



Raymond A. Mason
School of Business
WILLIAM & MARY

Predicting Future Sales

Team 20 - You Wu, Vanessa Guzman, Frank Wan, Elie Baaklini



Agenda

1. Overview
2. Critique of Other Competition Submissions
3. Our Solution
4. Final Prediction



Overview

Objective: Predict total sales (quantity) for every product and store in the next month (November 2015)

Dataset: includes 6 csv files from 1C company (Russian software company)

- items:

| | A | B | C |
|---|--|---------|------------------|
| 1 | item_name | item_id | item_category_id |
| 2 | ДВД-диск "ДВД-диск" (DVD-диск) | 0 | 40 |
| 3 | Professional Edition Full (PC, DVD, CD, DVD) | 1 | 76 |
| 4 | Учебное пособие по информатике (UNV) | 2 | 40 |
| 5 | "Университетский" DVD-диск (Univ) | 3 | 40 |
| 6 | Учебное пособие по информатике (DVD) | 4 | 40 |
| 7 | Учебное пособие по информатике (DVD) | 5 | 40 |

- shops:

| | A | B |
|---|---|---------|
| 1 | shop_name | shop_id |
| 2 | Магазин "Магазин" (Магазин, 56 Магазин) | 0 |
| 3 | Магазин "Магазин" (Магазин, 56 Магазин) | 1 |
| 4 | Магазин "Магазин" (Магазин, 56 Магазин) | 2 |
| 5 | Магазин "Магазин" (Магазин, 56 Магазин) | 3 |
| 6 | Магазин "Магазин" (Магазин, 56 Магазин) | 4 |
| 7 | Магазин "Магазин" (Магазин, 56 Магазин) | 5 |

- item categories:

| | A | B |
|---|--|------------------|
| 1 | item_category_name | item_category_id |
| 2 | PC - Магазин "Магазин" (Магазин, 56 Магазин) | 0 |
| 3 | Магазин "Магазин" (Магазин, 56 Магазин) | 1 |
| 4 | Магазин "Магазин" (Магазин, 56 Магазин) | 2 |
| 5 | Магазин "Магазин" (Магазин, 56 Магазин) | 3 |
| 6 | Магазин "Магазин" (Магазин, 56 Магазин) | 4 |
| 7 | Магазин "Магазин" (Магазин, 56 Магазин) | 5 |

Overview

- sales-train: training set

| | A | B | C | D | E | F |
|---|------------|----------------|---------|---------|------------|--------------|
| 1 | date | date_block_num | shop_id | item_id | item_price | item_cnt_day |
| 2 | 02.01.2013 | 0 | 59 | 22154 | 999 | 1 |
| 3 | 03.01.2013 | 0 | 25 | 2552 | 899 | 1 |
| 4 | 05.01.2013 | 0 | 25 | 2552 | 899 | -1 |
| 5 | 06.01.2013 | 0 | 25 | 2554 | 1709.05 | 1 |
| 6 | 15.01.2013 | 0 | 25 | 2555 | 1099 | 1 |
| 7 | 10.01.2013 | 0 | 25 | 2564 | 349 | 1 |

- sample: sample submission file in the correct format

| | A | B |
|----|----|----------------|
| 1 | ID | item_cnt_month |
| 2 | 0 | 0.5 |
| 3 | 1 | 0.5 |
| 4 | 2 | 0.5 |
| 5 | 3 | 0.5 |
| 6 | 4 | 0.5 |
| 7 | 5 | 0.5 |
| 8 | 6 | 0.5 |
| 9 | 7 | 0.5 |
| 10 | 8 | 0.5 |
| 11 | 9 | 0.5 |
| 12 | 10 | 0.5 |
| 13 | 11 | 0.5 |
| 14 | 12 | 0.5 |

- test: test set

| | A | B | C |
|---|----|---------|---------|
| 1 | ID | shop_id | item_id |
| 2 | 0 | 5 | 5037 |
| 3 | 1 | 5 | 5320 |
| 4 | 2 | 5 | 5233 |
| 5 | 3 | 5 | 5232 |
| 6 | 4 | 5 | 5268 |
| 7 | 5 | 5 | 5039 |

Critique of Other Competition Submissions

Random Forest

- **Pros:**
 - Suitable for large dataset
 - Easy data preparation comparing to other algorithms
- **Cons:**
 - Did not tune the hyperparameter
 - Computationally intense when number of trees is big
 - Will cause overfitting when noise is large
 - Lower output accuracy because it cannot guarantee the best tree
- **RMSE: 2.0182**

```
|:
#Random forest regressor model building
from sklearn.ensemble import RandomForestRegressor

RF_model = RandomForestRegressor()
RF_model.fit(X_prepared, y)
```

Critique of Other Competition Submissions

LightGBM (Light Gradient Boosting Machine)

- **Pros:**
 - Faster training speed, Light GBM uses a histogram-based algorithm i.e it buckets continuous feature values into discrete bins which fasten the training procedure
 - Replaces continuous values to discrete bins which results in lower memory usage
- **Cons:**
 - Light GBM split the tree leaf-wise which can lead to overfitting as it produces much complex trees
 - Did not do cross validation for hyperparameters
- **RMSE: 0.9610**

```
evals_result = {}  
gbm = lgb.train(  
    params,  
    lgb_train,  
    num_boost_round=3000,  
    valid_sets=(lgb_train, lgb_eval),  
    feature_name = feature_name,  
    #categorical_feature = categorical_features,  
    verbose_eval=5,  
    evals_result = evals_result,  
    early_stopping_rounds = 10)
```

Our Solution

XGBoost: eXtreme Gradient Boosting--"ALL in One" algorithm

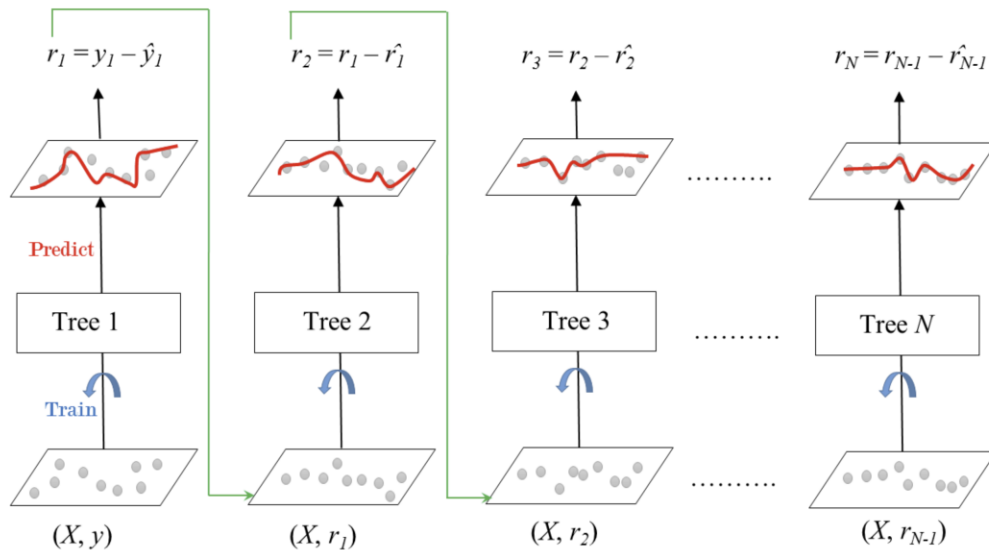
As a popular supervised-learning algorithm, XGBoost uses decision trees as base learners; combining many weak learners to make a strong learner to speed up and increase the performance of gradient boosted decision trees

- Regularization
- Parallel Processing
- Handling Missing Values
- Effective Tree Pruning



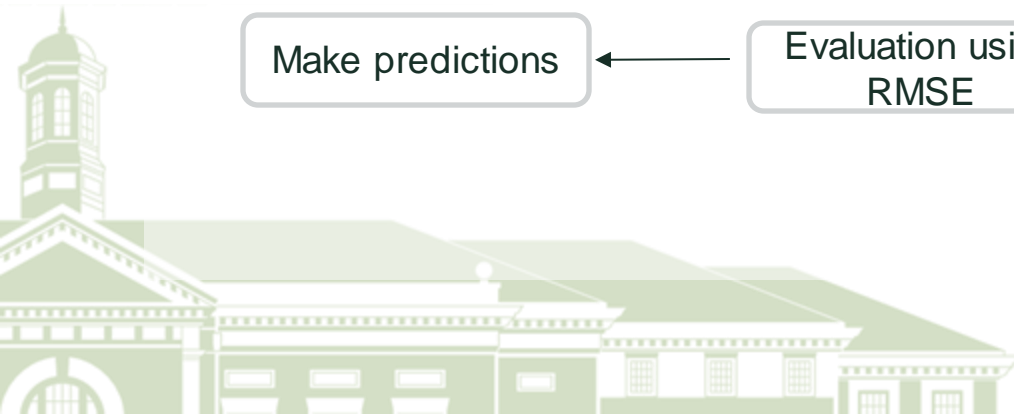
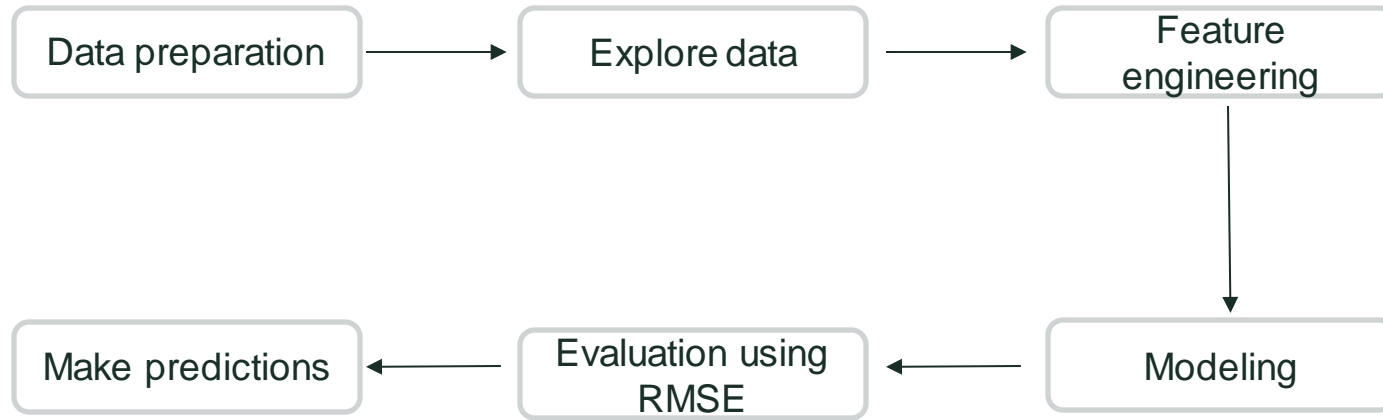
Our Solution

How does
XGBoost
work?



Our Solution

Flow chart



Final Prediction

Performance: RMSE=0.877

Final prediction results:

| ID | item_cnt_month |
|-------|----------------|
| 33050 | 0.01800462 |
| 32455 | 0.003886187 |
| 35412 | 0.019796828 |
| 32228 | 0.020295562 |
| 35385 | 0 |
| 32229 | 0.073716491 |
| 33048 | 0.007104043 |
| 32454 | 0 |
| 32446 | 0.036807135 |
| 33047 | 0.215048745 |
| 31748 | 0.042575739 |
| 33046 | 0.073391959 |
| 32370 | 0.073716491 |
| 34153 | 0.073716491 |
| 33045 | 0.027066618 |
| 30784 | 0.031092696 |
| 34150 | 0.073716491 |
| 31880 | 0.036855627 |
| 32286 | 0.121237382 |
| 30785 | 0.008510824 |
| 33044 | 0.015567578 |
| 32371 | 0.004501749 |

