

## Mercari Price Prediction Taskforce

- Data exploration / Data cleaning
- Feature engg / selection
- Performance metric (RMSE ?)
- Model eval & selection / optimization
- Conclusion / report doc

← Already did part →

To visualize text data  $\xrightarrow{\text{plot}}$  word clouds

- Sort data acc. to prices (low  $\rightarrow$  high)
- Divide data into 4 equal parts
- 1<sup>st</sup> Quarter  $\rightarrow$  products with lowest prices
- 4<sup>th</sup> Quarter  $\rightarrow$  products with highest prices.

Observt<sup>n</sup>: Sellers try to put good words (<sup>never used, brand new etc</sup>) to sell the products quickly <sup>commonly found in all quarters</sup>

— % rows with No-descript<sup>n</sup> yet.

- \* Boxplot to check if products are impacted when description is given or not.  
 $\hookrightarrow$  not much diff (median at both occasions is almost same).

Wht abt 75<sup>th</sup> percentile values in both? —

Features selected for modelling:

- |                     |                    |
|---------------------|--------------------|
| 1) brand-name       | 6) subcategory - 1 |
| 2) brand-name-given | 7) subcategory - 2 |

- 3) item-description-id
- 4) shipping
- 5) main-category
- 8) item-description

Note: While applying ONE-HOT, apply it to test data with respect to train data to avoid **data leakage**

To represent item-description feature:

- Let's apply Bag of words to test in simple.

*(similar to One-hot encoding)*

→ builds vocabulary of unique words in our dataset & associate unique column to each word in vocab.

• Each sentence is rep as list with length = # distinct words in vocabulary

② each col. in list, mark how many times given word appears in sentence. *(ignores order of words)*.

Evaluation metric:

\* RMSE (Root Mean Squared Error) is used for regression tasks.

becz of skew in target col. (price), RMSLE (log error) would be relevant.

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

RMSLE

$\epsilon \rightarrow$  RMSLE score  
 $n \rightarrow$  #total observations in public/private data set.  
 $p_i \rightarrow$  your prediction of price  
 $a_i \rightarrow$  actual price

Establishing a baseline:

To compare our results against.

For regression problems  $\rightarrow$  mean of price in train set.

Baseline error  $\rightarrow$  \_\_\_\_\_.

## Model selection / hyper-param tuning

- Try diff algos & see which works best

- Linear Regression
- Support vector Regression (?)
- Decision Tree Regression (?)
- Logistic Regression

\* Above models can be optimized by tuning their **Hyperparameters**. (eg: #trees in random forest)  
# neighbours in K-nearest

**Model parameters**  $\rightarrow$  which model learns during training (eg: weights in Linear Regression)

Control balance bet **Underfitting & Overfitting**

when model  $\leftarrow$  is not complex enough to learn map from features  $\rightarrow$  target  
**(High bias)**

corrected by making model more complex

soot of memorizes training data.  
**(High variance)**

corrected by - limiting complexity of model tho'

**Regularization**

④ There's no universal correct hyperparameters, every ML probelm will be diff.

Hyper-parameter Tuning technique:

→ Random search with K-fold cross validation.

↳ defines a grid, then randomly sample diff combinations & perform K-fold CV.

K-fold CV: Given dataset is split into K. no of sections / folds, where each fold is used as testing set at some point of time.

Test error from all iters are averaged for that specific hyper-parameter combination.



e.g: 5-fold cross-validation.

(one-fold used for test, rest all for training.)