

OLX Code & the Curious

Popularity Based Ad Recommender System



Akhil Punia

<https://www.linkedin.com/in/akhil-punia-0922bab6/>

30 July 2017

Part 1- Dataset Exploration

Comments on Dataset Quality & Errors Found

(Refer to OLX-Data-Analysis-Graphlab.ipynb)

#1. ads_data.csv

Most redundant variables were **latitude, longitude and product description**.

Most useful variables include **ad_id, category_id , seller_id , title , price and source**.

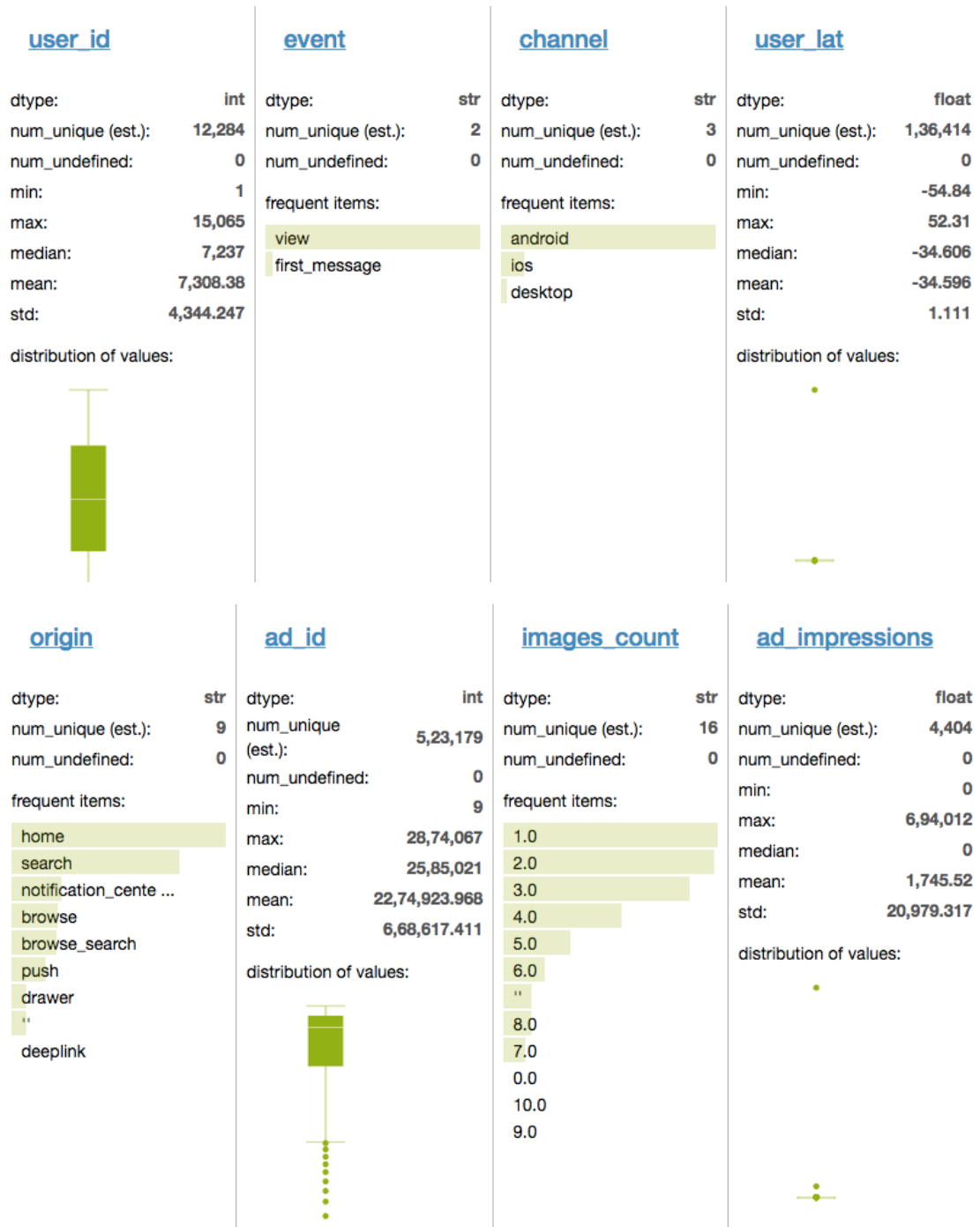
<u>ad_id</u>		<u>category_id</u>		<u>seller_id</u>		<u>creation_time</u>	
dtype:	int	dtype:	int	dtype:	int	dtype:	str
num_unique (est.):	5,70,262	num_unique (est.):	14	num_unique (est.):	1,00,956	num_unique (est.):	5,57,693
num_undefined:	0	num_undefined:	0	num_undefined:	0	num_undefined:	0
min:	0	min:	5	min:	2	frequent items:	
max:	29,22,041	max:	888	max:	6,46,815	No values appear with \geq	
median:	23,75,611	median:	815	median:	3,86,256	0.01% occurrence.	
mean:	21,08,882.133	mean:	810.952	mean:	3,66,124.074		
std:	7,57,628.493	std:	95.939	std:	1,88,590.59		

<u>title</u>		<u>description</u>		<u>price</u>		<u>lat</u>	
dtype:	str	dtype:	str	dtype:	int	dtype:	str
num_unique (est.):	4,30,873	num_unique (est.):	4,62,669	num_unique (est.):	2,158	num_unique (est.):	49,621
num_undefined:	0	num_undefined:	0	num_undefined:	0	num_undefined:	0
frequent items:		frequent items:		min:	0	frequent items:	
Campera		"		max:	1.127e+8	"	
Zapatos		Excelente estado		median:	300	-34.60373306	
Vestido		Nuevo		mean:	4,013.705	-34.60399999	
Zapatillas		Talle M		std:	1,77,488.926	0E-8	
Remera		Muy buen estado		distribution of values:		-34.59000015	
Vestido de fiesta		Talle S				-34.57099999	
Botas		Nueva				-34.70381546	
Cartera		Buen estado				-34.57114959	
Camisa		Sin uso				-34.65490729	
Jeans		Talle L				-34.60373450	
Campera de cuero		Impecable				-34.56208440	
Musculosa		En buen estado				-34.52812060	

#2. userdata.csv

Most redundant variables were **ad_messges**, **user_long** and **user_lat**.

Most useful variables include **user_id**, **event**, **channel**, **origin**, **ad_id**, **images_count** and **ad_views**.



Part 2- Data Handling

Summary of Dataset Preprocessing Steps

Step 1: Reformatting user_messages.csv

(Refer to OLX-User-Data-Restructuring.ipynb)

Step 2: Segmenting the User Message Data based on 10 Unique Categories

```
In [77]: user_messages.sort_values('category_id', inplace=True)
user_messages.head(5)
```

```
Out[77]:
```

	user_id	category_id	ad_id
12167	7597	362	1730320
19091	11938	362	2388804
11279	6452	362	1238103
11280	6452	362	2187575
11281	6452	362	2863136

Step 3: Finding the the Ads with the Highest Frequency in the Given Category ID and sorting them in Descending Order

```
In [109]: pred = list()
for i in xrange(10):
    df1 = user_messages[user_messages['category_id'] == category[i]]
    ads_all = list(df1['ad_id'])
    ads_uni = df1['ad_id'].unique()

    add_freq = list()

    for i in xrange(len(ads_uni)):
        c = ads_all.count(ads_uni[i])
        add_freq.append(c)

    d = {'Frequency': add_freq, 'ad_id': ads_uni}
    dl = df(data = d)

    dl.sort_values('Frequency', ascending = False, inplace = True)
```

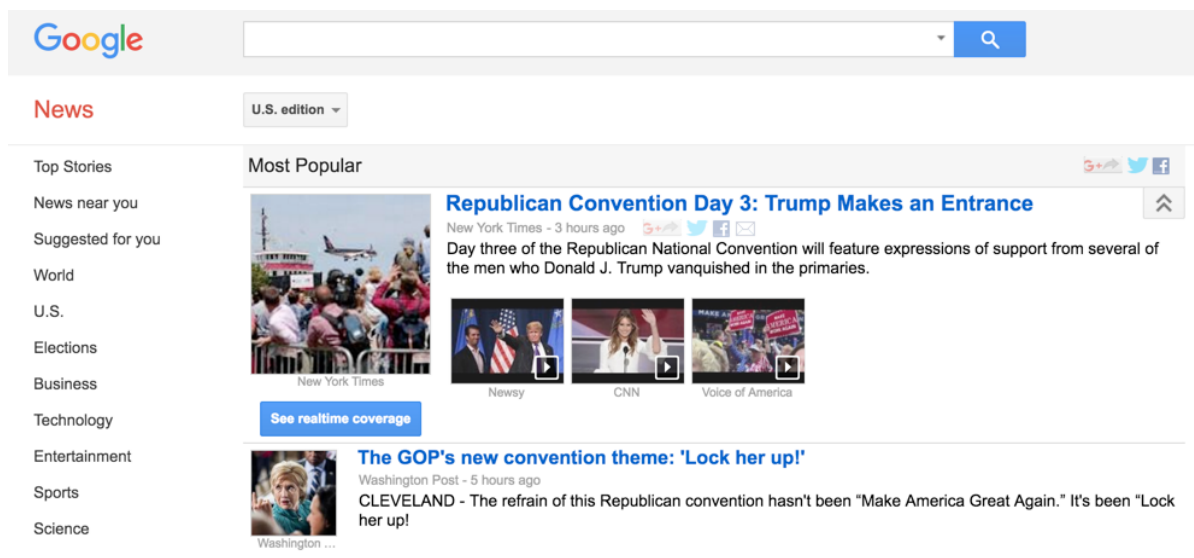
Part 3- Model Justification

Why I chose Popularity based recommender system model?

The **simplest approach** could be to recommend the ads which will be suitable for most number of users. This is a blazing fast and dirty approach and thus has a major drawback.

The thing is, there is **no personalisation involved** with this approach.

Basically, the 10 most popular ads within a particular category (1 out of 10 at a time) would be same for each user since, popularity is defined on the entire user pool. So everybody will see the same ads. It sounds like, 'a website recommends you to buy HP Laptop in the Electronics category just because it's been liked by other users and doesn't care if you are even interested in buying or not'.



Surprisingly, such approach still works in places like news portals. Whenever you login to say Google News where we will see a column of “Popular News” which is subdivided into sections and the most read articles of each section are displayed. This approach works in this case because:

- There is division by section so user can look at the section of his interest.
- At a time there are only a few hot topics and there is a high chance that a user wants to read the news which is being read by most others

Similarly, it can be used OLX dataset as the users are divided into 10 separate categories.