

kaggle™



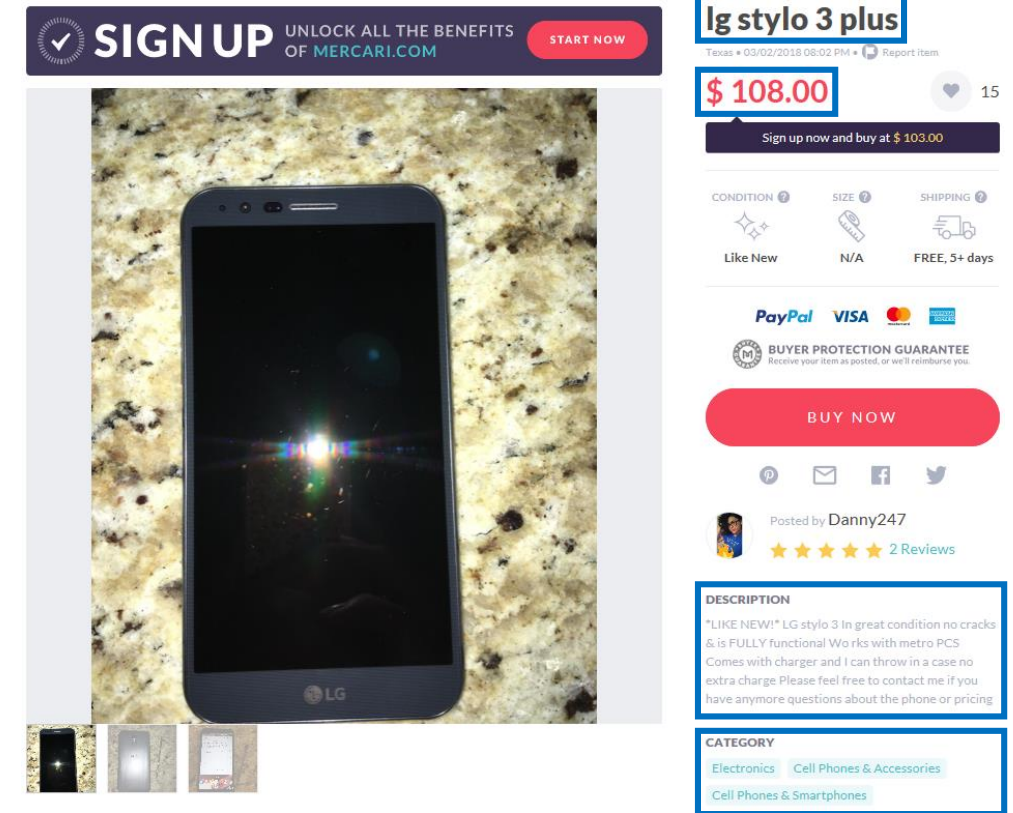
**mercari**

# MERCARI PRICE SUGGESTION CHALLENGE

46TH PLACE SOLUTION – BY THOMAS SELECK

# COMPETITION CONTEXT AND GOAL

- Mercari: Japan's biggest community-powered shopping app
- It can be hard to know how much something's really worth: small details can mean big differences in pricing.
- Goal of the competition: Build an algorithm that automatically suggests the right product prices using user-inputted text descriptions of their products, product category name, brand name, and item condition.

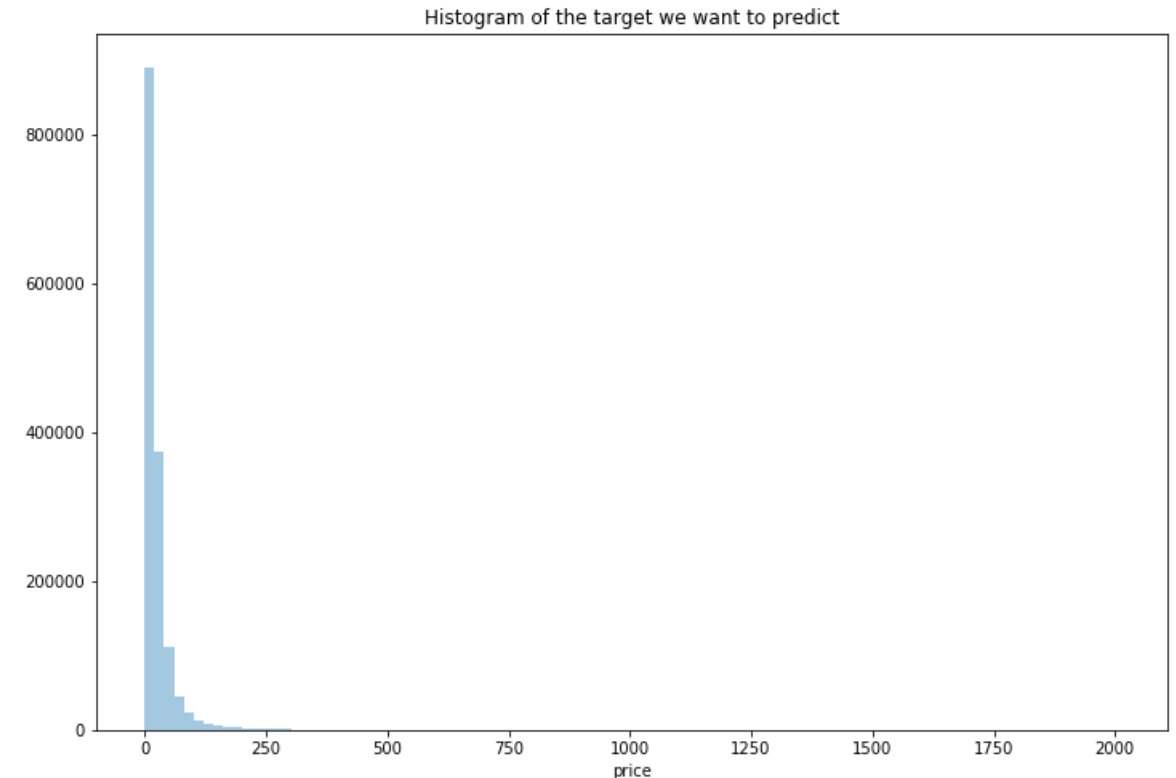


# EVALUATION METRIC AND DATA

- Evaluation metric: RMSLE (Root Mean Squared Logarithmic Error)

$$RMSLE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\log(p_i + 1) - \log(a_i + 1))^2}$$

- $n$  is the total number of observations in the data
  - $p_i$  is the predicted price
  - $a_i$  is the actual sale price for the  $i^{th}$  article
- Data:
    - 1,482,535 samples in train set
    - 693,359 samples in test set (phase 1)
    - 3,5 million samples in test set (phase 2)
    - 6 features: *name*, *item\_condition\_id*, *category\_name*, *brand\_name*, *shipping*, *item\_description*



# FEATURE ENGINEERING (1 / 2)

- 874 items had a zero price in train set  $\Rightarrow$  remove those items
- *brand\_name* feature: 42.67% of items in train set don't have a brand  $\Rightarrow$  Create dummy: 1 == no brand
- For *item\_condition\_id*: goes from 1 to 5; 1 = brand new and 5 = broken; reverse levels order (1 = 5, 2 = 4, ...)
- Compute average price for each category found in *category\_name*.
  - Group categories that have less than 10 samples to avoid overfitting.
  - Create 7 new features based on it using binning:  $\text{avg} \in [0;10[, [10;20[, [20;40[, [40;50[, [50;75[, [75;+\infty[$
- Extract first 3 levels out of 5 from *category\_name* and fill missing values with “missing” string
  - Last 2 levels are only used by 0.3% of samples
  - E.g. Women/Jewelry/Necklaces  $\Rightarrow$  *category\_1* = Women, *category\_2* = Jewelry, *category\_3* = Necklaces
- Fill missing values in *brand\_name*, *name* and *item\_description* with “missing” string
- Group least occurring brands, as 1,243 brands out of 4,810 only appears once in train set

# FEATURE ENGINEERING (2 / 2)

- Luxury brands are expensive  $\Rightarrow$  Dummy equal to 1 for 22 luxury brands (*Louis Vuitton, Rolex, Apple, ...*)
  - Do the same for cheapest brands
- Look for important keywords
  - “dust”  $\Rightarrow$  refers to “dust bag”, a women’s handbag accessory: brand new luxury bags have one, others not
  - “gold”  $\Rightarrow$  as gold is a precious metal, having gold in the item raises its price
  - “lularoe”  $\Rightarrow$  clothes brand that was one of the most important features for LightGBM
  - “bundle”  $\Rightarrow$  this means several objects to sell and can increase the price
- Compute some statistics on *name* and *item\_description*
  - Number of characters, tokens, words, numbers, letters, digits
- Add brand groups depending on *category\_1*
- Normalize the created features (subtract mean and divide by standard deviation)
- Use WordBatch on *name* and *item\_description* (text processing), LabelBinarizer for *brand\_name*, CountVectorizer for *category\_1*, *category\_2* and *category\_3*

# PREDICTIVE MODELS

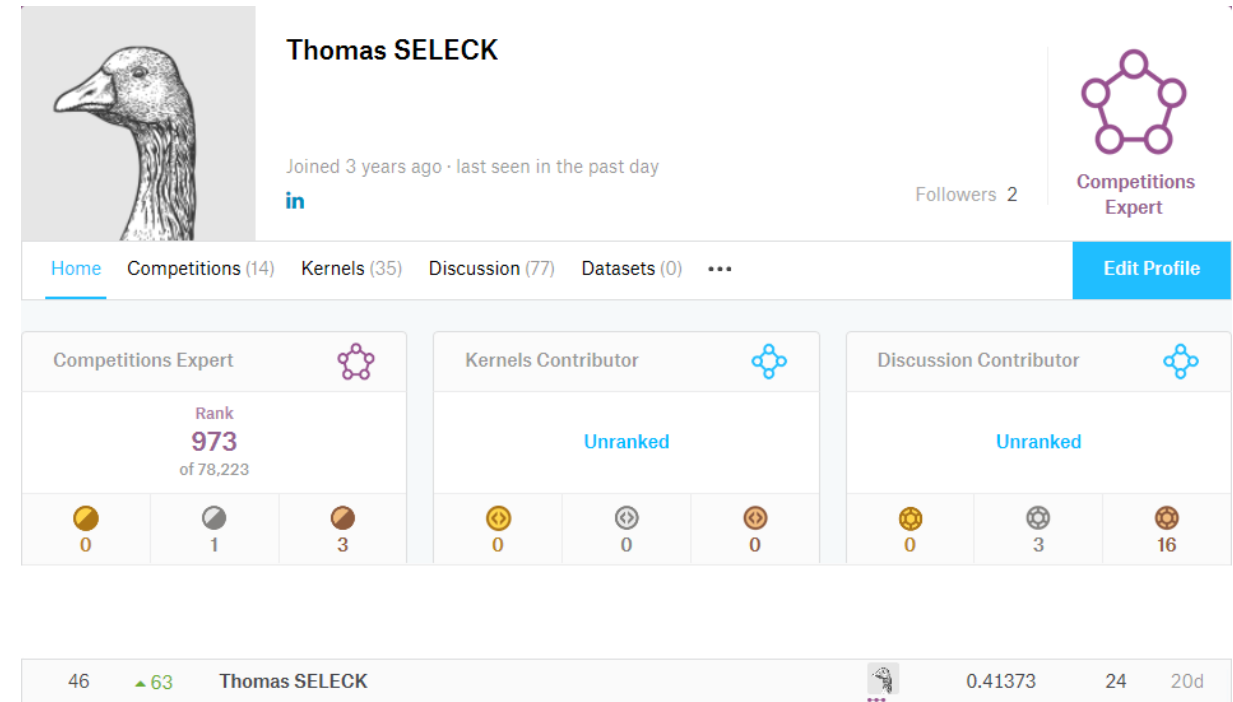
- FTRL (Follow The Regularized Leader): kind of adaptive-learning-rate sparse linear regression with efficient L1-L2-regularization
  - Comes from WordBatch package: <https://github.com/anttttti/Wordbatch>
  - Original paper: H. B. McMahan, “Follow-the-Regularized-Leader and Mirror Descent: Equivalence Theorems and L1 Regularization”: <https://static.googleusercontent.com/media/research.google.com/fr//pubs/archive/37013.pdf>
- FM\_FTRL: Factorization Machine with linear effects estimated with FTRL and factor effects estimated with adaptive SGD.
  - Comes from WordBatch package: <https://github.com/anttttti/Wordbatch>
- LightGBM: Gradient boosted trees library developed by Microsoft
  - More info here: <https://github.com/Microsoft/LightGBM>

# BLENDING STRATEGY

- Quick and Dirty way:
  - Split train set into 2 sets:  $X_{train}$  and  $X_{valid}$  ( $y_{train}$  and  $y_{valid}$  for target)
  - Group predictions in a DataFrame and split it by `category_1`
  - Fit a linear regression without intercept for each split, trying to predict price using predictions made by all three models.
  - Use linear regression coefficients as weights for blending
  - E.g. “Men”:  $[-0.01, 0.65, 0.36] \Rightarrow$  preds for “Men” split =  $-0.01 \times \text{FTRL} + 0.65 \times \text{FM\_FTRL} + 0.36 \times \text{LightGBM}$

# FINAL RESULTS AND CONCLUSION

- Public LB score: 0.41316
- Private LB score: 0.41373



Thomas SELECK

Joined 3 years ago · last seen in the past day

Followers 2

Competitions Expert

Home Competitions (14) Kernels (35) Discussion (77) Datasets (0) ... Edit Profile

Competitions Expert	Kernels Contributor	Discussion Contributor
Rank <b>973</b> of 78,223	Unranked	Unranked
0  1  3	0  0  0	0  3  16

46 ▲ 63 Thomas SELECK 0.41373 24 20d

- Code and slides available here:
  - [https://github.com/ThomasSELECK/Kaggle\\_Mercari\\_competition](https://github.com/ThomasSELECK/Kaggle_Mercari_competition)



# ANY QUESTIONS?

DON'T FORGET: [HTTPS://GITHUB.COM/THOMASSELECK/KAGGLE\\_MERCARI\\_COMPETITION](https://github.com/thomasseleck/kaggle_mercari_competition)