

Santander Product Recommendation

TEAM KWT Wen Li Yisong Tao Lydia Kan



01

Introduction

02

Data Cleaning and EDA

03

Feature Engineering

04

Training Model

05

Result and Finding

06

Future Steps

AGENDA



Introduction

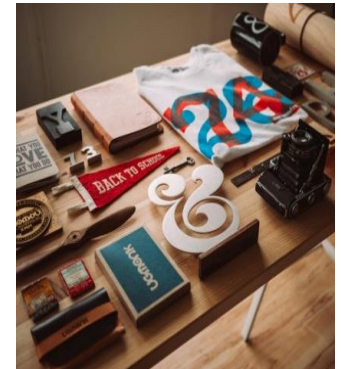
Project Description

Santander Bank offers their customers personalized product recommendations time to time, in order to meet the individuals needs and satisfaction.

This challenge seeks to improve the recommendation system by predicting which products their existing customers will use in the next month based on their past behavior.

Goal

Achieve top 5% ranking and MAP@7 score on Kaggle leader board



Introduction

01

Data Size

Training Set:
13,647,409

Test Set:

02

Input Features

Categorical: 45

Continuous: 3

Customer Info. :

1: 24

Product Purchased Info:

25:48

03

Output Features

Multi-Classifer

Recommended

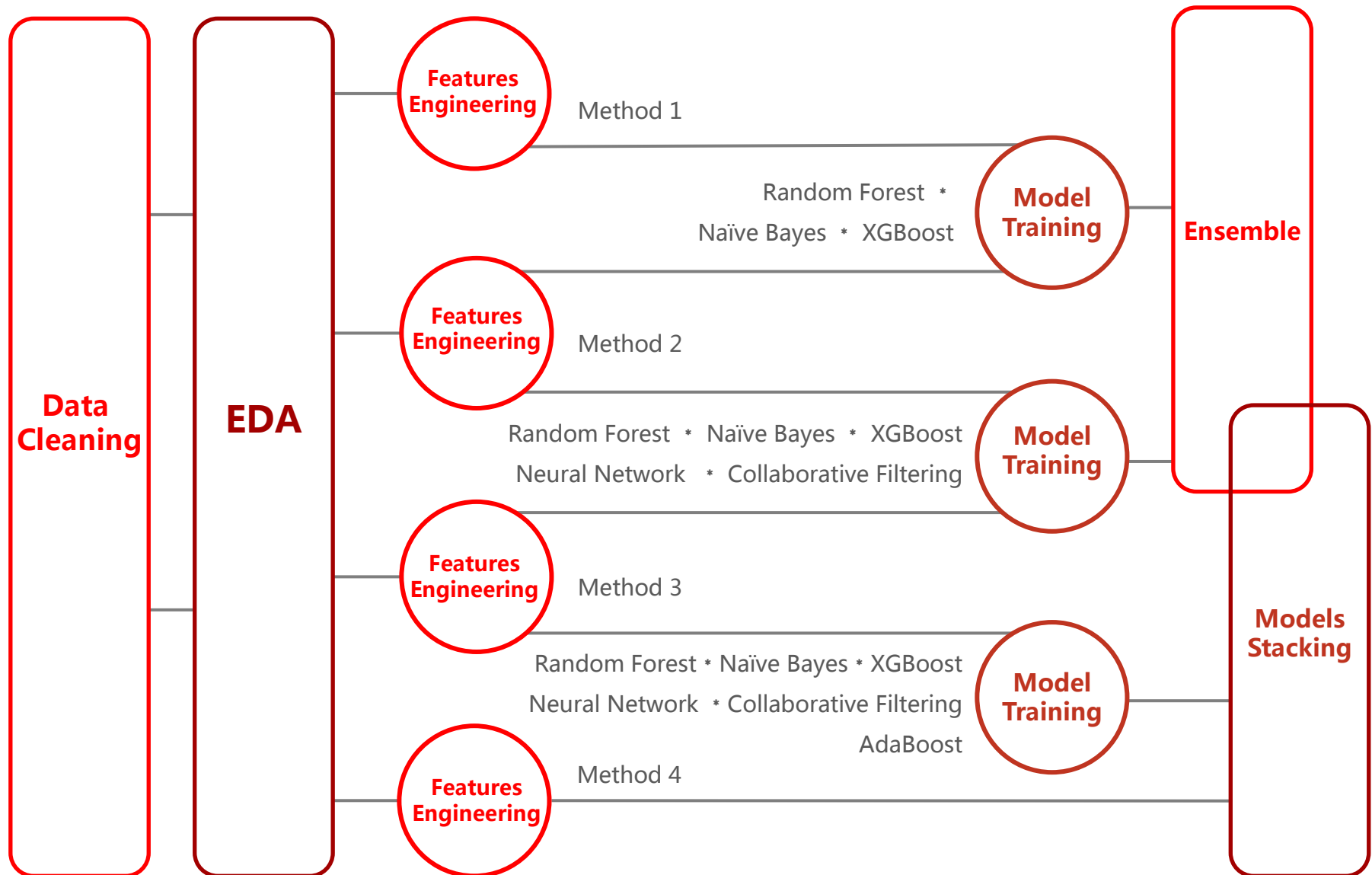
Products : 7

04

Evaluation

MAP@7

Workflow



Data Cleaning

Imputation

Contain Missing Values:

24 Features

Dropping Features

Drop 5 Features:

- **Having over 95% missing value**
- **Repetitive of other features**

Imputation



Unknown

- Sex
- Employee Index
- Country Residency
- Segmentation
- Residence Index
- Foreigner Index
- Channel to Join
- Primary
- Province Name



Common Type

- Customer Type
- Activity Index
- Rent



Others

- New Customer – New
- Seniority – Min
- Age – Scale, Mean
- Relationship Type – 'A'
- Deceased Index – 'N'

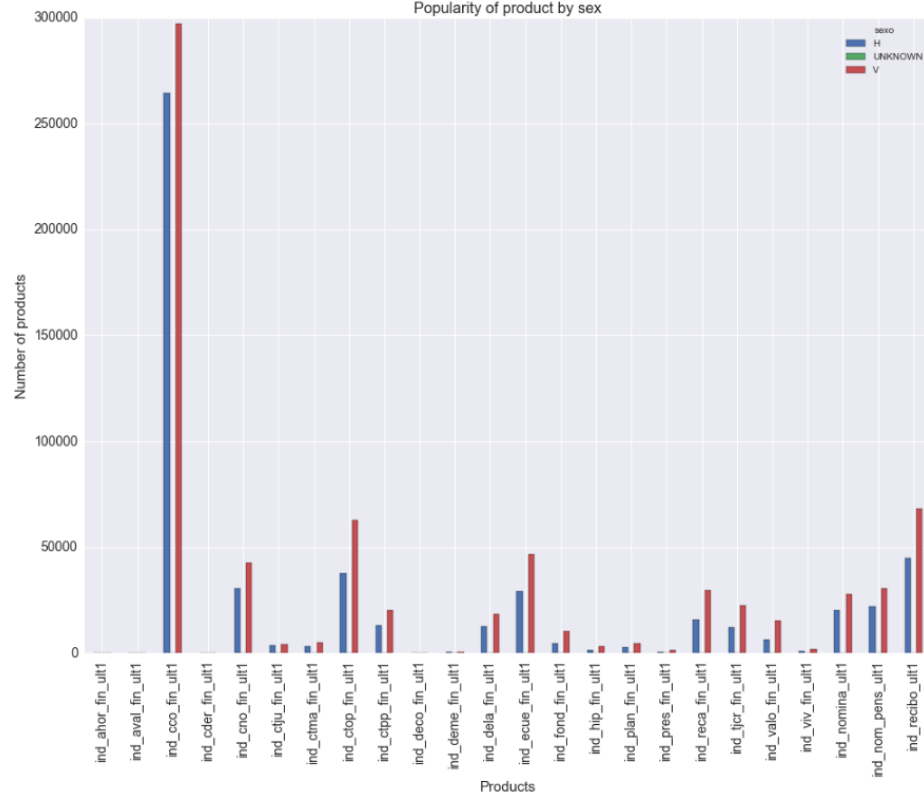


Products

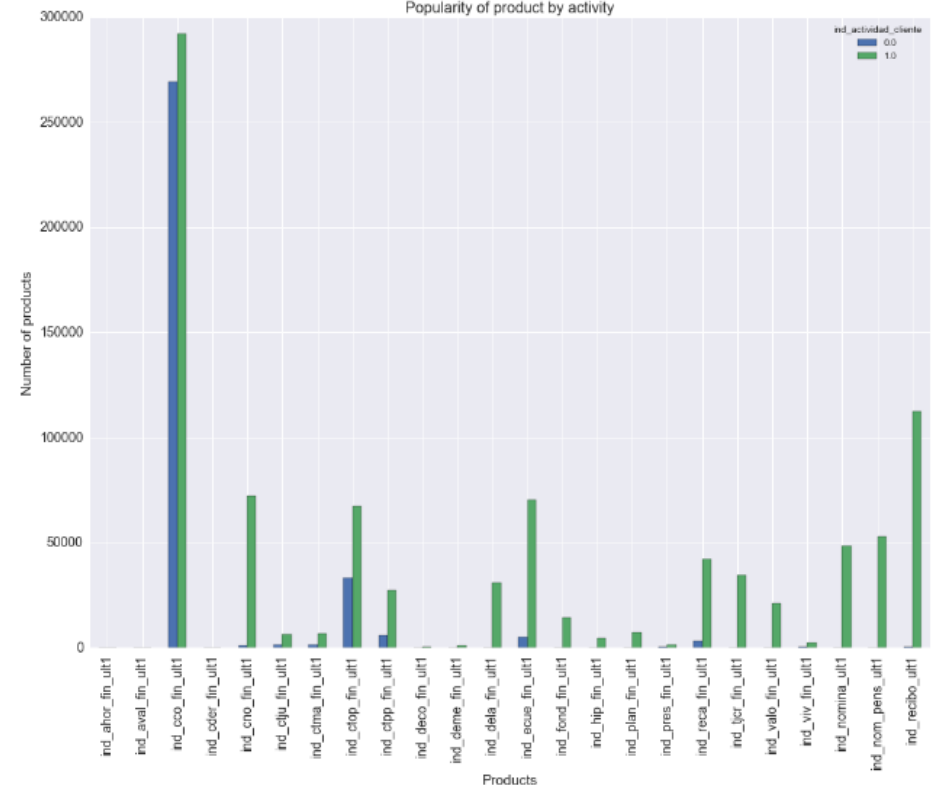
- Payroll - 0
- Pensions - 0

Product Sales Related to Customer's Info - 2016.5

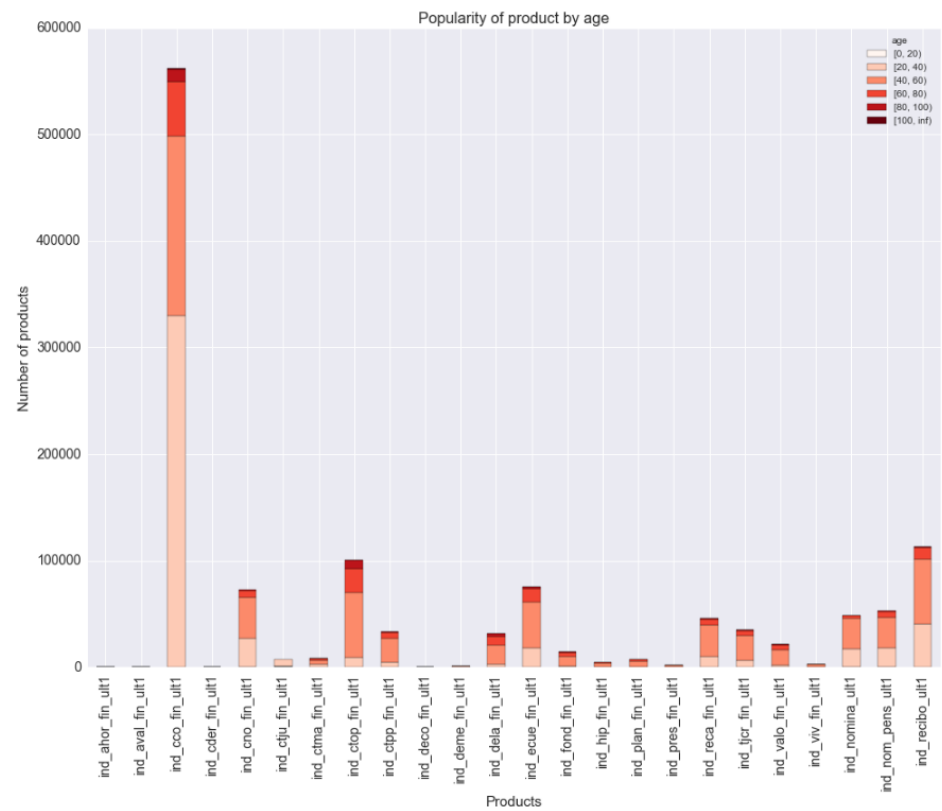
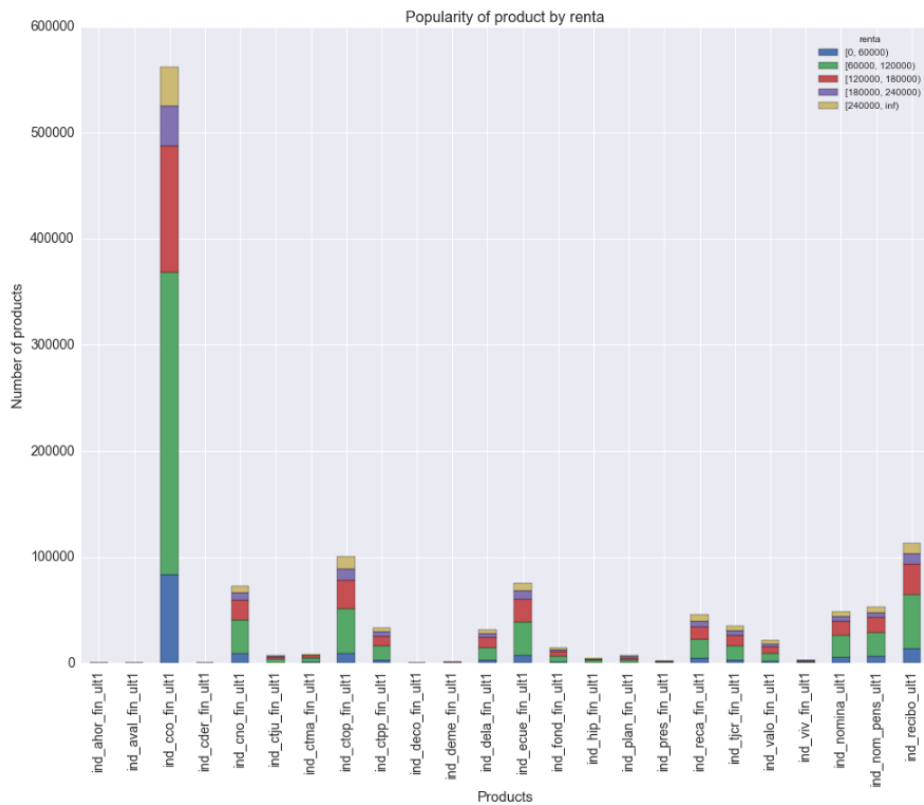
Popularity of product by sex



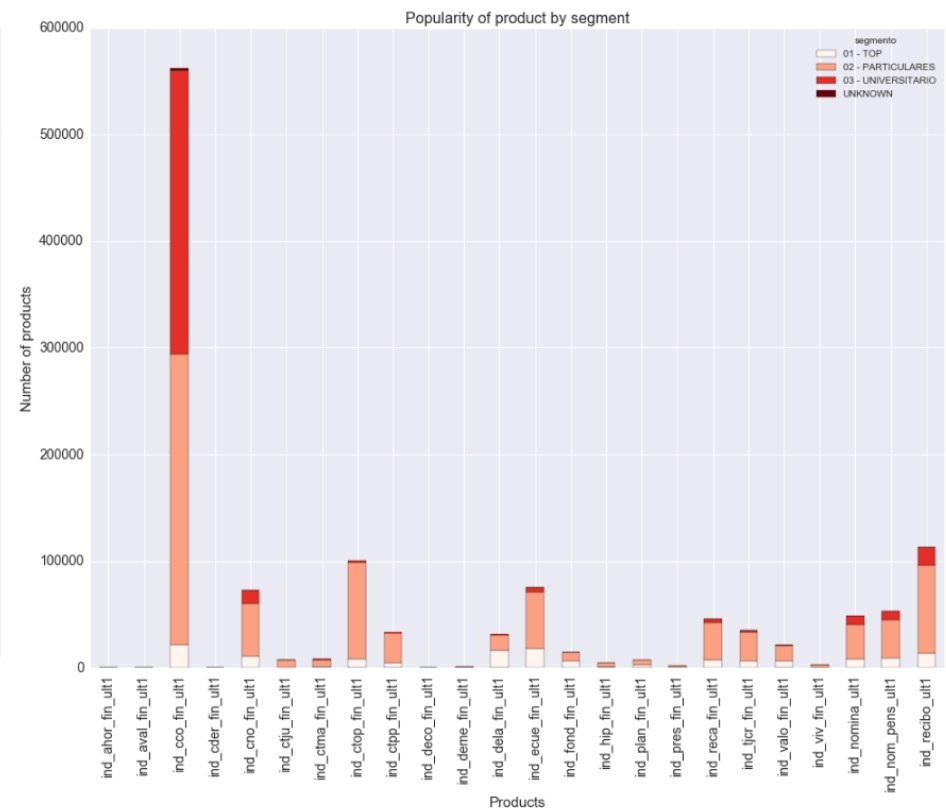
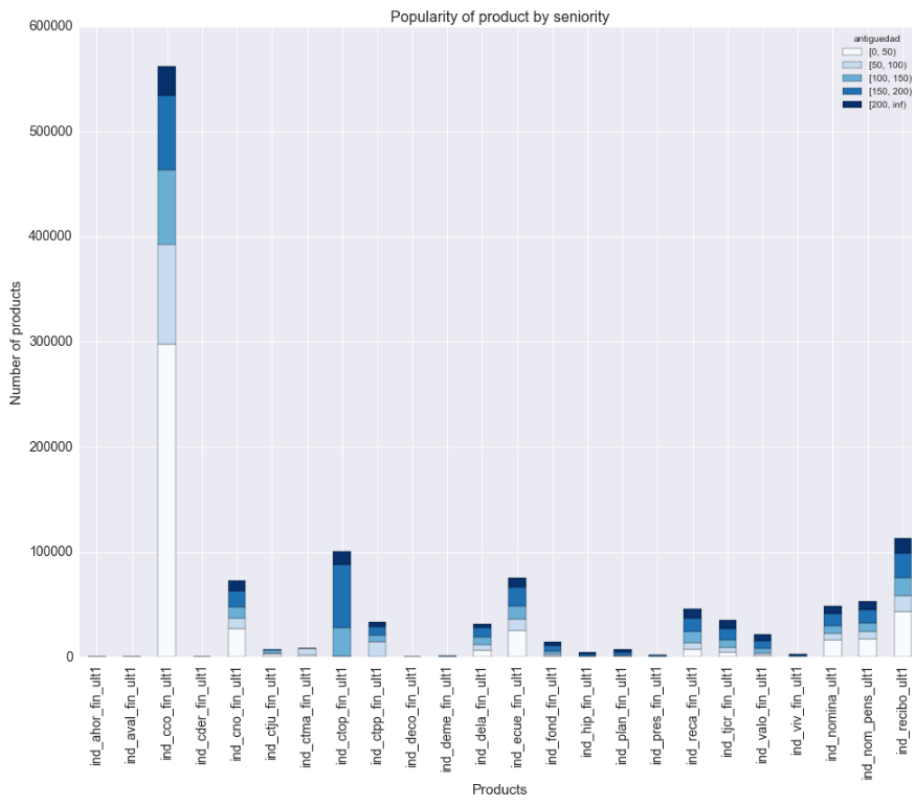
Popularity of product by activity



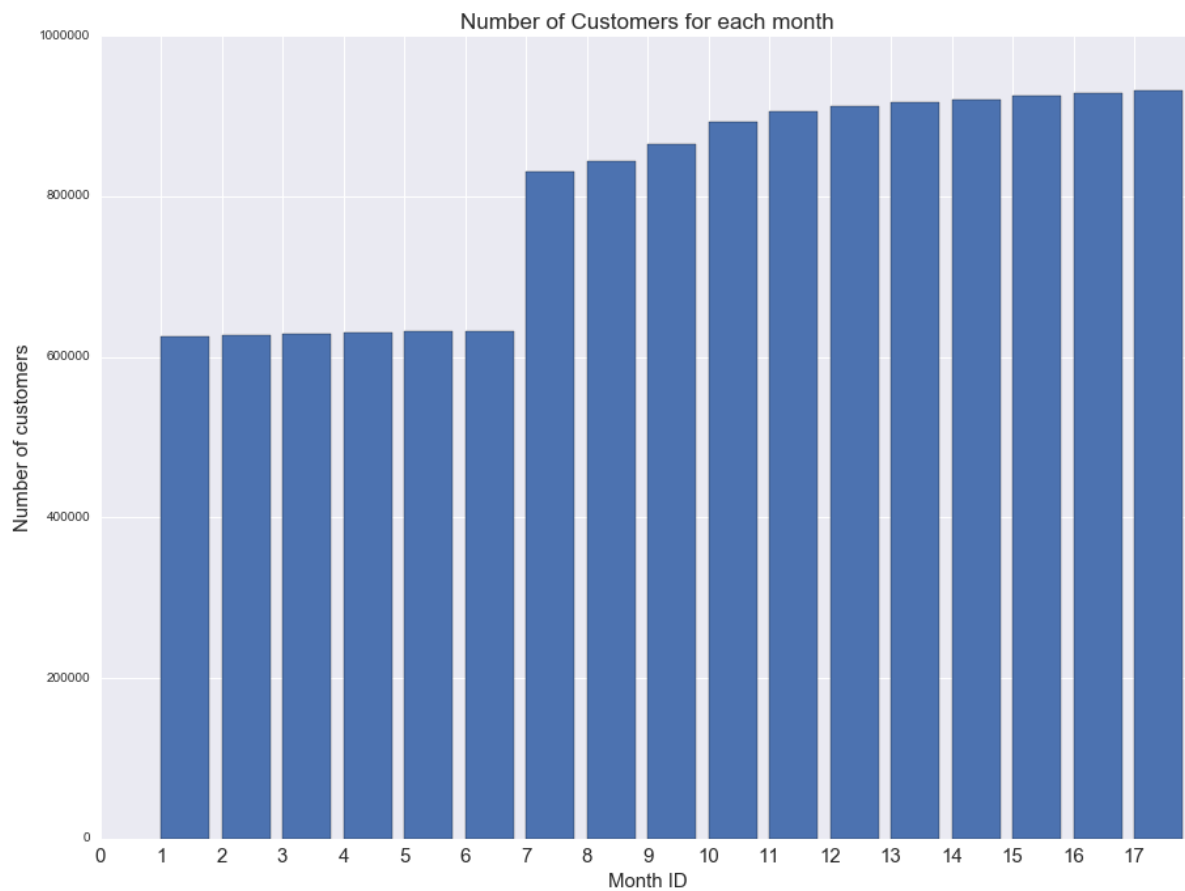
Product Sales Related to Customer's Info - 2016.5



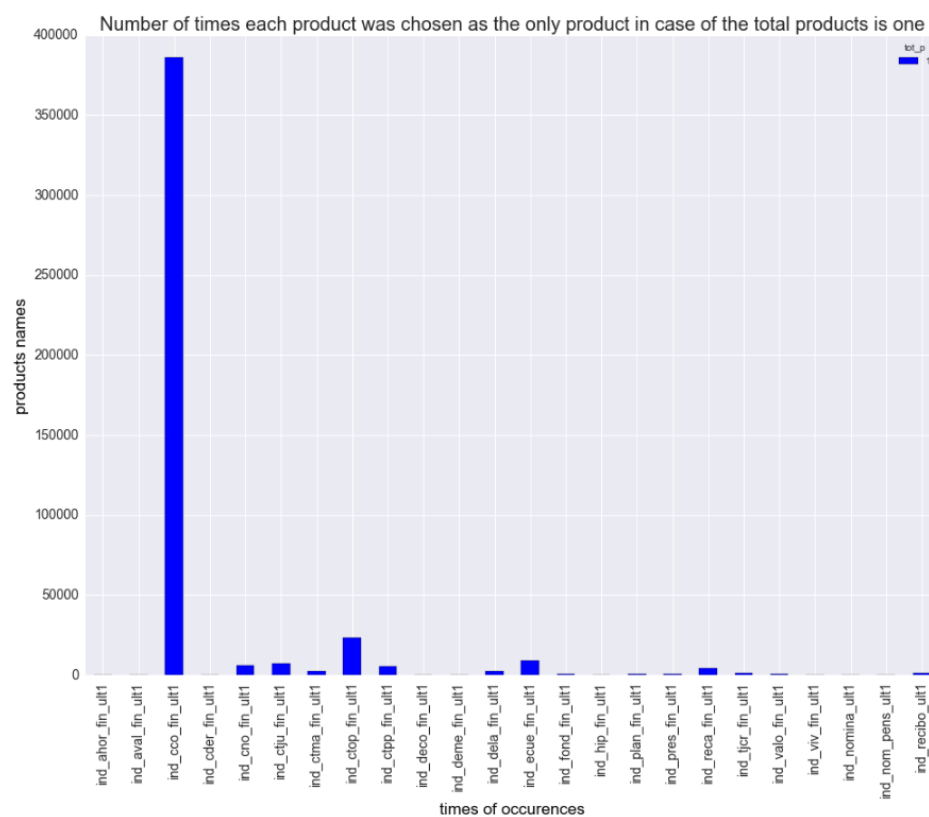
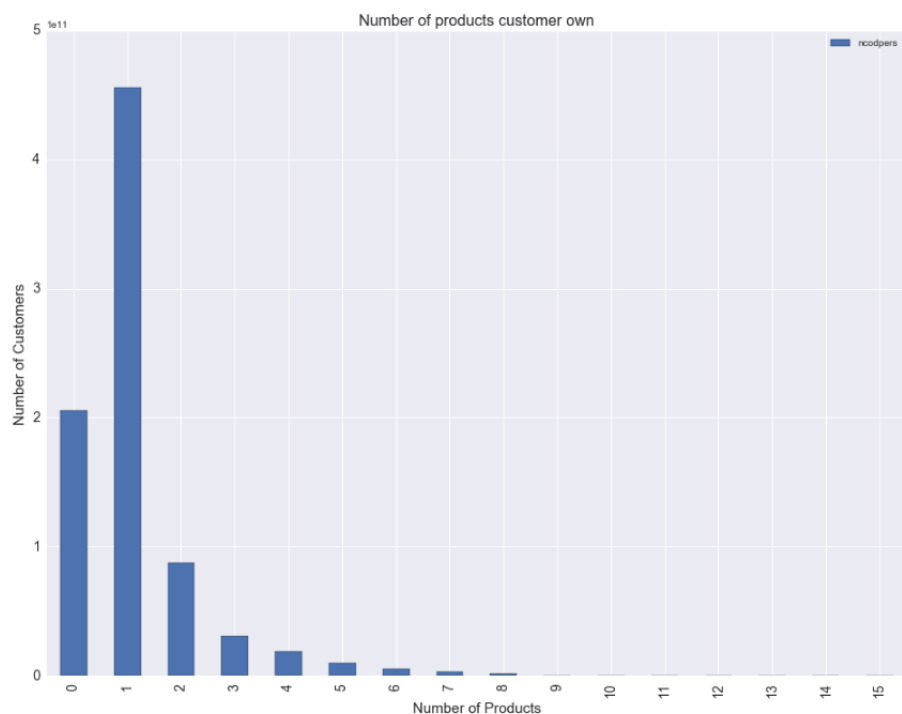
Product Sales Related to Customer's Info - 2016.5



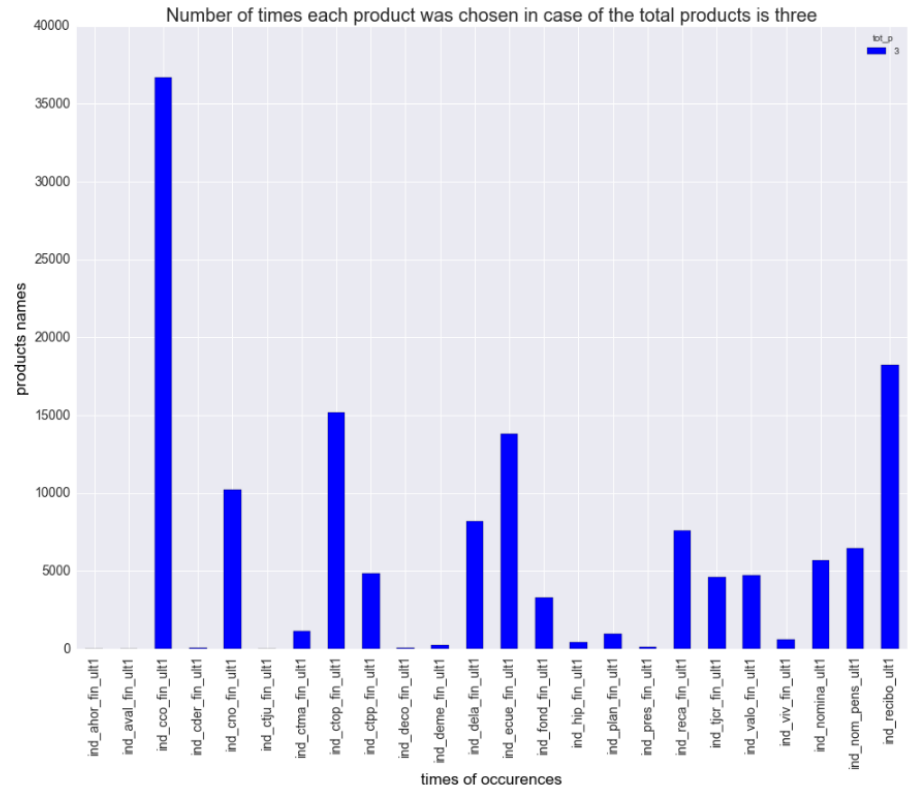
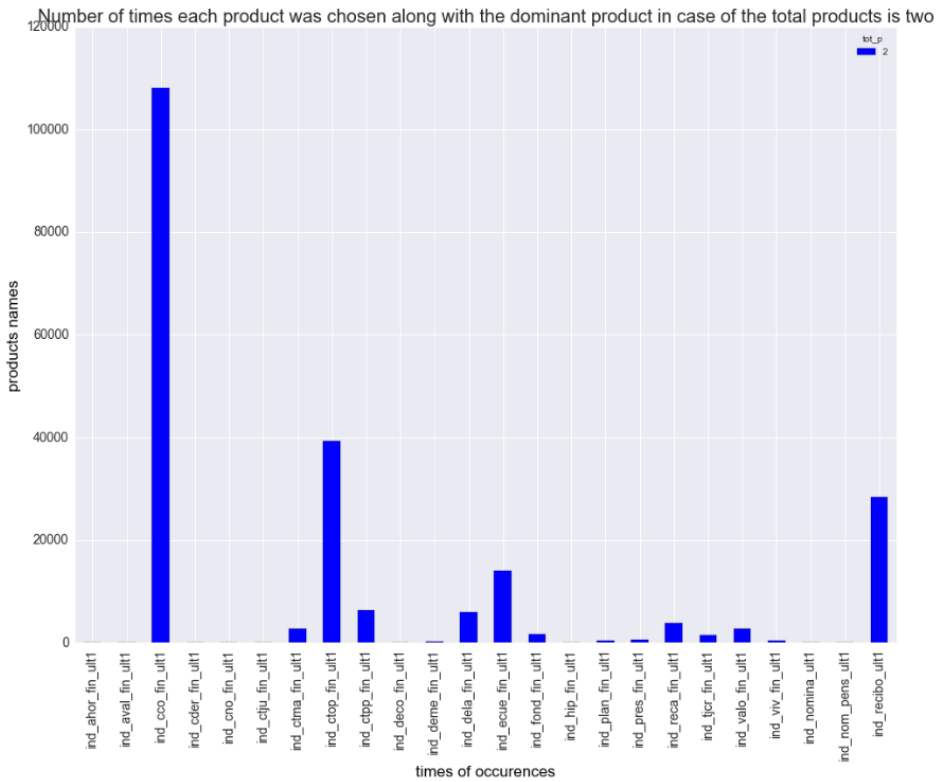
Number of Customers by Time



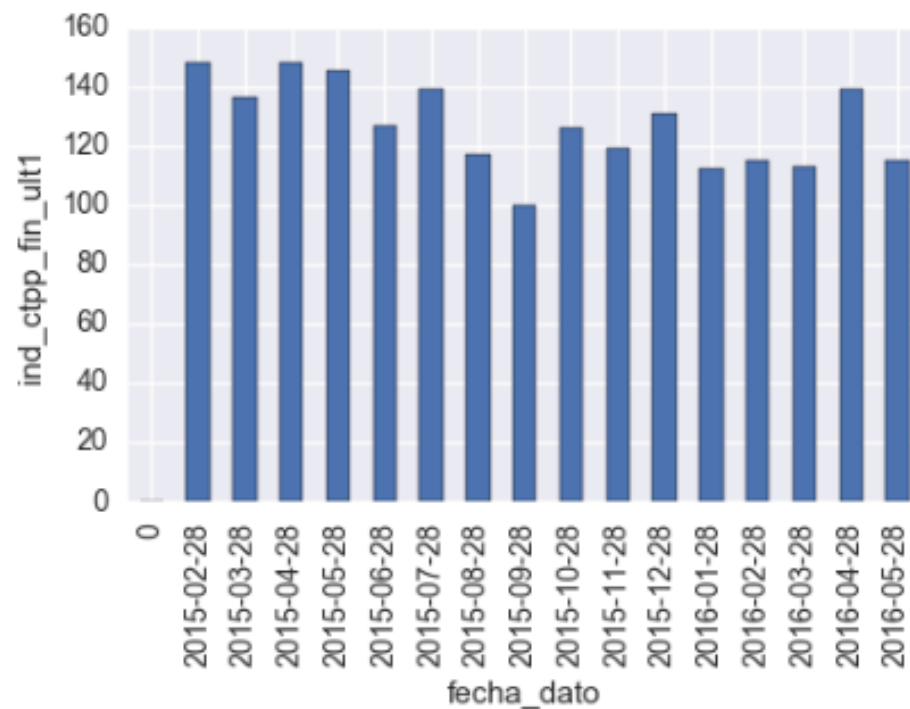
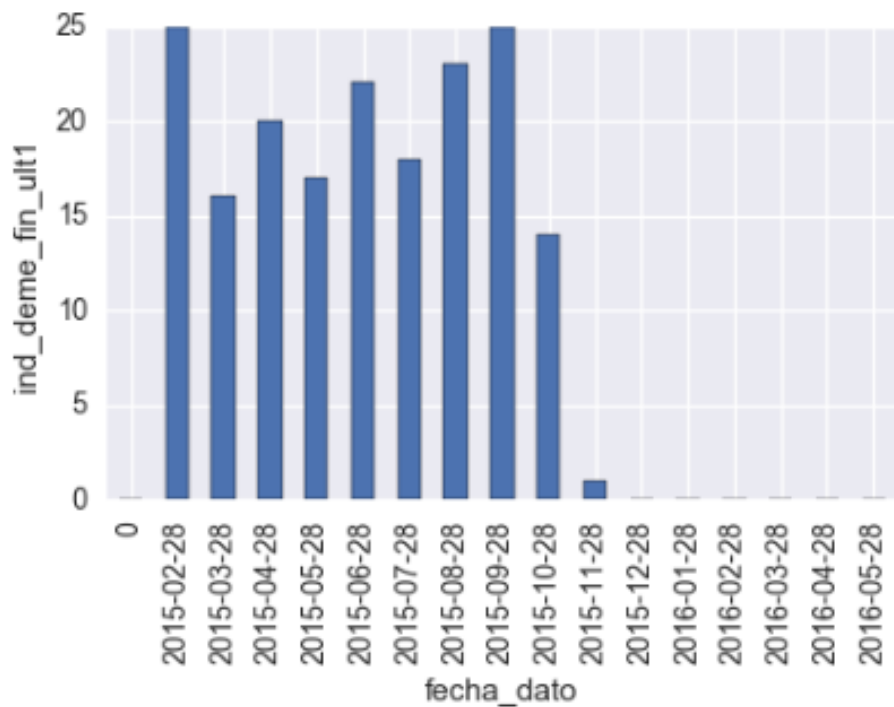
Number of Product Own - 2016.5



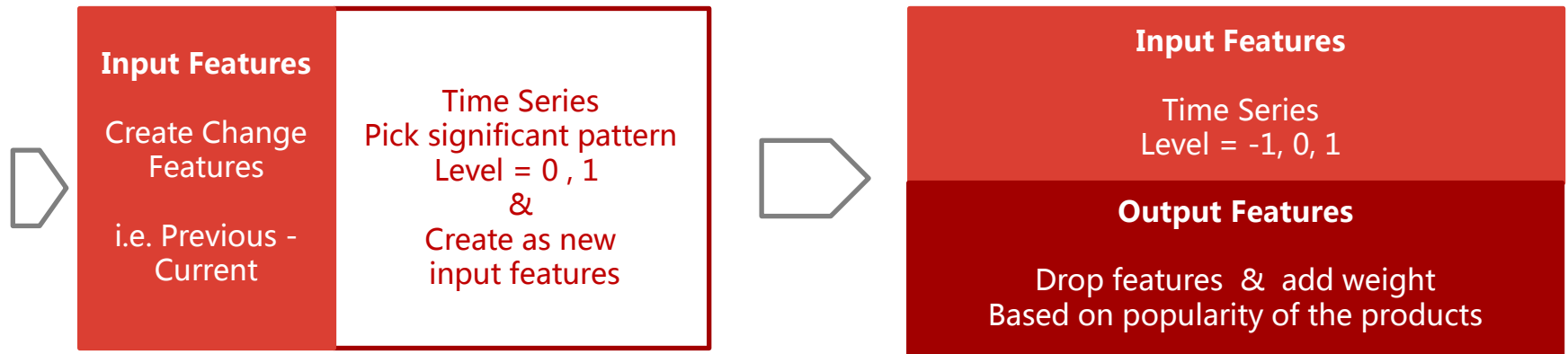
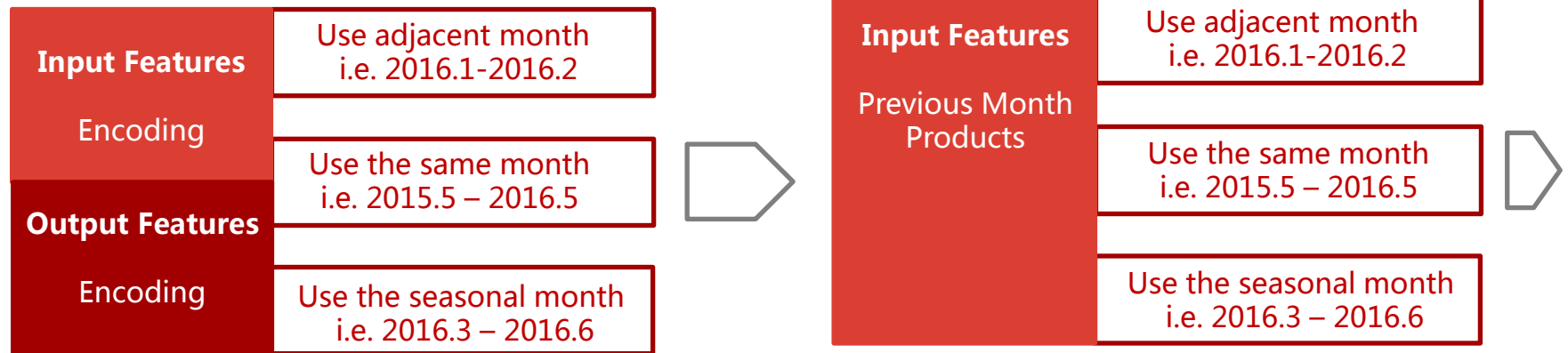
Number of Product Own - 2016.5



Number of Product Sales by Time

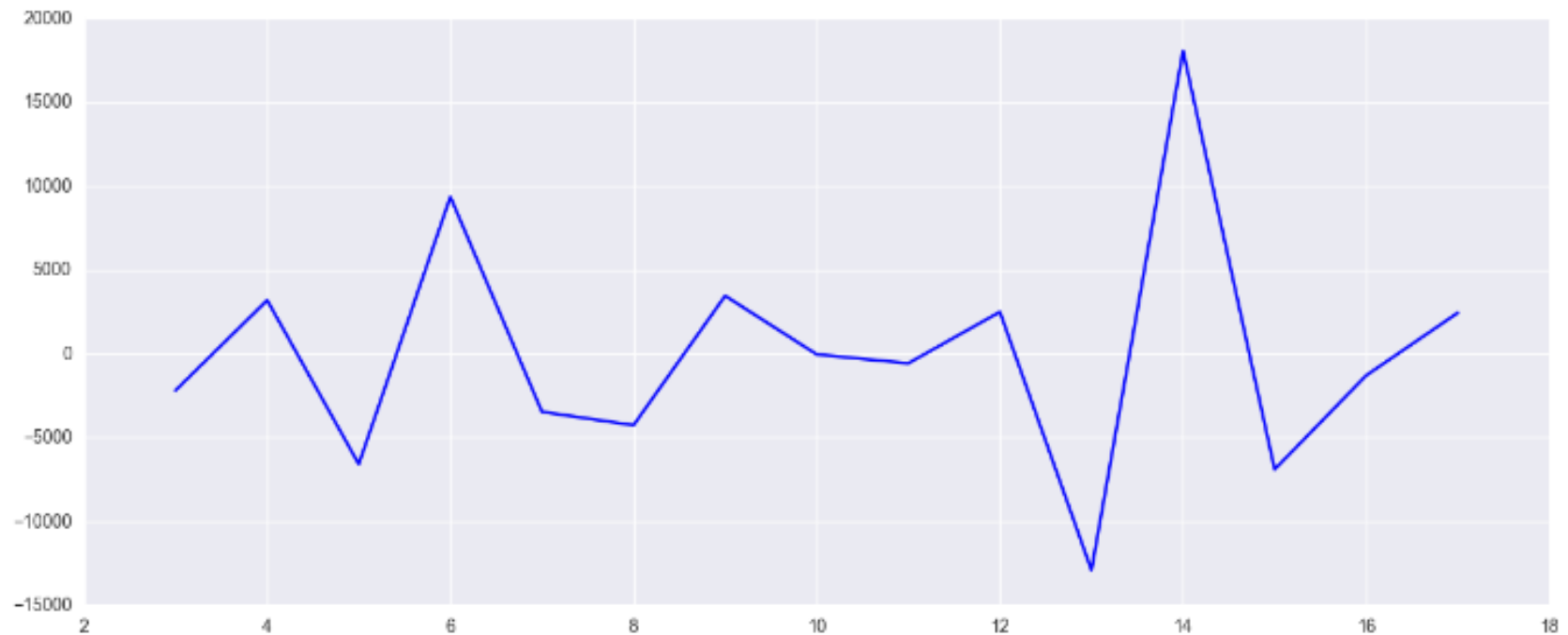


Feature Engineering



Make recommendation based on products' popularity

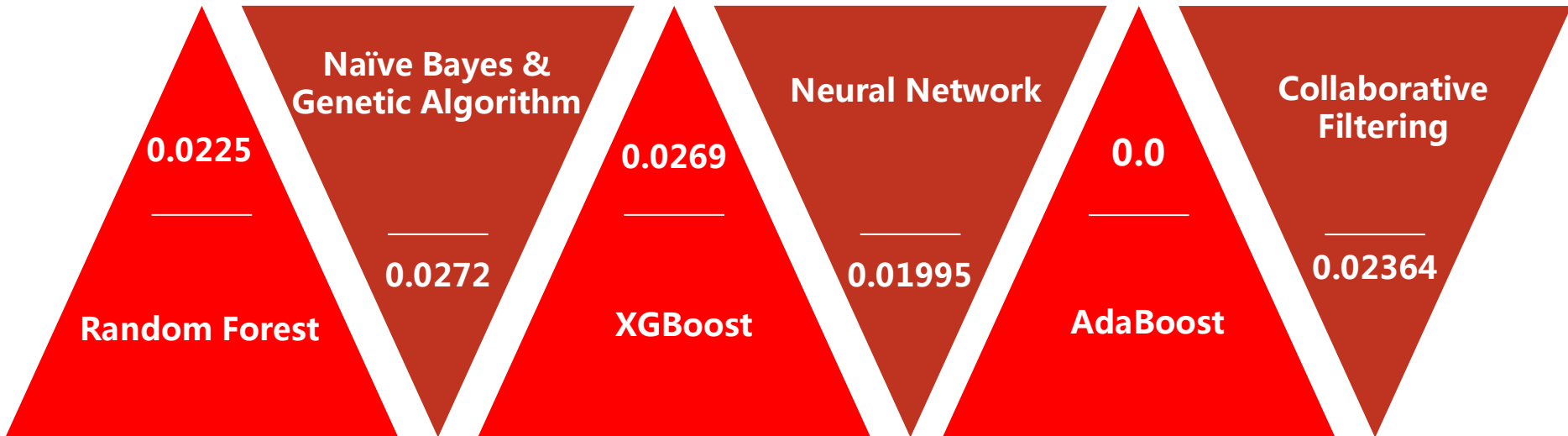
Time Series



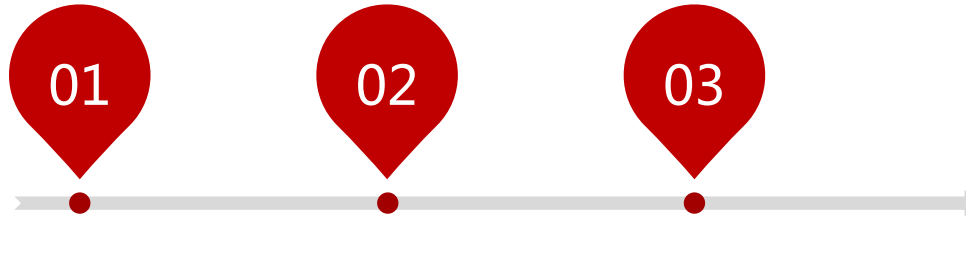
Results of Dickey - Fuller Test Product Pension

Test Statistic	-3.163039
p-value	0.022226
No. Lags Used	4.000000
Critical Value (5%)	-3.232950
Critical Value (1%)	-4.331573
Critical Value (10%)	-2.748700

Models Training



Ensemble - Voting



Popularity Collaborative Naïve Bayes
Filter



Popularity Collaborative Naïve Bayes XGBoost XGBoost Random Forest
Filter (new Features) (multiclass)



Result & Findings

Using the month from the previous year has better prediction than the previous month from the current year

Single Model, Naïve Bayes has the best performance

From the evaluation, it only penalized the false negative

Multiclass has better performance than multi-labels

Future Steps

- Using best models for model stacking
- Trying more features engineering
- More Ensemble and Stacking

