Predicting Ads Deal Probability

Danijel Kopčinović, IT Market

The Problem

When selling used goods online, a combination of tiny, nuanced details in a product description (e.g. well taken photos, nicely formatted copy etc.) can make a big difference in drumming up interest.

Avito, Russia's largest classified advertisements website, is deeply familiar with this problem. Sellers on their platform sometimes feel frustrated with both too little demand (indicating something is wrong with the product or the product listing) or too much demand (indicating a hot item with a good description was underpriced).

Our task is to predict demand for an online advertisement based on its full description (title, description, images, etc.), its context (geographically where it was posted, similar ads already posted) and historical demand for similar ads in similar contexts. With this information, Avito can inform sellers on how to best optimize their listing and provide some indication of how much interest they should realistically expect to receive.

The task was originally published on Kaggle: https://www.kaggle.com/c/avito-demand-prediction

The Data

The data we work with consists of:

- item id Ad id.
- user_id User id.
- region Ad region.
- city Ad city.
- parent category name Top level ad category as classified by Avito's ad model.
- category_name Fine grain ad category as classified by Avito's ad model.
- param_1 Optional parameter from Avito's ad model.
- param_2 Optional parameter from Avito's ad model.
- param_3 Optional parameter from Avito's ad model.
- title Ad title.
- description Ad description.
- price Ad price.
- item seq number Ad sequential number for user.
- activation_date- Date ad was placed.
- user_type User type.
- image Id code of image. Ties to a jpg file in train_jpg. Not every ad has an image.
- image top 1 Avito's classification code for the image.
- deal_probability The target variable. This is the likelihood that an ad actually sold something.
 It's not possible to verify every transaction with certainty, so this column's value can be any float from zero to one.

Finding Predictors for the Target Variable

After initial data processing, our main task is to find good (highly correlated) predictors for the deal_probability variable. The initial set of predictors are the initial data variables but we also include their "interactions"/combinations to cover as much variability in the target variable as possible.

Correlating all these predictors with deal_probability gives us the following table:

```
Predictor Correlation
                                     parent_category_name_Лсчные_вещс
                                                                           0.3249185
1 2 3 4
                             category_name_одесда__обувь__аксессуаdjы
                                                                           0.2147931
            parent_category_name_Услугс__X__description_length
category_name_пdједлосенсе_услуг__X__description_length
                                                                           0.2108996
                                                                           0.2108996
5
                                           parent_category_name_Услугс
                                                                           0.2105639
6
                                     category_name_Пdjедлосенсе_услуг
                                                                           0.2105639
                                 _X__category_name_Пdjедлосенсе_услуг
   parent_category_name_Услугс_
                                                                           0.2105639
8
                                 category_name_Детская_одесда_с_обувь
                                                                           0.1978882
9
                                                param_1_Сенская_одесда
                                                                           0.1965910
              parent_category_name_Услугс_X__user_type_Private
category_name_Пdjедлосенсе_услуг_X__user_type_Private
10
                                                                           0.1920129
11
                                                                           0.1920129
                     item_seq_number__x__parent_category_name_Услугс
12
                                                                           0.1895534
13
                item_seq_number__X__category_name_Пdjедлосенсе_услуг
                                                                           0.1895534
                          image_top_1__X__parent_category_name_Услугс
14
                                                                           0.1870032
15
                    image_top_1__X__category_name_Пdjедлосенсе_услуг
                                                                           0.1870032
16
                                price_x_parent_category_name_Услугс
                                                                           0.1700630
                                                                           0.1700630
17
                           price__X__category_name_Пdjедлосенсе_услуг
18
                                                        item_seq_number
                                                                           0.1690906
                                    parent_category_name_Heдвcccмость
19
                                                                           0.1646846
20
           parent_category_name_Heдвcccмocть__X__description_length
                                                                           0.1621378
21
                                                            image_top_1
                                                                           0.1606214
22
                                                                           0.1479474
                                                   param_1_Для_девочек
23
                          price_X_parent_category_name_недвсссмость
                                                                           0.1470145
24
                                     parent_category_name_TdjaHcnodjT
                                                                           0.1467983
25
             parent_category_name_TdjaHcnodjt__X__description_length
                                                                           0.1462346
                    image_top_1__X__parent_category_name_TdjaHcnodjT
26
                                                                           0.1441984
27
               item_seq_number__X__parent_category_name_Недвсссмость
                                                                           0.1435174
28
                           price_x_parent_category_name_TdjaHcnodjT
                                                                           0.1424918
29
                                   image_top_1__X__description_length
                                                                           0.1424293
30
                                                   param_1_C_ndjo6erom
                                                                           0.1418551
31
           parent_category_name_TdjaHcnodjt__X__param_1_C_ndjo6erom
                                                                           0.1418551
32
                    category_name_Автомобслс__X__param_1_C_ndjoбегом
                                                                           0.1418551
                          param_1_C_ndjo6erom__X__description_length
33
                                                                           0.1416224
              parent_category_name_TdjaHcnodjt__X_user_type_Private
34
                                                                           0.1404075
35
                                  image_top_1__X__param_1_C_ndjo6erom
                                                                           0.1397038
36
                                              category_name_Автомобслс
                                                                           0.1395902
37
      parent_category_name_TdjaHcnodjt__X__category_name_ABTOMO6CAC
                                                                           0.1395902
38
                                          price_x_description_length
                                                                           0.1392518
39
                     category_name_Aвтомобслс__X__description_length
                                                                           0.1392317
40
                                        price_x_param_1_c_ndjo6erom
                                                                           0.1387086
41
                             image_top_1_
                                          _X__category_name_Автомобслс
                                                                           0.1375089
42
                                                                  price
                                                                           0.1368594
                                   price__X__category_name_Автомобслс
43
                                                                           0.1363071
44
                                              category_name_кваdjтcdjы
                                                                           0.1342821
45
     parent_category_name_недвсссмость__X__category_name_кваdjтcdjы
                                                                           0.1342821
46
                     category_name_кваdjтcdjы__X__description_length
                                                                           0.1331503
47
                      category_name_Aвтомобслс__X_user_type_Private
                                                                           0.1323884
                            param_1_C_ndjo6erom__X_user_type_Private
48
                                                                           0.1323248
49
             parent_category_name_Heдвcccмость__X_user_type_Private
                                                                           0.1320766
              param_1_Djeмонт__ctdjocteльство__X__description_length
50
                                                                           0.1295317
```

To better understand the meaning of the predictors, here is a brief translation of some words from Russian to English:

```
* ЛСЧНЫЕ ВЕЩС : Laser products
 одесда обувь аксессуаdjы : clothes accessories shoes
 Услуг : Services
* Пиједлосенсе услуг : Paid services
 детская одесда с обувь : Children's clothes with shoes
 Сенская одесда : Women's clothes
* Недвсссмость : Undecided
 Для девочек : For girls
 Tdjaнcnodjт : Transport
 С пdjoбегом : With mileage
 Автомобслс : Automotive
 Djeмонтстdjостельство : Demonstration
* сдам : Rent
* кваdjтcdjы : Quads
* Пdjодам : Sell
* Tdjaнcnodjт пеdjевозкс : Transit train
```

The "__X__" sequence of characters means the interaction/combination of two variables.

We can see that the top (mostly correlated) predictor for the deal_probability is the parent category of "Laser products". This means that if the parent category is "Laser products", then the deal probability increases. Coefficient 0.3249185 is the Spearman coefficient. Its magnitude indicates solid positive correlation between the two variables. Similarly, we can see that:

- the deal_probability increases for the category "clothes accessories shoes"
- the deal_probability increases with the ad text description length for the parent category "Services"
- the deal_probability increases for "Paid services" within the "Services" category etc.

Let's have a look at some predictors negatively correlated with the deal_probability:

```
Predictor Correlation
4198 parent_category_name_Лсчные_вещс__X__user_type_Private -0.2575794
4199 parent_category_name_Лсчные_вещс__X__description_length -0.2625350
4200 image_top_1__X_parent_category_name_Лсчные_вещс -0.2800965
4201 price__X_parent_category_name_Лсчные_вещс -0.2906887
4202 item_seq_number__X_parent_category_name_Лсчные_вещс -0.3131590
```

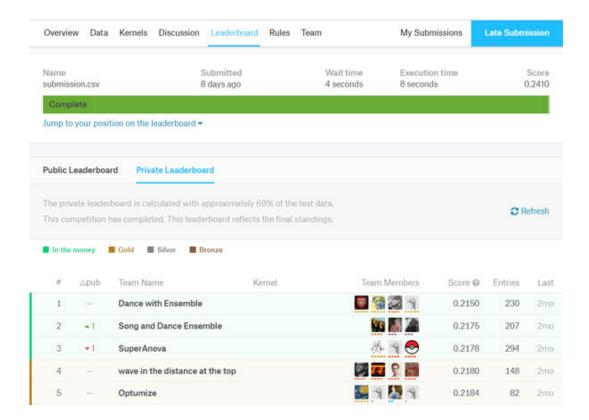
We can see that the deal_probability decreases in the category "Laser products" for "Private" users, for greater "description length" (this is a bit surprising, but it's what the data suggests), for greater price and for greater item_seq_number. Greater item_seq_number indicates greater number of ads by the same seller, so the conclusion would be that if a seller within "Laser products" publishes many ads, he/she actually decreases the probability of the deal/sale. The image_top_1 is an internal Avito image categorization so we cannot describe it verbally, but is also decreases the deal_probability for the "Laser products".

Sometimes it's not easy to interpret the interactions and combinations of predictors, but if they do provide an important impact on the target variable (indicated by the correlation coefficient) then they should be included in the model to increase its quality.

Regression Models

Once we've set up the most important predictors, we build the regression models for the deal_probability. We use 3 different methods: linear regression, xgboost and neural network. All the methods gave us roughly 20% R-squared value meaning that they managed to explain 20% of the deal_probability variations. This is not too good and it indicates that there are other variables influencing the deal_probability variable. Some of those would probably be the quality of the ad images that we didn't take into account.

All in all, we predicted the target variable for the test set in the challenge and scored a solid result of 0.2410 RMSE (root mean squared error). The top results were around 0.21 RMSE.



Conclusion

Data science provides many different tools and approaches to analyze data and make educated business decisions.

Are you analyzing your data?;)

For further information please contact danijel.kopcinovic@itmarket.hr.