

# Predict Future Sales

Final Report

Kaicheng Jia

# CONTENTS

**01**

**INTRO**

Why important ?

**02**

**DATA**

EDA  
Feature Engineering

**03**

**ALGORITHM**

Model Ensemble

**04**

**RESULTS**

Results Analysis

**05**

**CONCLUSION**

Project conclusion

# 01 INTRODUCTION

# Background



- In this project we will work with a challenging timeseries dataset consisting of daily sales data, kindly provided by one of the largest Russian software firms - 1C Company.
- The task is to forecast the total amount of products sold in every shop for the test set in the next month.

# 02 DATA

# Files

## File descriptions

- `sales_train.csv` - the training set. Daily historical data from January 2013 to October 2015.
- `test.csv` - the test set. You need to forecast the sales for these shops and products for November 2015.
- `sample_submission.csv` - a sample submission file in the correct format.
- `items.csv` - supplemental information about the items/products.
- `item_categories.csv` - supplemental information about the items categories.
- `shops.csv` - supplemental information about the shops.

## Basic info of train

```
sales_train.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 2935849 entries, 0 to 2935848  
Data columns (total 6 columns):  
date                object  
date_block_num      int64  
shop_id             int64  
item_id             int64  
item_price          float64  
item_cnt_day        float64  
dtypes: float64(2), int64(3), object(1)  
memory usage: 134.4+ MB
```

	date	date_block_num	shop_id	item_id	item_price	item_cnt_day
0	02.01.2013	0	59	22154	999.00	1.0
1	03.01.2013	0	25	2552	899.00	1.0
2	05.01.2013	0	25	2552	899.00	-1.0
3	06.01.2013	0	25	2554	1709.05	1.0
4	15.01.2013	0	25	2555	1099.00	1.0

**date** - date in format dd/mm/yyyy

**date\_block\_num** - a consecutive month number, used for convenience. January 2013 is 0

**shop\_id** - unique identifier of a shop

**item\_id** - unique identifier of a product

**item\_price** - current price of an item

**item\_cnt\_day** - number of products sold. predicting a monthly amount of this measure

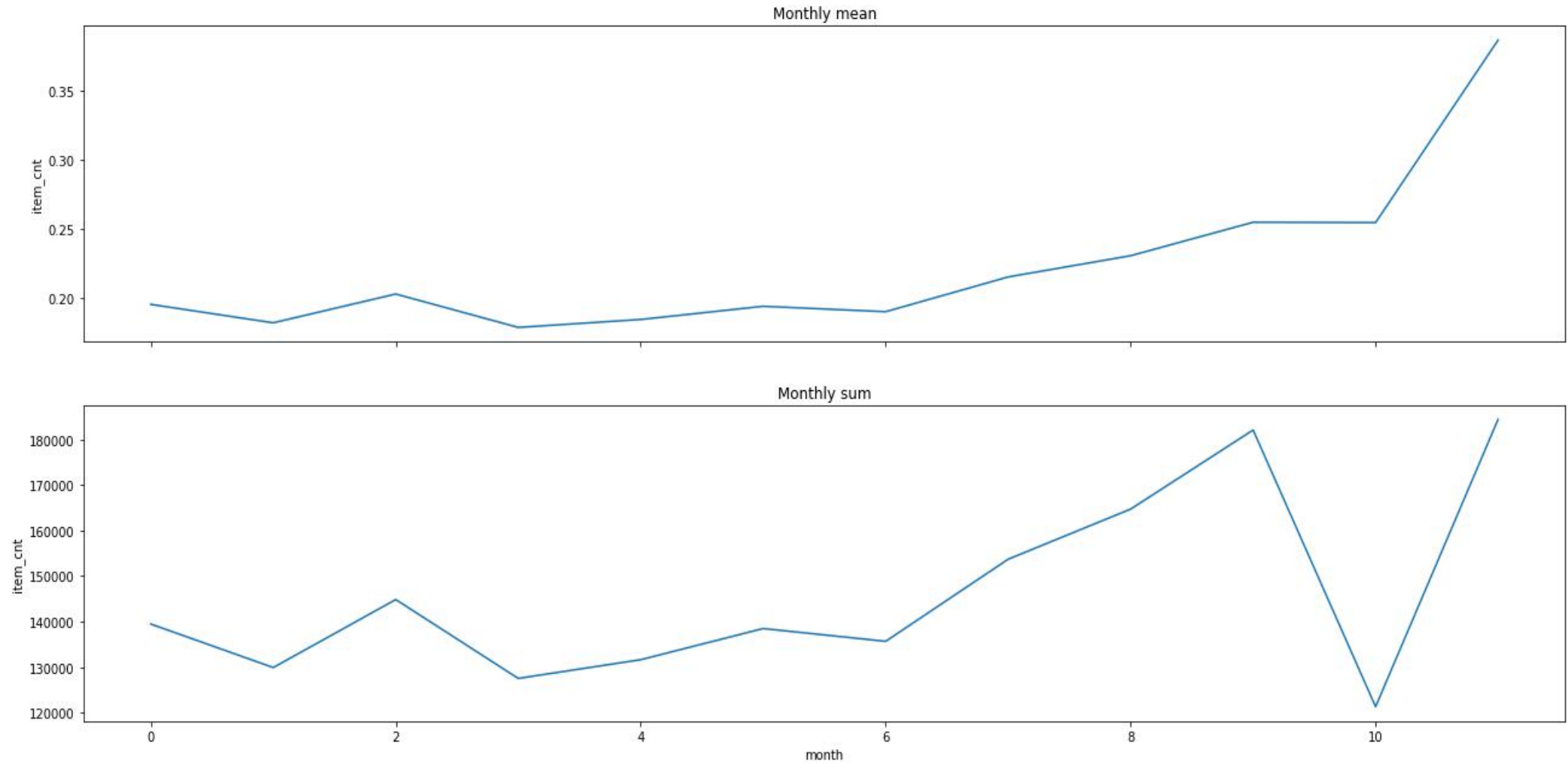


## Duplicated And Missing values

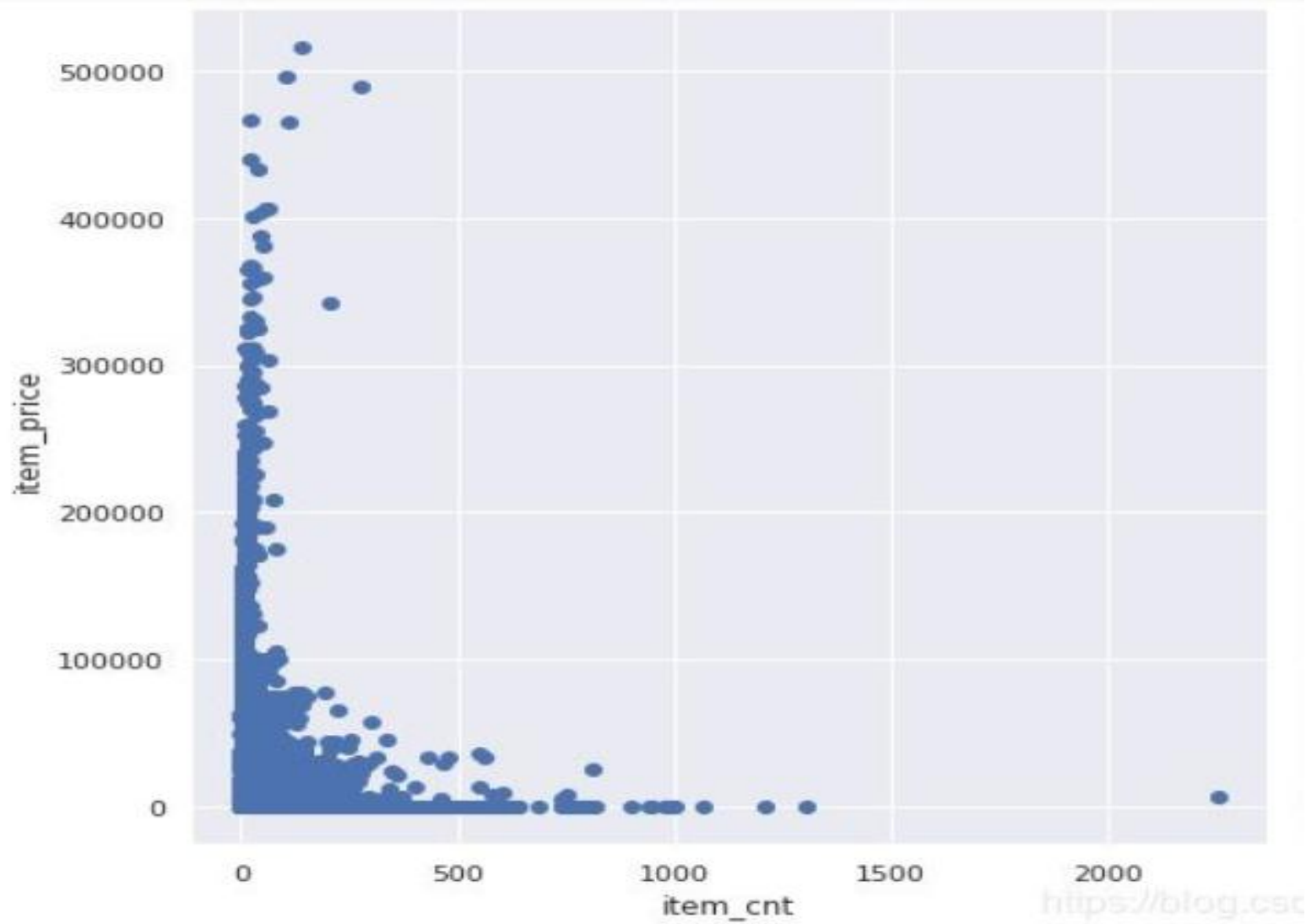
```
↳ duplicated lines in sales_train is 5  
missing value in sales_train is  
date                0  
date_block_num      0  
shop_id             0  
item_category_id    0  
item_id             0  
item_price          0  
item_cnt_day        0  
dtype: int64
```



# Item\_cnt month trend



## Item\_price VS item\_cnt



## Basic info of shops

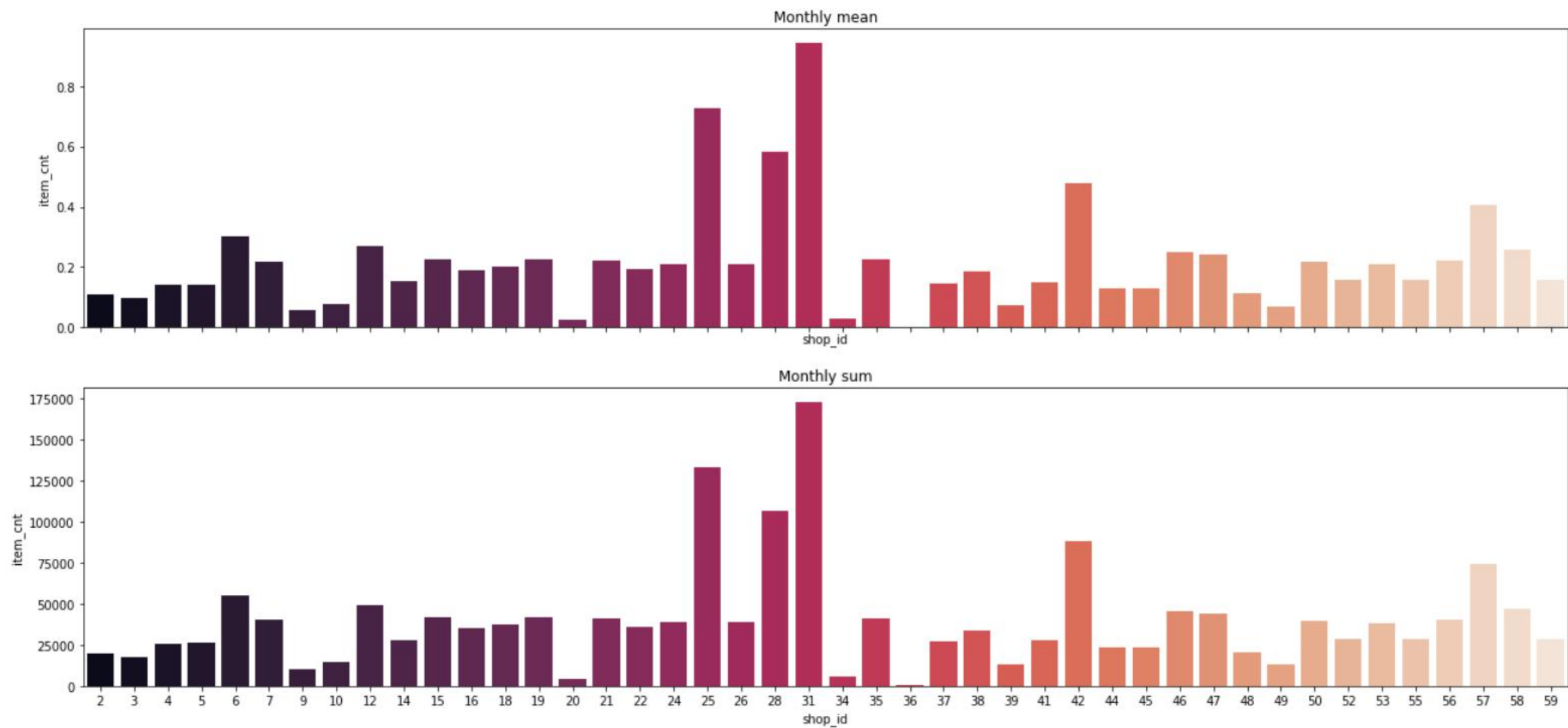
	shop_name	shop_id
0	!Якутск Орджоникидзе, 56 фран	0
1	!Якутск ТЦ "Центральный" фран	1
2	Адыгея ТЦ "Мега"	2
3	Балашиха ТРК "Октябрь-Киномир"	3
4	Волжский ТЦ "Волга Молл"	4

Shop\_name: City|type|name

Волжский ТЦ "Волга Молл"

Volga shopping center "Volga Mall"

# Shops VS item\_cnt



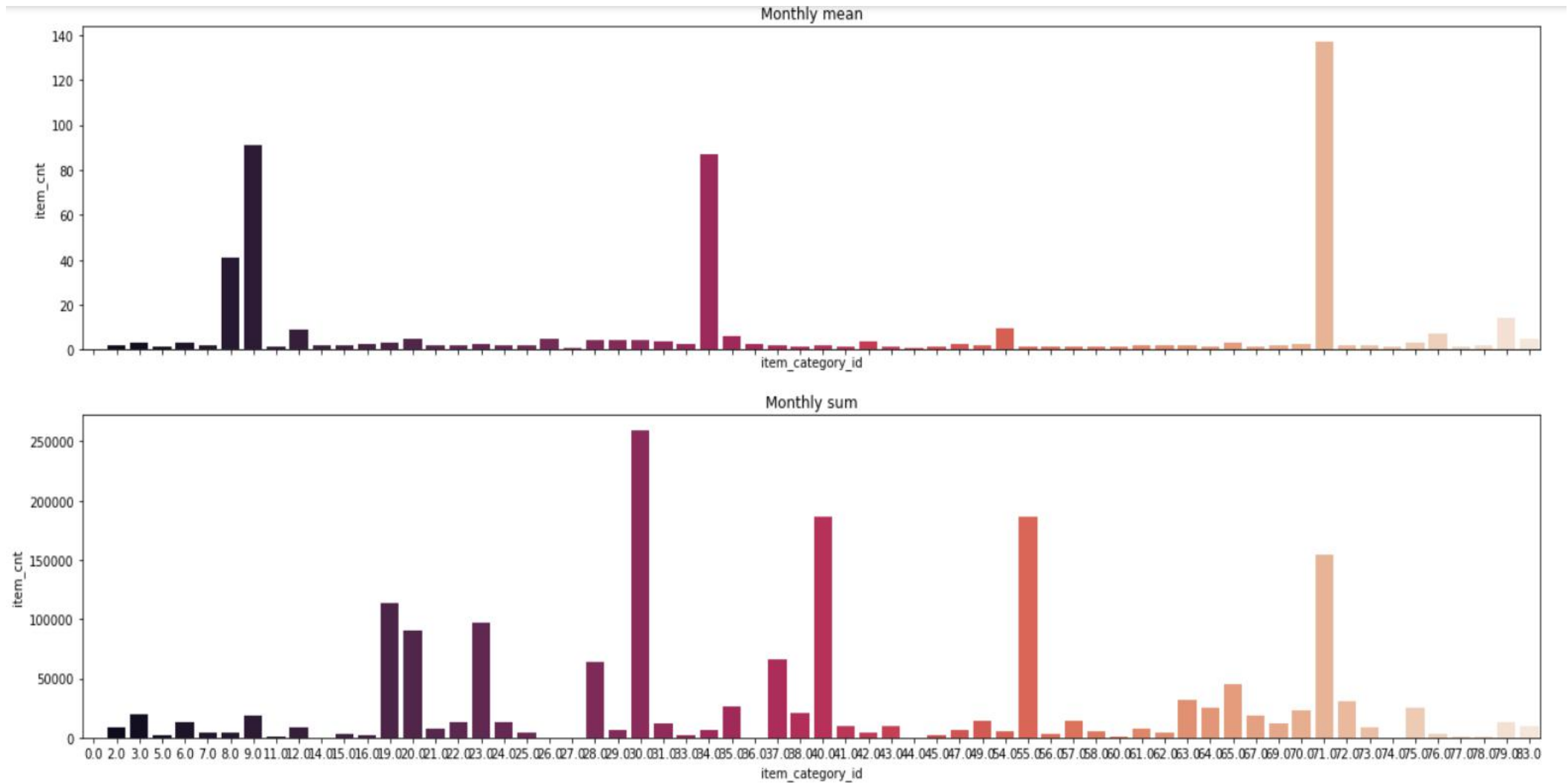
## Basic info of item\_category

	item_category_name	item_category_id
0	PC - Гарнитуры/Наушники	0
1	Аксессуары - PS2	1
2	Аксессуары - PS3	2
3	Аксессуары - PS4	3
4	Аксессуары - PSP	4

Item\_category\_name: category|sub\_cat

Служебные - Билеты	Service - Tickets
Чистые носители (шпиль)	Net carriers (spire)

# item\_category VS item\_cnt



## Basic info of items

	item_name	item_id	item_category_id
0	! ВО ВЛАСТИ НАВАЖДЕНИЯ (ПЛАСТ.) D	0	40
1	!ABBY FineReader 12 Professional Edition Full...	1	76
2	***В ЛУЧАХ СЛАВЫ (UNV) D	2	40
3	***ГОЛУБАЯ ВОЛНА (Univ) D	3	40
4	***КОРОБКА (СТЕКЛО) D	4	40

Item\_name: name|sub\_type1|sub\_type2

! ВО ВЛАСТИ НАВАЖДЕНИЯ (ПЛАСТ.) D	! VO VLASTI NAVAZHDENIYA (PLAST.) D
-----------------------------------	-------------------------------------



## Basic info of test

	ID	shop_id	item_id
<b>0</b>	0	5	5037
<b>1</b>	1	5	5320
<b>2</b>	2	5	5233
<b>3</b>	3	5	5232
<b>4</b>	4	5	5268

how many lines in train set: (2935849, 6)

unique items in train set: 21807

unique shops in train set: 60

how many lines in test set: (214200, 3)

unique items in test set: 5100

unique shops in test set: 42

- ID - an Id that represents a (Shop, Item) tuple within the test set

# Feature Engineering

- Data Preprocessing ( won't be mentioned here )
- Split features & Label Encoding
- Lag features & Mean Encoding
- Dummy variable & Other features
- Merge features

# Split features & Label Encoding

item\_cnt\_day  
item\_price

item\_name  
shop\_name  
item\_category\_name

ID  
item\_id  
shop\_id  
item\_category\_id  
date  
date\_block\_num

sum  
mean



item\_cnt\_month  
item\_price\_month

split name  
+  
Label Encoding

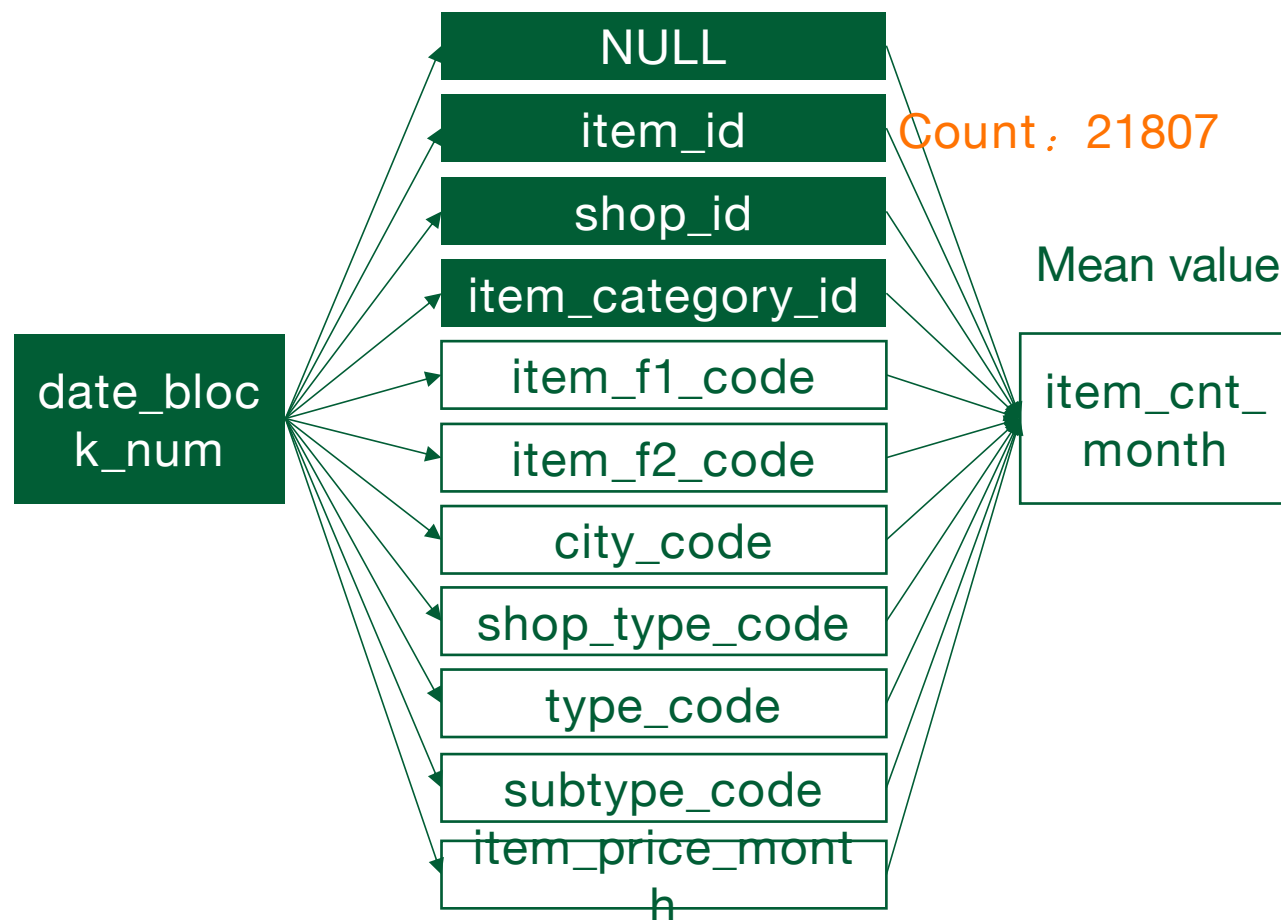


item\_f1\_code  
city\_code  
type\_code

item\_f2\_code  
shop\_type\_code  
subtype\_code

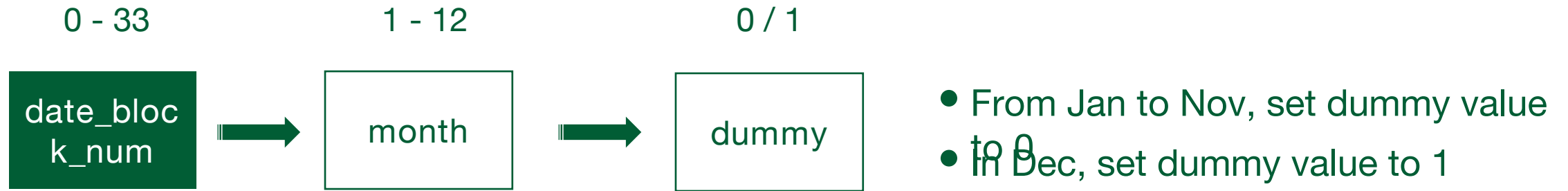
city_code	shop_type_code	item_category_id	item_f1_code	item_f2_code	type_code	subtype_code
0	5	40	729	138	10	4
0	5	19	729	138	5	10
0	5	30	729	138	8	53
0	5	23	729	138	5	16
0	5	40	729	138	10	4
...	...	...	...	...	...	...
20	5	55	729	138	12	2
20	5	64	729	138	13	41
20	5	55	729	138	12	2
20	5	40	853	138	10	4
20	5	37	364	138	10	1

# Lag features & Mean Encoding



- Mean Encoding:  
Sum “item\_cnt\_month” of the same “data\_block\_num” and the features in the center column, then divide its frequency.

# Dummy variable & Other features



hist\_min\_item\_price

- The lowest daily price per item in history

hist\_max\_item\_price

- The highest daily price per item in history

item\_first\_sale

- The first daily sales of an item in history

item\_shop\_first\_sale

- The first daily sales of a pair in history

holidays

- The number of Russian official holidays in every month

# Merge features

- Original features: 4
- Split features & Label Encoding: 8
- Lag features & Mean Encoding: 19
- Dummy variable & Other features: 7

	date_block_num	item_cnt_month	item_id	shop_id	city_code	shop_type_code	item_category_id	item_f1_code
0	0	0.0	19	2	0	5	40	729
1	0	1.0	27	2	0	5	19	729
2	0	0.0	28	2	0	5	30	729
3	0	0.0	29	2	0	5	23	729
4	0	0.0	32	2	0	5	40	729
...	...	...	...	...	...	...	...	...
11056272	34	0.0	18454	45	20	5	55	729
11056273	34	0.0	16188	45	20	5	64	729
11056274	34	0.0	15757	45	20	5	55	729
11056275	34	0.0	19648	45	20	5	40	853
11056276	34	0.0	969	45	20	5	37	364

11056276 rows x 39 columns

# 03 ALGORITHM



# ALGORITHMs

## LinearRegression

A linear approach to modeling the relationship between a scalar response and one or more explanatory variables. The case of one explanatory variable is called simple linear regression. For more than one explanatory variable, the process is called multiple linear regression.

**Not employed**

## RandomForest

An ensemble learning method for classification, regression and other tasks that operates by constructing a multitude of decision trees at training time and outputting the class that is the mode of the classes or mean prediction (regression) of the individual trees.

**employed**

## XGBoost

XGBoost implements machine learning algorithms under the Gradient Boosting framework. It provides a parallel tree boosting that solve many data science problems in a fast and accurate way.

**employed**

## LightGBM

A fast, distributed, high performance gradient boosting (GBT, GBDT, GBRT, GBM or MART) framework based on decision tree algorithms, used for ranking, classification and many other machine learning tasks

**employed**

# Prediction without Parameters Tuning

## EVALUATION METRIC :

Root Mean Square Error均方根误差

$$\text{RMSE}(X, h) = \sqrt{\frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2}$$

	Train_rmse	Valid_rmse	Test_rmse
lightGBM	0.69266	0.89392	0.91685
XGBoost	0.68734	0.89816	0.91812

# Parameters Tuning

## XGBoost Parameters:

The overall parameters have been divided into 3 categories by XGBoost authors:

1. General Parameters: Guide the overall functioning
2. Booster Parameters: Guide the individual booster (tree/regression) at each step
3. Learning Task Parameters: Guide the optimization performed



# Parameters Tuning

## 8 Booster Parameters:

eta (default=0.3, alias: learning\_rate)

Step size shrinkage used in update to prevents overfitting. After each boosting step, we can directly get the weights of new features, and eta shrinks the feature weights to make the boosting process more conservative.  
range: [0,1]

gamma (default=0, alias: min\_split\_loss)

Minimum loss reduction required to make a further partition on a leaf node of the tree. The larger gamma is, the more conservative the algorithm will be.  
range: [0,∞)

max\_depth (default=6)

Maximum depth of a tree. Increasing this value will make the model more complex and more likely to overfit. 0 is only accepted in lossguided growing policy when tree\_method is set as hist and it indicates no limit on depth. Beware that XGBoost aggressively consumes memory when training a deep tree.  
range: [0,∞)

min\_child\_weight (default=1)

Minimum sum of instance weight needed in a child. If the tree partition step results in a leaf node with the sum of instance weight less than min\_child\_weight, then the building process will give up further partitioning. In linear regression task, this simply corresponds to minimum number of instances needed to be in each node. The larger min\_child\_weight is, the more conservative the algorithm will be.  
range: [0,∞)

# Parameters Tuning

## 8 Booster Parameters:

subsample (default=1)

Subsample ratio of the training instances. Setting it to 0.5 means that XGBoost would randomly sample half of the training data prior to growing trees. and this will prevent overfitting.

Subsampling will occur once in every boosting iteration.

range: (0,1]

colsample\_bytree (default=1)

The subsample ratio of columns when constructing each tree. Subsampling occurs once for every tree constructed.

range: (0, 1], the default value of 1, and specify the fraction of columns to be subsampled.

alpha (default=0, alias: reg\_alpha)

L1 regularization term on weights. Increasing this value will make model more conservative.

lambda (default=1, alias: reg\_lambda)

L2 regularization term on weights. Increasing this value will make model more conservative.

## 1 other Parameters:

n\_estimators (alias: num\_boosting\_rounds)

# Parameters Tuning

## Tuning Methods:

- Early Stopping

```
model = xgb.train(params, d_train, 10000, watchlist, early_stopping_rounds=100, verbose_eval=10)
```

Train until valid-rmse hasn't improved in 100 rounds.

- [sklearn.model\\_selection.GridSearchCV](#) from Scikit-learn

CV :

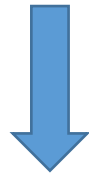
```
cv_params = {'n_estimators': [550, 575, 600, 650, 675]}
```

```
other_params = {'learning_rate': 0.1, 'max_depth': 5, 'min_child_weight': 1,  
                'colsample_bytree': 0.8, 'gamma': 0,  
                'reg_alpha': 0, 'reg_lambda': 1}
```

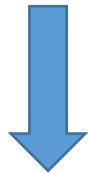
# Parameters Tuning

## Final Values:

```
params['eta'] = 0.01  
params['min_child_weight'] = 2  
params['colsample_bytree'] = 0.8  
params['subsample'] = 0.8  
params['max_depth'] = 5  
params['gamma'] = 0  
params['lambda'] = 0.3  
params['alpha'] = 0.6
```



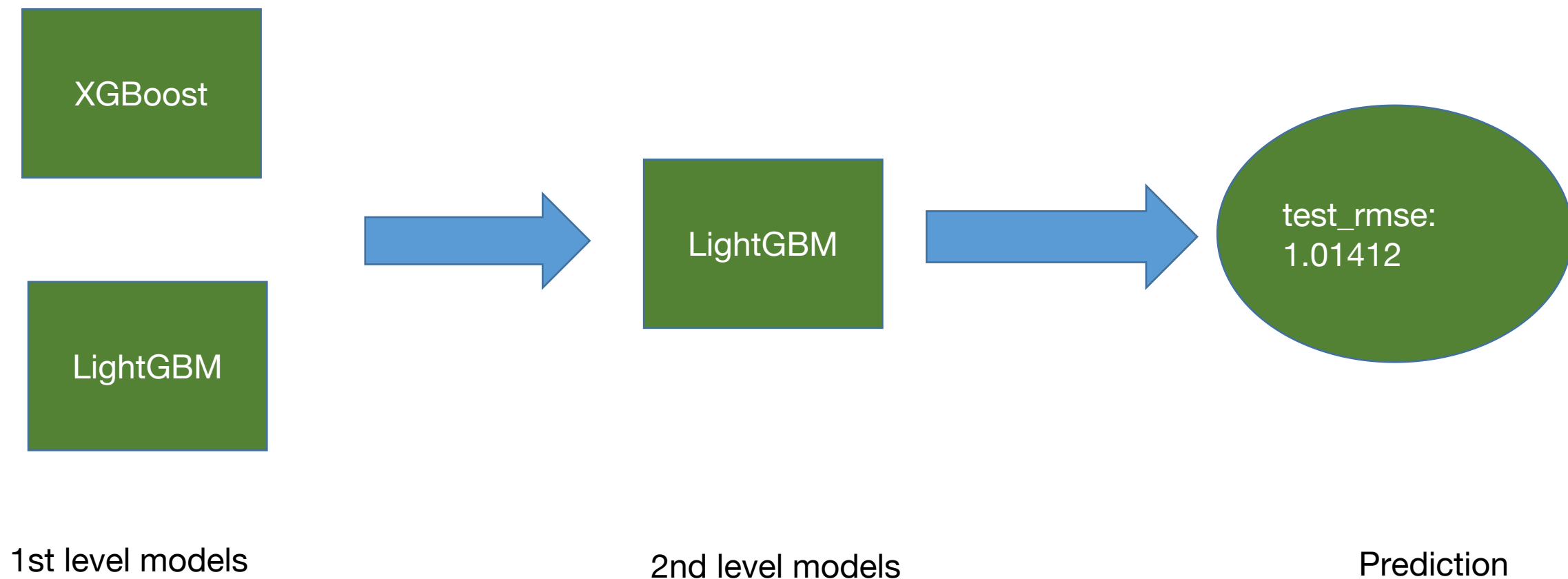
Best iteration:[850]



test\_rmse:0.90912



# Model Ensembling



# 04 RESULTS



date\_item\_avg\_item\_cnt\_lag\_1  
date\_cat\_item\_f1\_avg\_item\_cnt\_lag\_1  
item\_id  
hist\_min\_item\_price  
hist\_max\_item\_price  
item\_first\_sale  
delta\_price\_lag

# 05 CONCLUSION

THANKS FOR WATCHING