

# Mercari price suggestion challenge

Ankit Bhadu  
2018CSB1073

Adarsh Kumar  
2018CSB1066

Indian Institute of Technology Ropar

Guided By :  
Dr Subramanian

---

## Abstract

Mercari is Japan's biggest community powered shopping app where people can sell and buy a variety of brand new and used products of different brands, from sweaters to smartphones. Mercari would like to suggest the correct prices to the sellers but this is tough because their sellers are enabled to put just about anything, or any bundle of things on Mercari's marketplace. Our objective is to build a model for Mercari that automatically suggests the right product prices to the sellers.

## 1 Introduction

To start with we try out various regression techniques as this problem lies under the category of supervised regression machine learning.

### 1.1 Exploratory Data Analysis

#### 1. Contents :-

Train id — the id of the product

Name — the title of the product

Item condition id — the condition of the product provided by the sellers

Category name — category of the product

Brand name — the product's brand name

Shipping — 1 if shipping fee is paid by seller and 0 if shipping fee is paid by buyer

Item description — the full description of the product

Price(Target Variable) — the price that the product was sold for

#### 2. Price Distribution:-

We plot a histogram to observe how the price of items is distributed and come to know that most items are priced between [15,25]. Since the data is highly right skewed we take logarithm. Since price can have value 0 we take  $\log(\text{price}+1)$  to avoid log of 0. After this transformation the plot seems to be well distributed as shown in FIG-1(b). We can observe that now most of the price are between 2-4.

#### 3. How bearer of shipping charges affect price

We plot another histogram(Figure-2) to see how it changes price of an item when

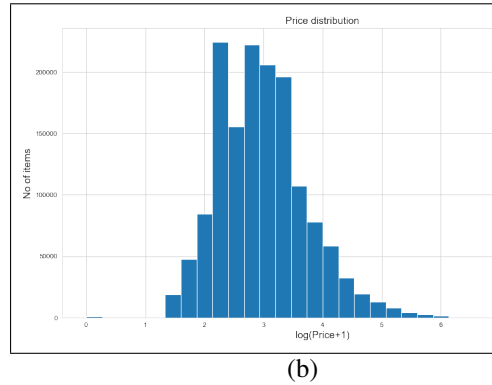
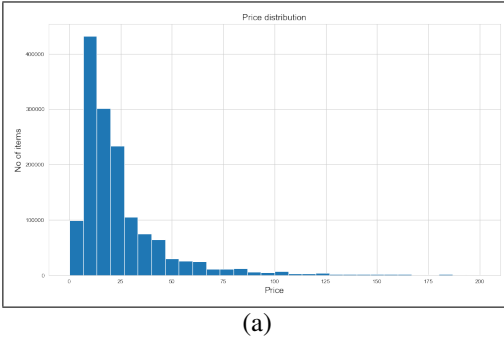


Figure 1: price distribution (a) original (b) logarithmic transform

shipping charges are paid of by the buyer vs when they are paid by the seller. We can observe shipping charges are more often paid often by the seller as opposed for the higher priced products where its paid by the buyers. It is probably due to profit reasons for the lower priced products however as we can see there is lot of overlap too between the two categories.

#### 4. Category of product

Category column can be divided into 3 based on '/' say main category and then sub categories. We found that data contains 1288 categories and in main category section women had maximum frequency and in sub category Athletic apparel was the top one. We plotted several plots for distribution according to categories and observing relations between them but it was of very less use to the main problem,

#### 5. Description length

A scatter plot[figure 4] of Distribution length vs price showed that higher priced items required less description as compared to lower priced ones. It shows that branded items which are usually higher priced do not need much description, i.e. they sell themselves.

#### 6. Correlation

The correlation matrix[figure 5] makes the picture altogether clear as the relations drawn until now are supported by it. We can see that description length has a fair correlation with the price of an item and hence we have selected it as a feature in addition to item\_description.

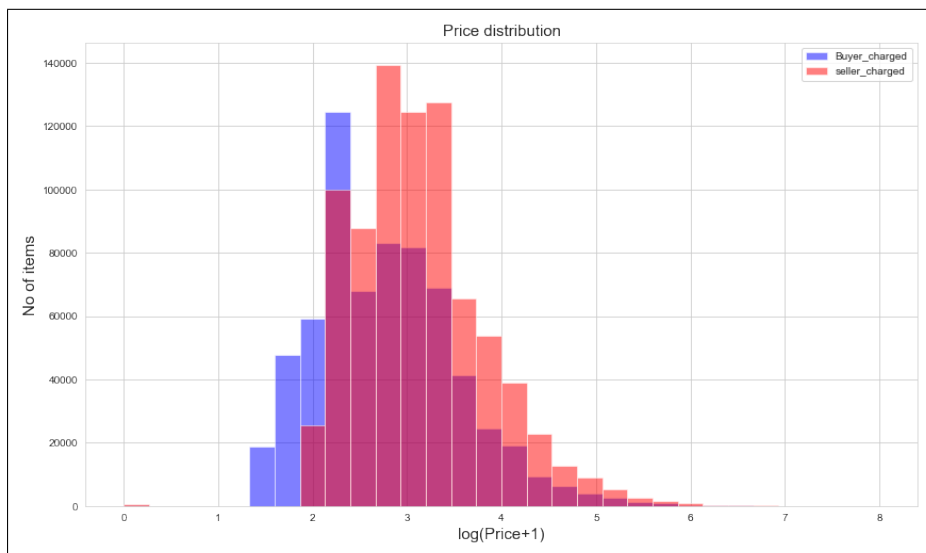
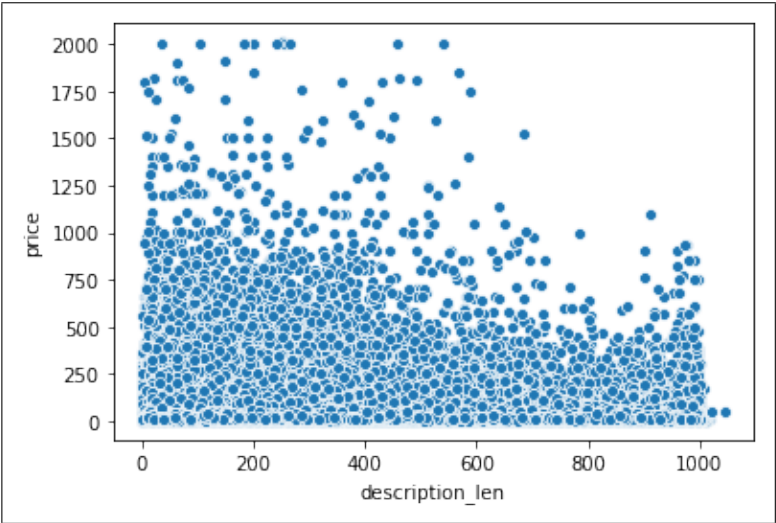


Figure 2: How bearer of shipping charges affect price

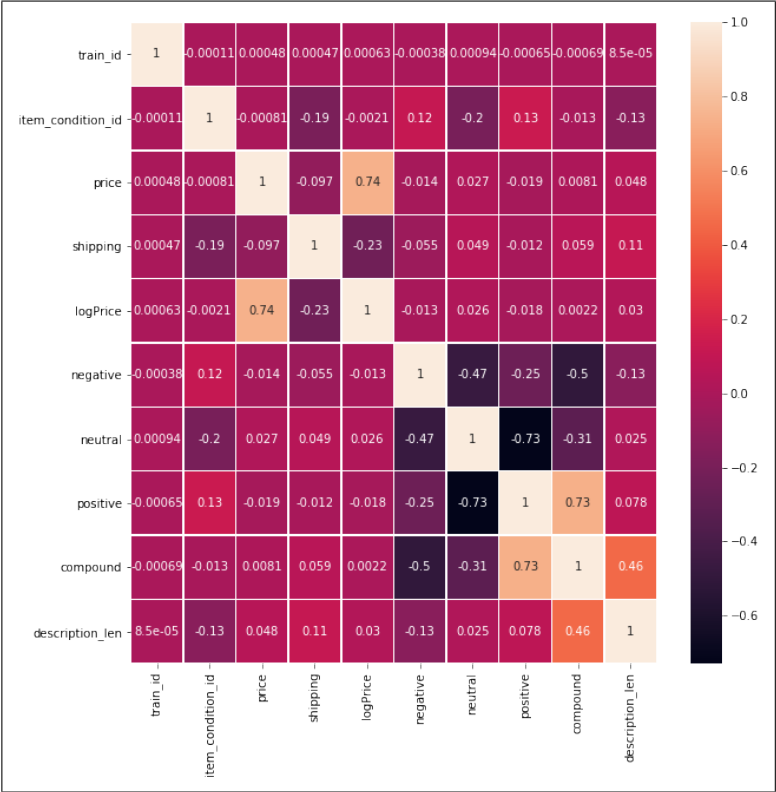


Figure 3: visualization of frequently used words in description through wordcloud



(a)

Figure 4: Price vs Description length



(a)

Figure 5: Correlation between different features

## 1.2 Pre-processing

Firstly we change all short forms of words like aren't, we've, we'll to their complete forms to avoid separate vectorization of say aren't and "are not". Then we remove all stop words which won't convey any meaning. All stopwords have been taken from "<https://gist.github.com/sebleier/554280>" which contains NLTK's all stop words.

Next we perform one hot encoding to convert our categorical values to numerical values. Now to represent the column 'item\_description' numerically we use the technique of "bag of words". In this model, a text is represented as the bag ( multiset ) of its words, disregarding grammar and even word order but keeping multiplicity. We ignore punctuation and frequently occurring meaningless words to deal with the large set of vocabulary.

## 2 Literature Survey

### 2.1 Existing Models

Since this is a Natural Language Processing problem, existing models have used Recurrent Neural Networks and Gated Recurrent Units, such models have provided RMSLE of 0.43. Some ensemble models which combined Ridge Regression and RNNs also exist, having RM-SLE of 0.44.

Some models have also used the Gradient Boost Algorithm (paper link:<https://arxiv.org/pdf/1603.02754>) which hasn't given good results because it tends to overfit on the training data when principal features aren't identified correctly.

### 2.2 Regression

Regression techniques like logistic regression, SGD have long been applied to NLP based classification problems. This article ("<https://medium.com/natural-language-processing-machine-learning/nlp-for-beginners-how-simple-machine-learning-model-compete-with-the-complex-neural-network-on-b9f7f93c79e6>") shows how a simple machine learning model like Naive-Bayes logistic regression can result in a performance as good as an RNN (LSTM) models. Regression models have been successfully applied to large-scale and sparse machine learning problems often encountered in text classification and natural language processing. Thus we experiment in here with various regression based models for this challenge.

### 2.3 Recurrent Neural Networks

Such kind of neural nets are used when there is a temporal component in the input data. For example speech recognition or language processing tasks like sentiment analysis. They were based on David Rumelhart's work in 1986. Hopfield networks - a special kind of RNN - were discovered by John Hopfield in 1982. There are various kinds of RNNs like LSTM or Long Short Term Memory networks, also improved large-vocabulary speech recognition and text-to-speech synthesis and was used in Google Android. In 2015, Google's speech recognition reportedly experienced a dramatic performance jump of 49% through the use of LSTM networks, which was used by Google voice search.

## 3 Performance Metric and Models Used

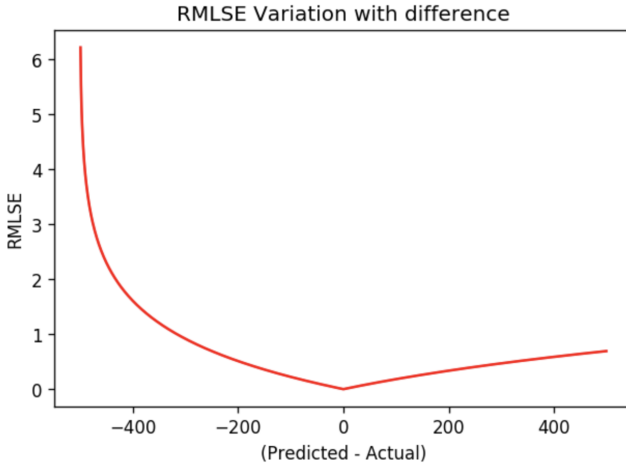
### 3.1 Performance Metric

Since we had to predict the price of a product and provide it to a seller, if the predicted price is higher than the actual price, it is of not much concern to the seller, but if the predicted price is lesser than the actual price, then the seller might suffer a loss. Hence, while training the model, we should incur a larger penalty if the predicted value is lesser than the actual value. Hence the performance metric that we chose is Root Mean Squared Log Error. We didn't choose Root Mean Squared Error because it incurs an equal penalty for over estimation and under estimation.

The formula for RMSLE is:

$$\text{RMSLE} = \sqrt{\frac{1}{N} \sum_{i=1}^N (\log(y_i + 1) - \log(\hat{y}_i + 1))^2}$$

This plot shows the penalty incurred by RMSLE vs the difference in actual and predicted values.



Hence we have applied the algorithms after changing the actual and predicted prices  $p$  to  $\log(p+1)$ .

### 3.2 Models

We have experimented with a lot of machine learning models of various types such as:

#### 1. Linear Models

- Linear Regression
- Stochastic Gradient Descent Regression

- Support Vector Regression
2. Non-Linear Models
    - Decision Tree Regression
    - Light Gradient Boosting Machine Regression
  3. Ensemble Models
    - XGBoost Regression

To choose the hyperparameters for these models, we assigned them random values and then used 5-fold cross validation to find out which value gives the best results.

## 4 Results

The table shows the value of RMSLE for each of the models used:

Model	RMSLE Value
Linear Regression	0.47
SGD Regression	0.48
SVM Regression	0.52
Tree Regression	0.5
LGBM Regression	0.48
XGB Regression	0.59

We can see that the least RMSLE value is for Linear Regression, hence indicating that the dataset is linearly separable to some extent.

Also the Gradient Boost regression method has performed poorly because it isn't good at identifying the principal features, hence it tends to overfit the data.

## 5 Tuning the hyperparameters

Hyperparameters are important because they directly control the behaviour of the training algorithm and have a significant impact on the performance of the model is being trained. Hyperparameter tuning plays an important role in predictions as we can then avoid under fitting and overfitting.

We use GridSearch with 3 fold cross validation for this purpose. Grid search trains the algorithm for all combinations by using the two set of hyperparameters (learning rate and number of layers) and measures the performance using "Cross Validation" technique. This validation technique gives assurance that our trained model got most of the patterns from the dataset. It suffers when the data is high dimensional which is not the case in our data and hence we can use it.

### 5.1 SGD Regressor

For SGD Regressor we found learning rate ="adaptive" and alpha= 1e-09 as the best values. Using these hyperparameters RMSLE reduced to **0.46** from .48.

### 5.2 LGBM Regressor

LGBM has many hyperparameters to tune in and the best set of these was found to be ‘boosting\_type’: ‘gbdt’, ‘n\_estimators: 200’, ‘learning\_rate’: 0.1, ‘max\_depth: 15, , ‘num\_leaves: 75’. Using these set of hyperparameters RMSLE reduced to **0.45** from 0.48.

## 6 Ensembling

Ensembling helps to reduce noise, bias and variance in the models. Earlier we used XG-Boost for ensemble but the results were not as expected and therefore now we used Adaptive boosting. After ensembling the top regression models with weights decided by their RMSLE score, the final RMSLE was found to be **.44**

## 7 Using CNNs for NLP problems

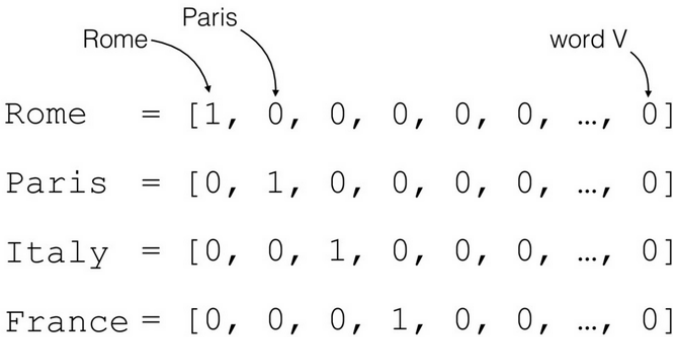
### 7.1 Introduction

Convolutional Network Models(CNNs) are generally used for images, since convolutions can be used to find patterns in images and dense layers can then be used to classify the input based on these patterns.

But CNNs can also be used for Natural Language Processing(NLP) problems, by converting written text like descriptions and names into images. The so called images have each row representing one word, words can be converted into vectors by several ways.

### 7.2 Ways of converting words to vectors

The simplest way to convert a particular word into a vector is by one hot encoding. We use a set which contains all the unique words that appear in the document and we assign each word a vector of the size of the set. For each word one bit of the vector is 1 and the rest are 0.





Another way is to use TF-IDF which stands for term frequency-inverse document frequency, in this case also each word is a vector of size equal to the number of distinct words in the document, but each word has an importance associated with it, the importance is calculated by multiplying two numbers, the first is the number of times the word occurs in this particular sentence and the second is the logarithm of the total number of documents divided by the number of sentences in which this word occurs. The process is shown here:

$$w_{i,j} = tf_{i,j} \times \log \left( \frac{N}{df_i} \right)$$

$tf_{i,j}$  = number of occurrences of  $i$  in  $j$

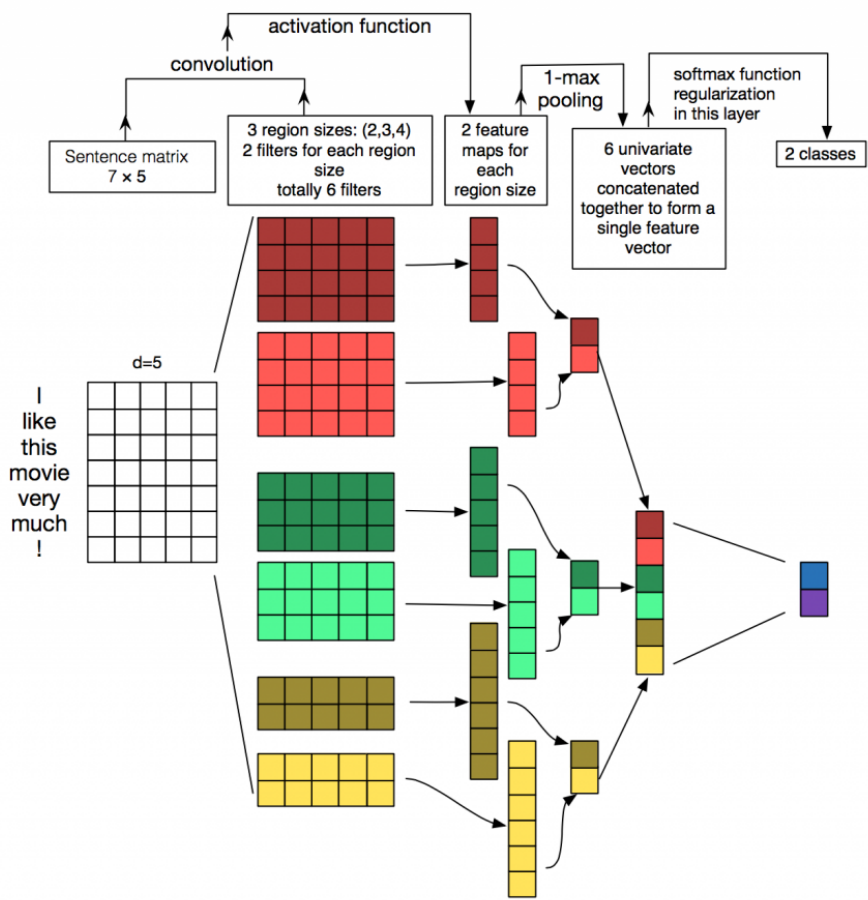
$df_i$  = number of documents containing  $i$

$N$  = total number of documents

The problem with one hot encoding is that it assigns each word a very different vector even if the words have very similar meanings, it basically treats every word as independent from other words. But if we want similar words to have similar vectors we need to take into account the dependence between them, this is called the **cosine dependence** or **cosine similarity**. If we have a vector for every word, then similar words should have smaller angles between their vectors.

To facilitate the formation of such word embeddings, neural networks are used which assign each word a vector. Several pre-trained models exist for creating such word embeddings, like **Word2Vec** and **GloVE**.

### 7.3 Using CNNs for NLP problems:An Example



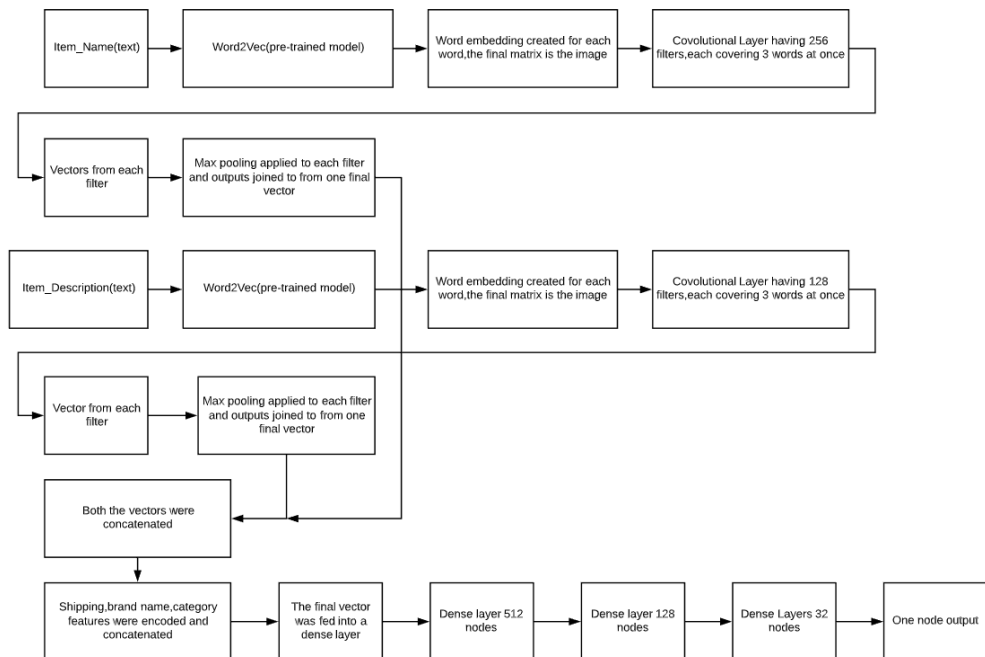
In the above network we have the input image which is created using word embedding models, for each sentence. This image is convolved with 6 filters, these filters find the pattern in the sentences' words and their values are trained. The first pair of filters considers 4 words at once, the second 3 words at once and the last pair 2 words at once. After convolution vectors are formed from each filter that was used. Pooling is applied on these vectors, pooling is of two types average pooling and maximum-pooling, in average pooling the output is the average of all the values of the vector and in maximum pooling the max is taken. After max-pooling has been applied, the output of each filter has been concatenated and fed into the softmax classification layer which outputs the probabilities of the input being in each of the classes. The class with the max-probability is taken and outputted.

## 8 Mercari Price Suggestion Model

In our case, the input had the following features, item name, item description, item category broken into three parts, item brand, item condition, shipping cost paid by seller or

buyer(bool). The first two features i.e. the item name and the item description are textual and hence were converted into vectors using the pre-trained **Word2Vec** model by using transfer learning. These vectors were then convolved using a similar way as shown in the above example and the final output of both were concatenated. The item condition, brand name, category and shipping were encoded using one-hot-encoding and concatenated to the initial vectors. The final vector was fed into a dense layer having 512 nodes, which was then fed into another layer having 128 nodes and another layer having 32 nodes, and finally the last one node layer which predicted the final price.

The final model architecture was:



## 9 Results

The CNN model gave an RMSLE of 0.4155 which is better than what the ML algorithms have been able to achieve.

We applied a lot of different algorithms on the dataset like regressors of all kinds such tree, boosted and Support Vector Machines. We tuned the hyperparameters of these algorithms to achieve their best performance and also used ensembling which lowered the RMSLE from 0.48 to 0.44.

The CNN model worked better than these algorithms because the filters were better able to identify the patterns between the descriptions and the prices, and because the pre-trained Word2Vec model helped create very strong word-embeddings, which were then fed into three dense layers which identified the pattern between the word vectors and the prices.

---

## 10 References

For understanding the workings of CNN and neural networks:

DeepLearning.ai Andrew Ng's course

For understanding the application of CNN on NLP problems:

<http://www.wildml.com/2015/11/understanding-convolutional-neural-networks-for-nlp/>