

# Machine Learning Engineer Nanodegree

## Capstone Project

### Avito Demand Prediction Challenge

Sachin Lamba

Mar, 2019

#### I. Definition

##### Project Overview

Today, any website we visit always tried to recommend something depending on our taste (or we can say history). Machine Learning is a field to check from data what history can say about our future tastes.

In India, Some companies like OLX are helping their users to sell their products online. But What Sellers (Customer for company) don't know about it is that what details they are filling, can make their items price go up or down. Machine Learning can help in predicting the likelihood of the price depending on the info provided by the seller and directly suggest that seller what need to be updated to make more effective product sell.

Similar to this situation, Avito (Russian based company) created a challenge on kaggle (<https://www.kaggle.com/c/avito-demand-prediction>) to predict demand for an online advertisement based on its full description, its context and historical demand for similar ads in similar contexts.

I want to check how provided details of a product can affect its demand in live market. It is like a recommended system, except that the user will get a prediction of how well their product is and how well he defined it for a profit.

##### Problem Statement

When selling used goods online, a combination of tiny, nuanced details in a product description can make a big difference in drumming up interest. Details like:

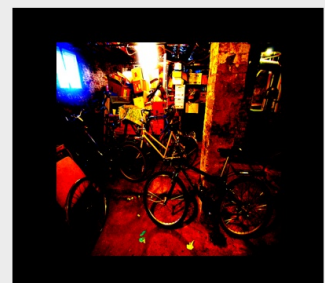
#### Well-Taken, Authentic Photos



Too Glossy



Authentic



Poor Quality

#### Believable and Informative Description Copy

**Description:**

\*\*\*AMAZING WATCH  
FOR SALE!!!!\*\*\*

DON'T MISS THIS  
DEAL. IT'S THE DEAL  
OF THE CENTURY!!

Unlikely

**Description:**

I have an adjustable  
Chaleur D'Animale  
Watch for sale.

It's never been worn  
and still in the original  
box. Battery included.

Informative

**Description:**

fancy watch for sale

no low ball offers, cash  
and carry

Poor Quality

And, even with an optimized product listing, demand for a product may simply not exist—frustrating sellers who may have over-invested in marketing. 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). In this Kaggle competition, Avito is challenging 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.

I want to take a subset of their data to test and verify how prediction by details help a user to get best price of their product.

I will be treating this problem as Regression one. As the value can be in [0,1] both inclusive, as we can never be sure 100% of ad demand.

## Metrics

Metrics are used to evaluate how our algorithms are working. As per the Kaggle challenge, I will be using the following metric:

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y - \hat{y}_i)^2}{n}}$$

Where

predicted Value  $\Rightarrow \hat{y}_i$

original Value  $\Rightarrow y$

## II. Analysis

In [3]:

```
import numpy as np
import pandas as pd

import seaborn as sns
import matplotlib.pyplot as plt
```

## Data Exploration

The dataset is composed of multiple CSV files but I will be mainly using train.csv & test.csv file. It has been obtained from a [Kaggle Competition provided by Avito](#).

The most important file is the train.csv file which has 18 columns containing user id, category, region, price, image flag etc along with the target variable deal\_probability and has 1503424 rows. The test.csv is similar to the previous file discussed but does not have our target variable and we have to use these to predict the destination and has 508438 rows. We have a good amount of data to work with to produce meaningful models.

- train.csv - Train data
  - 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
- test.csv - Test data. Same schema as the train data, minus deal\_probability

In [4]:

```
all_train_dataset = pd.read_csv("../inputs/train.csv")
all_train_dataset.head()
```

Out[4]:

	item_id	user_id	region	city	parent_category_name	category_name	param_1
0	b912c3c6a6ad	e00f8ff2eaf9	Свердловская область	Екатеринбург	Личные вещи	Товары для детей и игрушки	Постельные принадлежности
1	2dac0150717d	39aeb48f0017	Самарская область	Самара	Для дома и дачи	Мебель и интерьер	Другое
2	ba83aefab5dc	91e2f88dd6e3	Ростовская область	Ростов-на-Дону	Бытовая электроника	Аудио и видео	Видео, DVD и Blu-ray плееры
3	02996f1dd2ea	bf5cccea572d	Татарстан	Набережные Челны	Личные вещи	Товары для детей и игрушки	Автомобильные кресла
4	7c90be56d2ab	ef50846afc0b	Волгоградская область	Волгоград	Транспорт	Автомобили	С пробегом

In [5]:

```
all_test_dataset = pd.read_csv("../inputs/test.csv")
all_test_dataset.head()
```

Out[5]:

	item_id	user_id	region	city	parent_category_name	category_name	param_1	param
0	6544e41a8817	db73ad6e4b5	Волгоградская область	Волгоград	Личные вещи	Детская одежда и обувь	Для мальчиков	Обувь
1	65b9484d670f	2e11806abe57	Свердловская область	Нижняя Тура	Хобби и отдых	Велосипеды	Дорожные	NaN
2	8bab230b2ecd	0b850bbebb10	Новосибирская область	Бердск	Бытовая электроника	Аудио и видео	Телевизоры и проекторы	NaN

	item_id	user_id	region	city	parent_category_name	category_name	param_1	param_2
3	8e348601fefc	5f1d5c3ce0da	Саратовская область	Саратов	Для дома и дачи	Бытовая техника	Для кухни	Вытяжка
4	8bd2fe400b89	23e2d97bfc7f	Оренбургская область	Бузулук	Личные вещи	Товары для детей и игрушки	Детские коляски	NaN

## Exploratory Visualization

### Abstract Data:

- Training Data have about 1503424 rows.
- Testing data have 508438.
- **Splitting our dataset in train:validation bins in ratio of 0.75:0.25 0.85:0.15** . As I go through my work, I updated my ratio for final run to this value.
- Deal\_probability is the targeted column.

With multiple runs, I found differnt column names which have NaN values inside, so that I can fill those in my preprocessing step with empty strings or zeros.

Check different paramters and dtypes for columns for better understanding of data types.

In [12]:

```
# all_train_dataset.info()
for col in train_input.columns:
    if (col == "param_1" or col == "param_2" or col == "param_3" or col == "description" or col == "image"): # as NA exist in these columns
        print("Column Name:", col, "    #Unique Values(+1 extra):", len(np.unique(train_input[col].fillna("missing"))))
    else:
        print("Column Name:", col, "    #Unique Values:", len(np.unique(train_input[col])))
```

```
Column Name: item_id      #Unique Values: 1503424
Column Name: user_id      #Unique Values: 771769
Column Name: region       #Unique Values: 28
Column Name: city         #Unique Values: 1733
Column Name: parent_category_name  #Unique Values: 9
Column Name: category_name    #Unique Values: 47
Column Name: param_1         #Unique Values(+1 extra): 372
Column Name: param_2         #Unique Values(+1 extra): 272
Column Name: param_3         #Unique Values(+1 extra): 1220
Column Name: title          #Unique Values: 788377
Column Name: description     #Unique Values(+1 extra): 1317103
Column Name: price          #Unique Values: 102368
Column Name: item_seq_number    #Unique Values: 28232
Column Name: activation_date   #Unique Values: 21
Column Name: user_type       #Unique Values: 3
Column Name: image          #Unique Values(+1 extra): 1390837
Column Name: image_top_1     #Unique Values: 115650
```

In [6]:

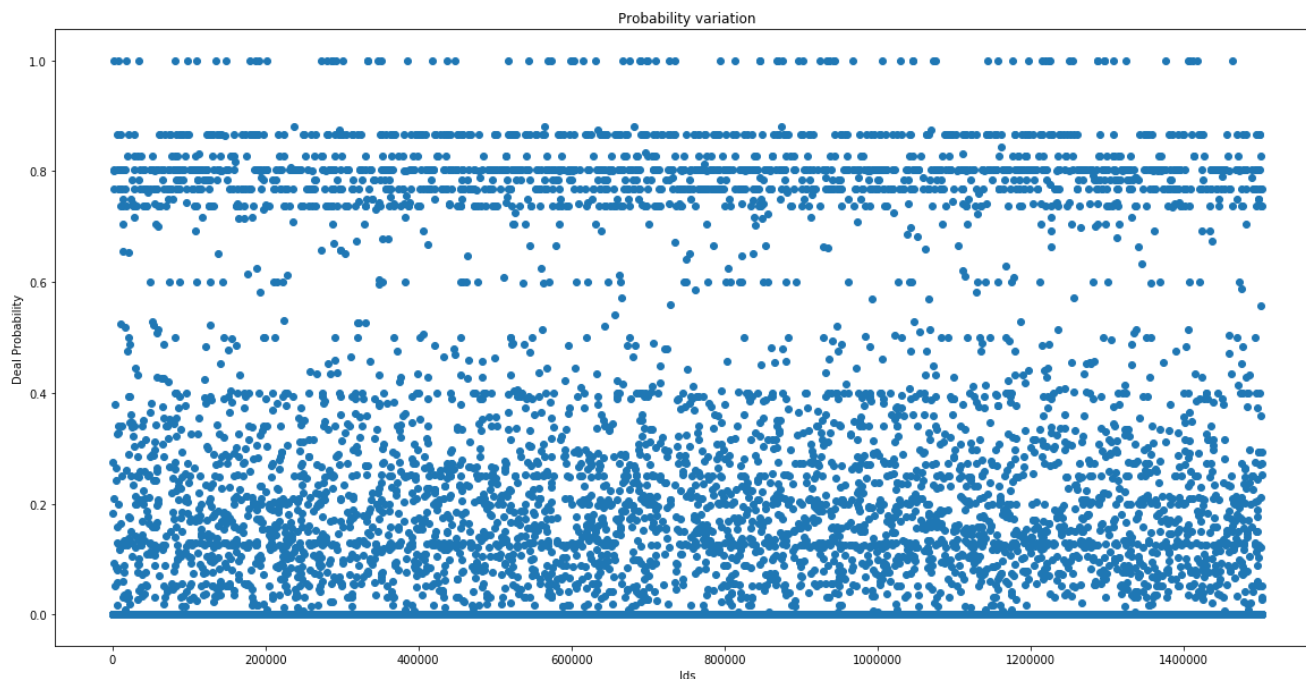
```
train_output = all_train_dataset['deal_probability'].astype('float32')
train_input = all_train_dataset.drop(['deal_probability'], axis=1)
# del all_train_dataset['deal_probability'] ##### can be used
test_input = all_test_dataset.copy()
```

To check how the dataset vary with different deal\_probability. In below graph, we can check it with 1% of the dataset that its more towards low range in deal\_probability.

In [13]:

```
plt.figure(figsize=(20, 10))
dataToShow = 0.01 # percentage
indexes = np.random.randint(0, len(train_output), size=int(dataToShow * len(train_output)))
x_select = train_output.iloc[indexes]
plt.xlabel("Ids")
plt.ylabel("Deal Probability")
plt.title("Probability variation")

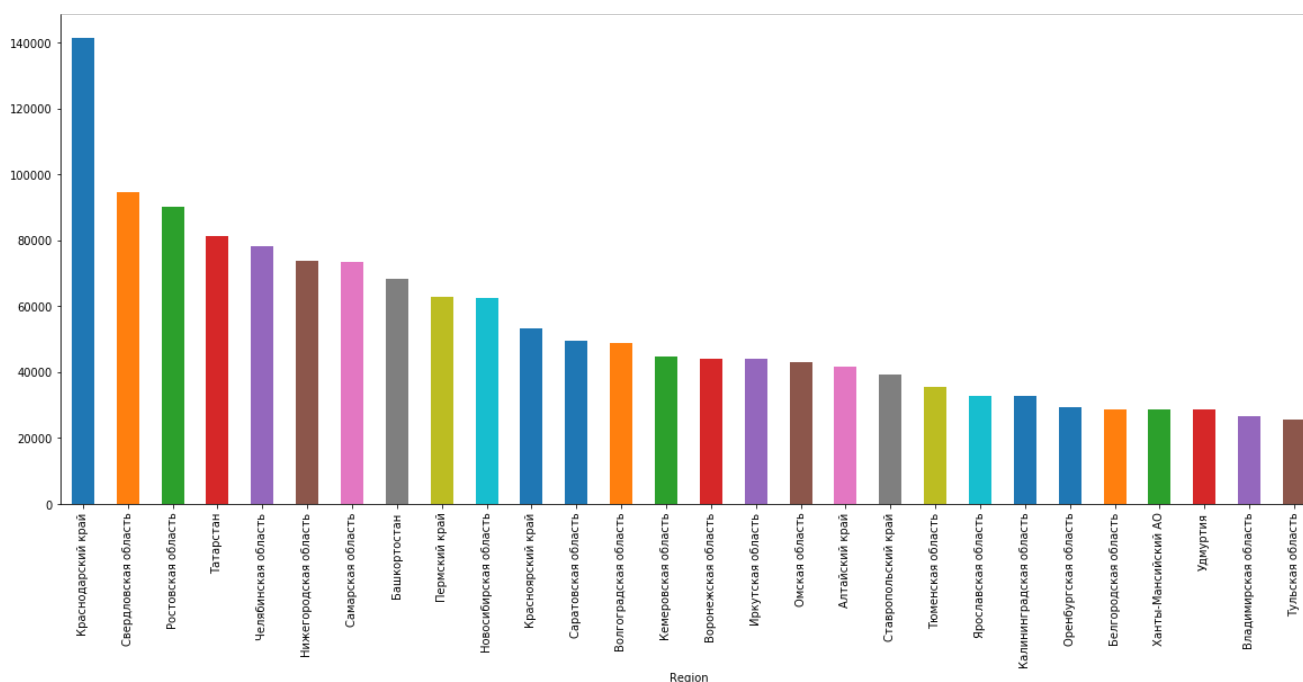
plt.scatter(indexes, x_select)
plt.show()
```



Below 3 graphs shows that region, Category Name and Parent Category Name have few unique values which can be encoded to numeric data for model training.

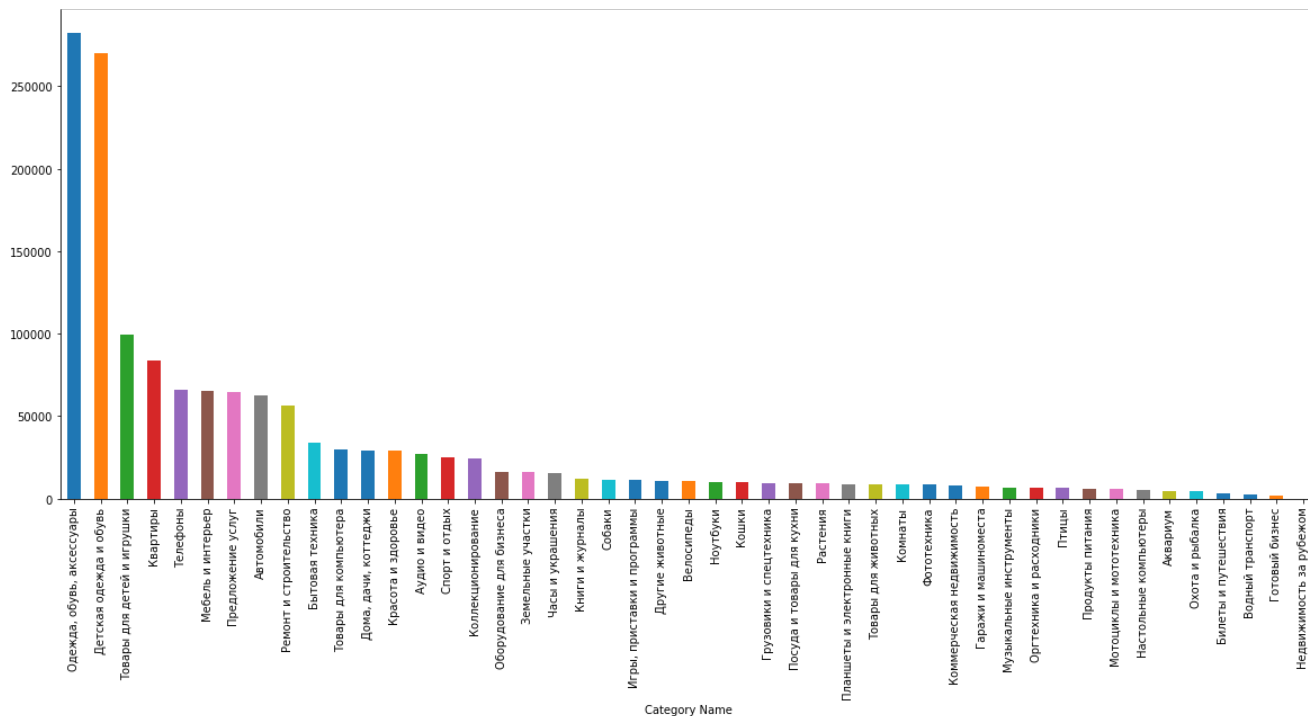
In [7]:

```
plt.figure(figsize=(20, 8))
all_train_dataset.region.value_counts(dropna=False).plot(kind='bar', rot=0)
plt.xlabel('Region')
plt.xticks(rotation=90)
sns.despine()
```



In [8]:

```
plt.figure(figsize=(20, 8))
all_train_dataset.category_name.value_counts(dropna=False).plot(kind='bar', rot=0)
plt.xlabel('Category Name')
plt.xticks(rotation=90)
sns.despine()
```

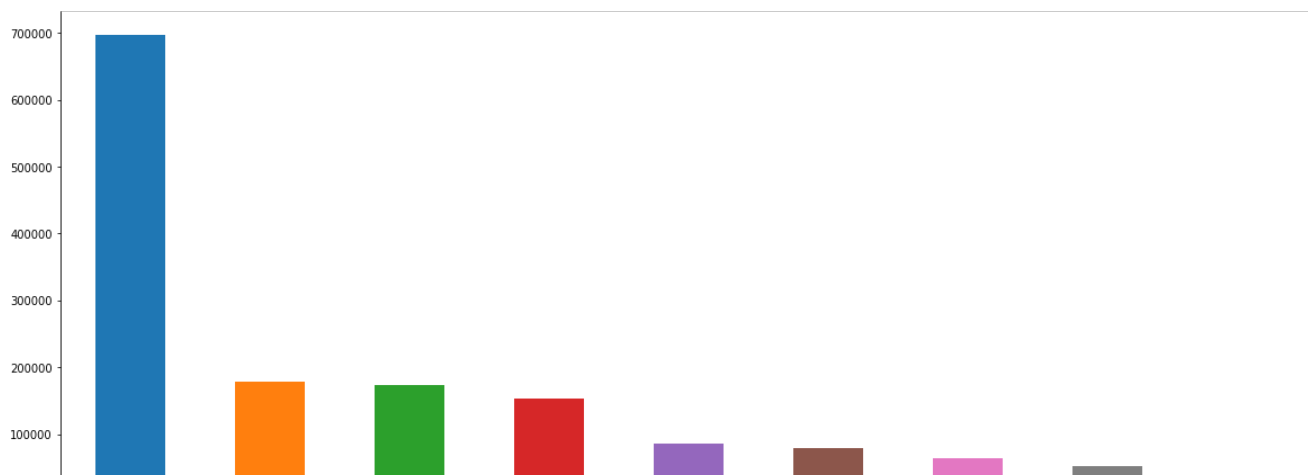


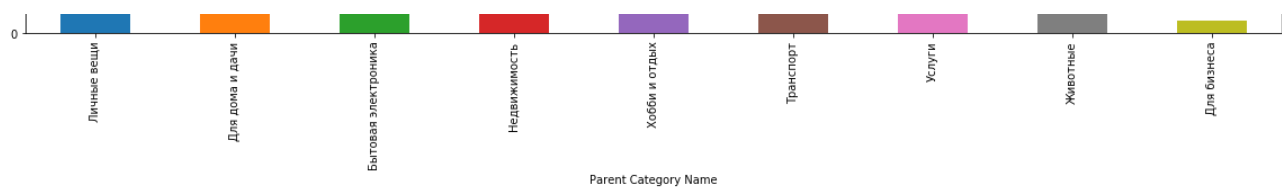
In [9]:

```
plt.figure(figsize=(20, 8))
all_train_dataset.parent_category_name.value_counts(dropna=False).plot(kind='bar', rot=0)
plt.xlabel('Parent Category Name')
plt.xticks(rotation=90)
sns.despine()

# plt.figure(figsize=(20, 8))
# all_train_dataset.city.value_counts(dropna=False).plot(kind='bar', rot=0)
# plt.xlabel('city')
# plt.xticks(rotation=90)
# sns.despine()

# plt.figure(figsize=(20, 8))
# all_train_dataset.price.value_counts(dropna=False).plot(kind='bar', rot=0)
# plt.xlabel('city')
# plt.xticks(rotation=90)
# sns.despine()
```





## Algorithms and Techniques

As this is a supervised learning problem with a regression solution with range  $[0, 1]$ ; I want to use different ensemble methods with tuning of hyper-parameters for the best model later for improvements. I tried with following algorithms:

- The **XGBoostRegressor** algorithm was chosen after research into Kaggle Competitions and it was found out that it proves extremely effective in such arenas. XGBoost is an optimized distributed gradient boosting library designed to be highly efficient, flexible and portable. It produces an ensemble of weak decision tree learners via additive training (boosting). XGBoost is short for Extreme Gradient Boosting. This is based on Gradient boosted trees. Boosted trees are basically an ensemble of decision trees which are fit sequentially so that each new tree makes up for errors in the previously existing set of trees. The model is "boosted" by focusing new additions on correcting the residual errors of the last version of the model. Then you take an approximate step in the gradient direction by training a model to predict the gradient given the data. XGBoost algorithm tuning is a tricky process. Heavy computation power is required for such level of tuning.
- **LGBMRegressor** & **CatBoostRegressor** used because of similar reasons. As these were designed later, they are faster than XGBoost.

More similar algorithms I will try are :

- GradientBoostingRegressor
- AdaBoostRegressor

After getting intuition, I used GridSearchCV for further analysis.

## Benchmark

As this is like a recommended system, I thought of using Decision tree as with GridSearchCV as hyperparameter tuning to save its best model, but later used `RandomForestRegressor` for comparison to other models.

**Random Forest** Classifier fits number of decision trees on subsamples of a dataset and averages the results.

A random forest is a collection of random decision trees. In which at each node you will randomly draw a subset of features and the decision tree will predict the classification. Then the same is done with several trees and bagged. This ensemble method will reduce overfitting and provide good classification. The bias-variance trade-off for this algorithm is good and therefore the possibility of overfitting is drastically reduced.

**RandomForestRegressor** have **RMSE** error of **0.10** in training data but have **0.24** in validation data, meaning its overfit the data in training set.

## III. Methodology

### Data Preprocessing

- Data preprocessing/cleaning
  - Text data one-hot encoded for the required text columns (except for item\_id, user\_id, title, description).
  - As I will be working with csv data, Image related columns converted to bool type.
  - For Outliers detection, I will be using Interquartile range process.
  - filled the missing values with some constant.
  - Split data into training & validation set

**Note :** I didn't use [word vector](#) algorithms (as I am not familiar for those topics right now much.)

### Preprocess steps functions

1. `processData` function is used to process the NaN values in data and also to update different data types as required by different boosting algorithms.

One e.g. I got something like `**ValueError: DataFrame.dtypes for data must be int, float or bool. Did not expect the data types in fields item_id, user_id, activation_date, user_type, title_description.**`

As These columns have text data with varying lengths, I need to drop them for my algorithm training by boosting algorithm methods.

1. `uniqueLabelsExtract` function is using `LabelEncoder` to find different category types for different data in train & test datasets. [Read more about Label Encoders](#).

1. `encode2Labels` function will encode our text data to their respective category for specific columns.

I used log for price column to adjust range of variation here. ( *Using a logarithmic transformation significantly reduces the range of values caused by outliers.* )

Also, for image column, I didn't do processing on those files of zip as per in dataset, But I want to make use of this column to make algorithm check if an image in description effect the outcome or not.

## Implementation

- Evaluate Algorithm

### 1. Build Models (try in between AdaBoost, GradientBoost, Lightgbm, Xgboost, Catboost)

## Training step for evaluating algorithms

1. `train_predict` function used to save different algorithms *train time*, *validation time*, *training error* and *validation error*.\*

As We can check: `XGBRegressor`, `CatBoostRegressor`, `LGBMRegressor` outperforms `GradientBoostingRegressor`, `AdaBoostRegressor`. Let us view the charts in more details:

- `RandomForestRegressor` have RMSE error of 0.10 in training data but have 0.24 in validation data, meaning its overfit the data in training set.
- `AdaBoost` speed is fine, but as per RMSE error, it lags behind other algorithms.
- `GradientBoostingRegressor` is fine also but I want to explore first 3 in more details, i didn't use it further :)
- Best one in all for now is `LGBMRegressor` & `CatBoostRegressor` as similar range of errors help me build my intuition.

## Refinement

- Model tuning to improve results
  - Selected best model with tuning hyper-parameters. First I will choose a range of parameters in some step increments & also some randomized values for covering large parameter space also.

### 2. Select best model

#### Training step with hyperparameters

`useStaticParameters` is used to set static parameters for (algorithms I decided above ) `GridSearchCV`.

`useDynamicParameters` is used to set dynamic parameters generation for (algorithms I decided above ) `GridSearchCV`.

### 3. Make predictions on validation set



modelEvaluator is training each model/learner with GridSearchCV for finding better paramters based on **RMSE metric**

## IV. Results

### Model Evaluation and Validation

Submission @Kaggle of csv file.

10 submissions for [LambaG](#)

Sort by 

Most recent

All

Successful

Selected

Submission and Description	Private Score	Public Score	Use for Final Score
<a href="#">submissionLGBM_S.csv</a> 5 minutes ago by <a href="#">Sachin Lamba</a> Non-Negative LGBM with static hyperparamters variations.	0.2355	0.2319	<input type="checkbox"/>
<a href="#">submissionCat_D.csv</a> 7 minutes ago by <a href="#">Sachin Lamba</a> Non-Negative CatBoost with Dynamic generated hyperparamters variations.	0.2366	0.2330	<input type="checkbox"/>
<a href="#">submissionXGB_S.csv</a> 9 minutes ago by <a href="#">Sachin Lamba</a> Non-Negative XGBoost with static hyperparamters variations.	0.2349	0.2312	<input type="checkbox"/>
<a href="#">submission.csv</a> 19 hours ago by <a href="#">Sachin Lamba</a> Different Boosting algo checker.	0.2343	0.2307	<input type="checkbox"/>

As I can compare, XGB outperform other by small margin. But I need to do more training for much better results. Top score in this competition is at 0.21 which I am not even close. I will be reading more about word embedding so that i can use them in my training set and not to drop them.

### Justification

Unseen data can tell us about a model scalability, how much good it is. My benchmark model failed on this step only. As I checked my model submission on kaggle private dataset, I am getting accuracy to what I expect with the limited dataset I used here. Because of reduction of dataset, Image processing & text to word vector conversion got skiped which made the model fast to predict to comparable best models exist for this problems at kaggle. Hence, These models seems to be robust. Alongside this, the model seems to be trustable and aligns well with expected solutions outcomes.

## V. Conclusion

### Free-Form Visualization

In [10]:

```
print(np.sum(train_output > 0.5))
print(np.sum(train_output <= 0.5))
```

```
177754
1325670
```

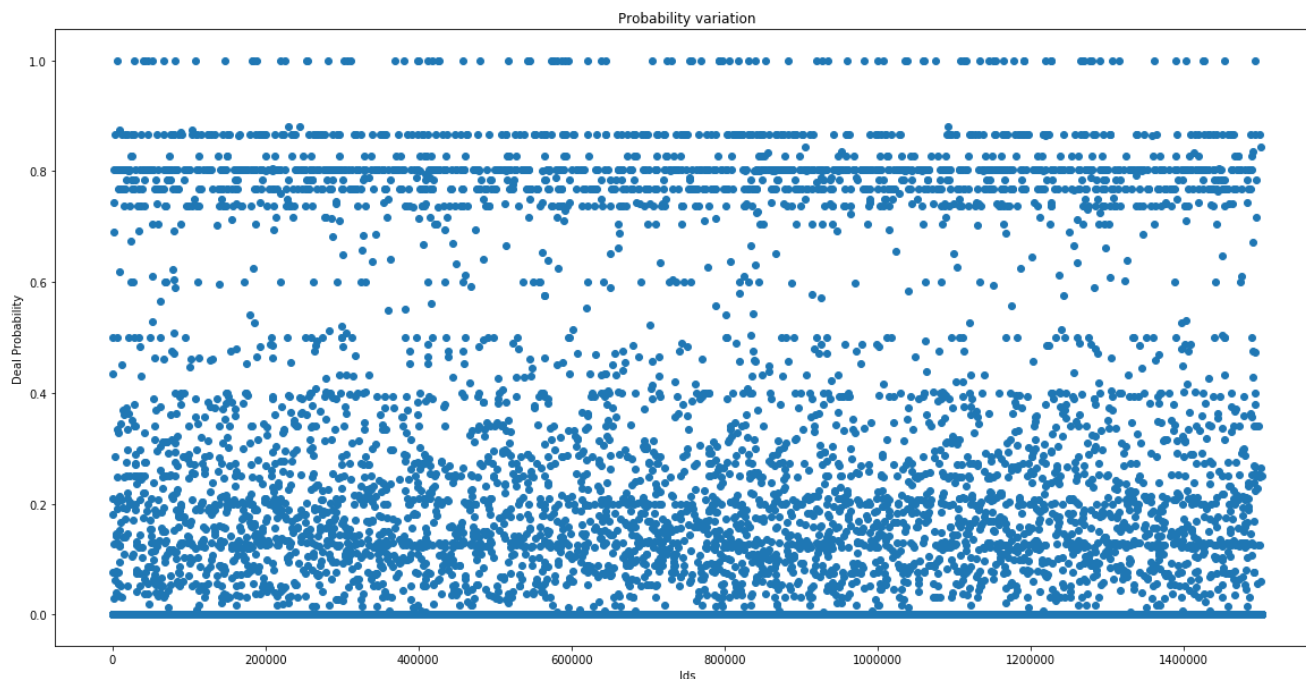
To check how the dataset vary with different deal\_probability. In below graph, we can check it with 1% of the dataset that its more towards low range in deal\_probability.

In [11]:

```
!!! [!!!]
```

```
plt.figure(figsize=(20, 10))
dataToShow = 0.01 # percentage
indexes = np.random.randint(0, len(train_output), size=int(dataToShow * len(train_output)))
x_select = train_output.iloc[indexes]
plt.xlabel("Ids")
plt.ylabel("Deal Probability")
plt.title("Probability variation")

plt.scatter(indexes, x_select)
plt.show()
```



## Reflection

As I find that I need to learn fasttext for text conversion, I tried to do that with some research. I was able to produce word vectors for the text columns too. But later I found that these columns as per word vector will consist of nested list of numeric data(experiments/Capstone Solution.ipynb) for each column which any boost algo didn't expect in its dataset. So, After this searching, I had to drop those columns for this projects to complete it and later learn more about them so that I can do better in NLP part.

## Improvement

Their can be various improvements try:

- One can be to use fasttext for text conversion to word vector. I need to read about more about NLP for that so that those text columns can be utilized by boosting algorithms.
- Use Computer vision to have images tell about product sentiments. It can help more in deducing the probability of buyer interest to some extent.