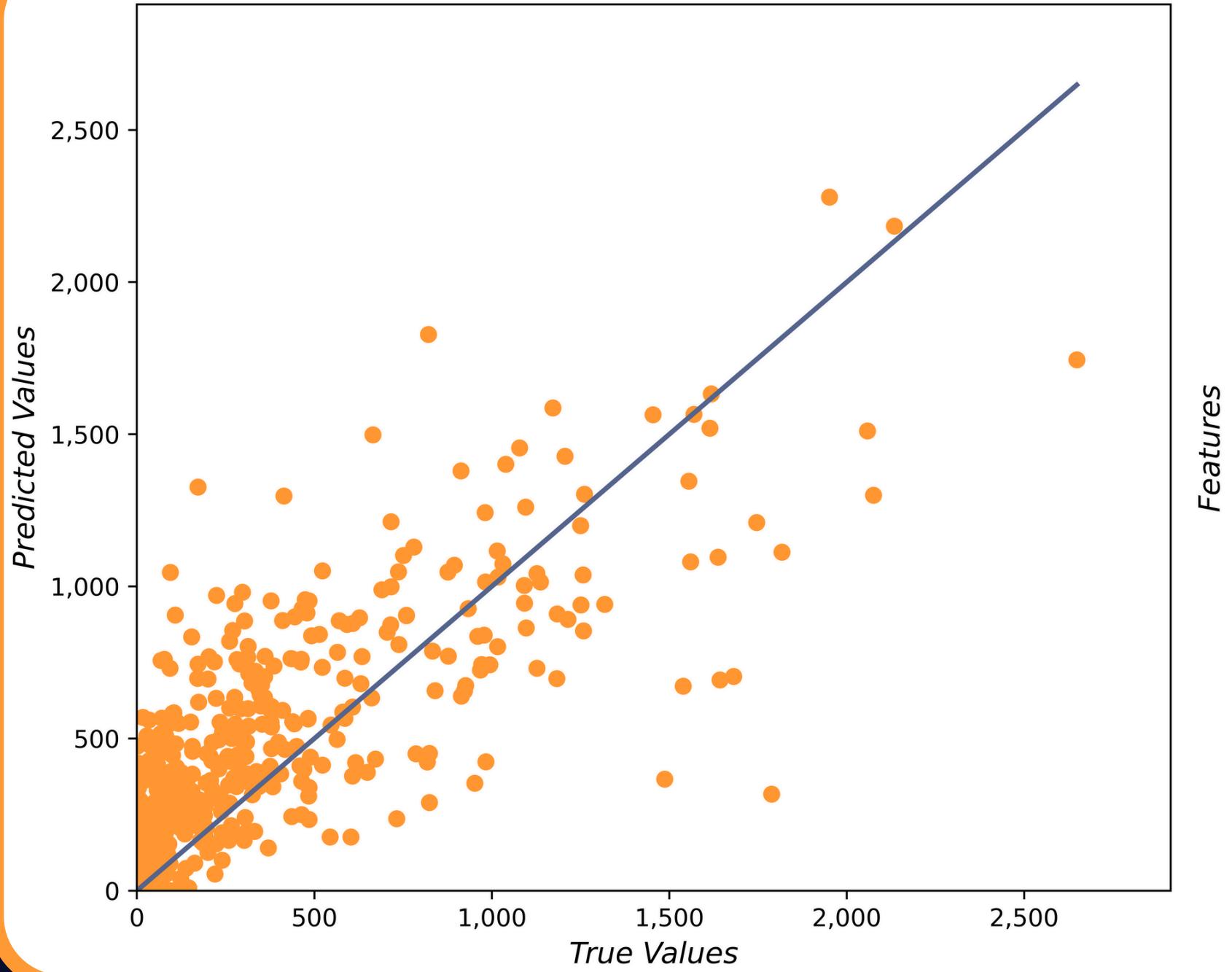
RMSE comparison among models for Smartphone



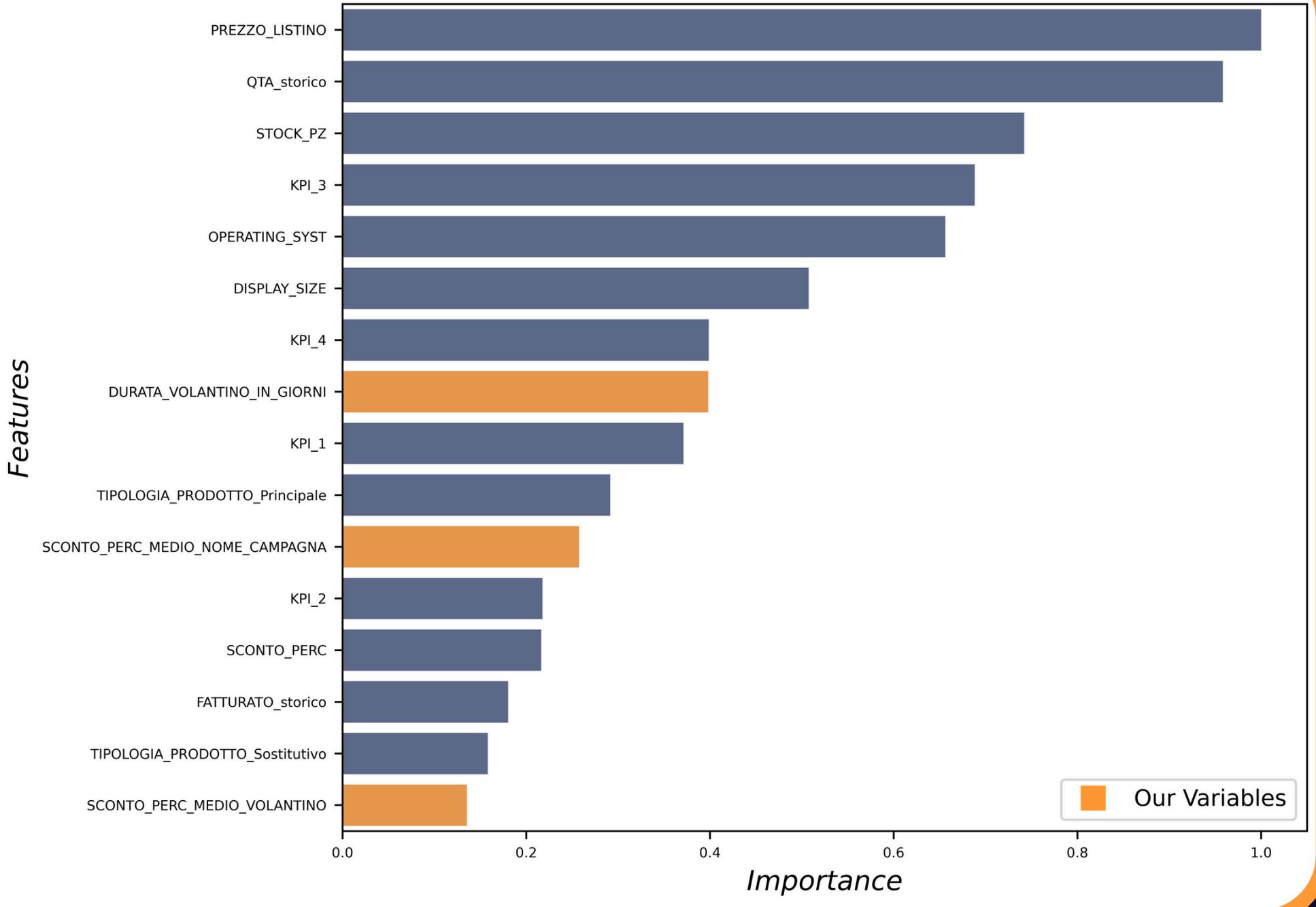
Smartphone - Linear Regression



Predictions



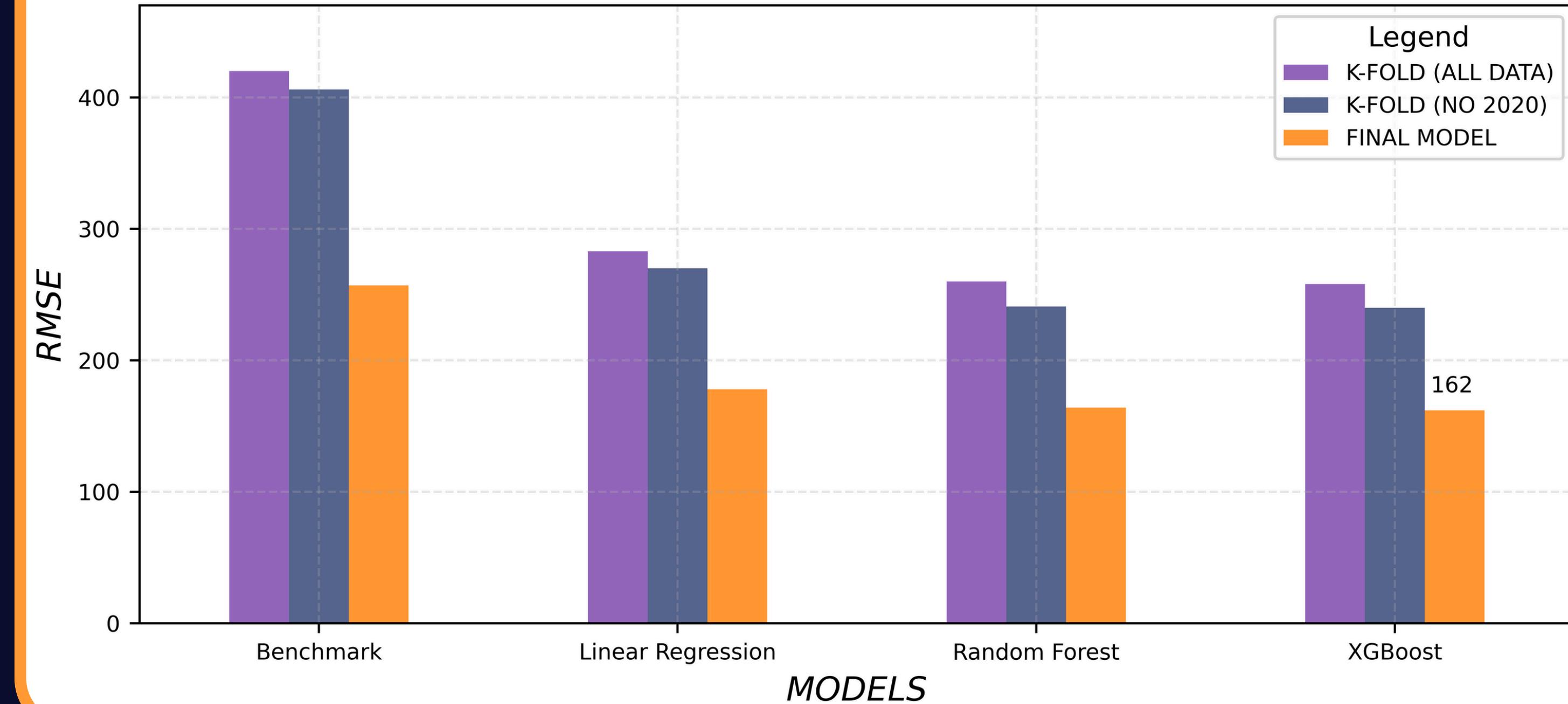
Feature Importance



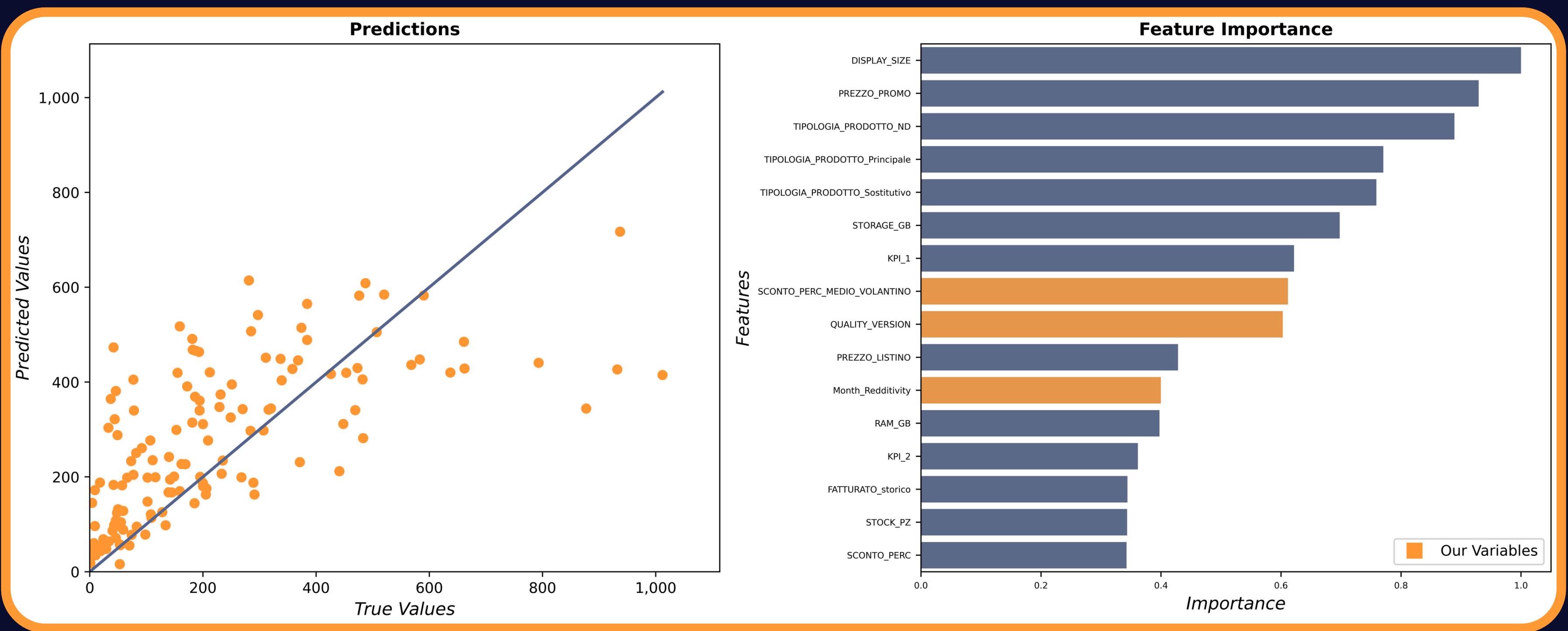
PC



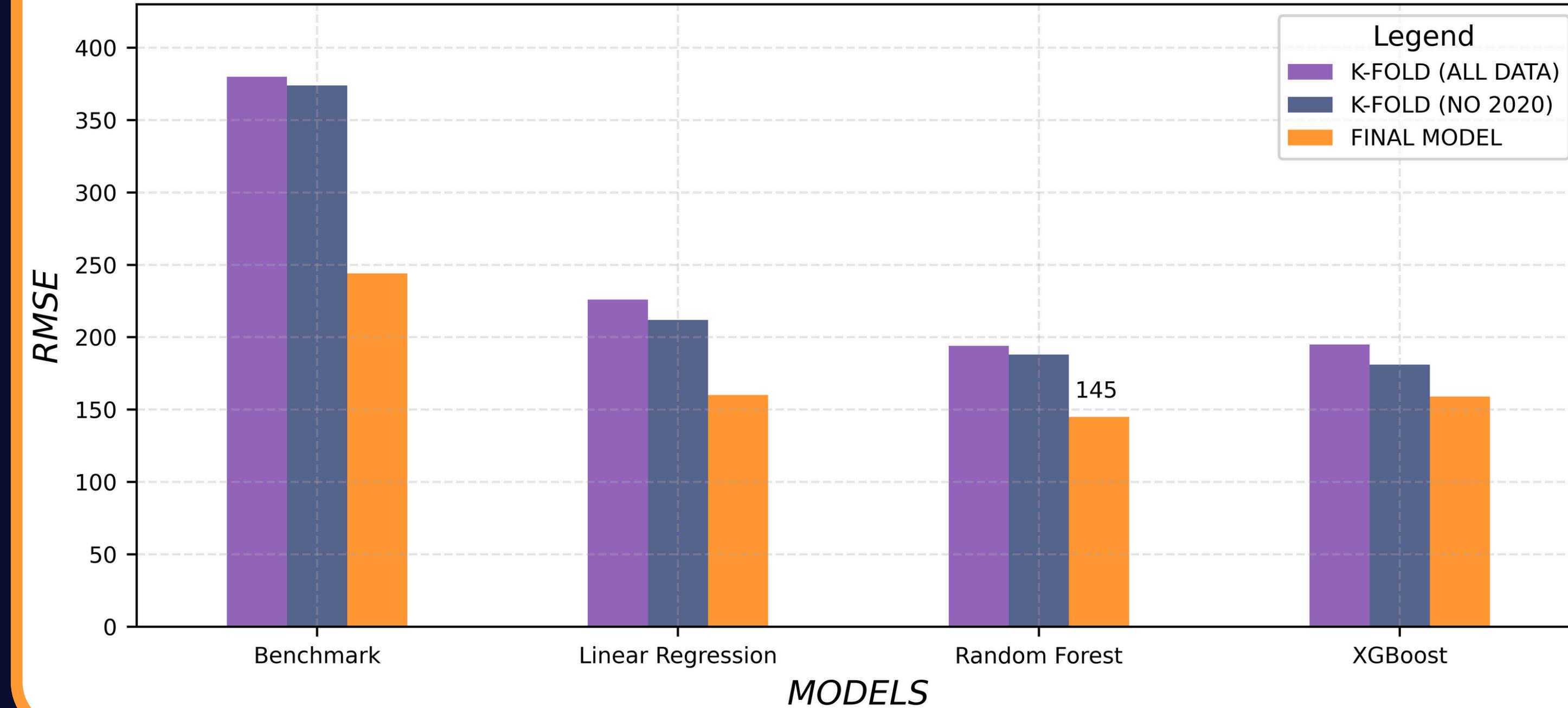
RMSE comparison among models for PC



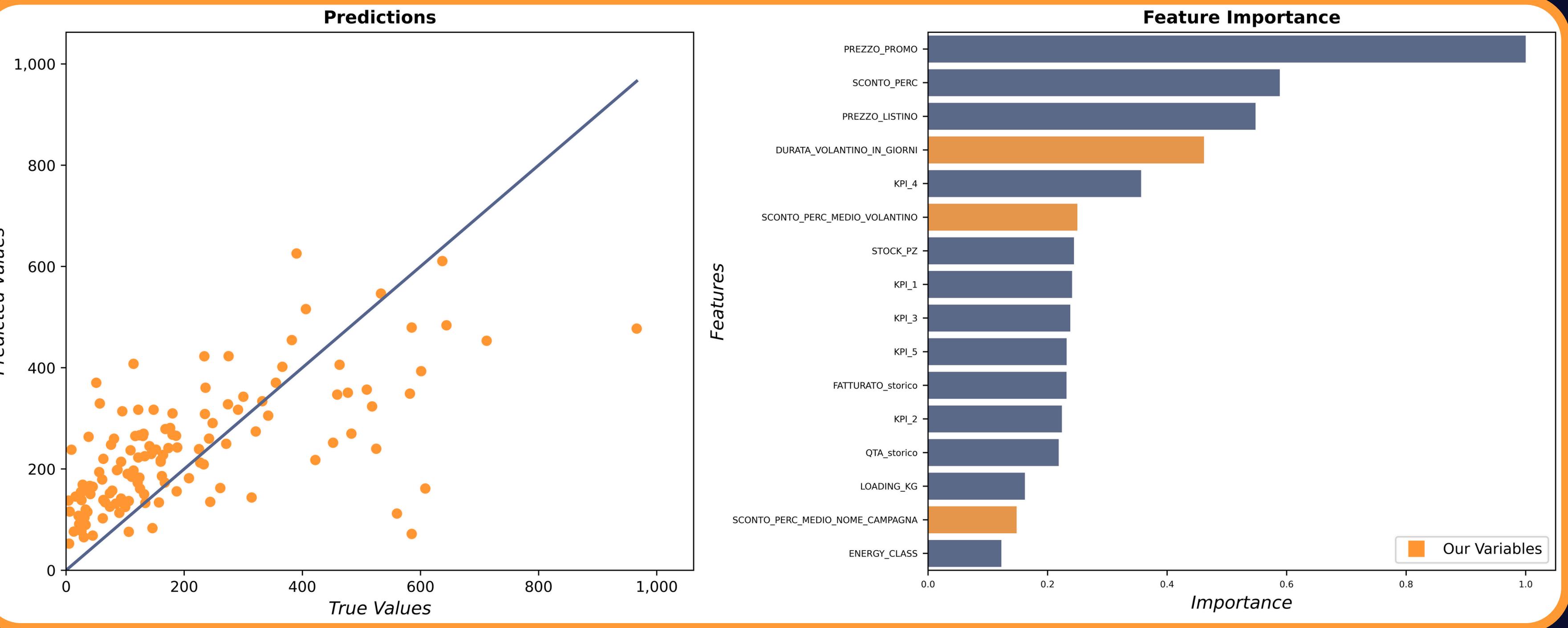
PC - XGBoost



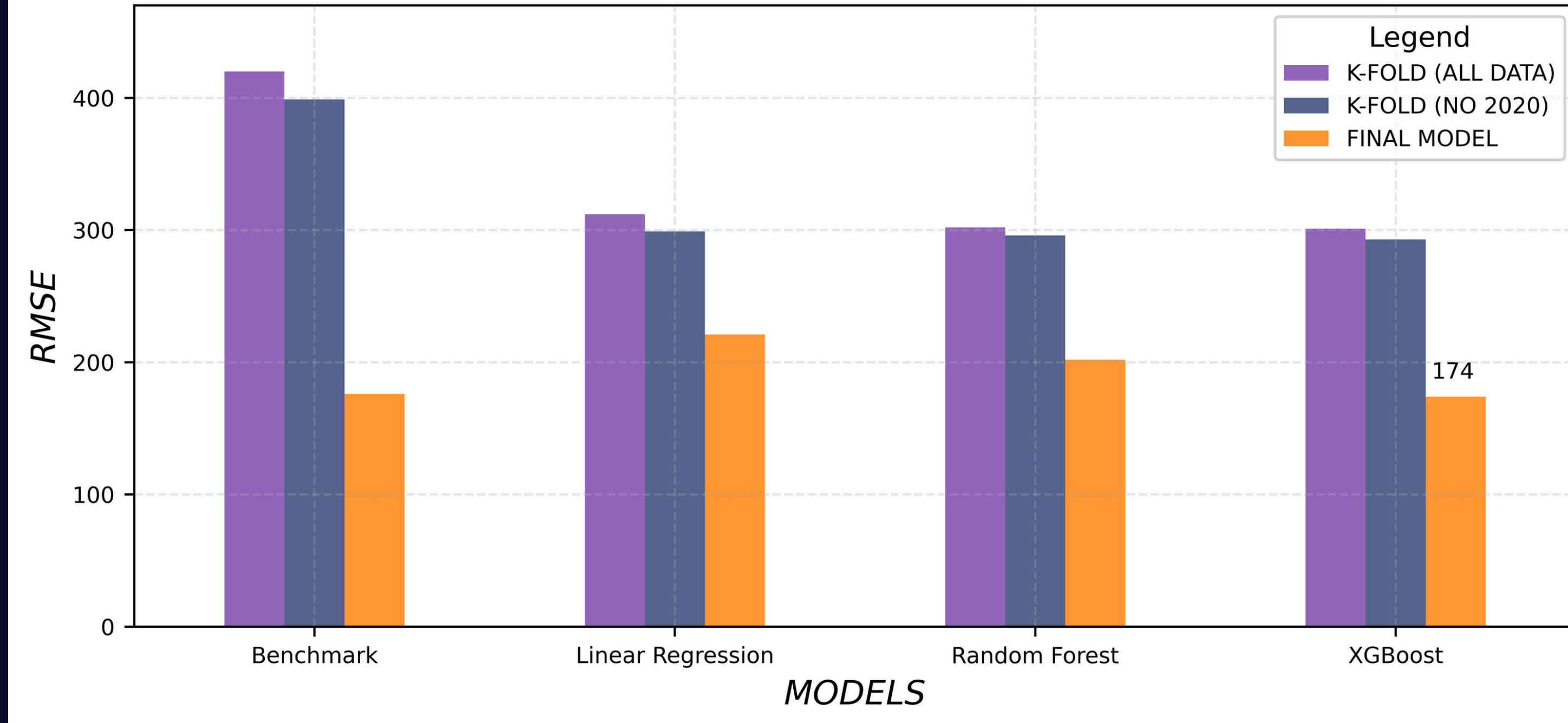
RMSE comparison among models for Wash



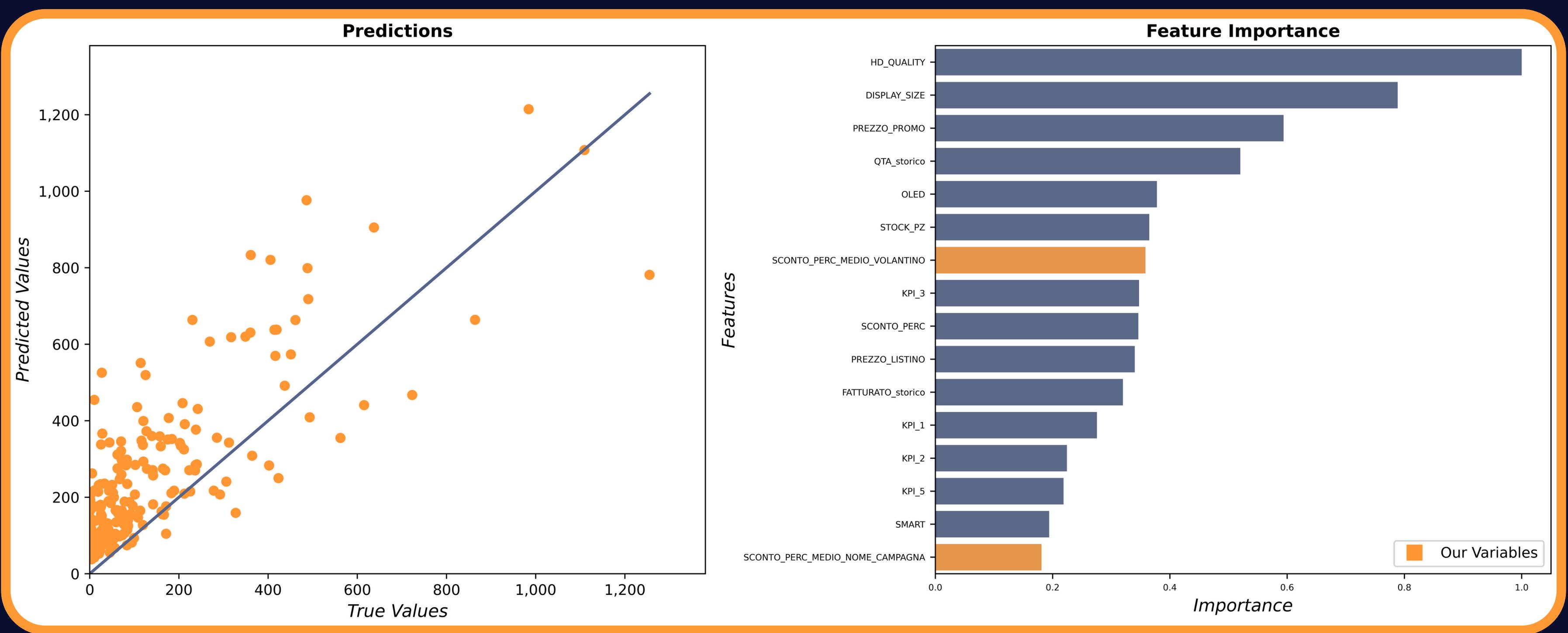
Wash - Random Forest



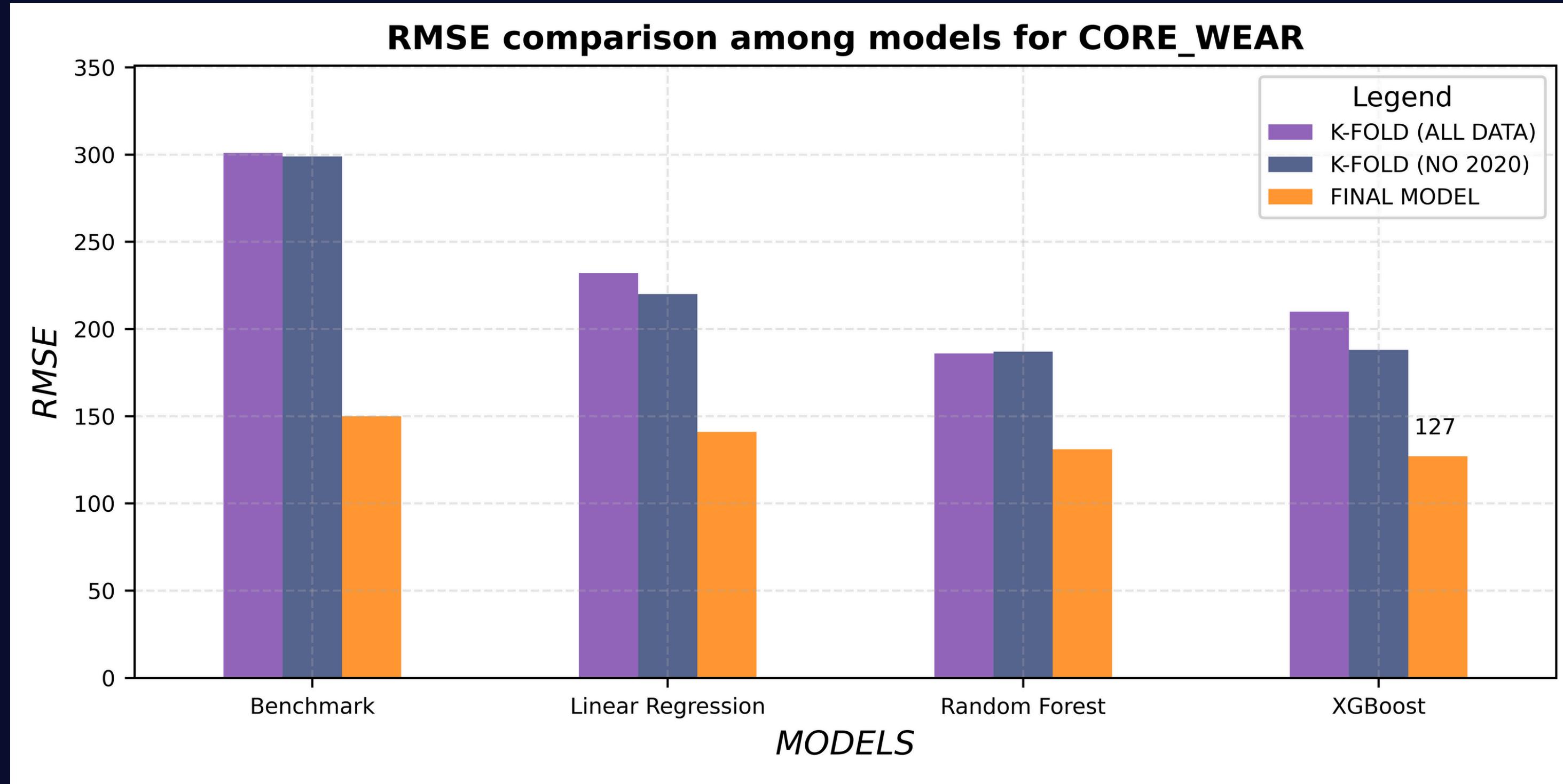
RMSE comparison among models for TV



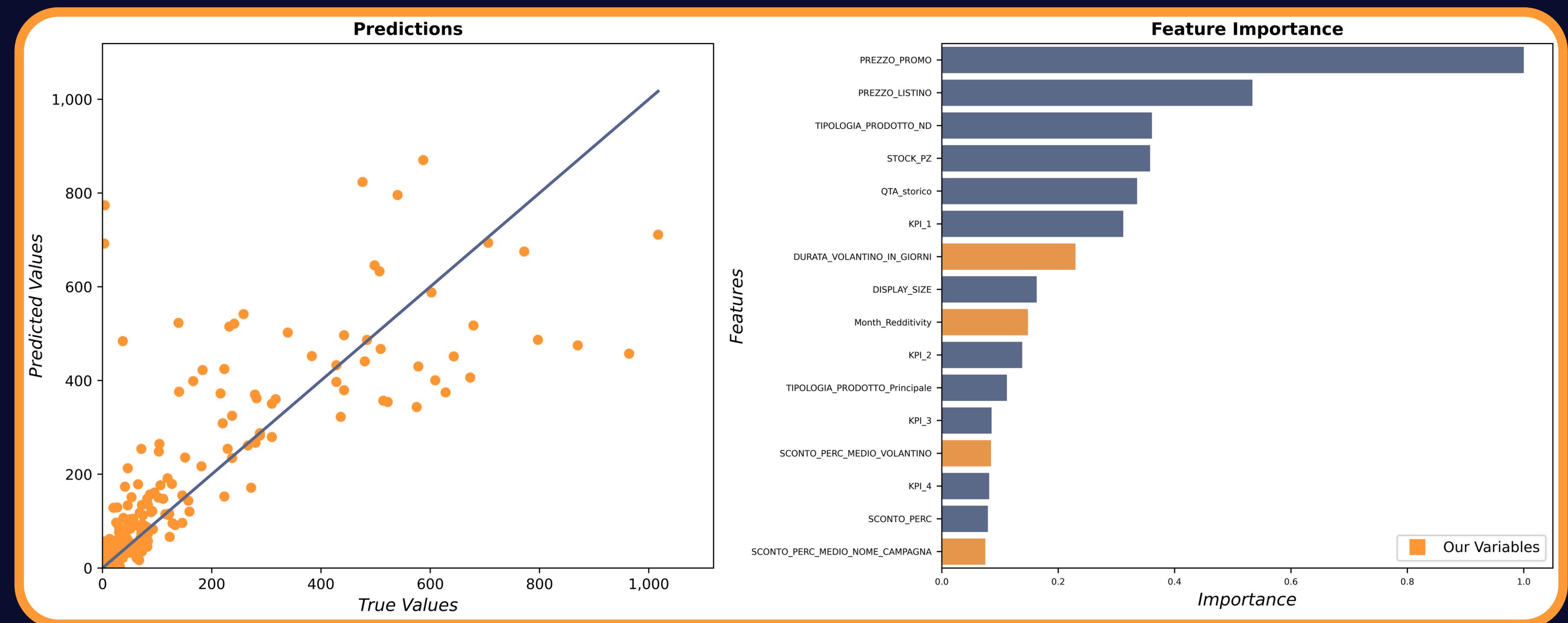
TV - XGBoost



Core Wearables



Core Wearables - XGBoost





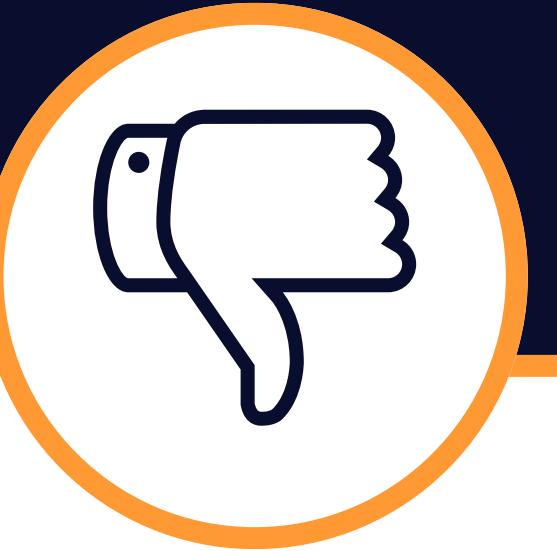
6

Discussion

Discussion



Preprocessing;
Feature engineering;
Computational time;
Good prediction for
core wear, wash, TV
and PC;



Hyperparametertuning;
Implementation of new
variables;
Poor prediction for
smartphone;
Dimensionality reduction;



Thank You!
for the Attention

Vincenzo Camerlengo 773731

Riccardo Paolantoni 773691

Raffaele Torelli 775831