

# Predicting the “Usefulness” of Customer Reviews

**Anonymous**

Affiliation

Address

*email*

## Abstract

This paper compares the performance of various machine learning regression techniques while predicting the usefulness of customer reviews. The dataset is obtained from Yelp, an Internet based business review service, as part of its recently announced dataset challenge on Kaggle.com. The paper suggests that the usefulness of a review depends on the contents of the review and the language it uses. During the experiment, it is discovered that feature selection and preparing and massaging the data for this task are very tedious processes. For this reason, the techniques used for feature selection and data massaging have also been described in detail in this paper. It is found that for the given dataset, a non-linear model best describes the model, and support vector machines using radial basis functions provide a good way to regress the number of useful votes for a review. It is also recommended to experiment with other non-linear regression techniques such as regression with Gaussian processes and neural networks.

## 1 Introduction

Online communities rely heavily on trust within the community. Building trust has always been a difficult task, and doing this in an online community is especially difficult. A commonly used technique that is often used to support such a community is to enable community powered ratings and reviews. This is commonly seen in Internet based crowd-sourced reviewing services such as Yelp and collaborative consumption based services such as Airbnb and GetAround.

However, the growing popularity of a community often results in a large amount of variance in the quality of customer reviews. This can have an adverse effect on the experience of the user, and it might impact the trust within the community. Therefore it becomes essential to control which reviews are displayed to users, and what order they are ranked in. Traditionally, community managers have done this manually. However, as the service scales, it becomes essential to find an automated way to rank these reviews. For this reason, Yelp gathers feedback on each user review from other reviewers. Yelp tries to measure each review with three community-powered metrics: useful, funny and cool. These metrics are then used to decide how the reviews should be ranked.

However, the problem is that gathering votes for how useful, funny or cool a review can be a time consuming process. It can often take weeks if not months. And the usefulness of a review also deteriorates with age, thus rendering reviews that would otherwise have been valuable useless. For this reason, it is very important to have an automated way to predict

42 how useful a review will be before the community finishes voting on the review. This will  
43 ensure that high quality reviews are ranked high regardless of the number of votes awarded  
44 by the community.

45 To be able to rank these reviews according to usefulness, it is important to understand what  
46 makes a review useful. This is where the concept of feature selection comes into play.  
47 Feature selection is a machine learning technique that helps create a set of relevant features  
48 for use in model construction. Since the reviews on Yelp are text-based, intuition would  
49 suggest that most of the relevant features would lie in the text of the reviews themselves.  
50 While one can also think of other non-textual features that would be relevant for such a task  
51 (such as number of past reviews by the same user), this paper exclusively focuses on the  
52 textual contents of the reviews.

53 Yelp has recently released a dataset of reviews from businesses in Phoenix, Arizona, and has  
54 hosted a competition on the data science competition website Kaggle.com. The goal of the  
55 competition is to predict the number of “useful” votes a review will get over its lifetime. The  
56 dataset consists of 220,000 training reviews. This paper aims to find the best technique that  
57 should be used to predict the number of “useful” votes for a review.

58

## 59 **2 Related Work**

60 As mentioned above, in this paper we are trying to predict the number of votes a review will  
61 get. This means that we are trying to predict how important the review is, and how high up it  
62 should be ranked when retrieved amongst other reviews for a specific business. Since the  
63 task at hand is to train a model for a ranking task, the problem can be classified as a learning  
64 to rank problem.

65 Since review ranking is a problem faced by several online services, it is very likely that  
66 many such services use their own machine learning techniques to rank their reviews.

67 Several researchers have previously explored the problem of learning to rank. A Microsoft  
68 paper from 2005 introduced RankNet- an implementation of gradient descent methods  
69 applied to learning ranking functions using a neural network to model the underlying  
70 function [1].

71 Learning to rank can be employed in a number of applications such as document search,  
72 definition search, information retrieval, key phrase extraction, collaborative filtering,  
73 document summarization and machine translation [2].

74 While we have found some prior work related to ranking methods using neural networks, we  
75 have not found any prior work where an item is ranked based on the regression of its textual  
76 contents using natural language processing. We believe that this might be an effective  
77 method of predicting the rank of a review, given labeled training data where reviews have  
78 already a qualitative indicator.

79

## 80 **3 Background**

81 The Yelp dataset consists of four different kinds of data- data about businesses such as their  
82 names, neighborhoods, addresses and categories, data about users such as their name and the  
83 number of times they have reviewed or rated businesses, review data such as the text of the  
84 review along with the useful, cool and funny votes assigned, and check-in data about how  
85 many users checked into a business in a specific period of time. The data that is most  
86 interesting for our task is the review data itself. The review data is what will define the  
87 feature set for our task. The figure below describes the contents of the review data.

88 The review data provided by Yelp is all text based. To use this data effectively, we must use  
89 some natural language processing techniques. Natural language processing is a field of  
90 artificial intelligence concerned with the interactions between computers and natural human  
91 languages. Several open source tools have been built to help process natural language text.  
92 One such tool called nltk (Natural Language Toolkit) was used for processing and massaging  
93 the data.

94

## 95 4 Feature Selection Method

96 While creating a feature set, we decided to use only the words within the textual contents of  
 97 the review as features. We did not choose contextual information or meta-data such as the  
 98 type of business, geographic location of the business, and reviewer information. This is  
 99 because we do not think that the contents of reviews are geographically dependent, and we  
 100 believe that the quality of the review does not change with type of business. While it can be  
 101 argued that restaurant reviews would be viewed by more users than say jewelry store  
 102 reviews, the usefulness of reviews does not vary with the type of business. Since each  
 103 review is to be ranked before the number of votes for the review starts to plateau, the test  
 104 reviews are assumed to be fresh. For this reason, review freshness is not used as a feature.

105 Yelp’s reviews are user generated. Each user can write reviews from Yelp’s mobile and web  
 106 apps after visiting a business. Yelp employees do not moderate reviews before they are  
 107 posted online. For this reason, there is lots of noise in the data. Words are often misspelt,  
 108 and sentences are often badly constructed. Several reviews don’t make grammatical sense.  
 109 Variance in these reviews is also high- some reviews comprise of multiple paragraphs,  
 110 whereas some are only a few words long. Some reviews describe the quality of services or  
 111 products offered, whereas some only talk about prices or wait times.

112 Intuition suggests that reviews with certain descriptive words or phrases will be more useful  
 113 than others. Words such as “awesome”, “clean” and “fresh” are clear indicators that the  
 114 review is positive in nature, and phrases such as “try the Calamari” are high in information  
 115 gain about what is good (or bad) at a business.

116 We started by trying to create a feature set of each unique word in all the reviews. This  
 117 resulted in a feature set of over 120,000 words. We discovered that This feature set, while  
 118 exhaustive was too large. Fitting the training data and predicting results with such a large  
 119 feature set would take hours, if not days. For this reason, it became essential that we filtered  
 120 according to parts of speech. To do this, we used nltk.

121 Adjectives, adverbs and nouns are the most descriptive and provide the highest amount of  
 122 information gain for our task. The table below lists the parts of speech that were used for  
 123 this project. Other parts of speech such as conjunctions, prepositions, verbs and modals were  
 124 not included in this list and were filtered out so that they would not be included in the  
 125 feature set.

126

127 Table 1: Parts of speech used to create feature set

128

FW	Foreign Word	Tandoori
JJ	Adjective	Big
JJR	Adjective, comparative	Bigger
JJS	Adjective, superlative	Biggest
NN	Noun, singular or mass	Door
NNP	Noun, plural	Doors
NNS	Proper noun, singular	John
NNPS	Proper noun, plural	Vikings
RB	Adverb	Good
RBR	Adverb, comparative	Better
RBS	Adverb, superlative	Best

129

130 Another important part of feature selection was the use of punctuation and emoticons in user  
131 generated content. Since reviews use informal English, it is not uncommon to see users' use  
132 of punctuation and emoticons to express tone in a review. For example, a smiley face would  
133 indicate that the user was pleased with the service/product, where as an exclamation would  
134 indicate an extremely positive or negative experience. For this reason, we included certain  
135 emoticons and punctuations were part of the feature list. A table in the appendix details the  
136 list of emoticons and punctuations that were included.

137 We encountered another problem while preparing a feature set that was filtered by parts of  
138 speech. The process of creating the feature set involved iterating through each of the  
139 220,000 reviews, tokenizing them and separating them according to the various parts of  
140 speech, filtering by part of speech, and checking whether the words of the review are already  
141 in the feature list. This was a tedious process, and when repeated over 220,000 times, it  
142 became difficult even for powerful machines. For this reason, we decided to cut down the  
143 number of reviews that we would create the feature set out of. We cut them down from  
144 220,000 reviews to just 5,000 reviews, to create a feature set of 11,000 words. Figure 1  
145 below shows a snippet of code that uses parts of speech tagging and emoticon/punctuation  
146 based filtering to create a feature set.

```
for review in reviews:
    review_index += 1
    if (review_index >= feature_set_size):
        break
    current_review = review.replace(".", " ")
    for each_pair in nltk.pos_tag(nltk.word_tokenize(current_review)):
        if each_pair[1] in ['FW', 'JJ', 'JJR', 'JJS', 'NN', 'NNS', 'NNP', 'NNPS', 'POS', 'RB', 'RBR', 'RBS',] \
            or each_pair[0] in [':', ':D', ':P', '=)', 'XD', '=D', ':o', '=]', 'D:', ';D', '!']:
            if not apriori.has_key(each_pair[0]):
                apriori[each_pair[0]] = index
                index += 1
```

147

148 Figure 1: Code snippet to demonstrate feature selection using parts of speech tagging

149

## 150 5 Experiments

151 We conducted several experiments to determine the best model for this dataset. We  
152 experimented with linear as well as non-linear regression.

153 We used the scikit-learn module to experiment with each method. This greatly reduced  
154 implementation and prototyping time, as the algorithms we used were already implemented.  
155 Most of the work now lay in massaging the data and cross-validating the parameters.

156 We used the metric of mean squared errors to measure the success of the model. Below is a  
157 summary of the results of each technique each we experimented with.

158

### 159 5.1 Ordinary Least Squares

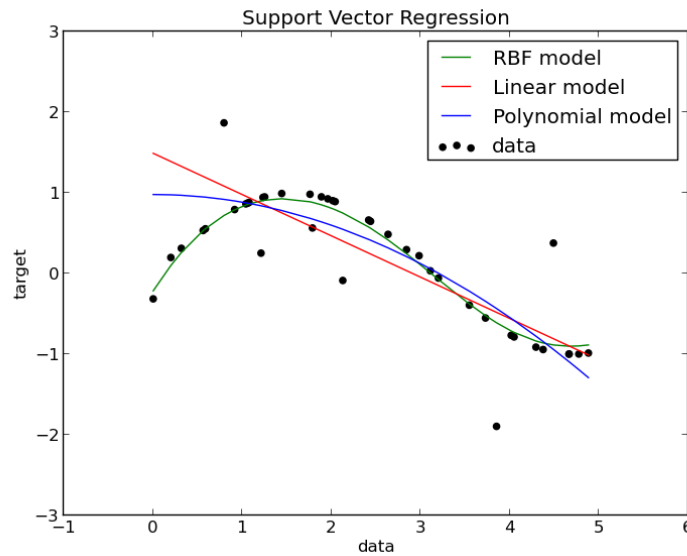
160 We started with a simple linear regression model using ordinary least squares. While this  
161 method was easy to implement, it also resulted in the highest mean squared error. This  
162 technique was discarded early on due to the poor initial results, and we moved from linear to  
163 non-linear regression after this.

164

### 165 5.2 Support Vector Machines

166 Next we decided to experiment with non-linear regression using Support Vector Machines,  
167 as they are known to be able to efficiently perform non-linear regression tasks.

168 While it is not possible to plot the result of support vector regression with a feature set of  
169 over 11,000 features, Figure 2 below shows an example of what the results of using support  
170 vector regression can look like with data from a cosine function.



171

172

173

174

Figure 2: Example of using regression with Support Vector Machines

### 5.2.1 SVR with polynomial kernels

175

176

177

178

179

We started experimenting with Support Vector Regression by using polynomial kernels of degree two and higher. Regressing with these models gave much better results than regressing with a linear model. This gave us an early indicator that linear models in fact are unsuitable for this data, and that further experimentation with non-linear models is a better approach.

180

181

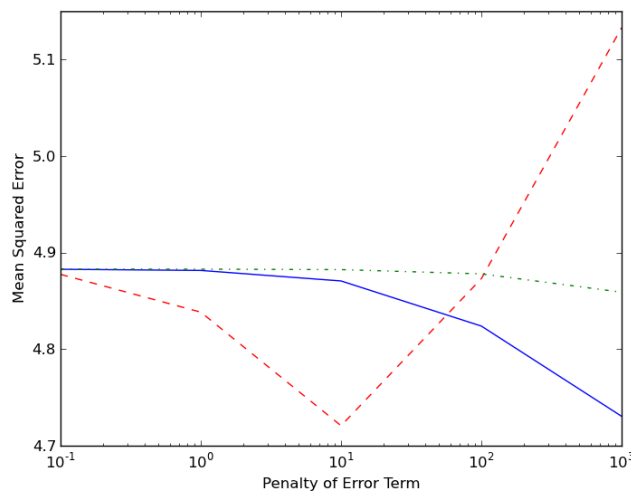
182

183

184

185

We cross-validated with different polynomials as well as other parameters such as the penalty of the error term. Figure 3 below shows the results of this experimentation. It can be seen that while the results do not vary much, on average, polynomials with degree three (blue) and two (red) performed slightly better than polynomials with higher degrees (green). As mentioned previously, mean squared error has been chosen as the metric for measuring the success of the regression.



186

187

188

189

Figure 3: Variation in mean square error with penalty of error term for polynomials with different degrees. Here, the red line is for second-degree polynomials, the blue line is for third degree polynomials, and the green line is for fourth degree polynomials.

190

### 191 5.2.2 SVR with radial basis function kernels

192 After experimenting with polynomial kernels, we decided to try using support vector  
193 machines with radial basis function kernels next.

194 We cross-validated with different values for penalty of error terms and kernel coefficients.  
195 As it can be seen in Table 2 below, the results show that using radial basis functions as  
196 kernels for support vector regression provides us with much better results than polynomial  
197 kernels. The lowest mean squared error we found through cross validation was 4.197 for  
198 penalty of error term 10, and kernel coefficient 0.01.

199

200 Table 2: Results of Support Vector Regression using radial basis functions

201

Penalty of Error term	Kernel coefficient	Mean squared error
1000	0.1	4.573830217
100	0.1	4.495839661
10	0.1	4.47455495
1	0.1	4.750245158
0.1	0.1	4.869590974
1000	0.01	4.536311987
100	0.01	4.364271694
<b>10</b>	<b>0.01</b>	<b>4.196941828</b>
1	0.01	4.431219925
0.1	0.01	4.574464125
1000	0	4.998818106
100	0	4.418512673
10	0	4.423488496
1	0	4.507002748
0.1	0	4.68913493

202

## 203 6 Conclusions and Recommendations

204 In this paper, we explored linear and non-linear regression techniques to predict the  
205 usefulness of user generated reviews. While we have determined that for this dataset non-  
206 linear regression with Support Vector Machines (SVMs) with radial basis functions as  
207 kernels are better than linear regression with original least squares, and SVMs with  
208 polynomial kernels, it is essential that other techniques are also investigated. Given the short  
209 time frame of the project, we were unable to experiment with other promising techniques  
210 such as regression with Gaussian processes and neural nets. As suggested by Burges [1],  
211 RankNet could also be used. In fact, deep nets could be used to learn a ranking function  
212 without the need for having user feedback for training the model.

213 Through cross-validation we found the best hyper parameters to tune the model. While this  
214 is a widely used and reliable technique, a better way to do tune the model would have been  
215 to use Bayesian optimization techniques.

216 On submitting our predictions of the test set to Kaggle.com, we found that our regression  
217 technique with support vector machines pushed us into the top 50 ranks on the leaderboard.

218 With further cross validation and experimentation, we are confident that this technique will  
219 push us higher up onto the leaderboard.

220 While we have only used a set of words as features, a more exhaustive feature set might  
221 provide better results. Instead of individual words, bigrams and collocations would be much  
222 more powerful as indicators of useful reviews. Examples include phrases such as “seasoned  
223 perfectly”, “money’s worth” and “no complaints”, which would be much more powerful  
224 indicators than the individual words that they consist of. Other meta-data and contextual  
225 information such as the reviewing experience of the reviewer are also likely to serve as  
226 useful features.

## 227 **References**

228 [1] Burges, C., Shaked, T., Renshaw, E., Lazier, A., Deeds, M., Hamilton, N., et al. (2005) Learning to Rank  
229 using Gradient Descent. In *Synthesis Lectures on Human Language Technologies*.

230 [2] Li, H. (2011) A Short Introduction to Learning to Rank. In *IEICE TRANS INF. & SYST., E94-*  
231 *D(10)*,1854-2862.