

# H1 Outline

## 1. Loading Data

## 2. Exploratory Data Analysis

- 2.1. Distribution of Target
- 2.2. Distribution of Demand by Region
- 2.3. Deal Probability by City
- 2.4. Deal Probability by Parent Category
- 2.5. Activation Date and Deal Probability
- 2.6. Duplicates in Train and Test

## 3. Pipeline

- 3.1. Cross-Decomposition of TF-IDF Vectors with BiGrams
- 3.2. Discretized Vector Cross-Decomposition
- 3.3. Discretized Vector Sums
- 3.4. Sentiment Analysis
- 3.5. Categorical Features
- 3.6. Other Features
- 3.7. Scaling Features

## 4. Evaluation

- 4.0. Baseline Linear Regression
- 4.1. PLSR 50 Components
- 4.2. Light Gradient Boosting
- 4.3. LGB with PLSR 50 Components
- 4.4. Feature Selection
- 4.5. LGB with 30 Ridge Features
- 4.6. Optimal N Features
- 4.7. PLSR by Original Feature
- 4.8. LGB with PLSR by Original Feature
- 4.9. Multi Layer Perceptron
- 4.10. Keras
- 4.11. RandomForestRegressor

## 5. Product

# H1 1. Loading Data

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4 import numpy as np
5 import cyrtranslit
```

```

6  from sklearn import preprocessing, model_selection, metrics,
   feature_selection, ensemble, linear_model, cross_decomposition,
   feature_extraction, decomposition, compose,neural_network
7  import lightgbm as lgb
8  from scipy import stats
9  import time
10 from sklearn.externals import joblib
11 import pickle
12 import os
13
14 import tensorflow as tf
15 import keras
16 from keras.models import Sequential
17 from keras.layers import Dense, Dropout, Flatten, Conv2D,
   MaxPooling2D
18 from keras.layers import LSTM, Input, TimeDistributed
19 from keras.models import Model
20 from keras.optimizers import RMSprop
21
22 # Import the backend
23 from keras import backend as K
24
25 color = sns.color_palette()
26 %matplotlib inline

```

```

1  Using TensorFlow backend.

```

```

1  train = pd.read_csv('../train.csv.zip',compression='zip',parse_dates=
   ['activation_date'])
2  test = pd.read_csv('../test.csv.zip',compression='zip',parse_dates=
   ['activation_date'])

```

## H1 2. Exploratory Data Analysis

[Back to Outline](#)

### Summary

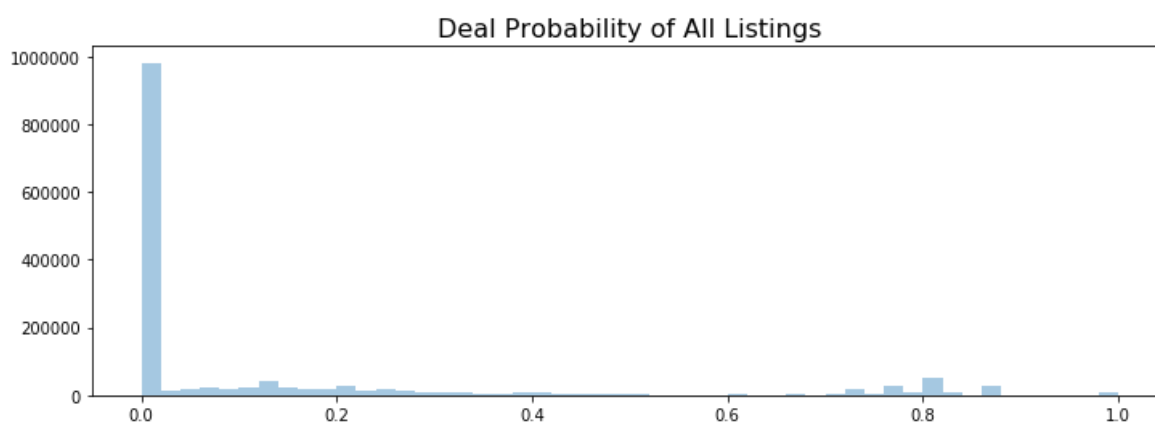
- Data is big.
- It's in Russian with Cyrillic alphabet.
- Outcome variable is very non-normal. There's three distinct groups. (Zero Range, Lower Range, Upper Range)
- Titles and descriptions offer endless NLP opportunities.
- Plenty of categorical data to binarize.

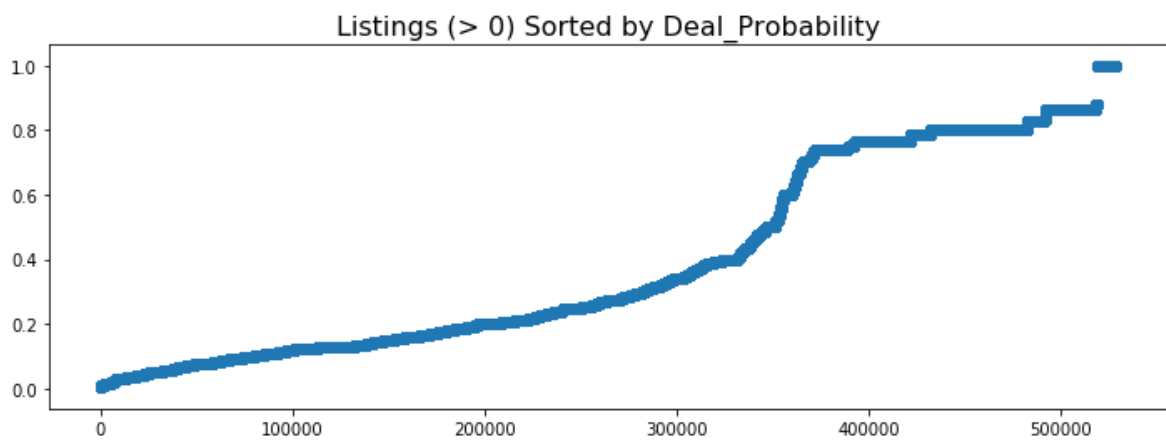
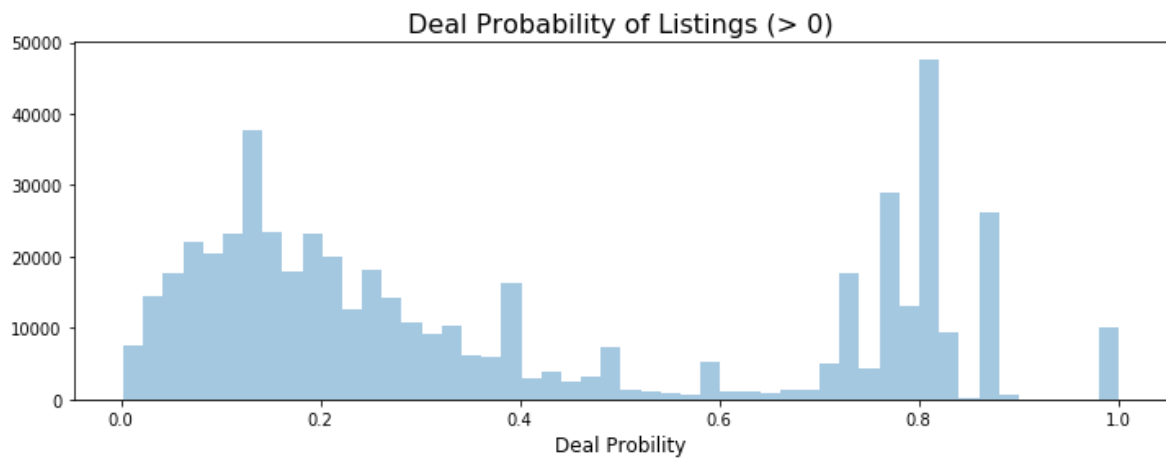
## H2 2.1. Distribution of Target

[Back to Outline](#)

There's about a million listings with zero demand, and a much smaller number with varying chance of selling.

```
1 # Define deal probabilities and those over zero
2 probs = train["deal_probability"].values
3 probs_no0 = probs[probs>0]
4
5 # Plot probability histogram
6 plt.figure(figsize=(12,4))
7 sns.distplot(probs, kde=False)
8 plt.title("Deal Probability of All Listings", fontsize=16)
9 plt.show()
10
11 # Probabilities > 0 hist
12 plt.figure(figsize=(12,4))
13 sns.distplot(probs_no0, kde=False)
14 plt.xlabel('Deal Probability', fontsize=12)
15 plt.title("Deal Probability of Listings (> 0)", fontsize=16)
16 plt.show()
17
18 # Scatter of sorted probs
19 plt.figure(figsize=(12,4))
20 plt.scatter(range(probs_no0.shape[0]), np.sort(probs_no0))
21 plt.title("Listings (> 0) Sorted by Deal_Probability", fontsize=16)
22 plt.show()
```





## H2 2.2. Distribution of Demand by Region

[Back to Outline](#)

First let's translate the regions from Russian using the `cyrtranslit` package. Then visualize the distribution of each region.

```

1  # Get unique regions in cyrilic
2  cyrilic_regs = train.region.unique().tolist()
3  # Get unique translations
4  latin_regs = [cyrtranslit.to_latin(reg, 'ru') for reg in cyrilic_regs]
5
6  # Put regions in a dictionary
7  reg_dict = {}
8  for cyr, lat in zip(cyrilic_regs, latin_regs):
9      reg_dict[cyr] = lat
10
11 # Create a translated list of each region in the dataset
12 en_list = []
13 for reg in train.region:
14     en_list.append(reg_dict[reg])
15
16 # Add english list as column
17 train['region_en'] = en_list

```

```

18
19 print('Translation of Russian Regions')
20 pd.DataFrame(latin_regs[:10], index=cyrilic_regs[:10], columns=
    ['Translations'])

```

1 Translation of Russian Regions

	Translations
Свердловская область	Sverdlovskaja oblast'
Самарская область	Samarskaja oblast'
Ростовская область	Rostovskaja oblast'
Татарстан	Tatarstan
Волгоградская область	Volgogradskaja oblast'
Нижегородская область	Nizhegorodskaja oblast'
Пермский край	Permskij kraj
Оренбургская область	Orenburgskaja oblast'
Ханты-Мансийский АО	Hanty-Mansijskij AO
Тюменская область	Tjumenskaja oblast'

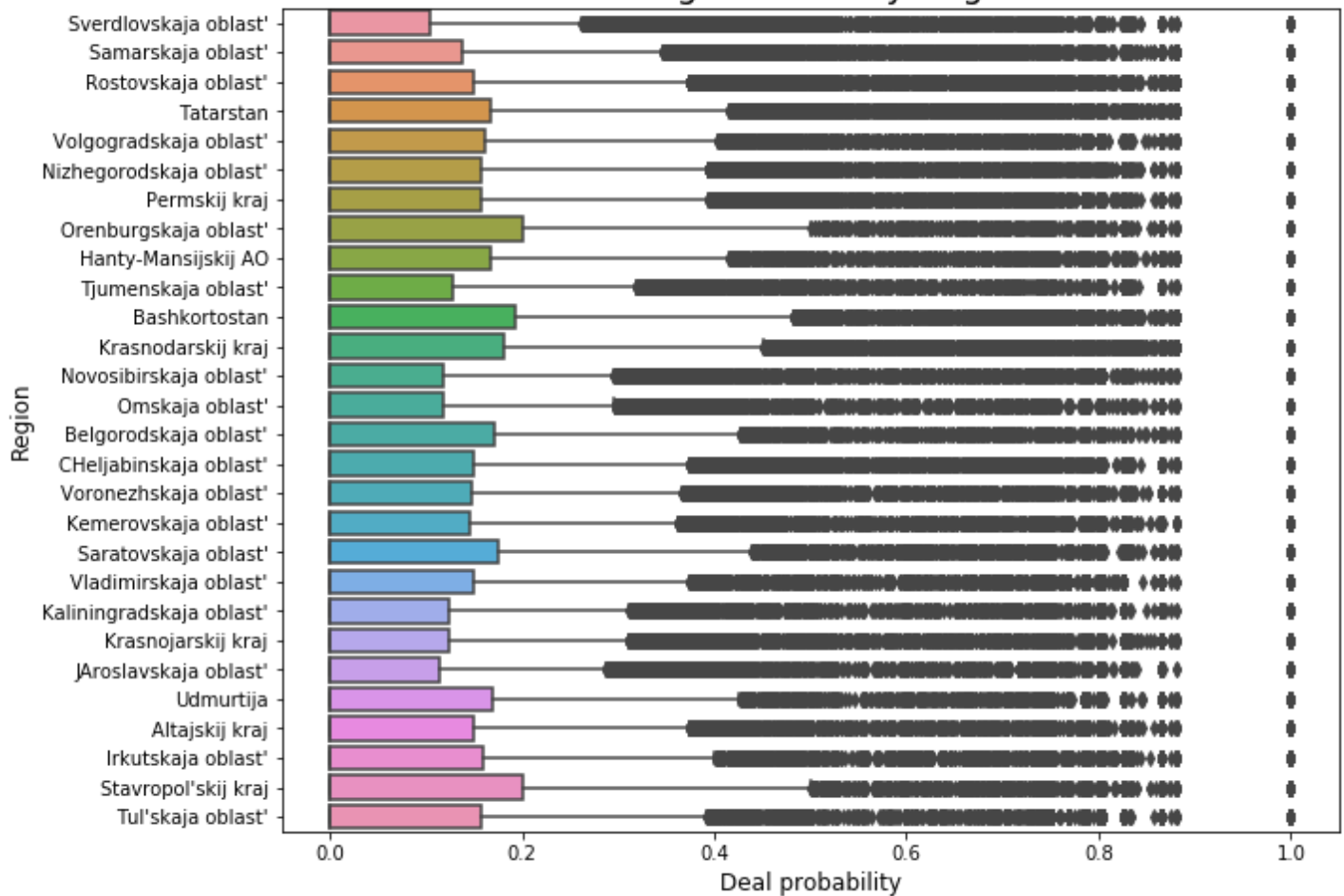
- There is a tremendous amount of outliers. Not surprising due to the large size of the data.
- Anything with over 50% chances of selling is an outlier. Based on regional data, expecting a sale is an exception rather than the norm.

```

1 # Boxplots by Region
2 plt.figure(figsize=(10,8))
3 sns.boxplot(x="deal_probability", y="region_en", data=train)
4 plt.xlabel('Deal probability', fontsize=12)
5 plt.ylabel('Region', fontsize=12)
6 plt.title("Listing Demand by Region", fontsize=18)
7 plt.show()

```

Listing Demand by Region



Percentage of Listings per Region

```

1 # Get region group counts, sort and divide by N of listings
2 region_perc =
  train.groupby('region_en').count().item_id.sort_values(ascending=False)/len(train)
3
4 # Top 5 regions
5 print('Percentage of Listings in Top Regions\n')
6 print(np.round(region_perc*100,2)[:5])
  
```

```

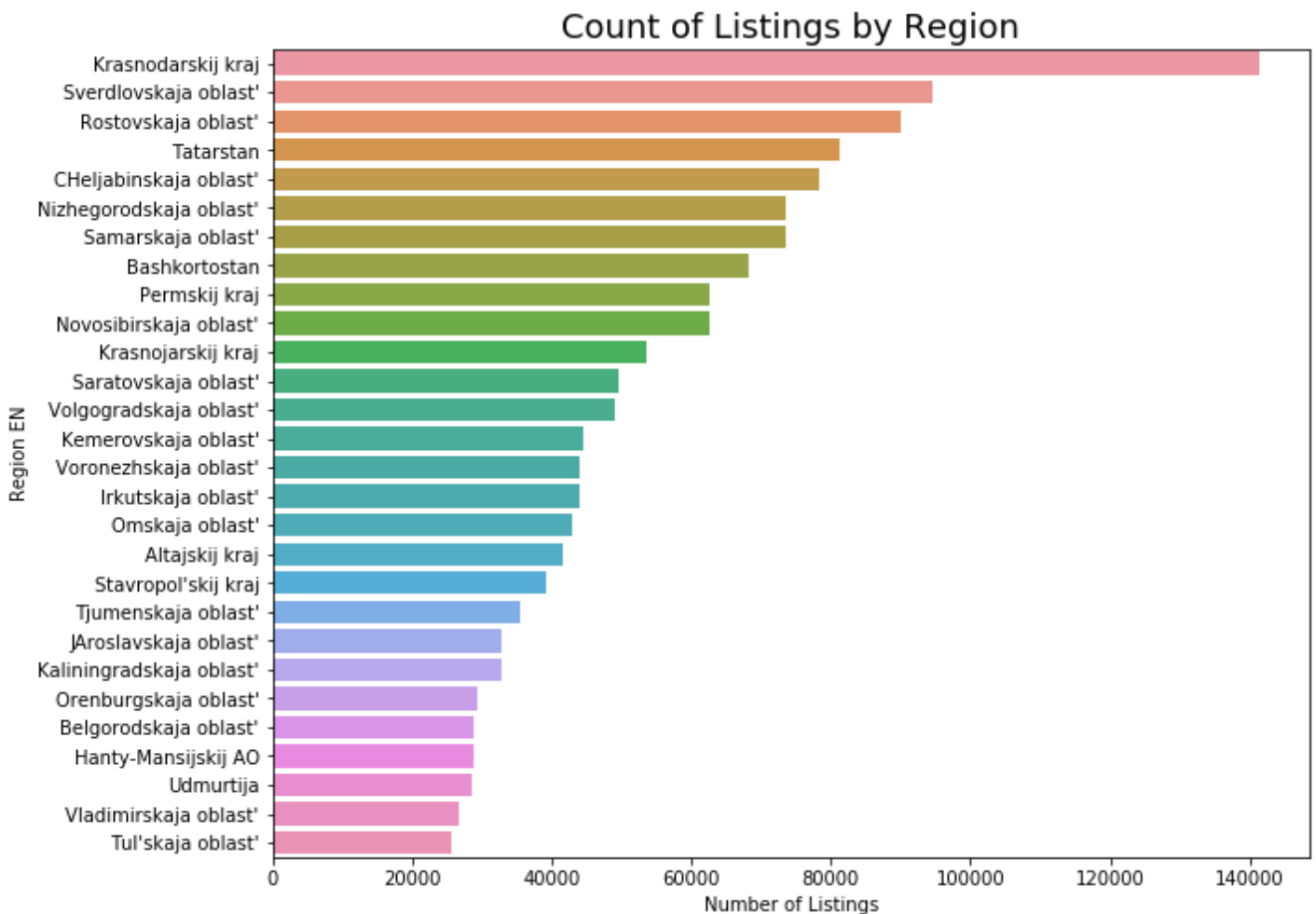
1 Percentage of Listings in Top Regions
2
3 region_en
4 Krasnodarskij kraj          9.41
5 Sverdlovskaja oblast'      6.28
6 Rostovskaja oblast'        5.99
7 Tatarstan                  5.41
8 CHeljabinskaja oblast'     5.21
9 Name: item_id, dtype: float64
  
```

Count of Listings by Region

```

1 # Visualize sorted listing counts by region
2 plt.figure(figsize=(10,8))
3 sns.countplot(data=train,y='region_en',order=region_perc.keys())
4 plt.title('Count of Listings by Region',fontsize=18)
5 plt.ylabel('Region EN')
6 plt.xlabel('Number of Listings')
7 plt.show()

```



## 2.3. Deal Probability by City

[Back to Outline](#)

```

1 # Get unique cities in cyrillic
2 cyrillic_cits = train.city.unique().tolist()
3 # Get unique translations
4 latin_cits = [cyrtranslit.to_latin(cit,'ru') for cit in cyrillic_cits]
5
6 # Put regions in a dictionary
7 cit_dict = {}
8 for cyr, lat in zip(cyrillic_cits,latin_cits):
9     cit_dict[cyr]=lat
10
11 # Create a translated list of each region in the dataset
12 en_list = []
13 for cit in train.city:

```

```

14     en_list.append(cit_dict[cit])
15
16 # Add english list as column
17 train['city_en'] = en_list
18
19 print('Translation of Russian Cities (First 10)')
20 pd.DataFrame(latin_cits[:10], index=cyrilic_cits[:10], columns=['Translations'])

```

```

1 Translation of Russian Cities (First 10)

```

	Translations
Екатеринбург	Ekaterinburg
Самара	Samara
Ростов-на-Дону	Rostov-na-Donu
Набережные Челны	Naberezhnye Chelny
Волгоград	Volgograd
Чистополь	CHistopol'
Нижний Новгород	Nizhnij Novgorod
Пермь	Perm'
Оренбург	Orenburg
Ханты-Мансийск	Hanty-Mansijsk

## Percentage of Listings per City

```

1 # Get city group counts, sort and divide by N of listings
2 city_perc =
  train.groupby('city_en').count().item_id.sort_values(ascending=False)/len(train)
3
4 # Top cities
5 print('Percentage of Listings in Top Cities\n')
6 print(np.round(city_perc*100,2)[:10])
7

```

```

1 Percentage of Listings in Top Cities
2
3 city_en
4 Krasnodar          4.23
5 Ekaterinburg       4.23
6 Novosibirsk        3.79
7 Rostov-na-Donu     3.48
8 Nizhnij Novgorod   3.46
9 CHeljabinsk        3.22

```



```

10 Perm' 3.11
11 Kazan' 3.10
12 Samara 2.79
13 Omsk 2.75
14 Name: item_id, dtype: float64

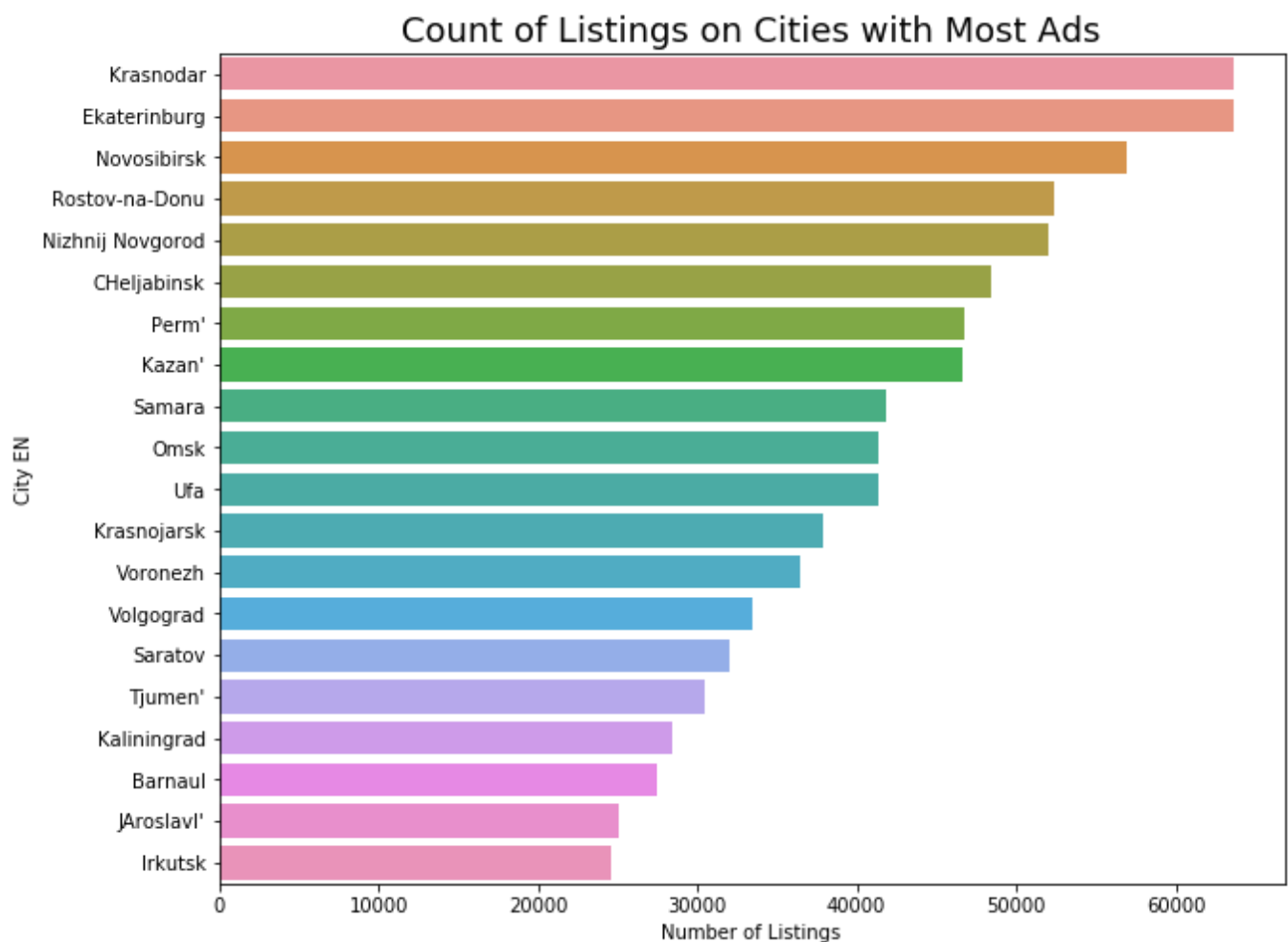
```

## Count of Listings by City

```

1 top20_cities = city_perc[:20].keys()
2
3 plot_data = train[train['city_en'].isin(top20_cities)]
4
5 # Visualize sorted listing counts by region
6 plt.figure(figsize=(10,8))
7 sns.countplot(data=plot_data,y='city_en',order=top20_cities)
8 plt.title('Count of Listings on Cities with Most Ads',fontsize=18)
9 plt.ylabel('City EN')
10 plt.xlabel('Number of Listings')
11 plt.show()

```



## 2.4. Deal Probability by Parent Category

[Back to Outline](#)

```

1 from io import StringIO

```

```

2 temp_data = StringIO("""
3 parent_category_name,parent_category_name_en
4 Личные вещи,Personal belongings
5 Для дома и дачи,For the home and garden
6 Бытовая электроника,Consumer electronics
7 Недвижимость,Real estate
8 Хобби и отдых,Hobbies & leisure
9 Транспорт,Transport
10 Услуги,Services
11 Животные,Animals
12 Для бизнеса,For business
13 """)
14
15 temp_df = pd.read_csv(temp_data)
16 train = pd.merge(train, temp_df, on="parent_category_name", how="left")
17 test = pd.merge(test, temp_df, on="parent_category_name", how="left")

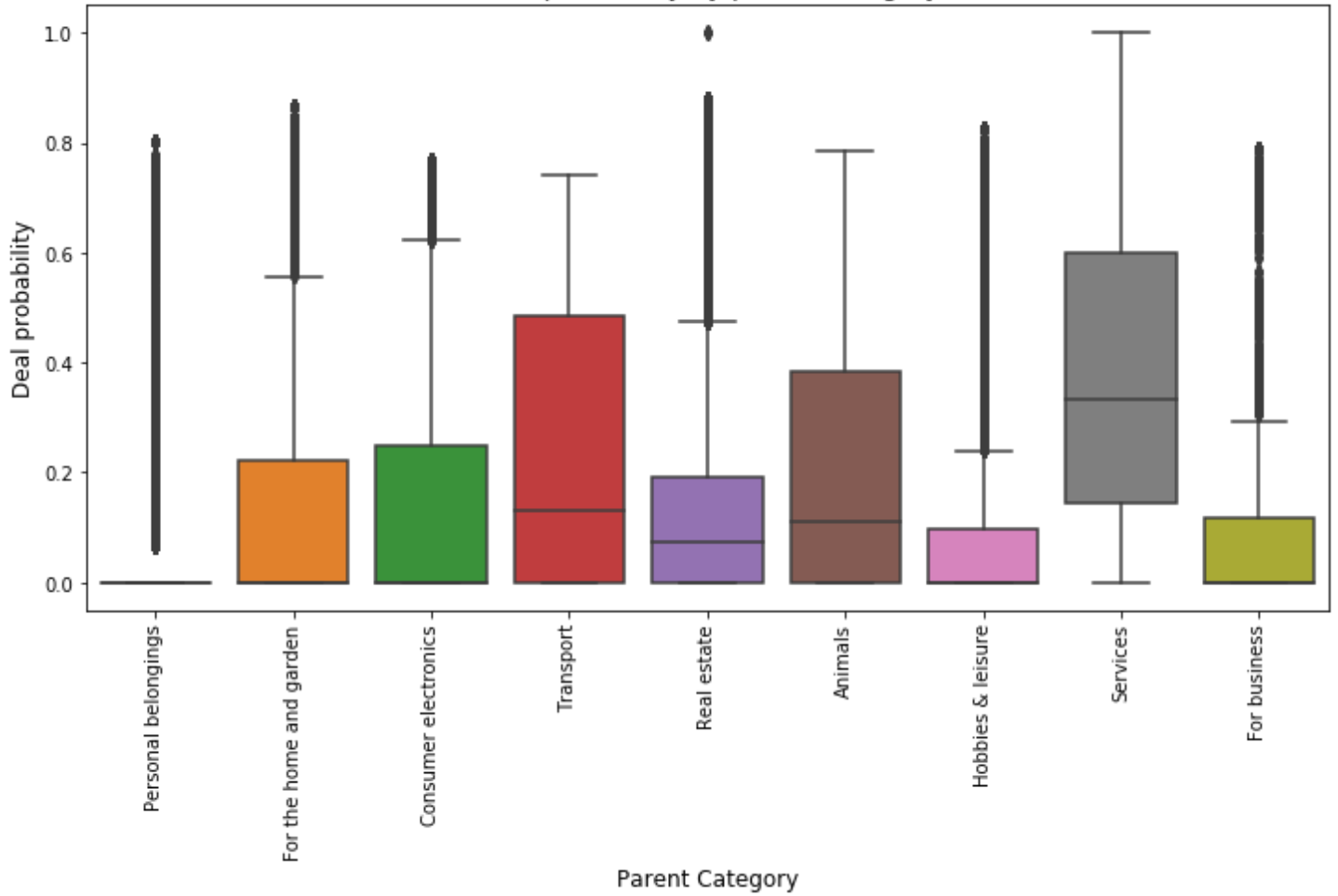
```

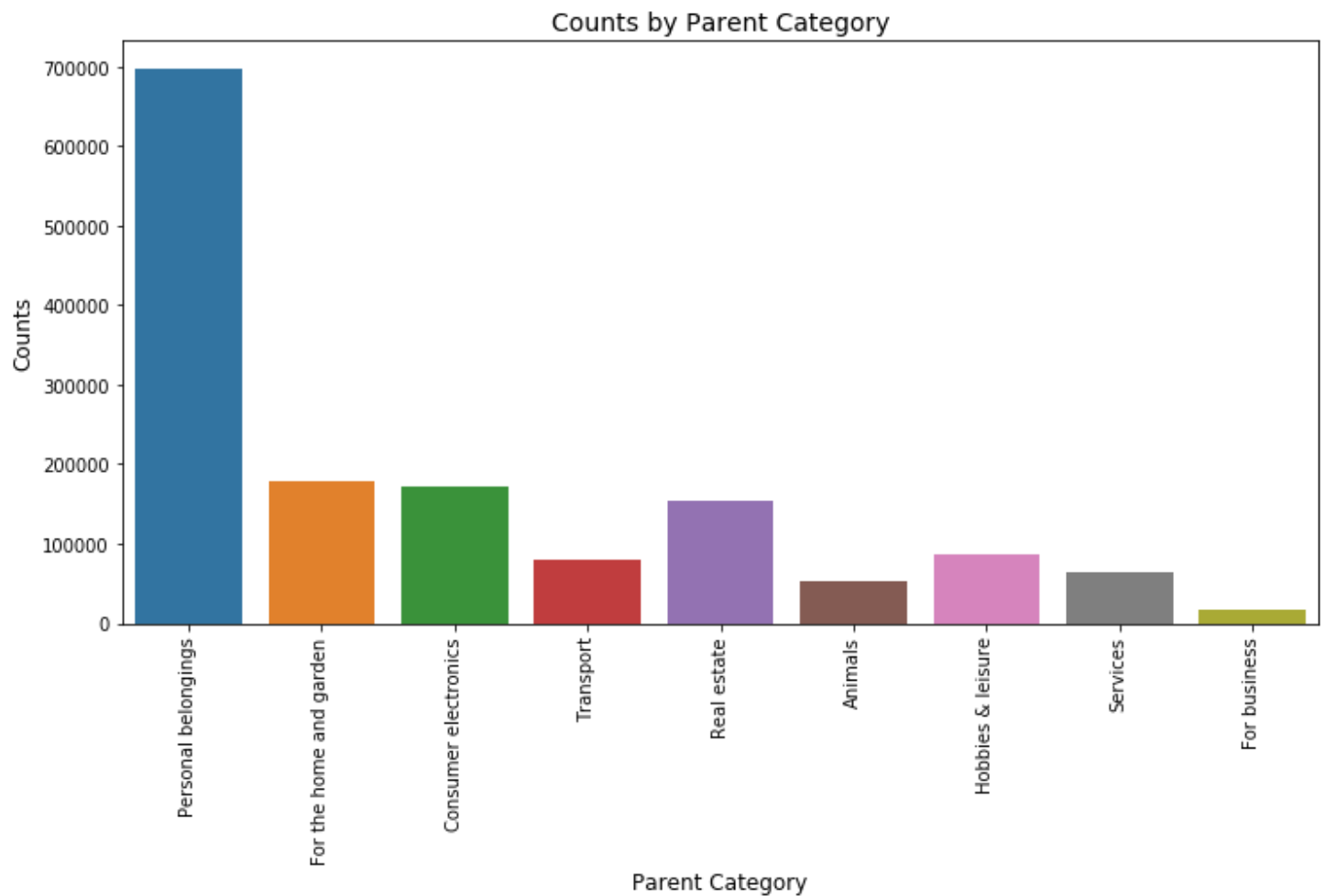
```

1 plt.figure(figsize=(12,6))
2 sns.boxplot(x="parent_category_name_en", y="deal_probability", data=train)
3 plt.ylabel('Deal probability', fontsize=12)
4 plt.xlabel('Parent Category', fontsize=12)
5 plt.title("Deal probability by parent category", fontsize=14)
6 plt.xticks(rotation='vertical')
7 plt.show()
8
9 plt.figure(figsize=(12,6))
10 sns.countplot(x="parent_category_name_en",data=train)
11 plt.ylabel('Counts', fontsize=12)
12 plt.xlabel('Parent Category', fontsize=12)
13 plt.title("Counts by Parent Category", fontsize=14)
14 plt.xticks(rotation='vertical')
15 plt.show()

```

Deal probability by parent category





## 2.5. Activation Date and Deal Probability

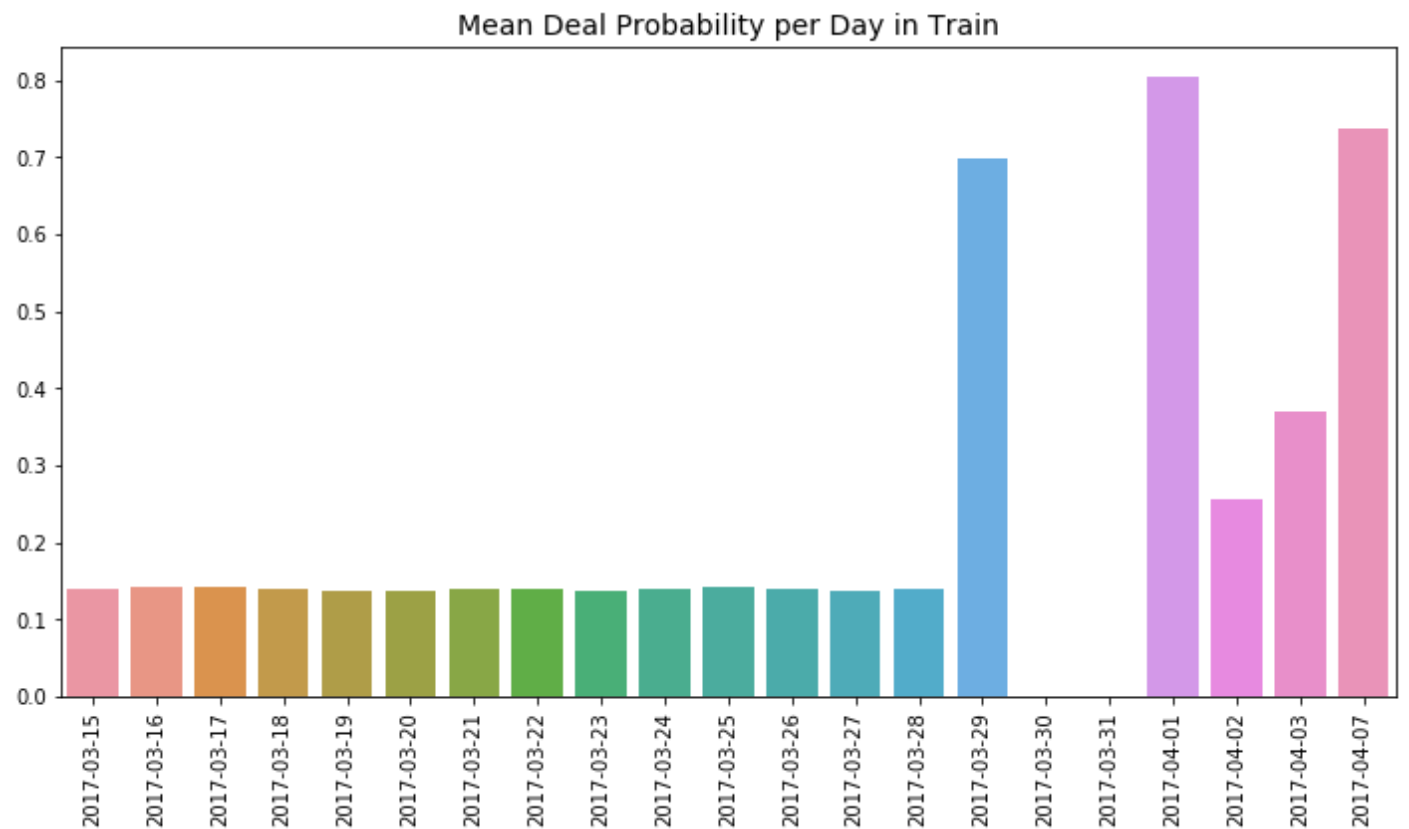
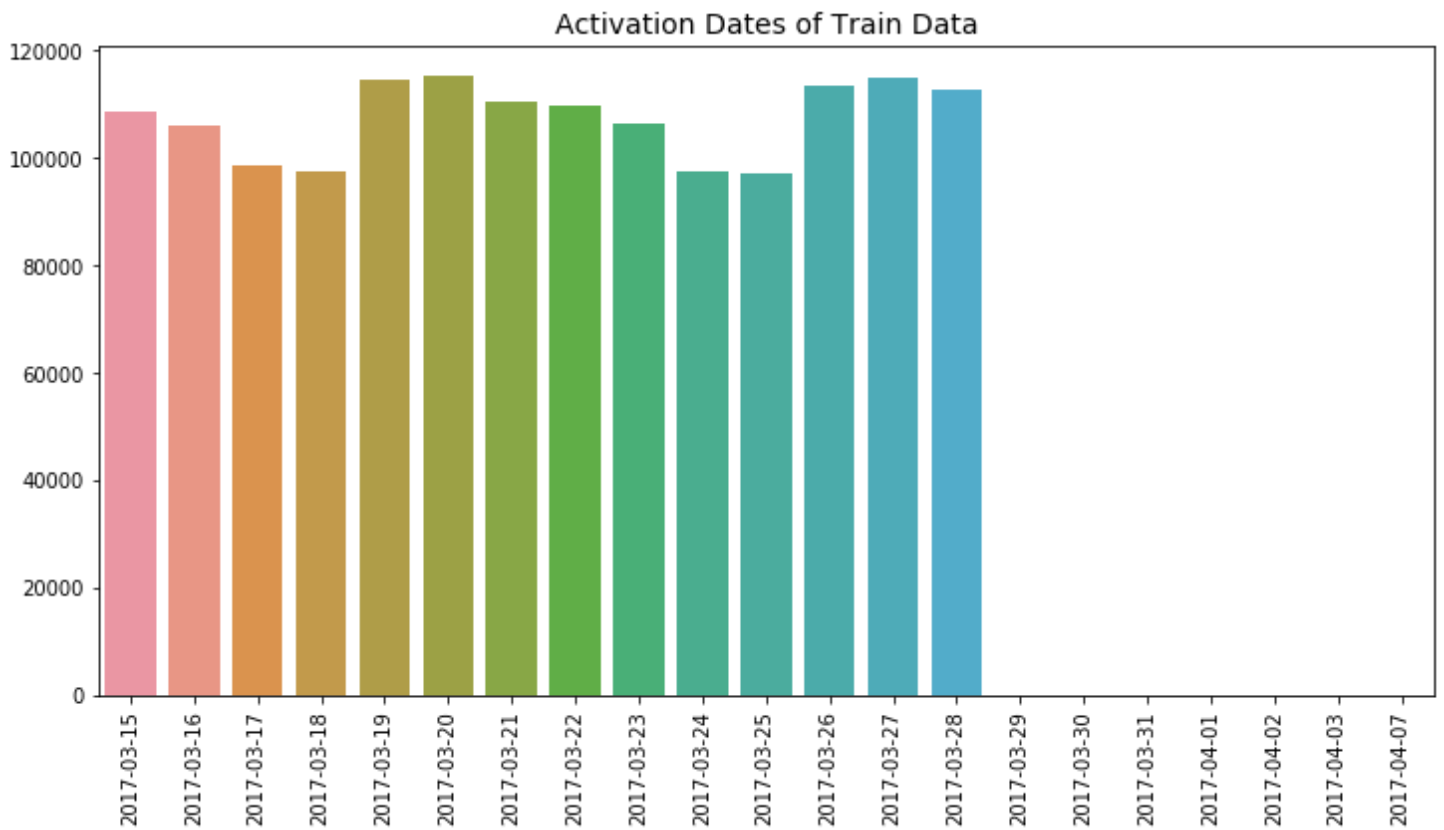
[Back to Outline](#)

```
1 srs = train.activation_date.value_counts().sort_index()
2
3 plt.figure(figsize=(12,6))
4 sns.barplot(x=srs.index.date,y=srs.values)
5 plt.xticks(rotation=90)
6 plt.title('Activation Dates of Train Data',fontsize=14)
7 plt.show()
8
9 srs = train.groupby('activation_date')['deal_probability'].mean()
10
11 plt.figure(figsize=(12,6))
12 sns.barplot(x=srs.index.date,y=srs.values)
13 plt.xticks(rotation=90)
14 plt.title('Mean Deal Probability per Day in Train',fontsize=14)
15 plt.show()
16
17 srs = test.activation_date.value_counts().sort_index()
18
19 plt.figure(figsize=(12,6))
20 sns.barplot(x=srs.index.date,y=srs.values)
21 plt.xticks(rotation=90)
```

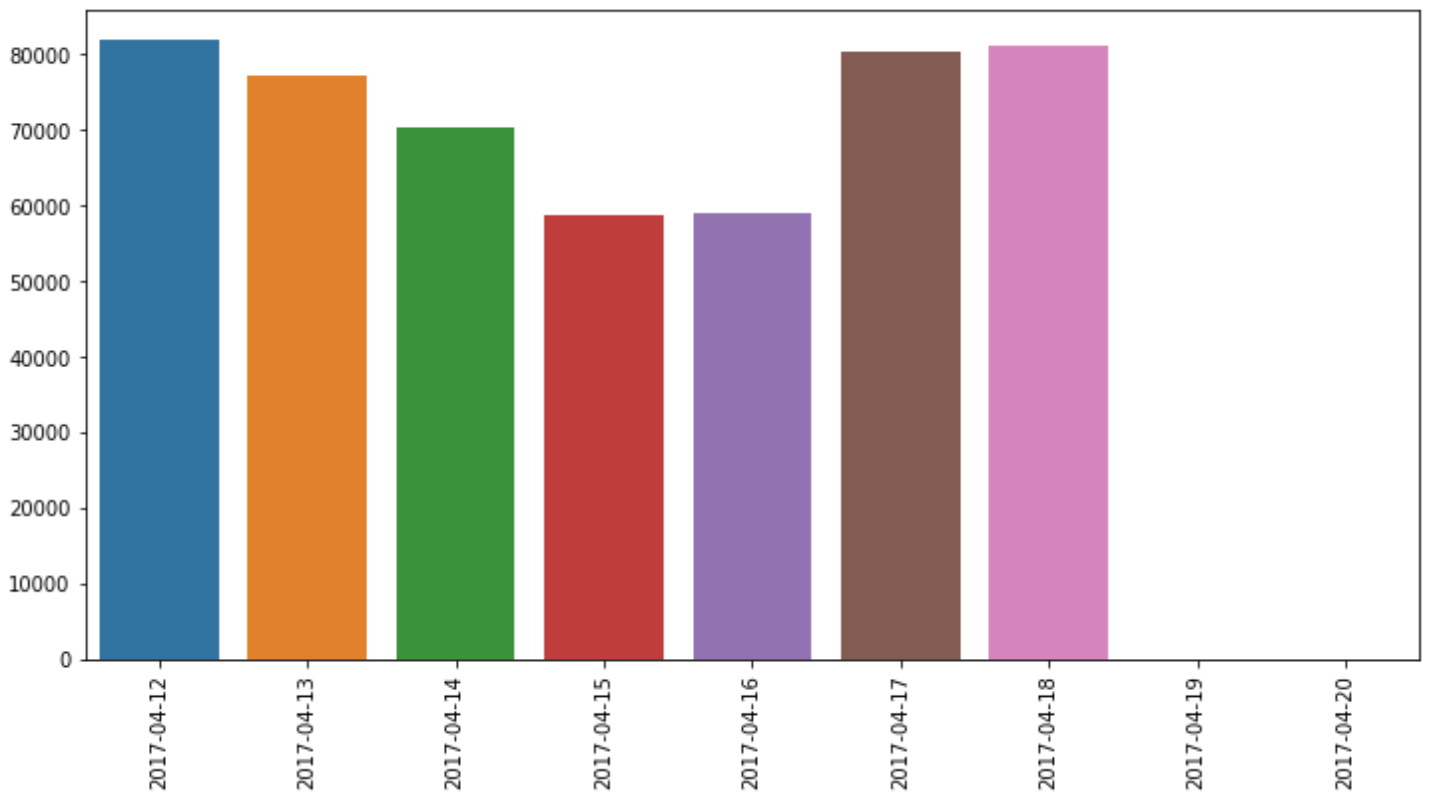
```

22 plt.title('Activation Dates of Test Data', fontsize=14)
23 plt.show()
24
25

```



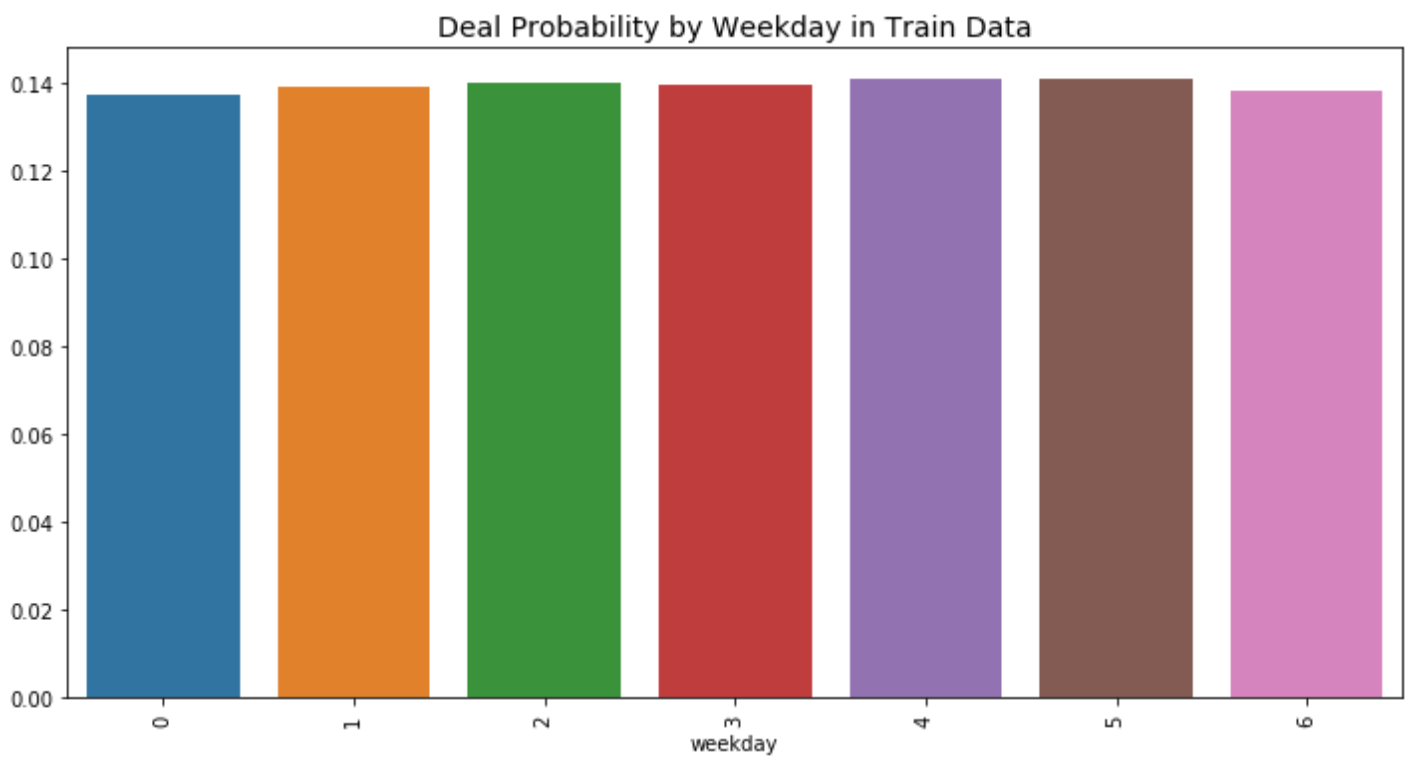
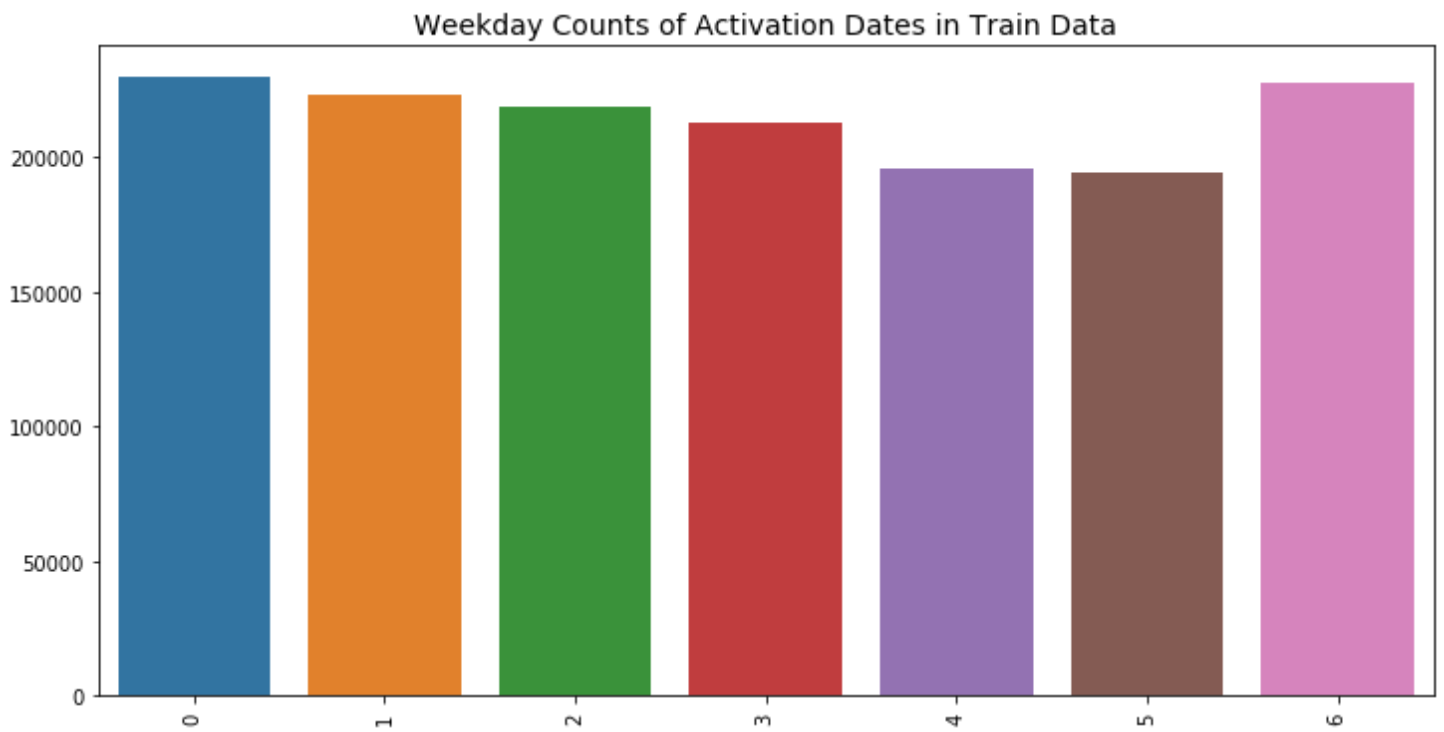
Activation Dates of Test Data



```

1 srs = train.activation_date.dt.weekday.value_counts().sort_index()
2
3 plt.figure(figsize=(12,6))
4 sns.barplot(x=srs.index,y=srs.values)
5 plt.xticks(rotation=90)
6 plt.title('Weekday Counts of Activation Dates in Train Data',fontsize=14)
7 plt.show()
8
9 train['weekday'] = train.activation_date.dt.weekday
10 srs = train.groupby('weekday').deal_probability.mean()
11
12 plt.figure(figsize=(12,6))
13 sns.barplot(x=srs.index,y=srs.values)
14 plt.xticks(rotation=90)
15 plt.title('Deal Probability by Weekday in Train Data',fontsize=14)
16 plt.show()

```



## 2.6. Duplicates in Train and Test

[Back to Outline](#)

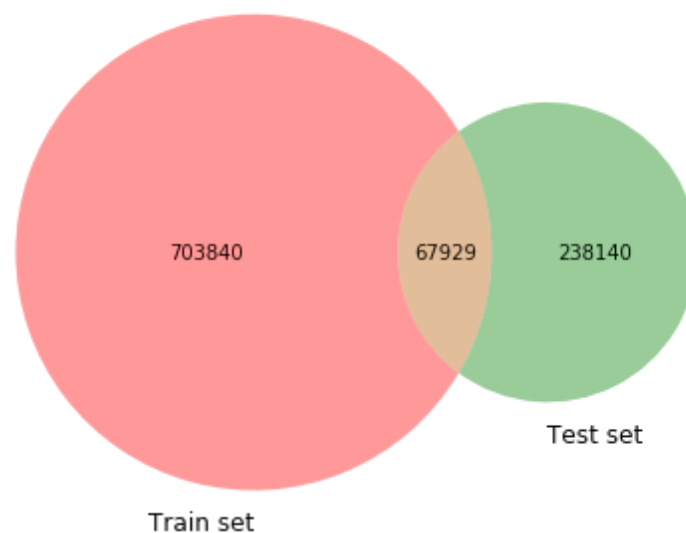
```
1 from matplotlib_venn import venn2
2
3 plt.figure(figsize=(7,6))
4 venn2([set(train.user_id), set(test.user_id)], set_labels = ('Train set', 'Test set'))
```

```

5 plt.title("Duplicate Users in Train/Test Data", fontsize=15)
6 plt.show()
7
8 plt.figure(figsize=(7,6))
9 venn2([set(train.title), set(test.title)], set_labels = ('Train set', 'Test set'))
10 plt.title("Duplicate Titles in Train/Test Data", fontsize=15)
11 plt.show()
12
13 plt.figure(figsize=(7,6))
14 venn2([set(train.item_id), set(test.item_id)], set_labels = ('Train set', 'Test set'))
15 plt.title("Duplicate Items in Train/Test Data", fontsize=15)
16 plt.show()
17
18 plt.figure(figsize=(7,6))
19 venn2([set(train.city), set(test.city)], set_labels = ('Train set', 'Test set'))
20 plt.title("Duplicate Cities in Train/Test Data", fontsize=15)
21 plt.show()
22
23 plt.figure(figsize=(7,6))
24 venn2([set(train.param_1), set(test.param_1)], set_labels = ('Train set', 'Test set'))
25 plt.title("Duplicate Param_1 Values in Train/Test Data", fontsize=15)
26 plt.show()

```

Duplicate Users in Train/Test Data





### Duplicate Titles in Train/Test Data



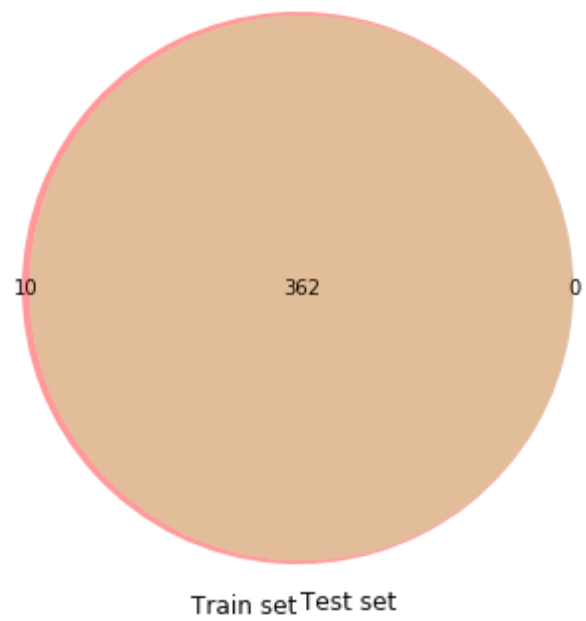
### Duplicate Items in Train/Test Data



### Duplicate Cities in Train/Test Data



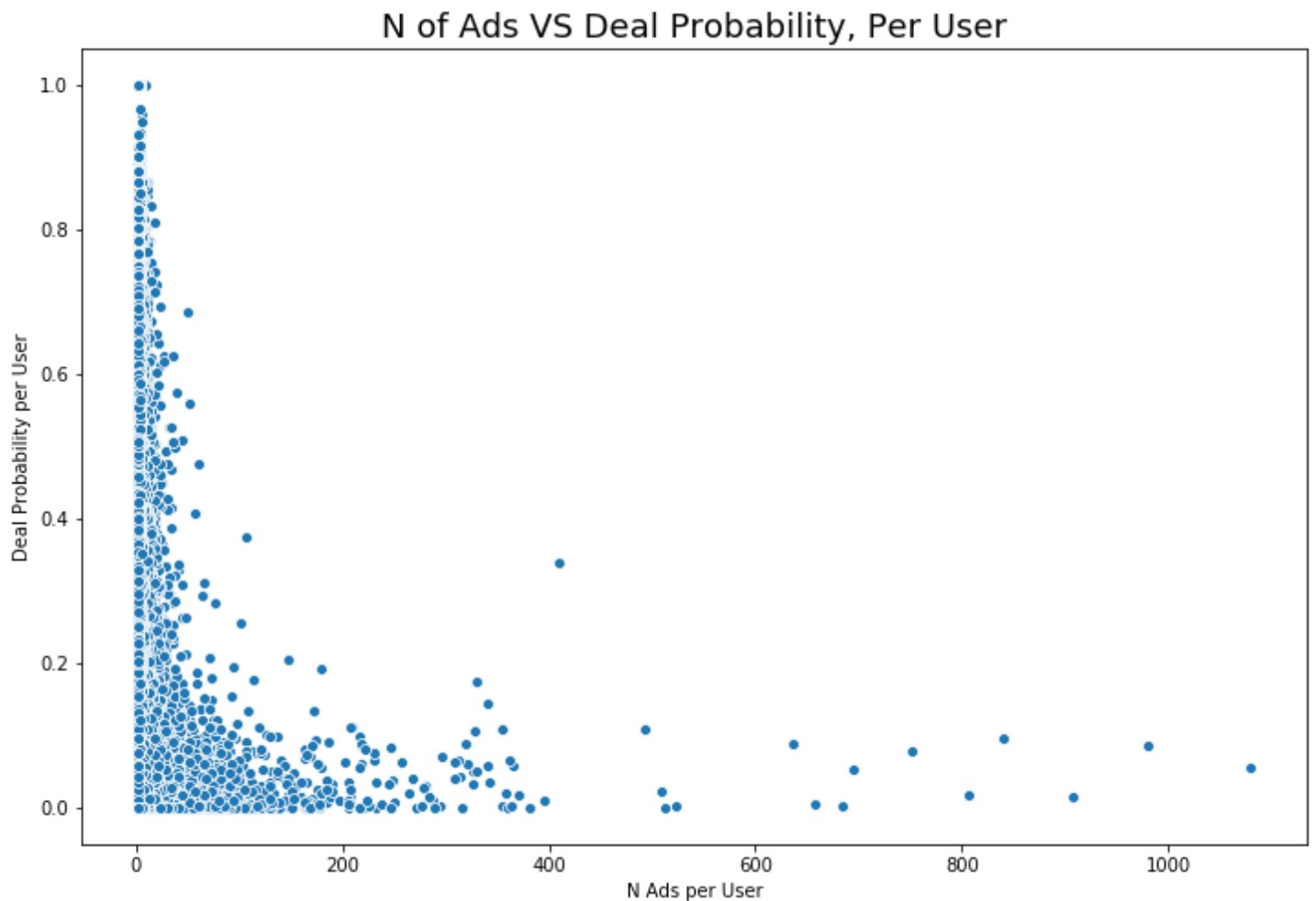
### Duplicate Param\_1 Values in Train/Test Data



```

1 user_prob = train.groupby('user_id').deal_probability.mean().sort_index()
2
3 user_items = train.user_id.value_counts().sort_index()
4
5 plt.figure(figsize=(12,8))
6 sns.scatterplot(user_items,user_prob)
7 plt.xlabel('N Ads per User')
8 plt.ylabel('Deal Probability per User')
9 plt.title('N of Ads VS Deal Probability, Per User',fontsize=18)
10 plt.show()

```



### 3. Pipeline

[Back to Outline](#)

```

1 # Load and describe engineered features
2 train_features = pd.DataFrame(index=train.index)
3 test_features = pd.DataFrame(index=test.index)

```

#### 3.1. Cross-Decomposition of TF-IDF Vectors with BiGrams

This process consists of extracting term frequency vectors using the text in each ad as a document. Tokens for unigrams and bigrams will be included in this stage. Lastly, the resulting matrix will be reduced to the smallest number of components that retain all potential predictive power. Perform onto both titles and descriptions and retain separate components for each.

```
1 path = 'feature_engineering/1.tfidf_ngrams/feature_dumps/'
2 print('Loading and joining feature sets:...')
3 for file in os.listdir(path):
4     print(file)
5     if file[:4] == 'test':
6         test_features = test_features.join(joblib.load(path+file))
7     else:
8         train_features = train_features.join(joblib.load(path+file))
```

```
1 Loading and joining feature sets:...
2 train_descr_idfngram.sav
3 train_title_idfngram.sav
4 test_descr_idfngram.sav
5 test_title_idfngram.sav
```

### 3.2. Discretized Vector Cross-Decomposition

This consists of splitting the dependent variable into discrete ranges and creating a vocabulary for each range. Then vectorize and cross-decompose each vocabulary independently. Resulting components for each vocabulary will reflect the presence of terms common in a certain discrete range of target.

```
1 path = 'feature_engineering/2.discrete-decomp/feature_dumps/'
2 print('Loading and joining feature sets:...')
3 for file in os.listdir(path):
4     print(file)
5     if file[:4] == 'test':
6         test_features = test_features.join(joblib.load(path+file))
7     else:
8         train_features = train_features.join(joblib.load(path+file))
```

```
1 Loading and joining feature sets:...
2 test_title_zeroidf.sav
3 test_title_lowcnt.sav
4 train_title_zeroidf.sav
5 train_title_lowcnt.sav
6 test_title_lowidf.sav
7 test_title_zerocnt.sav
8 train_title_zerocnt.sav
9 train_title_upidf.sav
10 train_title_upcnt.sav
11 train_title_lowidf.sav
12 test_title_upcnt.sav
13 test_title_upidf.sav
```

### 3.3. Discretized Vector Sums

Similar to previous procedure, vocabularies are created for discrete ranges of target. However instead of decomposing the vectors of those vocabularies, you simply sum their frequencies along the row axis of the term frequency matrix. This results in a single variable for each vocabulary, which represents the aggregate frequency of a vocabulary's terms per ad.

```
1 path = 'feature_engineering/3.vector-sums/feature_dumps/'
2 print('Loading and joining feature sets:...')
3 for file in os.listdir(path):
4     print(file)
5     if file[:4] == 'test':
6         test_features =
test_features.join(pd.read_pickle(path+file,compression='zip'))
7     else:
8         train_features =
train_features.join(pd.read_pickle(path+file,compression='zip'))
```

```
1 Loading and joining feature sets:...
2 test_sums.pkl
3 train_sums.pkl
```

## 3.4. Sentiment Analysis

An NLP library called `polyglot` offers multi-language tools, such as Sentiment-Analysis and Named-Entity-Recognition in Russian.

```
1 path = 'feature_engineering/4.sentiment/feature_dumps/'
2 print('Loading and joining feature sets:...')
3 for file in os.listdir(path):
4     print(file)
5     if file[:4] == 'test':
6         test_features = test_features.join(joblib.load(path+file))
7     else:
8         train_features = train_features.join(joblib.load(path+file))
```

```
1 Loading and joining feature sets:...
2 test_title_polarity.sav
3 train_title_polarity.sav
```

## 3.5. Categorical Features

### 3.5.1 Binary CountVectorizer

Several categorical variables in this data have thousands of unique values which would increase the dimensional space unreasonably if binarizing in dense format. A binary `CountVectorizer` does the heavy lifting of populating dummy counts in sparse format, and `PLSR` reduces the numerous columns to a few core components.

### 3.5.2. Target-Sorted Label Encodings

Additionally, a label encoder of each feature is made with particular considerations. Normally, label encoding isn't recommended for machine learning because the algorithm will interpret the code numbers as meaningful information. However, encodings can convey useful information if categorical values are sorted by their mean outcome value. This way, each label's code will represent an approximation of the target outcome.

```

1 path = 'feature_engineering/5.categorical/feature_dumps/'
2 print('Loading and joining feature sets:...')
3 for file in os.listdir(path):
4     print(file)
5     if file[:4] == 'test':
6         test_features = test_features.join(joblib.load(path+file))
7     else:
8         train_features = train_features.join(joblib.load(path+file))

```

```

1 Loading and joining feature sets:...
2 train_categorical.sav
3 test_categorical.sav

```

```

1 path = 'feature_engineering/5.categorical/feature_dumps_encoder/'
2 print('Loading and joining feature sets:...')
3 for file in os.listdir(path):
4     print(file)
5     if file[:4] == 'test':
6         test_features = test_features.join(joblib.load(path+file))
7     else:
8         train_features = train_features.join(joblib.load(path+file))

```

```

1 Loading and joining feature sets:...
2 train_codes.sav
3 test_codes.sav

```

```

1 dummies = train[['parent_category_name_en', 'user_type']]
2 dummies = pd.get_dummies(dummies)
3
4 train_features = train_features.join(dummies)
5
6 dummies = test[['parent_category_name_en', 'user_type']]
7 dummies = pd.get_dummies(dummies)
8
9 test_features = test_features.join(dummies)

```

### 3.6. Other Features

- Imputations
- Missing Indicators
- Day-of-Week dummies.

```

1 path = 'feature_engineering/6.other/feature_dumps/'
2 print('Loading and joining feature sets:...')
3 for file in os.listdir(path):
4     print(file)
5     if file[:4] == 'test':
6         test_features = test_features.join(joblib.load(path+file))
7     else:
8         train_features = train_features.join(joblib.load(path+file))

```

```

1 Loading and joining feature sets:...
2 train_othfeat.sav
3 test_othfeat.sa

```

### 3.7. Scaling Features

```

1 train_features = joblib.load('feature_dumps/train_features.sav')
2 test_features = joblib.load('feature_dumps/test_features.sav')

```

```

1 both = train_features.append(test_features, sort=False).reset_index(drop=True)

```

```

1 scaler = preprocessing.RobustScaler()

```

```

1 scaler.fit(both)

```

```

1 RobustScaler(copy=True, quantile_range=(25.0, 75.0), with_centering=True,
2           with_scaling=True)

```

## 4. Evaluation

[Back to Outline](#)

- Evaluation and comparison of multiple models via robust analysis of residuals and error.

```

1 train_scale = joblib.load('feature_dumps/train_scale.sav')
2 test_scale = joblib.load('feature_dumps/test_scale.sav')

```

```

1 def rmse(y, pred):
2     return metrics.mean_squared_error(y, pred)**0.5
3
4 score = metrics.make_scorer(rmse)
5 linear = linear_model.LinearRegression()

```

## 4.0. Baseline Linear Regression

```
1 preds = pd.DataFrame()
```

```
1 y = train.deal_probability
2 X_dev,X_val,y_dev,y_val = model_selection.train_test_split(train_scale,y)
```

```
1 # Baseline cv with linear regression
2 cv = model_selection.cross_val_score(
3     X=X_dev,y=y_dev,estimator=linear,
4     cv=10,scoring=score
5 )
```

```
1 print('CV Scores:\n',cv)
2 print('Mean CV score:\n',cv.mean())
```

```
1 CV Scores:
2 [0.2164809  0.21699175 0.21778808 0.21968502 0.21857161 0.21801418
3  0.21727649 0.21810469 0.21669987 0.21771957]
4 Mean CV score:
5 0.217733215497144
```

```
1 linear = linear.fit(X_dev,y_dev)
```

```
1 pred = linear.predict(X_val)
2 print('Less than zero:',sum(pred<0))
3 print('Over one:',sum(pred>1))
4 print('RMSE without modification:',metrics.mean_squared_error(y_val,pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:',metrics.mean_squared_error(y_val,pred)**0.5)
8 preds['lr'] = pred
```

```
1 Less than zero: 17953
2 Over one: 246
3 RMSE without modification: 0.2178849636631909
4 RMSE fit-to-range: 0.21752330734164146
```

```
1 pred = linear.predict(test_scale)
2 print('Less than zero:',sum(pred<0))
3 print('Over one:',sum(pred>1))
```

```
1 Less than zero: 12597
2 Over one: 163
```



```

1 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
2
3 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
4 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
5
6 sub.to_csv('predictions/sub_lr.csv', index=False)

```

```

1 !kaggle competitions submit -c avito-demand-prediction -f predictions/sub_lr.csv -m
  "Message"

```

```

1 Warning: Your Kaggle API key is readable by other users on this system! To fix this,
  you can run 'chmod 600 /home/user/.kaggle/kaggle.json'
2 100%|██████████████████████████████████████████████████████████████████████████| 15.7M/15.7M [00:29<00:00, 556kB/s]
3 Successfully submitted to Avito Demand Prediction Challenge

```

Private Score: 0.24664

Public Score: 0.24229

## 4.1. PLSR 50 Components

[Back to Outline](#)

```

1 plsr = cross_decomposition.PLSRegression(n_components=50)
2 plsr.fit(X_dev, y_dev)

```

```

1 PLSRegression(copy=True, max_iter=500, n_components=50, scale=True, tol=1e-06)

```

```

1 pred = plsr.predict(X_val)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
4 print('RMSE without modification:', metrics.mean_squared_error(y_val, pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:', metrics.mean_squared_error(y_val, pred)**0.5)
8 preds['plsr50'] = pred

```

```

1 Less than zero: [17953]
2 Over one: [247]
3 RMSE without modification: 0.21789044528359755
4 RMSE fit-to-range: 0.21752825004059145

```

```
1 pred = plsr.predict(test_scale)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
```

```
1 Less than zero: [12702]
2 Over one: [163]
```

```
1 pred = [p[0] for p in pred]
2 sub = pd.DataFrame({'item_id':test.item_id,'deal_probability':pred})
3
4 sub.loc[sub.deal_probability>1,'deal_probability'] = 1
5 sub.loc[sub.deal_probability<0,'deal_probability'] = 0
6
7 sub.to_csv('predictions/sub_plsr50.csv',index=False)
```

```
1 !kaggle competitions submit -c avito-demand-prediction -f predictions/sub_plsr50.csv -
m "Message"
```

```
1 Warning: Your Kaggle API key is readable by other users on this system! To fix this,
you can run 'chmod 600 /home/user/.kaggle/kaggle.json'
2 100%|████████████████████████████████████████| 15.7M/15.7M [00:29<00:00, 557kB/s]
3 Successfully submitted to Avito Demand Prediction Challenge
```

Private Score: 0.24661

Public Score: 0.24226

```
1 train_plsr50 = plsr.transform(train_scale)
2 test_plsr50 = plsr.transform(test_scale)
```

## 4.2. Light Gradient Boosting

[Back to Outline](#)

```
1 params = {
2     "objective" : "regression",
3     "metric" : "rmse",
4     "num_leaves" : 30,
5     "learning_rate" : 0.1,
6     "bagging_fraction" : 0.7,
7     "feature_fraction" : 0.7,
8     "bagging_frequency" : 5,
9     "bagging_seed" : 2342,
```

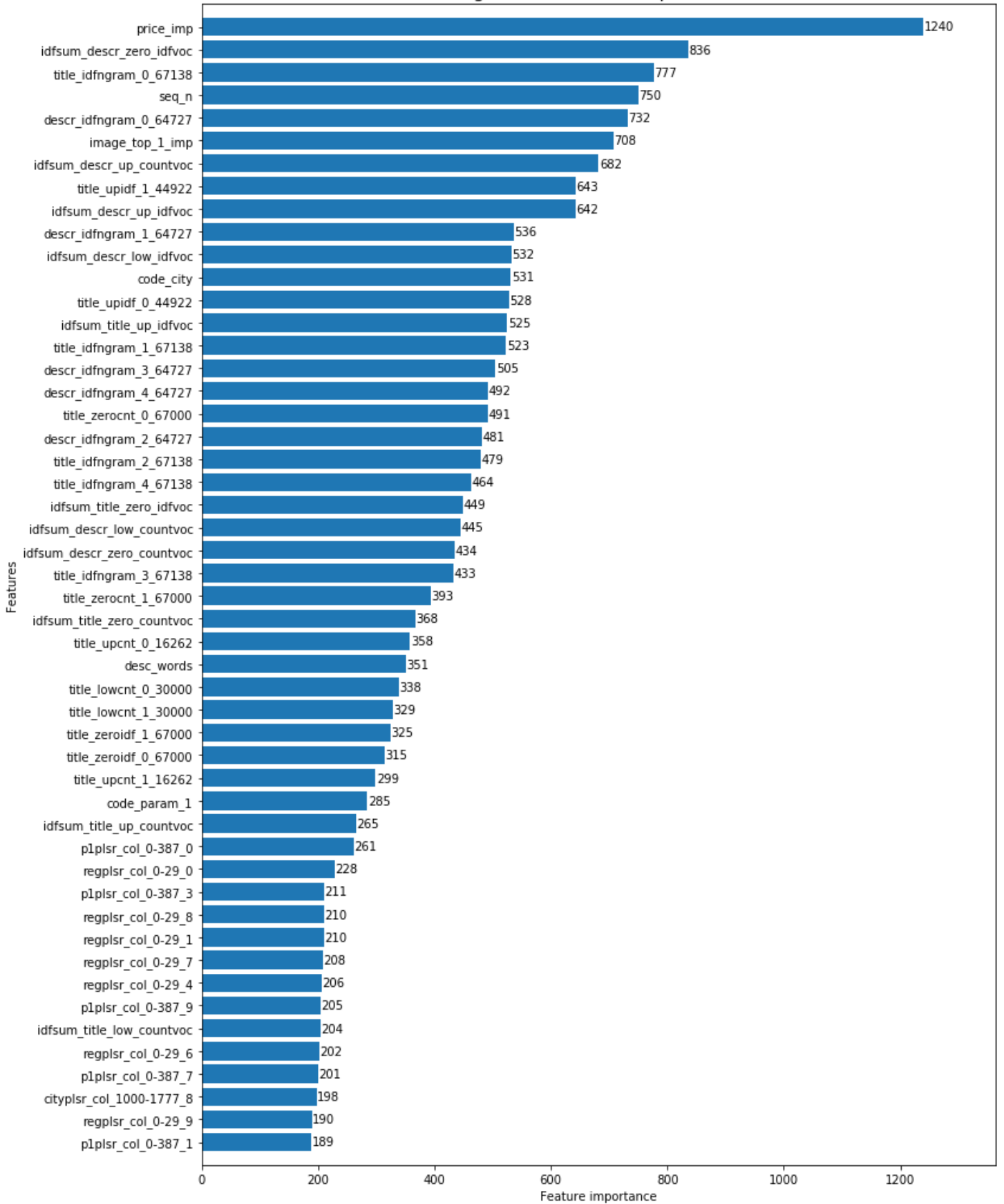
```
10     "verbosity" : -1
11 }
12
13 lgtrain = lgb.Dataset(X_dev, label=y_dev)
14 lgval = lgb.Dataset(X_val, label=y_val)
15 evals_result = {}
16
17 gb = lgb.train(params, lgtrain, 1000, valid_sets=[lgval],
18             early_stopping_rounds=100, verbose_eval=20,
19             evals_result=evals_result)
```

```
1 Training until validation scores don't improve for 100 rounds.
2 [20]    valid_0's rmse: 0.214165
3 [40]    valid_0's rmse: 0.210232
4 [60]    valid_0's rmse: 0.208477
5 [80]    valid_0's rmse: 0.207474
6 [100]   valid_0's rmse: 0.206755
7 [120]   valid_0's rmse: 0.206255
8 [140]   valid_0's rmse: 0.205848
9 [160]   valid_0's rmse: 0.205525
10 [180]   valid_0's rmse: 0.20528
11 [200]   valid_0's rmse: 0.205048
12 [220]   valid_0's rmse: 0.204862
13 [240]   valid_0's rmse: 0.20471
14 [260]   valid_0's rmse: 0.204591
15 [280]   valid_0's rmse: 0.204482
16 [300]   valid_0's rmse: 0.20438
17 [320]   valid_0's rmse: 0.2043
18 [340]   valid_0's rmse: 0.204212
19 [360]   valid_0's rmse: 0.204124
20 [380]   valid_0's rmse: 0.204032
21 [400]   valid_0's rmse: 0.203959
22 [420]   valid_0's rmse: 0.203919
23 [440]   valid_0's rmse: 0.203838
24 [460]   valid_0's rmse: 0.203796
25 [480]   valid_0's rmse: 0.20375
26 [500]   valid_0's rmse: 0.203702
27 [520]   valid_0's rmse: 0.203649
28 [540]   valid_0's rmse: 0.20362
29 [560]   valid_0's rmse: 0.203575
30 [580]   valid_0's rmse: 0.203512
31 [600]   valid_0's rmse: 0.203477
32 [620]   valid_0's rmse: 0.203442
33 [640]   valid_0's rmse: 0.203418
34 [660]   valid_0's rmse: 0.203381
35 [680]   valid_0's rmse: 0.203348
36 [700]   valid_0's rmse: 0.203299
37 [720]   valid_0's rmse: 0.203272
38 [740]   valid_0's rmse: 0.203259
39 [760]   valid_0's rmse: 0.203234
40 [780]   valid_0's rmse: 0.203218
41 [800]   valid_0's rmse: 0.203195
42 [820]   valid_0's rmse: 0.203176
```

```
43 [840] valid_0's rmse: 0.203156
44 [860] valid_0's rmse: 0.203136
45 [880] valid_0's rmse: 0.203109
46 [900] valid_0's rmse: 0.203081
47 [920] valid_0's rmse: 0.203049
48 [940] valid_0's rmse: 0.20304
49 [960] valid_0's rmse: 0.203008
50 [980] valid_0's rmse: 0.20299
51 [1000] valid_0's rmse: 0.202964
52 Did not meet early stopping. Best iteration is:
53 [1000] valid_0's rmse: 0.202964
```

```
1 fig, ax = plt.subplots(figsize=(12,18))
2 lgb.plot_importance(gb, max_num_features=50, height=0.8, ax=ax)
3 ax.grid(False)
4 plt.title("LightGBM - Feature Importance", fontsize=15)
5 plt.show()
```

LightGBM - Feature Importance



```

1 pred = gb.predict(X_val, num_iteration=gb.best_iteration)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
4 print('RMSE without modification:', metrics.mean_squared_error(y_val, pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:', metrics.mean_squared_error(y_val, pred)**0.5)
8 preds['gb'] = pred

```

```

1 Less than zero: 12674
2 Over one: 139
3 RMSE without modification: 0.20283926728961743
4 RMSE fit-to-range: 0.20276590248215942

```

```

1 pred = gb.predict(test_scale, num_iteration=gb.best_iteration)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))

```

```

1 Less than zero: 16640
2 Over one: 26

```

```

1 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
2
3 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
4 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
5
6 sub.to_csv('predictions/sub_gb.csv', index=False)

```

```

1 !kaggle competitions submit -c avito-demand-prediction -f predictions/sub_gb.csv -m
"Message"

```

```

1 Warning: Your Kaggle API key is readable by other users on this system! To fix this,
you can run 'chmod 600 /home/user/.kaggle/kaggle.json'
2 100%|████████████████████████████████████████| 15.7M/15.7M [00:20<00:00, 791kB/s]
3 Successfully submitted to Avito Demand Prediction Challenge

```

Private Score: 0.24261

Public Score: 0.23856

## 4.3. LGB with PLSR 50 Components

[Back to Outline](#)

```
preds = joblib.load('feature_dumps/preds.sav')
```

```
train_plsr50 = joblib.load('feature_dumps/train_plsr50.sav') test_plsr50 =  
joblib.load('feature_dumps/test_plsr50.sav')
```

```
1  params = {  
2      "objective" : "regression",  
3      "metric" : "rmse",  
4      "num_leaves" : 30,  
5      "learning_rate" : 0.1,  
6      "bagging_fraction" : 0.7,  
7      "feature_fraction" : 0.7,  
8      "bagging_frequency" : 5,  
9      "bagging_seed" : 2342,  
10     "verbosity" : -1  
11 }  
12  
13 lgtrain = lgb.Dataset(X_dev, label=y_dev)  
14 lgval = lgb.Dataset(X_val, label=y_val)  
15 evals_result = {}  
16  
17 gb = lgb.train(params, lgtrain, 1000, valid_sets=[lgval],  
18               early_stopping_rounds=100, verbose_eval=20,  
19               evals_result=evals_result)
```

```
1  Training until validation scores don't improve for 100 rounds.  
2  [20]    valid_0's rmse: 0.216184  
3  [40]    valid_0's rmse: 0.213355  
4  [60]    valid_0's rmse: 0.212656  
5  [80]    valid_0's rmse: 0.212232  
6  [100]   valid_0's rmse: 0.212004  
7  [120]   valid_0's rmse: 0.211837  
8  [140]   valid_0's rmse: 0.211719  
9  [160]   valid_0's rmse: 0.211625  
10 [180]   valid_0's rmse: 0.211549  
11 [200]   valid_0's rmse: 0.211446  
12 [220]   valid_0's rmse: 0.21136  
13 [240]   valid_0's rmse: 0.211308  
14 [260]   valid_0's rmse: 0.211228  
15 [280]   valid_0's rmse: 0.211167  
16 [300]   valid_0's rmse: 0.211116  
17 [320]   valid_0's rmse: 0.211073  
18 [340]   valid_0's rmse: 0.211042  
19 [360]   valid_0's rmse: 0.210998  
20 [380]   valid_0's rmse: 0.210972  
21 [400]   valid_0's rmse: 0.21095  
22 [420]   valid_0's rmse: 0.21093  
23 [440]   valid_0's rmse: 0.210885  
24 [460]   valid_0's rmse: 0.210848  
25 [480]   valid_0's rmse: 0.210811  
26 [500]   valid_0's rmse: 0.210783  
27 [520]   valid_0's rmse: 0.210758  
28 [540]   valid_0's rmse: 0.21074  
29 [560]   valid_0's rmse: 0.210708  
30 [580]   valid_0's rmse: 0.210696
```

```
31 [600] valid_0's rmse: 0.210675
32 [620] valid_0's rmse: 0.210657
33 [640] valid_0's rmse: 0.210635
34 [660] valid_0's rmse: 0.210606
35 [680] valid_0's rmse: 0.210583
36 [700] valid_0's rmse: 0.210565
37 [720] valid_0's rmse: 0.210558
38 [740] valid_0's rmse: 0.21054
39 [760] valid_0's rmse: 0.210525
40 [780] valid_0's rmse: 0.210518
41 [800] valid_0's rmse: 0.210499
42 [820] valid_0's rmse: 0.210491
43 [840] valid_0's rmse: 0.210483
44 [860] valid_0's rmse: 0.210474
45 [880] valid_0's rmse: 0.210475
46 [900] valid_0's rmse: 0.210453
47 [920] valid_0's rmse: 0.210449
48 [940] valid_0's rmse: 0.210443
49 [960] valid_0's rmse: 0.210436
50 [980] valid_0's rmse: 0.210428
51 [1000] valid_0's rmse: 0.210428
52 Did not meet early stopping. Best iteration is:
53 [994] valid_0's rmse: 0.210425
```

```
1 pred = gb.predict(X_val, num_iteration=gb.best_iteration)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
4 print('RMSE without modification:', metrics.mean_squared_error(y_val, pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:', metrics.mean_squared_error(y_val, pred)**0.5)
8 preds['gb_plsr50'] = pred
```

```
1 Less than zero: 4996
2 Over one: 50
3 RMSE without modification: 0.21042460691349882
4 RMSE fit-to-range: 0.2104147368480157
```

```
1 pred = gb.predict(test_plsr50)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
```

```
1 Less than zero: 4541
2 Over one: 14
```



```

1 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
2
3 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
4 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
5
6 sub.to_csv('predictions/sub_gb_plsr50.csv', index=False)

```

```

1 !kaggle competitions submit -c avito-demand-prediction -f
  predictions/sub_gb_plsr50.csv -m "Message"

```

Private Score: 0.24795

Public Score: 0.24364

## 4.4. Feature Selection

[Back to Outline](#)

```

1 def rmse(y, pred):
2     return metrics.mean_squared_error(y, pred)**0.5
3
4 score = metrics.make_scorer(rmse)
5 linear = linear_model.LinearRegression()

```

```

1 results = pd.DataFrame(index=
  ['Importances', 'CoefLasso', 'CoefRidge', 'FRegression'], columns=['Score', 'Selector'])

```

```

1 n_features = 30
2 data = train_scale
3 model = linear_model.LinearRegression()

```

```

1 # Score of SelectFromModel on Tree-based feature importances
2 selector1 = feature_selection.SelectFromModel(
3     ensemble.ExtraTreesRegressor(),
4     threshold=-np.inf,
5     max_features=n_features)
6
7 # Undersample dataset to reduce time
8 index = np.random.choice(len(train), size=int(4e5))
9 selector1.fit(data.iloc[index], train.iloc[index].deal_probability)
10 selection = data.iloc[:, selector1.get_support()]
11
12 # Score of these features
13 cv = model_selection.cross_val_score(model, selection, y, cv=5, scoring=score)
14 print('CV Scores:', cv)
15 print('Selection by Tree Importances:', np.mean(cv))
16 results.loc['Importances', 'Score'] = np.mean(cv)
17 results.loc['Importances', 'Selector']=selector1

```

```
1 /home/user/anaconda3/lib/python3.7/site-packages/sklearn/ensemble/forest.py:246:
FutureWarning: The default value of n_estimators will change from 10 in version 0.20
to 100 in 0.22.
2 "10 in version 0.20 to 100 in 0.22.", FutureWarning)
```

```
1 CV Scores: [0.21981017 0.22058937 0.21954811 0.22004812 0.22781175]
2 Mean CV Score: 0.2215615026676442
```

```
1 # Score of SelectFromModel from Lasso coeffs
2 selector2 = feature_selection.SelectFromModel(
3     linear_model.Lasso(),
4     threshold=-np.inf,
5     max_features=n_features)
6 selector2.fit(data, train.deal_probability)
7 selection = data.iloc[:, selector2.get_support()]
8
9 # Score of these features
10 cv = model_selection.cross_val_score(model, selection, y, cv=5, scoring=score)
11 print('CV Scores:', cv)
12 print('Selection by Lasso Coefs:', np.mean(cv))
13 results.loc['CoefLasso', 'Score'] = np.mean(cv)
14 results.loc['CoefLasso', 'Selector'] = selector2
```

```
1 CV Scores: [0.22494939 0.22541752 0.224815 0.22524057 0.23888665]
2 Selection by Coefs: 0.22786182658485582
```

```
1 # Score of SelectKBest from f_regression
2 selector3 = feature_selection.SelectKBest(
3     feature_selection.f_regression,
4     k=n_features)
5 selector3.fit(data, train.deal_probability)
6 selection = data.iloc[:, selector3.get_support()]
7 # Score of these features
8 cv = model_selection.cross_val_score(model, selection, y, cv=5, scoring=score)
9 print('CV Scores:', cv)
10 print('Selection by f_regression:', np.mean(cv))
11 results.loc['FRegression', 'Score'] = np.mean(cv)
12 results.loc['FRegression', 'Selector'] = selector3
```

```
1 CV Scores: [0.21925924 0.22004998 0.21892292 0.21950854 0.21973396]
2 Selection by f_regression: 0.21949492619756367
```

```
1 # Score of SelectFromModel from Ridge coeffs
2 selector4 = feature_selection.SelectFromModel(
```

```

3     linear_model.Ridge(),
4     threshold=-np.inf,
5     max_features=n_features)
6 selector4.fit(data, train.deal_probability)
7 selection = data.iloc[:, selector4.get_support()]
8
9 # Score of these features
10 cv = model_selection.cross_val_score(model, selection, y, cv=5, scoring=score)
11 print('CV Scores:', cv)
12 print('Selection by Ridge Coefs:', np.mean(cv))
13 results.loc['CoefRidge', 'Score'] = np.mean(cv)
14 results.loc['CoefRidge', 'Selector'] = selector4

```

```

1 CV Scores: [0.21923302 0.21989553 0.21890318 0.21946869 0.2197511 ]
2 Selection by Ridge Coefs: 0.2194503027900982

```

```

1 display(results.sort_values(by='Score', ascending=True))
2 best = results.sort_values(by='Score', ascending=True).iloc[0, 1]

```

	Score	Selector
<b>CoefRidge</b>	0.21945	SelectFromModel(estimator=Ridge(alpha=1.0, cop...
<b>FRegression</b>	0.219495	SelectKBest(k=30, score_func=<function f_regre...
<b>Importances</b>	0.221562	SelectFromModel(estimator=ExtraTreesRegressor(...
<b>CoefLasso</b>	0.227862	SelectFromModel(estimator=Lasso(alpha=1.0, cop...

```

1 # Keep the best features
2 train_sel30 = train_scale.loc[:, best.get_support()]
3 test_sel30 = test_scale.loc[:, best.get_support()]

```

```

1 joblib.dump(train_sel30, 'feature_dumps/train_sel30.sav')
2 joblib.dump(test_sel30, 'feature_dumps/test_sel30.sav')

```

## 4.5. LGB with 30 Ridge Features

[Back to Outline](#)

preds = joblib.load('feature\_dumps/preds.sav')

train\_plsr50 = joblib.load('feature\_dumps/train\_plsr50.sav') test\_plsr50 =  
joblib.load('feature\_dumps/test\_plsr50.sav')

```

1 params = {

```

```

2     "objective" : "regression",
3     "metric" : "rmse",
4     "num_leaves" : 30,
5     "learning_rate" : 0.1,
6     "bagging_fraction" : 0.7,
7     "feature_fraction" : 0.7,
8     "bagging_frequency" : 5,
9     "bagging_seed" : 2342,
10    "verbosity" : -1
11    }
12
13    lgtrain = lgb.Dataset(X_dev, label=y_dev)
14    lgval = lgb.Dataset(X_val, label=y_val)
15    evals_result = {}
16
17    gb = lgb.train(params, lgtrain, 1000, valid_sets=[lgval],
18                  early_stopping_rounds=100, verbose_eval=20,
19                  evals_result=evals_result)

```

```

1  Training until validation scores don't improve for 100 rounds.
2  [20]    valid_0's rmse: 0.217496
3  [40]    valid_0's rmse: 0.214649
4  [60]    valid_0's rmse: 0.213806
5  [80]    valid_0's rmse: 0.213349
6  [100]   valid_0's rmse: 0.213057
7  [120]   valid_0's rmse: 0.212834
8  [140]   valid_0's rmse: 0.212662
9  [160]   valid_0's rmse: 0.212531
10 [180]   valid_0's rmse: 0.212428
11 [200]   valid_0's rmse: 0.212326
12 [220]   valid_0's rmse: 0.212237
13 [240]   valid_0's rmse: 0.212158
14 [260]   valid_0's rmse: 0.212099
15 [280]   valid_0's rmse: 0.212044
16 [300]   valid_0's rmse: 0.211999
17 [320]   valid_0's rmse: 0.211952
18 [340]   valid_0's rmse: 0.21191
19 [360]   valid_0's rmse: 0.211871
20 [380]   valid_0's rmse: 0.211835
21 [400]   valid_0's rmse: 0.211789
22 [420]   valid_0's rmse: 0.211757
23 [440]   valid_0's rmse: 0.211718
24 [460]   valid_0's rmse: 0.211703
25 [480]   valid_0's rmse: 0.211689
26 [500]   valid_0's rmse: 0.211662
27 [520]   valid_0's rmse: 0.211644
28 [540]   valid_0's rmse: 0.211635
29 [560]   valid_0's rmse: 0.211613
30 [580]   valid_0's rmse: 0.211593
31 [600]   valid_0's rmse: 0.211576
32 [620]   valid_0's rmse: 0.211555
33 [640]   valid_0's rmse: 0.211521
34 [660]   valid_0's rmse: 0.211512

```

```
35 [680] valid_0's rmse: 0.211506
36 [700] valid_0's rmse: 0.211491
37 [720] valid_0's rmse: 0.211487
38 [740] valid_0's rmse: 0.211476
39 [760] valid_0's rmse: 0.21146
40 [780] valid_0's rmse: 0.211449
41 [800] valid_0's rmse: 0.211446
42 [820] valid_0's rmse: 0.211436
43 [840] valid_0's rmse: 0.211426
44 [860] valid_0's rmse: 0.211428
45 [880] valid_0's rmse: 0.211414
46 [900] valid_0's rmse: 0.211404
47 [920] valid_0's rmse: 0.211395
48 [940] valid_0's rmse: 0.211389
49 [960] valid_0's rmse: 0.211373
50 [980] valid_0's rmse: 0.211372
51 [1000] valid_0's rmse: 0.211368
52 Did not meet early stopping. Best iteration is:
53 [998] valid_0's rmse: 0.211368
```

```
1 pred = gb.predict(X_val, num_iteration=gb.best_iteration)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
4 print('RMSE without modification:', metrics.mean_squared_error(y_val, pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:', metrics.mean_squared_error(y_val, pred)**0.5)
8 preds['gb_sel30'] = pred
```

```
1 Less than zero: 5962
2 Over one: 49
3 RMSE without modification: 0.21136754065569408
4 RMSE fit-to-range: 0.21135017334460043
```

```
1 pred = gb.predict(test_sel30)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
```

```
1 Less than zero: 4939
2 Over one: 22
```

```

1 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
2
3 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
4 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
5
6 sub.to_csv('predictions/sub_gb_sel30.csv', index=False)

```

```

1 !kaggle competitions submit -c avito-demand-prediction -f predictions/sub_gb_sel30.csv
  -m "Message"

```

Private Score: 0.24551

Public Score: 0.24135

## 4.6. Optimal N Features

[Back to Outline](#)

```

1 train_scale = joblib.load('feature_dumps/train_scale.sav')
2 test_scale = joblib.load('feature_dumps/test_scale.sav')
3
4 def rmse(y, pred):
5     return metrics.mean_squared_error(y, pred)**0.5
6
7 score = metrics.make_scorer(rmse)
8
9 y = train.deal_probability
10 X_dev, X_val, y_dev, y_val = model_selection.train_test_split(train_scale, y)

```

```

1 n_features = np.arange(15, 41, 2).tolist()
2 data = train_scale
3 selector = feature_selection.SelectFromModel(
4     linear_model.Ridge(),
5     threshold=-np.inf)
6 model = linear_model.LinearRegression()

```

```

1 results = pd.DataFrame(columns=['Score'])

```

```

1 for n in n_features:
2     selector.max_features = n
3     selector.fit(data, y)
4     selection = data.iloc[:, selector.get_support()]
5     # Score of these features
6     cv = model_selection.cross_val_score(model, selection, y, cv=5, scoring=score)
7     print('n_features:', n)
8     print('CV Scores:', cv)
9     print('Mean CV Score:', np.mean(cv))
10    results.loc[n, 'Score'] = np.mean(cv)

```

```

1 n_features: 15

```

```

2 CV Scores: [0.2214024 0.22209101 0.22094013 0.22155788 0.22183177]
3 Mean CV Score: 0.22156463741803306
4 n_features: 17
5 CV Scores: [0.22100592 0.22173035 0.22052086 0.2211415 0.22143075]
6 Mean CV Score: 0.22116587668869325
7 n_features: 19
8 CV Scores: [0.22083071 0.22159906 0.2204022 0.22099091 0.22127407]
9 Mean CV Score: 0.22101939035838733
10 n_features: 21
11 CV Scores: [0.22041142 0.22113124 0.2200198 0.22056364 0.22084657]
12 Mean CV Score: 0.2205945303982421
13 n_features: 23
14 CV Scores: [0.21976962 0.22047235 0.21944153 0.21997114 0.22030295]
15 Mean CV Score: 0.21999151556366475
16 n_features: 25
17 CV Scores: [0.21946033 0.220143 0.21910807 0.21969886 0.22000164]
18 Mean CV Score: 0.2196823801454494
19 n_features: 27
20 CV Scores: [0.21935752 0.22001905 0.21901719 0.21959855 0.21987588]
21 Mean CV Score: 0.21957363990238027
22 n_features: 29
23 CV Scores: [0.219244 0.21990291 0.21890315 0.21946918 0.21975633]
24 Mean CV Score: 0.2194551130733152
25 n_features: 31
26 CV Scores: [0.21922796 0.21989394 0.21890329 0.21946472 0.21974498]
27 Mean CV Score: 0.2194469787959578
28 n_features: 33
29 CV Scores: [0.21909725 0.21974639 0.21874192 0.21930899 0.21957575]
30 Mean CV Score: 0.21929405938767052
31 n_features: 35
32 CV Scores: [0.21870135 0.21940735 0.2183705 0.21890897 0.21913743]
33 Mean CV Score: 0.2189051192649462
34 n_features: 37
35 CV Scores: [0.21866231 0.21936372 0.21832554 0.2188617 0.21909518]
36 Mean CV Score: 0.21886169117350002
37 n_features: 39
38 CV Scores: [0.21839662 0.21912002 0.21806015 0.21860995 0.21885036]
39 Mean CV Score: 0.21860741846370296

```

```

1 n_features = np.arange(42,81,5).tolist()
2 data = train_scale
3 selector = feature_selection.SelectFromModel(
4     linear_model.Ridge(),
5     threshold=-np.inf)
6 model = linear_model.LinearRegression()

```

```

1  for n in n_features:
2      selector.max_features = n
3      selector.fit(data,y)
4      selection = data.iloc[:,selector.get_support()]
5      # Score of these features
6      cv = model_selection.cross_val_score(model,selection,y,cv=5,scoring=score)
7      print('n_features:',n)
8      print('CV Scores:',cv)
9      print('Mean CV Score:',np.mean(cv))
10     results.loc[n,'Score'] = np.mean(cv)

```

```

1  n_features: 42
2  CV Scores: [0.21829902 0.21901327 0.21796659 0.21854189 0.21874101]
3  Mean CV Score: 0.21851235449796982
4  n_features: 47
5  CV Scores: [0.21815648 0.21885869 0.21780344 0.21836326 0.21858112]
6  Mean CV Score: 0.2183525972046391
7  n_features: 52
8  CV Scores: [0.21781718 0.21858665 0.21753391 0.2180219 0.21831257]
9  Mean CV Score: 0.21805444311090497
10 n_features: 57
11 CV Scores: [0.21767949 0.2184194 0.21738006 0.21788763 0.21818195]
12 Mean CV Score: 0.21790970823055159
13 n_features: 62
14 CV Scores: [0.21758837 0.21832427 0.21729401 0.21781125 0.21810756]
15 Mean CV Score: 0.21782509519951762
16 n_features: 67
17 CV Scores: [0.21755824 0.21828864 0.21725398 0.21777762 0.2180867 ]
18 Mean CV Score: 0.21779303464336142
19 n_features: 72
20 CV Scores: [0.21744609 0.21814406 0.21712654 0.21764411 0.21795778]
21 Mean CV Score: 0.21766371694471878
22 n_features: 77
23 CV Scores: [0.21740361 0.21811641 0.21709358 0.21761114 0.21792684]
24 Mean CV Score: 0.2176303186024358

```

## 4.7. PLSR by Original Feature

[Back to Outline](#)

```

1  train_plsrbf = pd.DataFrame(index=train.index)
2  test_plsrbf= pd.DataFrame(index=test.index)

```

### Title

```

1  f = 'title'
2  cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]

```



```

1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
  ['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
  for i in np.arange(n_comp)])

```

```

1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)

```

```

1 print(cv)
2 print(np.mean(cv))

```

```

1 [0.22706459 0.22759629 0.22707199 0.22745858 0.22764216]
2 0.22736672196238264

```

```

1 train_plsrbyf = train_plsrbyf.join(train_plsr)

```

```

1 test_plsrbyf = test_plsrbyf.join(test_plsr)

```

## Description

```

1 f = 'desc'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
  train_scale.columns.tolist()])]]

```

```

1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
  ['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
  for i in np.arange(n_comp)])

```

```

1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)

```

```

1 print(cv)
2 print(np.mean(cv))

```

```

1 [0.22920629 0.22972656 0.22875886 0.2294567 0.22963888]
2 0.2293574579086925

```

```

1 train_plsrbyf = train_plsrbyf.join(train_plsr)

```

```
1 test_plsrbyf = test_plsrbyf.join(test_plsr)
```

## Param\_1

```
1 f = 'p1'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]
```

```
1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
for i in np.arange(n_comp)])
```

```
1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)
```

```
1 print(cv)
2 print(np.mean(cv))
```

```
1 [0.23890726 0.23882191 0.23851118 0.23926176 0.23933843]
2 0.23896810819187678
```

```
1 train_plsrbyf = train_plsrbyf.join(train_plsr)
```

```
1 test_plsrbyf = test_plsrbyf.join(test_plsr)
```

## Param\_2

```
1 f = 'p2'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]
```

```
1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
for i in np.arange(n_comp)])
```

```
1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)
```

```
1 print(cv)
2 print(np.mean(cv))
```

```
1 [0.24378045 0.24348772 0.24340874 0.2440492 0.2441572 ]
2 0.24377666270113338
```

```
1 train_plsrbyf = train_plsrbyf.join(train_plsr)
```

```
1 test_plsrbyf = test_plsrbyf.join(test_plsr)
```

## Param\_3

```
1 f = 'p3'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]
```

```
1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
for i in np.arange(n_comp)])
```

```
1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)
```

```
1 print(cv)
2 print(np.mean(cv))
```

```
1 [0.25166324 0.25154532 0.25108586 0.25199132 0.25212277]
2 0.25168170236853965
```

```
1 train_plsrbyf = train_plsrbyf.join(train_plsr)
```

```
1 test_plsrbyf = test_plsrbyf.join(test_plsr)
```

## Region

```
1 f = 'reg'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]
```

```

1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
  ['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
  for i in np.arange(n_comp)])

```

```

1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)

```

```

1 print(cv)
2 print(np.mean(cv))

```

```

1 [0.2598769  0.25964108 0.25932434 0.26040616 0.2603358 ]
2 0.2599168552462979

```

```

1 train_plsrbyf = train_plsrbyf.join(train_plsr)

```

```

1 test_plsrbyf = test_plsrbyf.join(test_plsr)

```

## City

```

1 f = 'city'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
  train_scale.columns.tolist()])]]

```

```

1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
  ['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
  for i in np.arange(n_comp)])

```

```

1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)

```

```

1 print(cv)
2 print(np.mean(cv))

```

```

1 [0.25947971 0.25915577 0.25885413 0.25989806 0.25988651]
2 0.25945483436315697

```

```

1 train_plsrbyf = train_plsrbyf.join(train_plsr)

```

```
1 test_plsrbyf = test_plsrbyf.join(test_plsr)
```

## Category Name

```
1 f = 'catnam'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]]
```

```
1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
for i in np.arange(n_comp)])
```

```
1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)
```

```
1 print(cv)
2 print(np.mean(cv))
```

```
1 [0.2436163  0.24350094 0.24312993 0.24371419 0.24399855]
2 0.2435919807121306
```

```
1 train_plsrbyf = train_plsrbyf.join(train_plsr)
```

```
1 test_plsrbyf = test_plsrbyf.join(test_plsr)
```

## Parent

```
1 f = 'parent'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]]
```

```
1 n_comp = 5
2 plsr = cross_decomposition.PLSRegression(n_components=n_comp)
3 plsr.fit(train_scale[cols],y)
4 train_plsr = pd.DataFrame(plsr.transform(train_scale[cols]),columns=
['{}_{}'.format(f,i) for i in np.arange(n_comp)])
5 test_plsr = pd.DataFrame(plsr.transform(test_scale[cols]),columns=['{}_{}'.format(f,i)
for i in np.arange(n_comp)])
```

```
1 cv = model_selection.cross_val_score(model,train_plsr,y,cv=5,scoring=score)
```

```
1 print(cv)
2 print(np.mean(cv))
```

```
1 [0.24758108 0.24736126 0.24718824 0.24770739 0.2479346 ]
2 0.2475545146993571
```

```
1 train_plsrbyf = train_plsrbyf.join(train_plsr)
```

```
1 test_plsrbyf = test_plsrbyf.join(test_plsr)
```

## Encoders

```
1 f = 'code'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]]
```

```
1 cv = model_selection.cross_val_score(model, train_scale[cols], y, cv=5, scoring=score)
```

```
1 cv
```

```
1 array([0.23900483, 0.23894502, 0.23858882, 0.23932502, 0.23946458])
```

```
1 train_plsrbyf = train_plsrbyf.join(train_scale[cols])
```

```
1 test_plsrbyf = test_plsrbyf.join(test_scale[cols])
```

## Misses

```
1 f = 'miss'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]]
```

```
1 cv = model_selection.cross_val_score(model, train_scale[cols], y, cv=5, scoring=score)
```

```
1 cv
```

```
1 array([0.25628223, 0.25591782, 0.25567995, 0.25666355, 0.25671415])
```

```
1 train_plsrbyf = train_plsrbyf.join(train_scale[cols])
```

```
1 test_plsrbyf = test_plsrbyf.join(test_scale[cols])
```

## User Types

```
1 f = 'type'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]
```

```
1 cv = model_selection.cross_val_score(model, train_scale[cols], y, cv=5, scoring=score)
```

```
1 cv
```

```
1 array([0.259206 , 0.25890883, 0.25865925, 0.25973767, 0.25965589])
```

```
1 train_plsrbyf = train_plsrbyf.join(train_scale[cols])
```

```
1 test_plsrbyf = test_plsrbyf.join(test_scale[cols])
```

## Days

```
1 f = 'day'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]
```

```
1 cv = model_selection.cross_val_score(model, train_scale[cols], y, cv=5, scoring=score)
```

```
1 cv
```

```
1 array([0.26005697, 0.25978067, 0.25946028, 0.26059606, 0.26049315])
```

```
1 train_plsrbyf = train_plsrbyf.join(train_scale[cols])
```

```
1 test_plsrbyf = test_plsrbyf.join(test_scale[cols])
```

## Imputations

```
1 f = 'imp'
2 cols = train_scale.columns[[arg[0] for arg in np.argwhere([f in col for col in
train_scale.columns.tolist()])]]]
```

```
1 cv = model_selection.cross_val_score(model, train_scale[cols], y, cv=5, scoring=score)
```

```
1 cv
```

```
1 array([0.25578898, 0.25550334, 0.25520081, 0.25636177, 0.30268639])
```

```
1 train_plsrbf = train_plsrbf.join(train_scale[cols])
```

```
1 test_plsrbf = test_plsrbf.join(test_scale[cols])
```

## Seq\_n

```
1 train_plsrbf = train_plsrbf.join(train_scale.seq_n)
```

```
1 test_plsrbf = test_plsrbf.join(test_scale.seq_n)
```

```
1 cv = model_selection.cross_val_score(model, train_plsrbf, y, cv=5, scoring=score)
```

```
1 cv
```

```
1 array([0.21842824, 0.219155 , 0.21820018, 0.21869518, 0.21909233])
```

```
1 joblib.dump(train_plsrbf, 'feature_dumps/train_plsrbf.sav')
```

```
1 ['feature_dumps/train_plsrbf.sav']
```

```
1 joblib.dump(test_plsrbf, 'feature_dumps/test_plsrbf.sav')
```



```
1 ['feature_dumps/test_plsrbyf.sav']
```

## 4.8. LGB with PLSR by Original Feature

[Back to Outline](#)

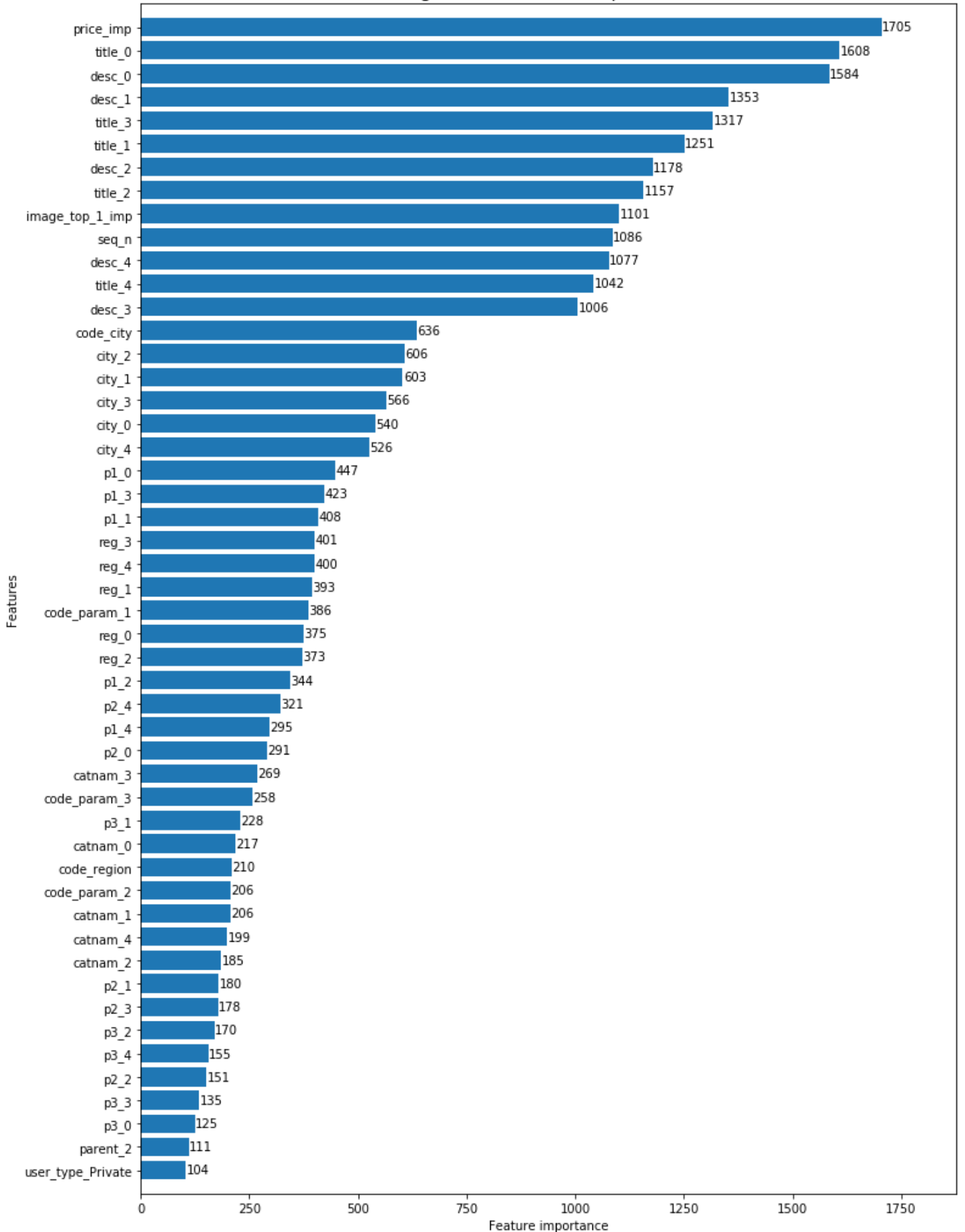
```
1 params = {
2     "objective" : "regression",
3     "metric" : "rmse",
4     "num_leaves" : 30,
5     "learning_rate" : 0.1,
6     "bagging_fraction" : 0.7,
7     "feature_fraction" : 0.7,
8     "bagging_frequency" : 5,
9     "bagging_seed" : 2342,
10    "verbosity" : -1
11 }
12
13 lgtrain = lgb.Dataset(X_dev, label=y_dev)
14 lgval = lgb.Dataset(X_val, label=y_val)
15 evals_result = {}
16
17 gb = lgb.train(params, lgtrain, 1000, valid_sets=[lgval],
18             early_stopping_rounds=100, verbose_eval=20,
19             evals_result=evals_result)
```

```
1 Training until validation scores don't improve for 100 rounds.
2 [20]    valid_0's rmse: 0.215517
3 [40]    valid_0's rmse: 0.21169
4 [60]    valid_0's rmse: 0.210382
5 [80]    valid_0's rmse: 0.209627
6 [100]   valid_0's rmse: 0.209132
7 [120]   valid_0's rmse: 0.20878
8 [140]   valid_0's rmse: 0.208521
9 [160]   valid_0's rmse: 0.208332
10 [180]   valid_0's rmse: 0.208156
11 [200]   valid_0's rmse: 0.208005
12 [220]   valid_0's rmse: 0.207872
13 [240]   valid_0's rmse: 0.207739
14 [260]   valid_0's rmse: 0.207636
15 [280]   valid_0's rmse: 0.207534
16 [300]   valid_0's rmse: 0.207445
17 [320]   valid_0's rmse: 0.207354
18 [340]   valid_0's rmse: 0.207287
19 [360]   valid_0's rmse: 0.207218
20 [380]   valid_0's rmse: 0.207147
21 [400]   valid_0's rmse: 0.207067
22 [420]   valid_0's rmse: 0.207019
23 [440]   valid_0's rmse: 0.20697
24 [460]   valid_0's rmse: 0.206927
25 [480]   valid_0's rmse: 0.206889
```

```
26 [500] valid_0's rmse: 0.206848
27 [520] valid_0's rmse: 0.20681
28 [540] valid_0's rmse: 0.206754
29 [560] valid_0's rmse: 0.206718
30 [580] valid_0's rmse: 0.206688
31 [600] valid_0's rmse: 0.206637
32 [620] valid_0's rmse: 0.206589
33 [640] valid_0's rmse: 0.206551
34 [660] valid_0's rmse: 0.206527
35 [680] valid_0's rmse: 0.206494
36 [700] valid_0's rmse: 0.20646
37 [720] valid_0's rmse: 0.206429
38 [740] valid_0's rmse: 0.206406
39 [760] valid_0's rmse: 0.206389
40 [780] valid_0's rmse: 0.206374
41 [800] valid_0's rmse: 0.20635
42 [820] valid_0's rmse: 0.206323
43 [840] valid_0's rmse: 0.206309
44 [860] valid_0's rmse: 0.206295
45 [880] valid_0's rmse: 0.206278
46 [900] valid_0's rmse: 0.206255
47 [920] valid_0's rmse: 0.206232
48 [940] valid_0's rmse: 0.206222
49 [960] valid_0's rmse: 0.206205
50 [980] valid_0's rmse: 0.206178
51 [1000] valid_0's rmse: 0.206158
52 Did not meet early stopping. Best iteration is:
53 [1000] valid_0's rmse: 0.206158
```

```
1 fig, ax = plt.subplots(figsize=(12,18))
2 lgb.plot_importance(gb, max_num_features=50, height=0.8, ax=ax)
3 ax.grid(False)
4 plt.title("LightGBM - Feature Importance", fontsize=15)
5 plt.show()
```

# LightGBM - Feature Importance



```

1 pred = gb.predict(X_val, num_iteration=gb.best_iteration)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
4 print('RMSE without modification:', metrics.mean_squared_error(y_val, pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:', metrics.mean_squared_error(y_val, pred)**0.5)
8 preds['gb_plsrbyf'] = pred

```

```

1 Less than zero: 8101
2 Over one: 53
3 RMSE without modification: 0.20615780720704704
4 RMSE fit-to-range: 0.20612838300391242

```

```

1 pred = gb.predict(test_plsrbyf)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))

```

```

1 Less than zero: 13258
2 Over one: 21

```

```

1 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
2
3 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
4 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
5
6 sub.to_csv('predictions/sub_gb_plsrbyf.csv', index=False)

```

```

1 !kaggle competitions submit -c avito-demand-prediction -f
   predictions/sub_gb_plsrbyf.csv -m "Message"

```

```

1 Warning: Your Kaggle API key is readable by other users on this system! To fix this,
   you can run 'chmod 600 /home/user/.kaggle/kaggle.json'
2 100%|████████████████████████████████████████| 15.8M/15.8M [00:16<00:00, 1.02MB/s]
3 Successfully submitted to Avito Demand Prediction Challenge

```

Private Score: 0.24004

Public Score: 0.23596

## 4.9. Multi Layer Perceptron

[Back to Outline](#)

```
1 mlp = neural_network.MLPRegressor(verbose=True,  
2                                     hidden_layer_sizes=(64,1),  
3                                     n_iter_no_change=10,  
4                                     tol=0.0001,  
5                                     learning_rate_init=0.1,  
6                                     )
```

```
1 mlp.fit(X_dev,y_dev)
```

```
1 Iteration 1, loss = 0.03414127  
2 Iteration 2, loss = 0.03412936  
3 Iteration 3, loss = 0.03412152  
4 Iteration 4, loss = 0.03412905  
5 Iteration 5, loss = 0.03412075  
6 Iteration 6, loss = 0.03409799  
7 Iteration 7, loss = 0.03412794  
8 Iteration 8, loss = 0.03412402  
9 Iteration 9, loss = 0.03412227  
10 Iteration 10, loss = 0.03412462  
11 Iteration 11, loss = 0.03411990  
12 Iteration 12, loss = 0.03411259  
13 Training loss did not improve more than tol=0.000100 for 10 consecutive epochs.  
    Stopping.
```

```
1 MLPRegressor(activation='relu', alpha=0.0001, batch_size='auto', beta_1=0.9,  
2               beta_2=0.999, early_stopping=False, epsilon=1e-08,  
3               hidden_layer_sizes=(64, 1), learning_rate='constant',  
4               learning_rate_init=0.1, max_iter=200, momentum=0.9,  
5               n_iter_no_change=10, nesterovs_momentum=True, power_t=0.5,  
6               random_state=None, shuffle=True, solver='adam', tol=0.0001,  
7               validation_fraction=0.1, verbose=True, warm_start=False)
```

```
1 pred = mlp.predict(X_val)  
2 print('Less than zero:',sum(pred<0))  
3 print('Over one:',sum(pred>1))  
4 print('RMSE without modification:',metrics.mean_squared_error(y_val,pred)**0.5)  
5 pred[pred>1] = 1  
6 pred[pred<0] = 0  
7 print('RMSE fit-to-range:',metrics.mean_squared_error(y_val,pred)**0.5)  
8 preds['mlp_plsrbf'] = pred
```

```
1 Less than zero: 0  
2 Over one: 0  
3 RMSE without modification: 0.2603847814321052  
4 RMSE fit-to-range: 0.2603847814321052
```

```
1 pred = mlp.predict(test_plsrbyf)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
```

```
1 Less than zero: 0
2 Over one: 0
```

```
1 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
2
3 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
4 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
5
6 sub.to_csv('predictions/sub_mlp_plsrbyf.csv', index=False)
```

```
1 !kaggle competitions submit -c avito-demand-prediction -f
   predictions/sub_mlp_plsrbyf.csv -m "Message"
```

```
1 Warning: Your Kaggle API key is readable by other users on this system! To fix this,
   you can run 'chmod 600 /home/user/.kaggle/kaggle.json'
2 100%|████████████████████████████████████████| 16.0M/16.0M [00:16<00:00, 1.01MB/s]
3 Successfully submitted to Avito Demand Prediction Challenge
```

Private Score: 0.27107

Public Score: 0.26633

## 4.10. Keras

[Back to Outline](#)

```
1 ker = Sequential()
2 ker.add(Dense(64, activation='relu', input_dim=72))
3 ker.add(Dense(1, activation='relu'))
4
5 ker.compile(loss='mean_squared_error', optimizer='adam', metrics=['accuracy'])
```

```
1 ker.fit(X_dev, y_dev, epochs=30, verbose=1)
```

```
1 WARNING:tensorflow:From /home/user/anaconda3/lib/python3.7/site-
   packages/tensorflow/python/ops/math_ops.py:3066: to_int32 (from
   tensorflow.python.ops.math_ops) is deprecated and will be removed in a future version.
2 Instructions for updating:
3 Use tf.cast instead.
4 Epoch 1/30
```

```
5 1127568/1127568 [=====] - 34s 31us/step - loss: 0.0730 - acc:
0.6466
6 Epoch 2/30
7 1127568/1127568 [=====] - 34s 30us/step - loss: 0.0455 - acc:
0.6480
8 Epoch 3/30
9 1127568/1127568 [=====] - 35s 31us/step - loss: 0.0448 - acc:
0.6480
10 Epoch 4/30
11 1127568/1127568 [=====] - 33s 29us/step - loss: 0.0447 - acc:
0.6479
12 Epoch 5/30
13 1127568/1127568 [=====] - 37s 33us/step - loss: 0.0447 - acc:
0.6479
14 Epoch 6/30
15 1127568/1127568 [=====] - 36s 32us/step - loss: 0.0446 - acc:
0.6479
16 Epoch 7/30
17 1127568/1127568 [=====] - 41s 37us/step - loss: 0.0446 - acc:
0.6479
18 Epoch 8/30
19 1127568/1127568 [=====] - 42s 37us/step - loss: 0.0445 - acc:
0.6479
20 Epoch 9/30
21 1127568/1127568 [=====] - 45s 40us/step - loss: 0.1481 - acc:
0.6480
22 Epoch 10/30
23 1127568/1127568 [=====] - 41s 36us/step - loss: 0.0445 - acc:
0.6480
24 Epoch 11/30
25 1127568/1127568 [=====] - 32s 28us/step - loss: 0.0444 - acc:
0.6480
26 Epoch 12/30
27 1127568/1127568 [=====] - 31s 28us/step - loss: 0.0445 - acc:
0.6480
28 Epoch 13/30
29 1127568/1127568 [=====] - 42s 38us/step - loss: 0.0444 - acc:
0.6480
30 Epoch 14/30
31 1127568/1127568 [=====] - 48s 43us/step - loss: 0.0444 - acc:
0.6480
32 Epoch 15/30
33 1127568/1127568 [=====] - 30s 26us/step - loss: 0.0444 - acc:
0.6480
34 Epoch 16/30
35 1127568/1127568 [=====] - 35s 31us/step - loss: 0.0444 - acc:
0.6480
36 Epoch 17/30
37 1127568/1127568 [=====] - 27s 24us/step - loss: 0.0444 - acc:
0.6480
38 Epoch 18/30
39 1127568/1127568 [=====] - 35s 31us/step - loss: 0.0443 - acc:
0.6481
40 Epoch 19/30
```

```

41 1127568/1127568 [=====] - 42s 37us/step - loss: 0.0445 - acc:
    0.6480
42 Epoch 20/30
43 1127568/1127568 [=====] - 43s 38us/step - loss: 0.0443 - acc:
    0.6480
44 Epoch 21/30
45 1127568/1127568 [=====] - 37s 33us/step - loss: 0.0444 - acc:
    0.6480
46 Epoch 22/30
47 1127568/1127568 [=====] - 36s 32us/step - loss: 0.0445 - acc:
    0.6480
48 Epoch 23/30
49 1127568/1127568 [=====] - 38s 33us/step - loss: 0.0443 - acc:
    0.6480
50 Epoch 24/30
51 1127568/1127568 [=====] - 38s 34us/step - loss: 0.0443 - acc:
    0.6480
52 Epoch 25/30
53 1127568/1127568 [=====] - 39s 34us/step - loss: 0.0443 - acc:
    0.6480
54 Epoch 26/30
55 1127568/1127568 [=====] - 39s 35us/step - loss: 0.0444 - acc:
    0.6480
56 Epoch 27/30
57 1127568/1127568 [=====] - 38s 34us/step - loss: 0.0446 - acc:
    0.6481
58 Epoch 28/30
59 1127568/1127568 [=====] - 36s 32us/step - loss: 0.0446 - acc:
    0.6480
60 Epoch 29/30
61 1127568/1127568 [=====] - 37s 33us/step - loss: 0.0443 - acc:
    0.6479
62 Epoch 30/30
63 1127568/1127568 [=====] - 38s 34us/step - loss: 0.0445 - acc:
    0.6480

```

```

1 pred = ker.predict(X_val)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
4 print('RMSE without modification:', metrics.mean_squared_error(y_val, pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:', metrics.mean_squared_error(y_val, pred)**0.5)
8 preds['ker_plsrbyf'] = pred

```

```

1 Less than zero: [0]
2 Over one: [308]
3 RMSE without modification: 0.21108059545026234
4 RMSE fit-to-range: 0.21100376236913634

```



```
1 pred = ker.predict(test_plsrbyf)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
```

```
1 Less than zero: [0]
2 Over one: [227]
```

```
1 pred = [p[0] for p in pred]
2 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
3
4 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
5 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
6
7 sub.to_csv('predictions/sub_ker_plsrbyf.csv', index=False)
```

```
1 !kaggle competitions submit -c avito-demand-prediction -f
   predictions/sub_ker_plsrbyf.csv -m "Message"
```

```
1 Warning: Your Kaggle API key is readable by other users on this system! To fix this,
   you can run 'chmod 600 /home/user/.kaggle/kaggle.json'
2 100%|██████████████████████████████████████████████████████████████████████████| 14.9M/14.9M [00:14<00:00, 1.07MB/s]
3 Successfully submitted to Avito Demand Prediction Challenge
```

Private Score: 0.24423

Public Score: 0.23975

## 4.11. RandomForestRegressor

[Back to Outline](#)

```
1 rf = ensemble.RandomForestRegressor(
2     n_estimators=100,
3     verbose=10,
4     random_state=33
5 )
```

```
1 rf.fit(X_dev, y_dev)
```

```
1 [Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
```

```
1 building tree 1 of 100
```

```
1 [Parallel(n_jobs=1)]: Done    1 out of    1 | elapsed:   58.9s remaining:    0.0s
```

```
1 building tree 2 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 2.0min remaining: 0.0s
```

```
1 building tree 3 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 3.0min remaining: 0.0s
```

```
1 building tree 4 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 4.1min remaining: 0.0s
```

```
1 building tree 5 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 5.2min remaining: 0.0s
```

```
1 building tree 6 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 6.1min remaining: 0.0s
```

```
1 building tree 7 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 7.0min remaining: 0.0s
```

```
1 building tree 8 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 8.0min remaining: 0.0s
```

```
1 building tree 9 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 9.1min remaining: 0.0s
```

```
1 building tree 10 of 100
2 building tree 11 of 100
3 building tree 12 of 100
4 building tree 13 of 100
5 building tree 14 of 100
6 building tree 15 of 100
7 building tree 16 of 100
8 building tree 17 of 100
9 building tree 18 of 100
10 building tree 19 of 100
11 building tree 20 of 100
12 building tree 21 of 100
```

13 building tree 22 of 100  
14 building tree 23 of 100  
15 building tree 24 of 100  
16 building tree 25 of 100  
17 building tree 26 of 100  
18 building tree 27 of 100  
19 building tree 28 of 100  
20 building tree 29 of 100  
21 building tree 30 of 100  
22 building tree 31 of 100  
23 building tree 32 of 100  
24 building tree 33 of 100  
25 building tree 34 of 100  
26 building tree 35 of 100  
27 building tree 36 of 100  
28 building tree 37 of 100  
29 building tree 38 of 100  
30 building tree 39 of 100  
31 building tree 40 of 100  
32 building tree 41 of 100  
33 building tree 42 of 100  
34 building tree 43 of 100  
35 building tree 44 of 100  
36 building tree 45 of 100  
37 building tree 46 of 100  
38 building tree 47 of 100  
39 building tree 48 of 100  
40 building tree 49 of 100  
41 building tree 50 of 100  
42 building tree 51 of 100  
43 building tree 52 of 100  
44 building tree 53 of 100  
45 building tree 54 of 100  
46 building tree 55 of 100  
47 building tree 56 of 100  
48 building tree 57 of 100  
49 building tree 58 of 100  
50 building tree 59 of 100  
51 building tree 60 of 100  
52 building tree 61 of 100  
53 building tree 62 of 100  
54 building tree 63 of 100  
55 building tree 64 of 100  
56 building tree 65 of 100  
57 building tree 66 of 100  
58 building tree 67 of 100  
59 building tree 68 of 100  
60 building tree 69 of 100  
61 building tree 70 of 100  
62 building tree 71 of 100  
63 building tree 72 of 100  
64 building tree 73 of 100  
65 building tree 74 of 100  
66 building tree 75 of 100

```
67 building tree 76 of 100
68 building tree 77 of 100
69 building tree 78 of 100
70 building tree 79 of 100
71 building tree 80 of 100
72 building tree 81 of 100
73 building tree 82 of 100
74 building tree 83 of 100
75 building tree 84 of 100
76 building tree 85 of 100
77 building tree 86 of 100
78 building tree 87 of 100
79 building tree 88 of 100
80 building tree 89 of 100
81 building tree 90 of 100
82 building tree 91 of 100
83 building tree 92 of 100
84 building tree 93 of 100
85 building tree 94 of 100
86 building tree 95 of 100
87 building tree 96 of 100
88 building tree 97 of 100
89 building tree 98 of 100
90 building tree 99 of 100
91 building tree 100 of 100
```

```
1 [Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 97.1min finished
```

```
1 RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
2                       max_features='auto', max_leaf_nodes=None,
3                       min_impurity_decrease=0.0, min_impurity_split=None,
4                       min_samples_leaf=1, min_samples_split=2,
5                       min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=None,
6                       oob_score=False, random_state=33, verbose=10, warm_start=False)
```

```
1 pred = rf.predict(X_val)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))
4 print('RMSE without modification:', metrics.mean_squared_error(y_val, pred)**0.5)
5 pred[pred>1] = 1
6 pred[pred<0] = 0
7 print('RMSE fit-to-range:', metrics.mean_squared_error(y_val, pred)**0.5)
8 preds['rf_plsrbyf'] = pred
```

```

1 [Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
2 [Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.3s remaining: 0.0s
3 [Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.5s remaining: 0.0s
4 [Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 0.8s remaining: 0.0s
5 [Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 1.1s remaining: 0.0s
6 [Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 1.3s remaining: 0.0s
7 [Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 1.6s remaining: 0.0s
8 [Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 1.9s remaining: 0.0s
9 [Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 2.1s remaining: 0.0s
10 [Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 2.4s remaining: 0.0s
11 [Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 24.9s finished

```

```

1 Less than zero: 0
2 Over one: 0
3 RMSE without modification: 0.21064608123386366
4 RMSE fit-to-range: 0.21064608123386366

```

```

1 pred = rf.predict(test_plsrbyf)
2 print('Less than zero:', sum(pred<0))
3 print('Over one:', sum(pred>1))

```

```

1 [Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
2 [Parallel(n_jobs=1)]: Done 1 out of 1 | elapsed: 0.3s remaining: 0.0s
3 [Parallel(n_jobs=1)]: Done 2 out of 2 | elapsed: 0.7s remaining: 0.0s
4 [Parallel(n_jobs=1)]: Done 3 out of 3 | elapsed: 1.0s remaining: 0.0s
5 [Parallel(n_jobs=1)]: Done 4 out of 4 | elapsed: 1.3s remaining: 0.0s
6 [Parallel(n_jobs=1)]: Done 5 out of 5 | elapsed: 1.6s remaining: 0.0s
7 [Parallel(n_jobs=1)]: Done 6 out of 6 | elapsed: 1.9s remaining: 0.0s
8 [Parallel(n_jobs=1)]: Done 7 out of 7 | elapsed: 2.2s remaining: 0.0s
9 [Parallel(n_jobs=1)]: Done 8 out of 8 | elapsed: 2.6s remaining: 0.0s
10 [Parallel(n_jobs=1)]: Done 9 out of 9 | elapsed: 2.9s remaining: 0.0s
11 [Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 30.2s finished

```

```

1 Less than zero: 0
2 Over one: 0

```

```

1 #pred = [p[0] for p in pred]
2 sub = pd.DataFrame({'item_id':test.item_id, 'deal_probability':pred})
3
4 sub.loc[sub.deal_probability>1, 'deal_probability'] = 1
5 sub.loc[sub.deal_probability<0, 'deal_probability'] = 0
6
7 sub.to_csv('predictions/sub_rf_plsrbyf.csv', index=False)

```

```

1 !kaggle competitions submit -c avito-demand-prediction -f
predictions/sub_rf_plsrbyf.csv -m "Message"

```

```
1 Warning: Your Kaggle API key is readable by other users on this system! To fix this,
  you can run 'chmod 600 /home/user/.kaggle/kaggle.json'
2 100%|████████████████████████████████████████| 14.7M/14.7M [00:15<00:00, 994kB/s]
3 Successfully submitted to Avito Demand Prediction Challenge
```

Private Score: 0.24491

Public Score: 0.24088

## 4.12. Error Analysis

```
1 train_scale = joblib.load('feature_dumps/train_scale.sav')
2 test_scale = joblib.load('feature_dumps/test_scale.sav')
```

```
1 y = train.deal_probability
2 X_dev, X_val, y_dev, y_val = model_selection.train_test_split(train_scale, y)
```

```
1 params = {
2     "objective" : "regression",
3     "metric" : "rmse",
4     "num_leaves" : 30,
5     "learning_rate" : 0.1,
6     "bagging_fraction" : 0.7,
7     "feature_fraction" : 0.7,
8     "bagging_frequency" : 5,
9     "bagging_seed" : 2342,
10    "verbosity" : -1
11 }
12
13 lgtrain = lgb.Dataset(X_dev, label=y_dev)
14 lgval = lgb.Dataset(X_val, label=y_val)
15 evals_result = {}
16
17 gb = lgb.train(params, lgtrain, 1000, valid_sets=[lgval],
18             early_stopping_rounds=100, verbose_eval=20,
19             evals_result=evals_result)
```

```
1 Training until validation scores don't improve for 100 rounds.
2 [20]  valid_0's rmse: 0.214165
3 [40]  valid_0's rmse: 0.210232
4 [60]  valid_0's rmse: 0.208477
5 [80]  valid_0's rmse: 0.207474
6 [100] valid_0's rmse: 0.206755
7 [120] valid_0's rmse: 0.206255
8 [140] valid_0's rmse: 0.205848
9 [160] valid_0's rmse: 0.205525
10 [180] valid_0's rmse: 0.20528
11 [200] valid_0's rmse: 0.205048
12 [220] valid_0's rmse: 0.204862
```

```
13 [240] valid_0's rmse: 0.20471
14 [260] valid_0's rmse: 0.204591
15 [280] valid_0's rmse: 0.204482
16 [300] valid_0's rmse: 0.20438
17 [320] valid_0's rmse: 0.2043
18 [340] valid_0's rmse: 0.204212
19 [360] valid_0's rmse: 0.204124
20 [380] valid_0's rmse: 0.204032
21 [400] valid_0's rmse: 0.203959
22 [420] valid_0's rmse: 0.203919
23 [440] valid_0's rmse: 0.203838
24 [460] valid_0's rmse: 0.203796
25 [480] valid_0's rmse: 0.20375
26 [500] valid_0's rmse: 0.203702
27 [520] valid_0's rmse: 0.203649
28 [540] valid_0's rmse: 0.20362
29 [560] valid_0's rmse: 0.203575
30 [580] valid_0's rmse: 0.203512
31 [600] valid_0's rmse: 0.203477
32 [620] valid_0's rmse: 0.203442
33 [640] valid_0's rmse: 0.203418
34 [660] valid_0's rmse: 0.203381
35 [680] valid_0's rmse: 0.203348
36 [700] valid_0's rmse: 0.203299
37 [720] valid_0's rmse: 0.203272
38 [740] valid_0's rmse: 0.203259
39 [760] valid_0's rmse: 0.203234
40 [780] valid_0's rmse: 0.203218
41 [800] valid_0's rmse: 0.203195
42 [820] valid_0's rmse: 0.203176
43 [840] valid_0's rmse: 0.203156
44 [860] valid_0's rmse: 0.203136
45 [880] valid_0's rmse: 0.203109
46 [900] valid_0's rmse: 0.203081
47 [920] valid_0's rmse: 0.203049
48 [940] valid_0's rmse: 0.20304
49 [960] valid_0's rmse: 0.203008
50 [980] valid_0's rmse: 0.20299
51 [1000] valid_0's rmse: 0.202964
52 Did not meet early stopping. Best iteration is:
53 [1000] valid_0's rmse: 0.202964
```

```
1 pred = gb.predict(X_val, num_iteration=gb.best_iteration)
```

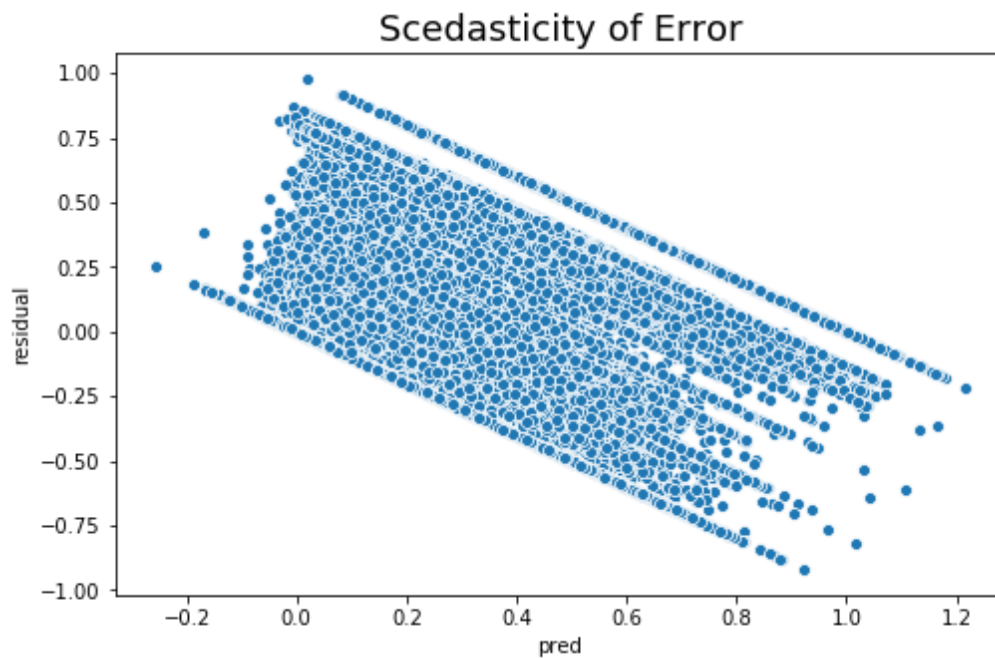
```
1 residual = y_val - pred
```

```
1 resdf = pd.DataFrame({'y':y_val, 'pred':pred, 'residual':residual})
```

```

1 plt.figure(figsize=(8,5))
2 sns.scatterplot(x='pred',y='residual',data=resdf)
3 plt.title('Scedasticity of Error',fontsize=18)
4 plt.show()

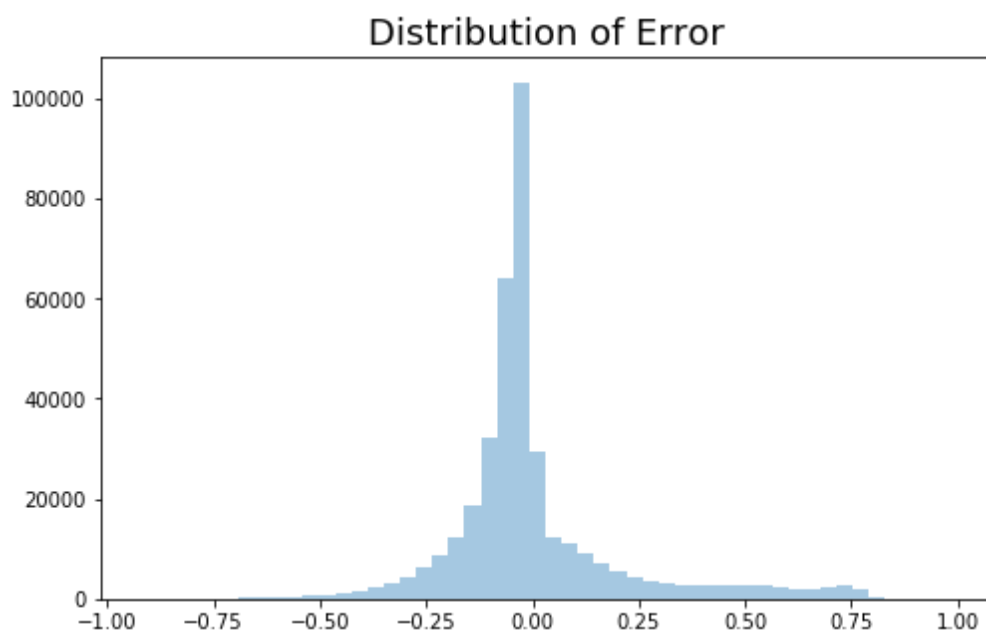
```



```

1 plt.figure(figsize=(8,5))
2 sns.distplot(residual.values,kde=False)
3 plt.title('Distribution of Error',fontsize=18)
4 plt.show()

```



## 5. Product



[Back to Outline](#)

Based on several comparisons, the above-described feature-engineering pipeline followed by LightGBM is the best approach at predicting the deal probability of an online ad.

### **Why it works?**

It works due to the robustness and variety of feature-extraction and decomposition techniques. While feature extraction can generate a lot of data, this is only useful in a reduced dimensional space. Decomposing large CSR matrices produces predictively powerful components.

### **What problem it solves?**

Online platforms for selling used goods rely on regular people selling their belongings online to achieve high-traffic. These sellers blindly sell their things with erroneous expectations and bad listing practices, thus becoming frustrated with online sales. Helping sellers understand the demand of their listings contributes in several ways: Informed sellers can optimize their listings for maximum deal probability and also optimize their choice of goods to sell, based on the deal probability of particular categories.

### **How will it work in production?**

In a production environment this model would learn from the sales information of a historic time window in order to predict the deal probability of new ads. Necessary maintenance might involve adjusting some of the feature-engineering procedures to ensure they are capturing the most valuable information.