

Avito demand prediction

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).

Business objective

To predict the deal probability of a product given 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.

Data Information

1. train.csv

- 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.

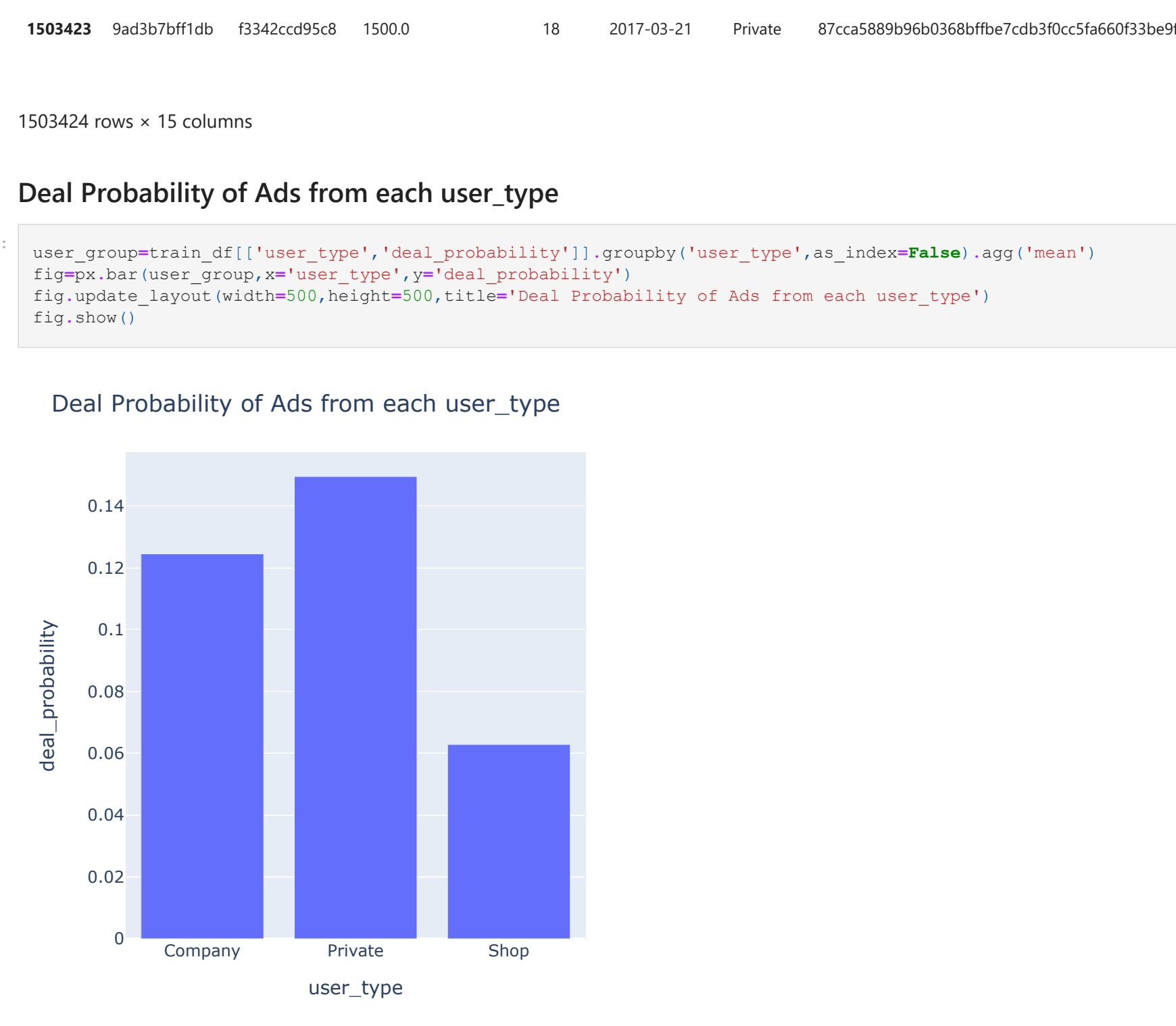
2. test.csv

- Test data. Same schema as the train data, minus deal_probability.

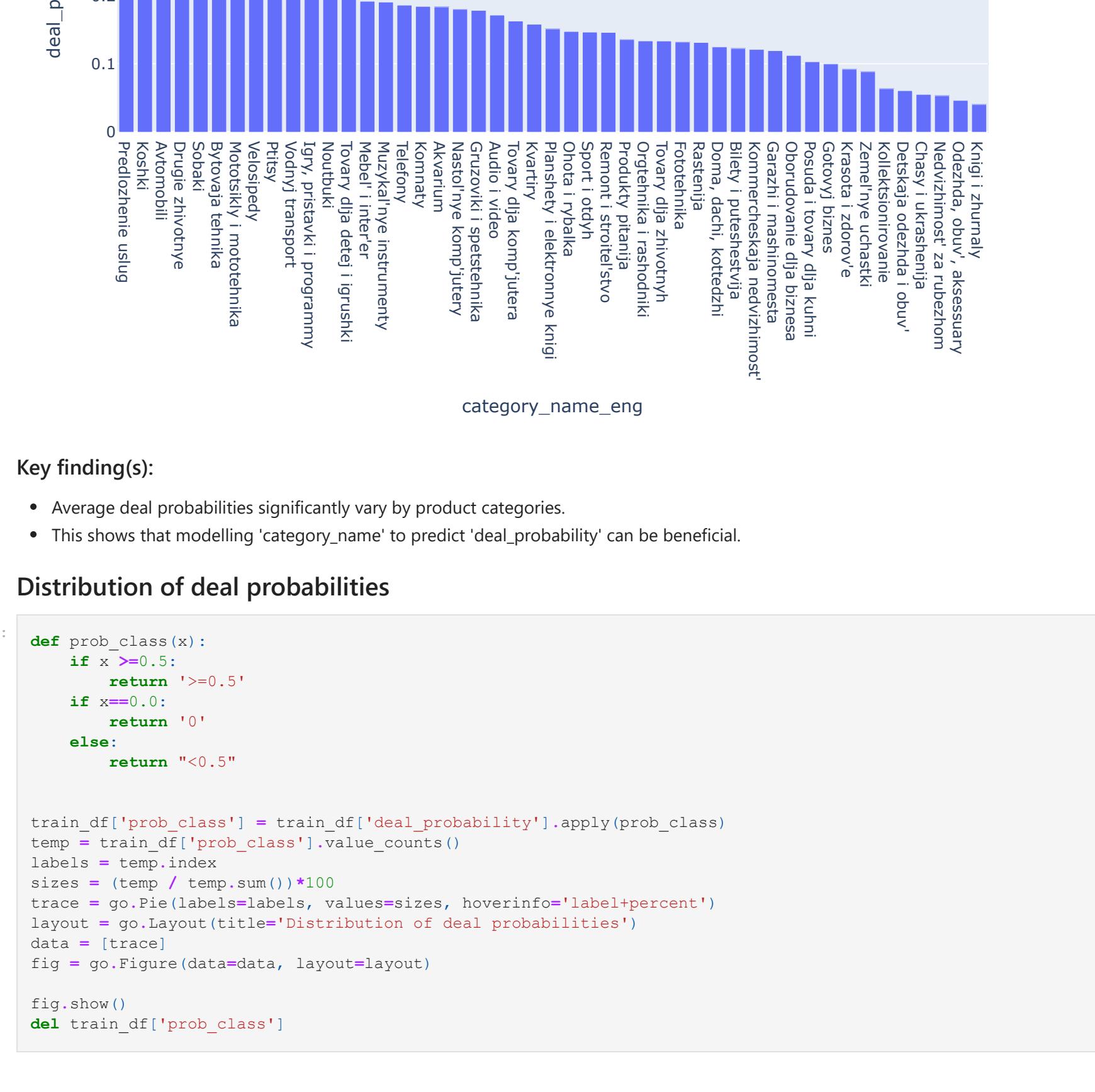
3. Ad images

- train.jpg.zip - Images from the ads in train.csv.
- test.jpg.zip - Images from the ads in test.csv.

Well-Taken, Authentic Photos



Believable and Informative Description Copy



3. Other files

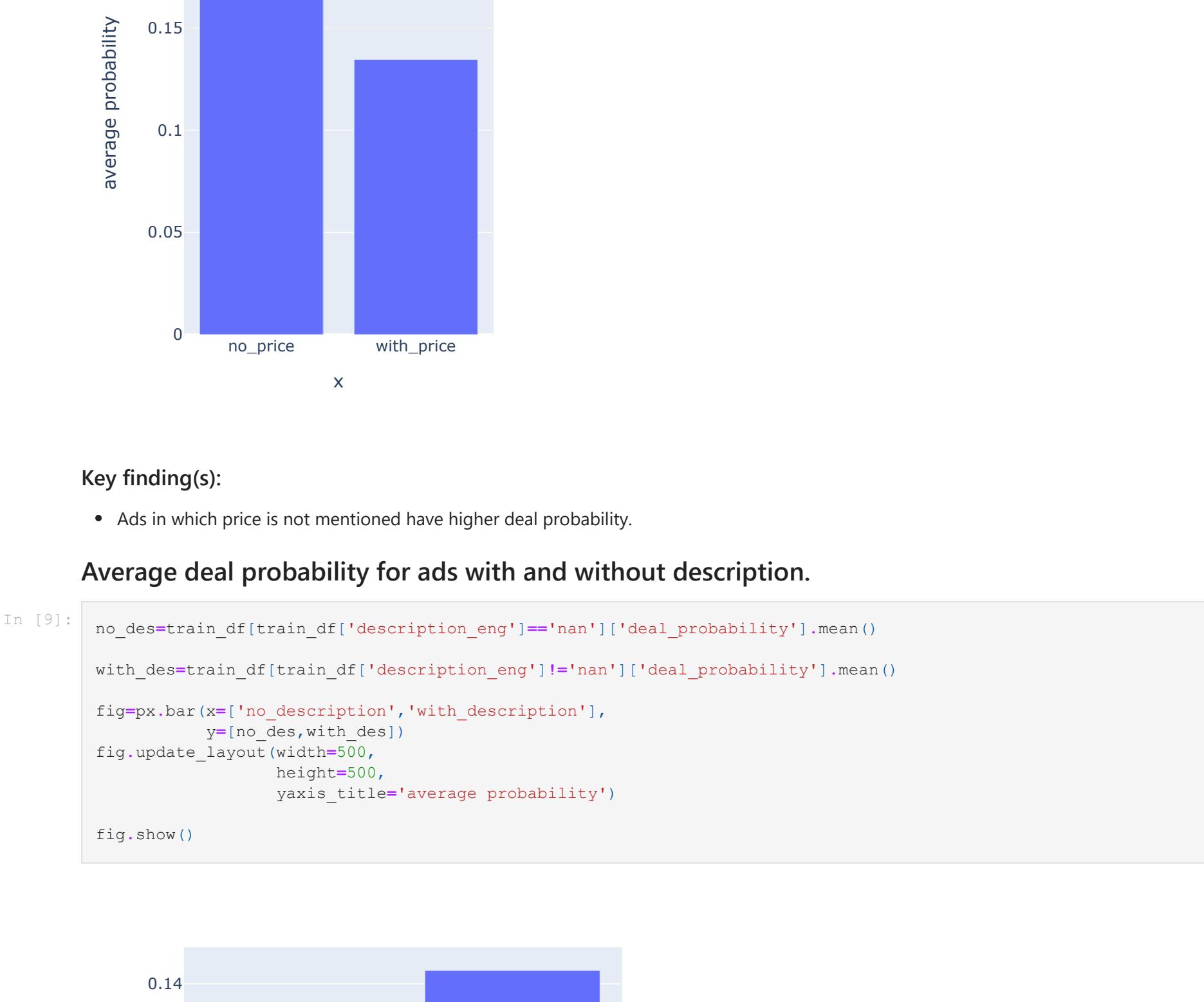
- train_active.csv - Supplemental data from ads that were displayed during the same period as train.csv. Same schema as * the train data minus deal_probability, image, and image_top_1.
- test_active.csv - Supplemental data from ads that were displayed during the same period as test.csv. Same schema as the * train data minus deal_probability, image, and image_top_1.
- periods_train.csv - Supplemental data showing the dates when the ads from train_active.csv were activated and when they were displayed.
 - item_id - Ad id. Maps to an id in train_active.csv. IDs may show up multiple times in this file if the ad was renewed.
 - activation_date - Date the ad was placed.
 - date_from - First day the ad was displayed.
 - date_to - Last day the ad was displayed.
- periods_test.csv - Supplemental data showing the dates when the ads from test_active.csv were activated and when they were displayed. Same schema as periods_train.csv, except that the item_ids map to an ad in test_active.csv.
- sample_submission.csv - A sample submission in the correct format.

Evaluation metric

- RMSE (Root mean squared error)

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (Predicted_i - Actual_i)^2}{N}}$$

- Where,
 - predicted = The predicted value for the ith observation.
 - actual = The observed(actual) value for the ith observation
 - N = Total number of observations.



Key finding(s):

- Ads from user type 'Shop' have the lowest probability to get a deal.
- Ads from user type 'Private' and 'Company' have the comparatively much higher deal probability than that of 'Shop'.
- This shows that modelling 'user_type' to predict 'deal_probability' can be beneficial.



Key finding(s):

- Average deal probabilities significantly vary by product categories.
- This shows that modelling 'category_name' to predict 'deal_probability' can be beneficial.

Distribution of deal probabilities



Key finding(s):

- More than 65% of ads have zero deal probability.
- More than 1% of ads have less than 0.5 deal probability.
- Only 20.25% of ads have more than 0.5 deal probability.

Key finding(s):

- Number of ads posted drastically vary region to region.

Average deal probability for ads with and without mentioned price.

Key finding(s):

- Ads in which price is not mentioned have higher deal probability.

Average deal probability for ads with and without description.

Key finding(s):

- Ads with product description have higher deal probability.

Checking corrupted images

Key finding(s):

- Ads in which no image have higher deal probability.

Average deal probability for ads with and without product image.

Key finding(s):

- Ads with product image have higher deal probability.

Checking corrupted images

