
CMPE 257 - Project

Herbarium and Store sales forecasting using Time Series

GROUP 14: Kashyap T, Mahesh Chandra Mareedu, Army P, Mehulkumar K, Wen-Hao Tseng

Herbarium : Identify novel plant species

- **About Data**

- Number of images: 1.05 M
- Number of classes : 15501
- Data size :160 GB (train)

- **Challenges involved**

- Cannot load data on to Kaggle RAM
- Requires high processing GPU (Kaggle provides only 36hrs).
- Classification models cannot be developed due to high dimensional target (15k+ classes).
- Data Augmentation is not possible as it requires 2X space (160GB+160GB)
- Neural network output layer becomes sparse due to high number of target variables
- Alternate approach takes 9 hrs to run each time
 - Faiss + ResNET

Transfer learning with ResNet+faiss

- Removed last layer of ResNet and used it to generate embeddings of image each image of dimension 512.
- Similarly for all the images belonging to 15k+ classes, we have composed image embeddings for 839772 training samples.
- We fed all the image embeddings to *faiss* (Facebook AI Similarity Search) model to create indexes ([Paper](#))
- The index file stores information of all the image vectors with their corresponding image ids.

```
print("Shape of image embedding",vec.shape)
print(vec[0:10])
plt.imshow(img)
```

```
Shape of image embedding (512,)
[0.48710352 0.3259773 0.8979322 0.28265658 0.3696866 0.3621263
 2.0340765 1.2006575 0.949871 0.11798392]
; <matplotlib.image.AxesImage at 0x7fe3137c1390>
```

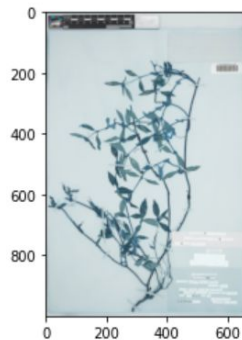
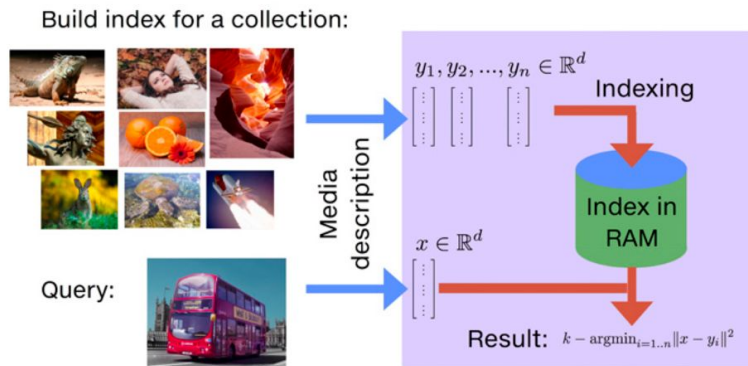


Image classification using Faiss

- Faiss clusters images embedding to P possible clusters.
- When ever we want find similar images for a given image, the module checks target image embeddings distance with the centroid of each clusters.
- The cluster centroid with minimum distance with target image is selected.
- Now, distance between the target image and images in that cluster are composed.
- The closes image ids are returned by the model.
- Runtime complexity comes down from $O(N) \Rightarrow O(P)$



Sample predictions

- As we are using non-annotated data, the cluster based classification model uses the non important features on images as well.
- The first two predictions seem right, but the remaining three have high cosine similarity due to the color block (a non important feature)



9507-Muhlenbergia(0.99)



9507-Muhlenbergia(0.99)



1081-Aristida(0.89)



12555-Rhynchospora(0.88)



6159-Euphorbia(0.88)



Sales Forecasting

- Motivation - current market is all about predicting real time growth for the business
- Objective - Forecast store sales using time-series on data from a large Ecuadorian-based grocery retailer
- Evaluation - Submission evaluated based on RMSLE (Root Mean Squared Logarithmic Error)

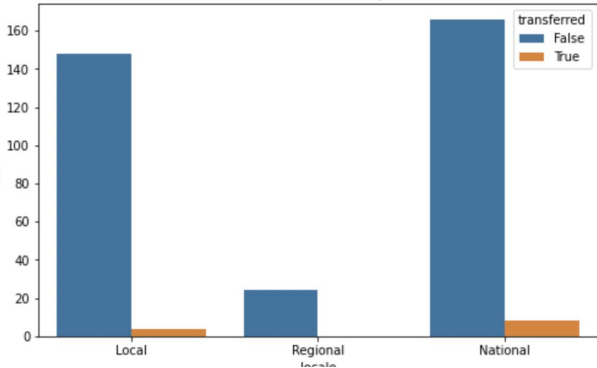


Dataset

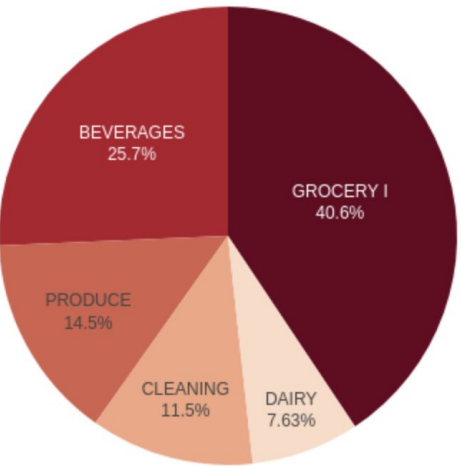
- Six files that contains information about
 - Holidays and events
 - Transactions
 - Stores
 - Oil prices over time
 - Train
 - Test
- Explored each of the dataset in deep to study the hidden patterns
- Examined the combination of each dataset with the other dataset to get more information on the data provided
- The data collected is from 2013-01-01 to 2017-08-31(combining both train and test)

Data sources, cleansing, validation, transformation, visualization.

Transferred Holidays



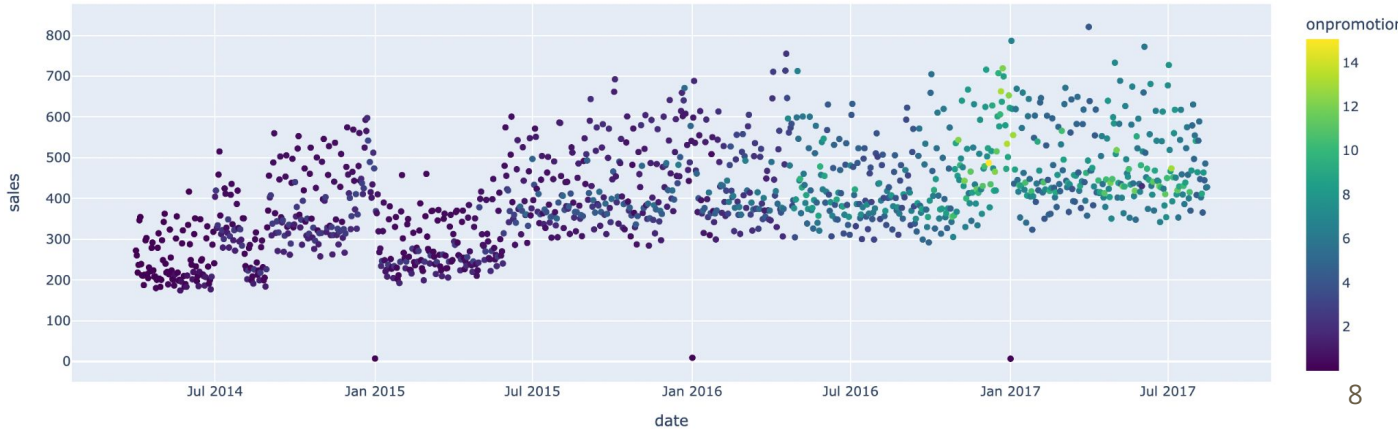
5 categories that make the most sales

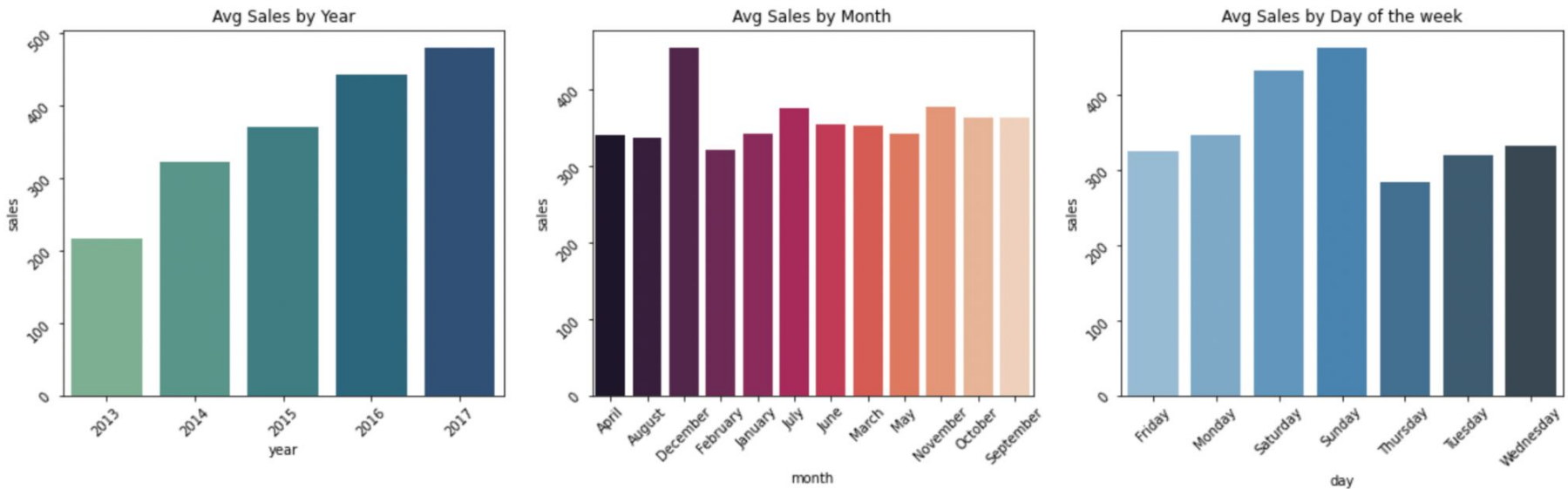


Average sales over time



Total average sales for items on promotion

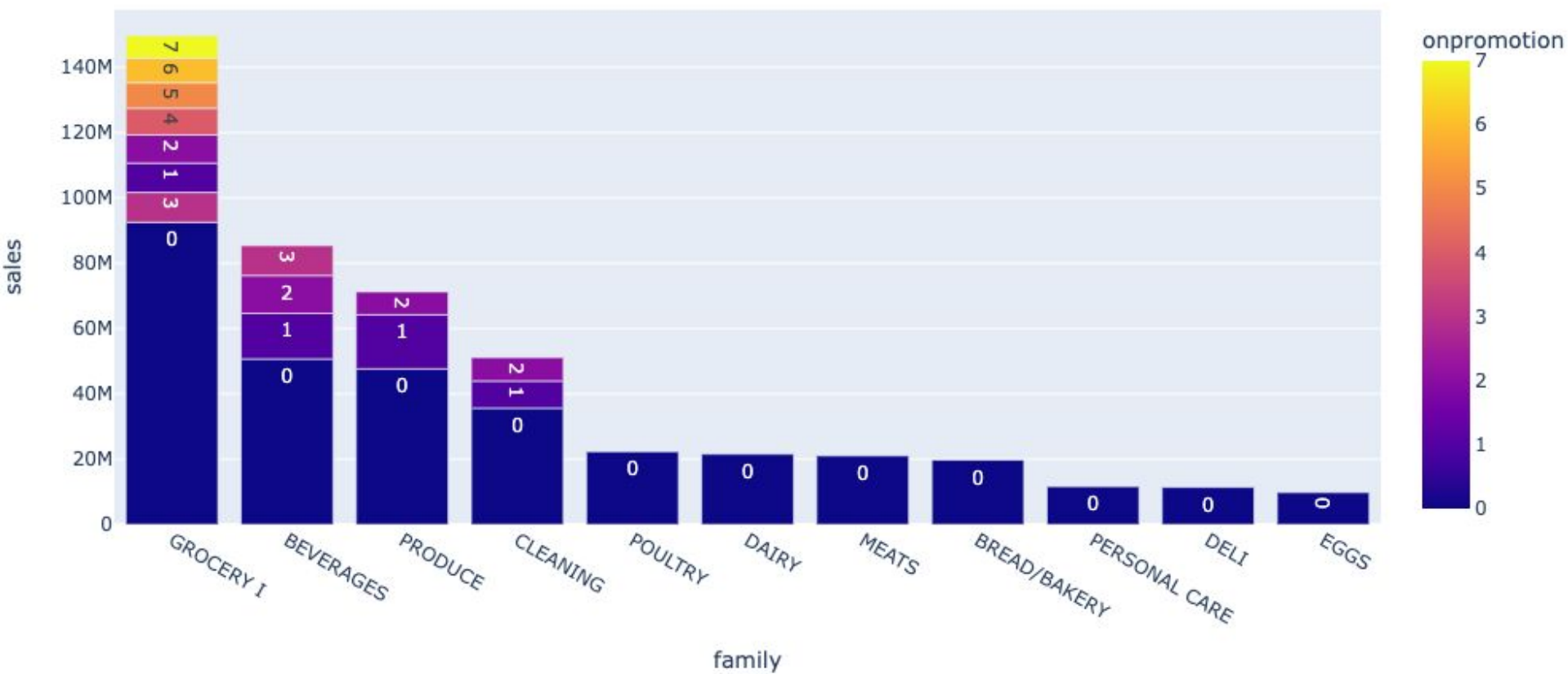




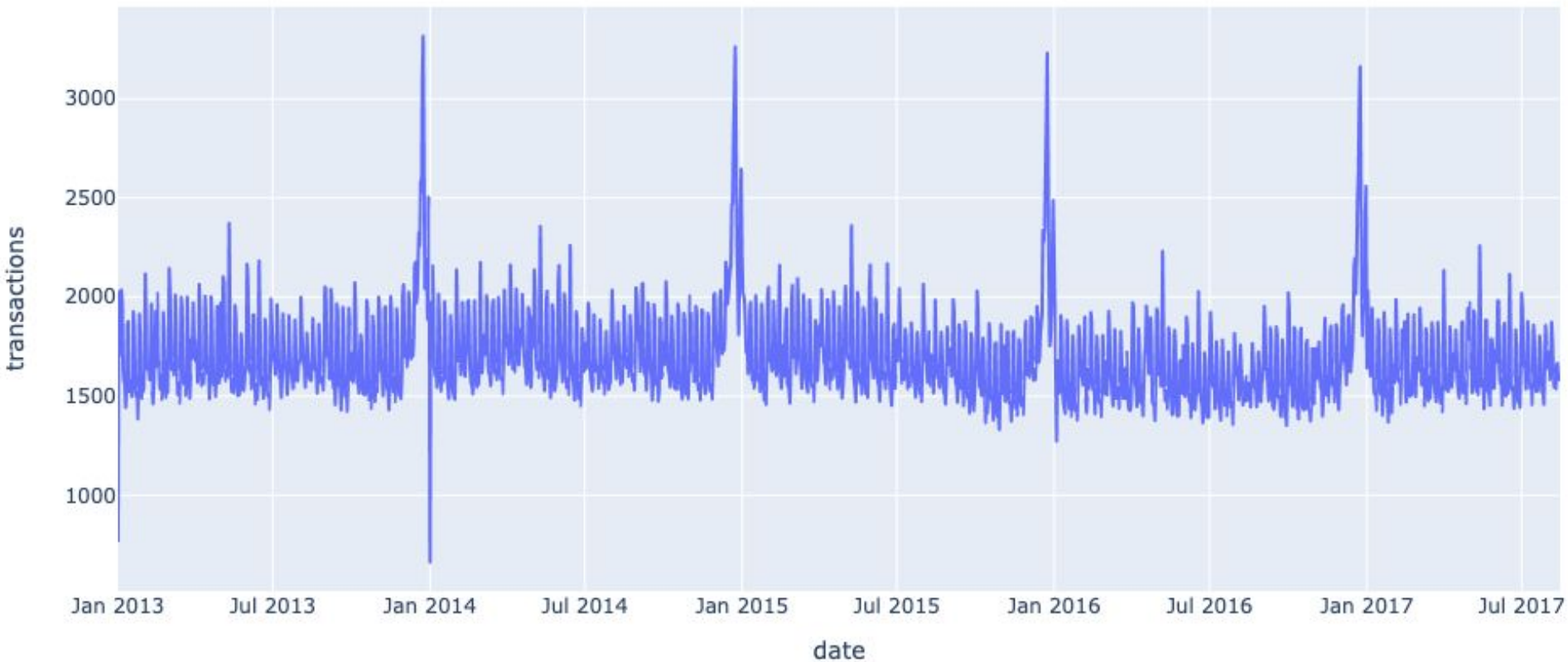
Oil prices over time



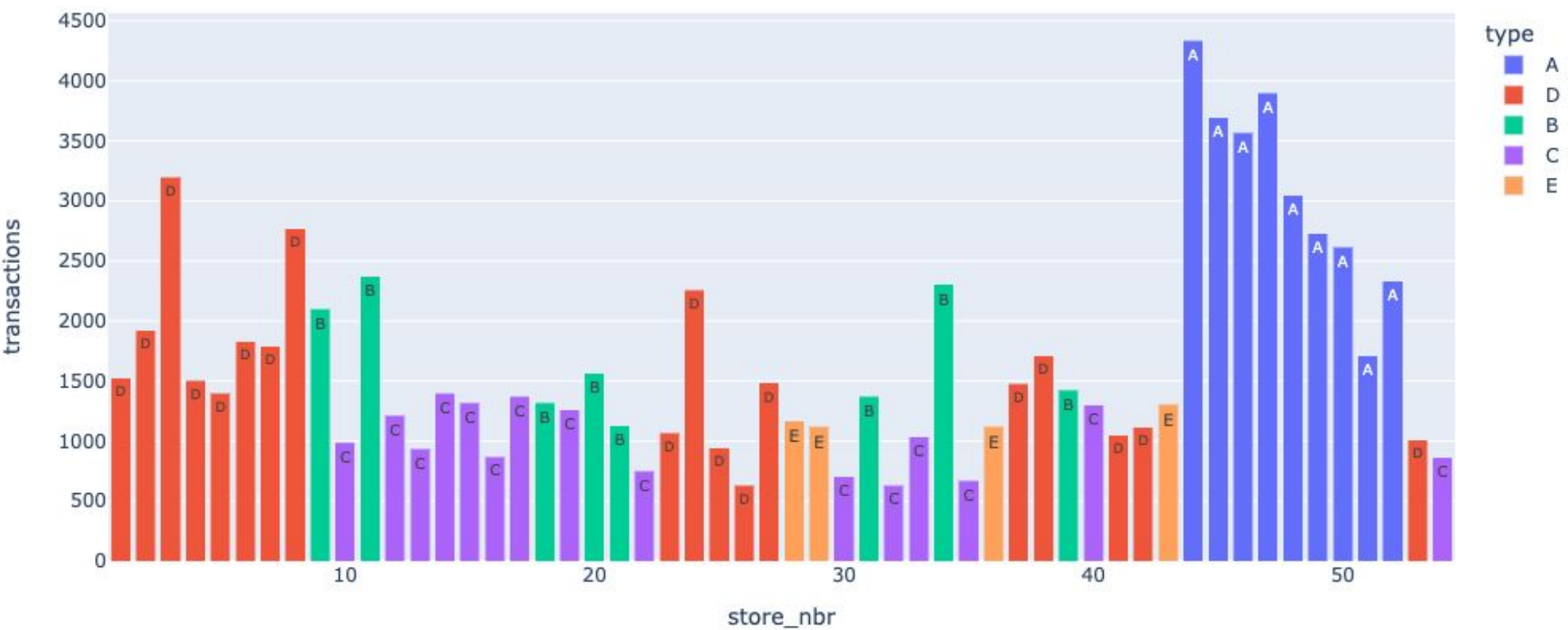
Total number of items sold in each family of products



Average transactions over time



Total number of transactions for each store

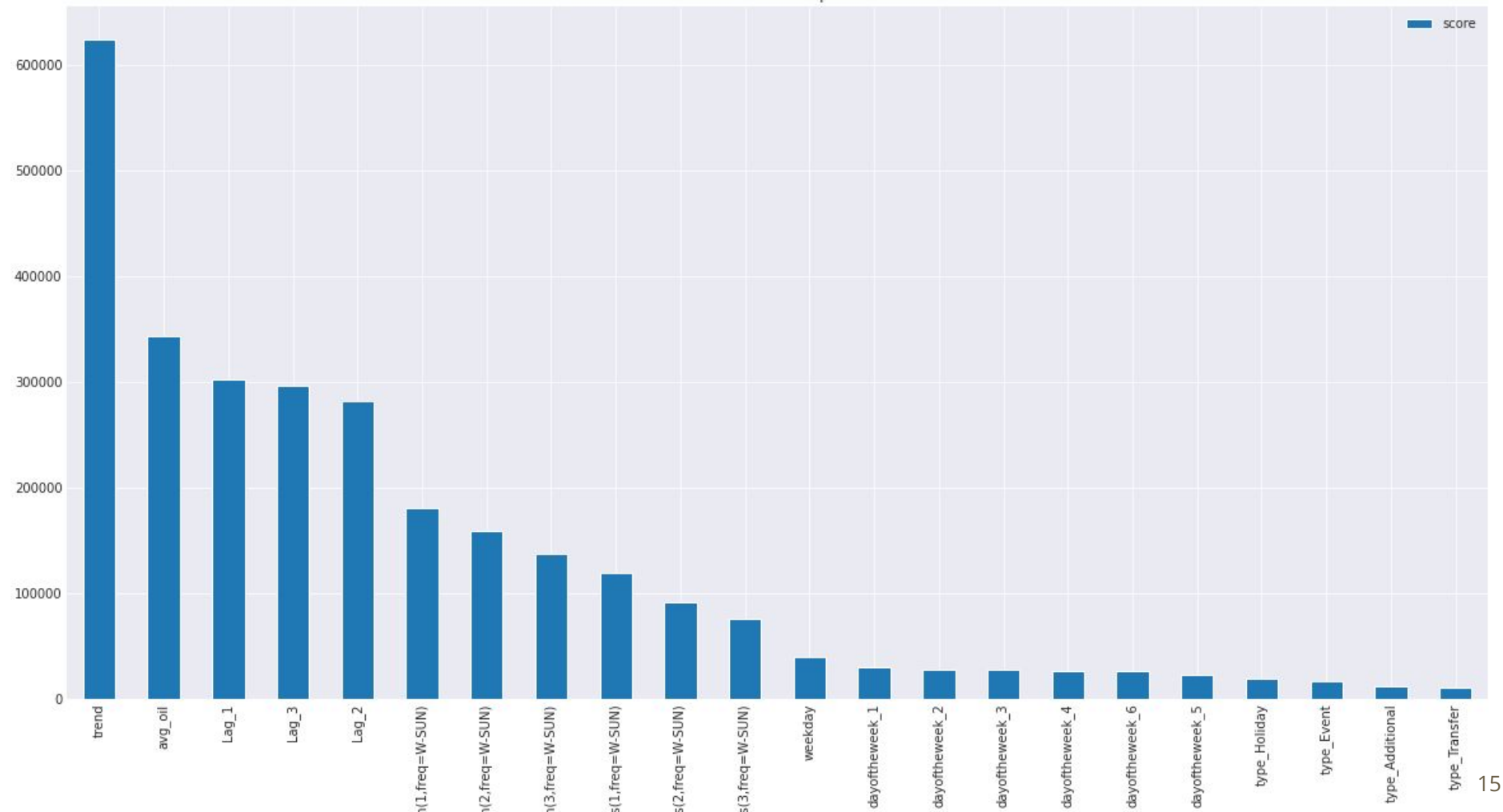


Modeling

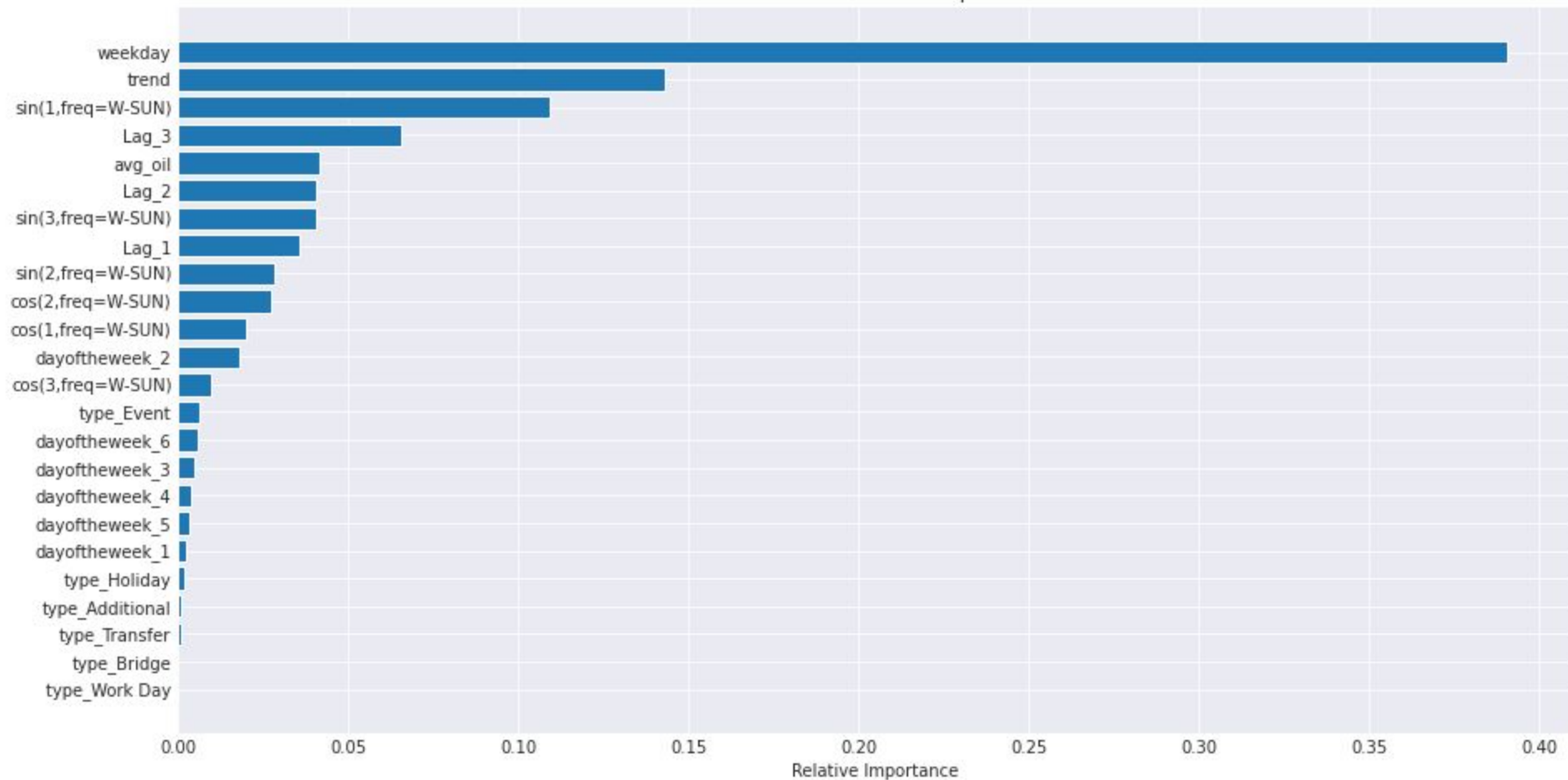
- Regression techniques
 - Linear Regression
 - Ridge
 - Lasso
 - Random Forest
 - SVR
 - XGboost
- To consider time dependency
 - ARIMA
 - SARIMA
- Blending and Custom Regressor (0.40408)
 - Used Linear model as base model
 - Built a custom regressor with Random Forest, Extra Tree Regressor, Ridge regression and svr.

Code demo

XGboost Feature importance



Random forest Feature Importances



Kaggle Leaderboard

Overview	Data	Code	Discussion	Leaderboard	Rules	Team	My Submissions	Submit Predictions	...
6	henry savier						0.40285	33	2d
7	wantAccepted						0.40299	15	17d
8	lyjps11						0.40299	11	21d
9	SOADM2						0.40300	1	2mo
10	ylquan						0.40303	28	19d
11	SaltyFishhh						0.40303	2	19d
12	HampnieHambart						0.40366	9	12d
13	曹纓						0.40371	16	1mo
14	RVK						0.40402	1	2mo
15	Kashyap Tamakuwala						0.40405	17	1s

Your Best Entry!
Your submission scored 0.40480, which is not an improvement of your previous score. Keep trying!

Total Participants: 900

Our team is in top: 2 %

App demo