

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

Evaluate Helpfulness in Amazon Reviews Using Deep Learning

Bobby Nguy
Stanford University
cnguy002@stanford.edu

Abstract

Understanding how users perceive reviews' quality is crucial since it can reveal insights on what drives a user's purchasing decisions. We examine this problem by looking at Amazon review data and its helpfulness rating. We proceed by casting it as a binary classification problem using deep learning models such as RNN and LSTM. Despite presence of preferential attachment and other social biases, our results suggest that users value keywords and relevancy over sentence structure and readability. In addition, it also seem to suggest that users may have a limited "attention span" for reading each review.

1 Introduction

Consumer reviews are abundant in today's e-commerce setting. Yet not all reviews are perceived equally. Some seem to be more influential than the others, and are more relevant to a user's purchasing decision. It will be helpful to understand what constitutes a helpful review. Different from sentiment analysis, this task can be difficult even for a human agent. Two reviews may both look informative, but often one was perceived as helpful and the other was not. We would like to evaluate the performance of NLP and deep learning on such task.

In our problem setting, we examine Amazon reviews data, which provides an option for helpfulness votes. If choose to, a user can upvote or downvote a review on its helpfulness. For example 3/5 would mean that three people thought the review was helpful and two thought otherwise. There might be many more who saw the review but did not express their views.

We approach the problem by casting it as a classification problem. Each review is labeled as "helpful" or "unhelpful"; we then predict the labels by various algorithms. Here we will use Bag-of-Words with Random Forest as a baseline. In addition, we ran simple-RNN and LSTM on the data and compare the outcome.

2 Related Work

Studies from Yun Wan et al [1] indicate that Amazon helpfulness ratings may suffer from sequential and self-selection bias, and might not be reflective of the "true quality of the review".

Danescu-Niculescu-Mizil et al [2] examined the phenomena from a social psychological perspective and considered various opinion distributions in review rating. They also confirmed that text quality is not the only explanatory factor in helpfulness rating.

Given these circumstances, compared to well-researched problems such as sentiment analysis, previous attempts to predict helpfulness ratings were only able to achieve moderate accuracy rate (approx. 70-80%). These attempts mostly use handcrafted features and apply them with Support Vector Machines. Many of these attempts however tend to discard the zero-vote population and hence cast

054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107

doubt on the generality of the results. We will further discuss this issue and our approach in the later section. Nonetheless, these attempts agree that review lengths and review texts are correlated with helpfulness perception. Little is revealed, however, as to in what ways review texts affect its perceived quality.

3 Approach

We will apply bag-of-words(BoW) with random forest, simple Recurrent Neural Network(RNN) and Long Short Term Memory (LSTM) with various hyperparameters to determine the best performing model.

The bag-of-words model is a simple model that keeps the frequency of words occurrence. In our experiment, we restrict our feature space to 5000 words that occur the most frequently. Random forest, a ensemble tree classifier, utilizes these BoW features and make subsequent predictions. On the other hand, RNN and LSTM are neural network models that are able to capture information from sequences of words.

In a RNN, each input will be fed forward sequentially, through a number of hidden states, and compute a probability output vector.

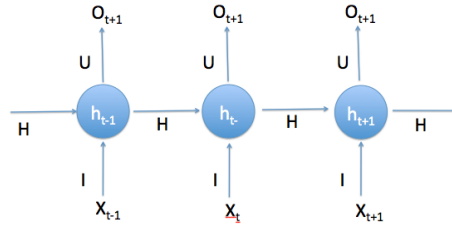


Figure 1: Flow diagram of RNN

The mathematical details are as follows [9]:

$$\begin{aligned}
 \mathbf{e}^{(t)} &= \mathbf{x}^{(t)} \mathbf{L} \\
 \mathbf{h}^{(t)} &= \text{sigmoid} \left(\mathbf{h}^{(t-1)} \mathbf{H} + \mathbf{e}^{(t)} \mathbf{I} + \mathbf{b}_1 \right) \\
 \hat{\mathbf{y}}^{(t)} &= \text{softmax} \left(\mathbf{h}^{(t)} \mathbf{U} + \mathbf{b}_2 \right) \\
 \hat{P}(\mathbf{x}_{t+1} = \mathbf{v}_j \mid \mathbf{x}_t, \dots, \mathbf{x}_1) &= \hat{y}_j^{(t)}
 \end{aligned}$$

LSTM is a variant of recurrent neural network with more complex units that include a memory cell [9]:

$$\begin{aligned}
 i_t &= \sigma \left(W^{(i)} x_t + U^{(i)} h_{t-1} \right) \\
 f_t &= \sigma \left(W^{(f)} x_t + U^{(f)} h_{t-1} \right) \\
 o_t &= \sigma \left(W^{(o)} x_t + U^{(o)} h_{t-1} \right) \\
 \tilde{c}_t &= \tanh \left(W^{(c)} x_t + U^{(c)} h_{t-1} \right) \\
 c_t &= f_t \circ c_{t-1} + i_t \circ \tilde{c}_t \\
 h_t &= o_t \circ \tanh(c_t)
 \end{aligned}$$

where i is the input gate, f is the forget gate, o the output gate and c the new memory cell.

108 The rationale behind these choices of models is that we can evaluate different aspects of review
109 text in our problem. Bag-of-words is a good baseline for keyword-only evaluation, while RNN
110 and LSTM also captures sequence information and can disambiguate various meanings of the same
111 words. Implementation is performed using TFLearn, a deep learning library built on top of Tensor-
112 Flow.

114 4 Experiment

116 4.1 Data and setup

118 Our experiments use Amazon review data collected by Stanford SNAP, publicly available at:

119 `https://snap.stanford.edu/data/web-Amazon.html`
120

121 The data comes from reviews for 24 product categories, where users provide comments or feedback
122 on their product experiences. Specifically, we used the reviews data from the *Digital Music*, *Musical*
123 *Instruments*, and *Patio, Lawn and Garden* categories, totaling to 88239 reviews. The reviews range
124 from 1 to 2055 words, with an average of 147.

125 Each review comes with a helpful rating $[x, y]$, directly relatable to the annotation “ x out of y people
126 found this review helpful” on Amazon website. Note that there might be many more who read the
127 review but did not vote. In many cases, reviews have never been voted and have helpful rating $[0,0]$.

128 Previous literature has focused on using helpfulness ratio, i.e. x/y , in making predictions and anal-
129 yses. There are several problems in this approach:

- 131 1. Many reviews ($>70\%$) are never rated and cannot be accounted for without correction. A
132 lot work simply disregards this bulk of the population in their analyses.
- 133 2. Viewing each $[x, y]$ pair as binomial (y, p_i) , variance of the helpfulness ratio is relatively
134 high with small y .
- 135 3. Preferential attachment is likely to be at force here for the helpfulness rating. We plot
136 the proportion of reviews with x upvotes (Figure 2) and its log-log plot in (Figure 3).
137 As suggested from the relatively straight line (until the tail portion, where variance is not
138 stabilized), in the log-log plot, the distribution of x follows the power law, which can be
139 generated from a network with preferential attachment. Preferential attachment is a result
140 of the richer-gets-richer effects, where reviews that happen to have already possessed some
141 upvotes are much likelier to gain more upvotes faster than the otherwise.

142 Since our focus is to discover any hints in review text that relates to helpfulness, we want to restrict
143 measuring the effect of preferential attachment. As a result, we decide to cast the problem as a binary
144 classification problem rather than a regression problem. To set up binary labels for the dataset, we
145 first perform correction on the data by adding one upvote and one downvote to each observation (i.e.
146 1 to x and 2 to y). Hence, all cases of $[0,0]$ will now have a helpfulness ratio of 0.5. We then label
147 reviews with ratio >0.5 as “helpful” and “unhelpful” otherwise.

148 We then perform NLP and deep learning on the review text and evaluate the model using cross
149 entropy loss. The score metrics that we are using are F1 scores and Accuracy:

$$151 F1 = 2 \times precision \times recall / (precision + recall)$$

$$152 Accuracy = (truePositives + trueNegatives) / totalEntries$$

154 We expect the model to explain a proportion of variances in the “helpful” and “unhelpful” text.
155 The neural network models are expected to perform better than the bag-of-words baseline since
156 LSTM captures both keywords and sentence structure, where readability is assumed to play a role in
157 indicating helpfulness. But because of various biases, including the effect of preferential attachment,
158 we expect none of the models to have a very high accuracy ($>90\%$).

160 4.2 Models

161 We performed classification using review text as features on the following models:

162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215

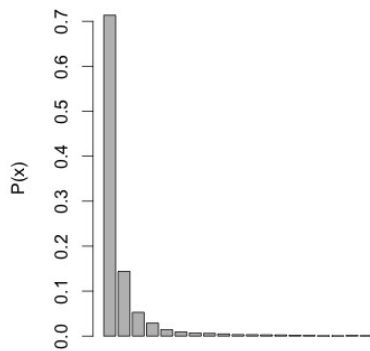


Figure 2: Proportion of reviews with upvotes = x.

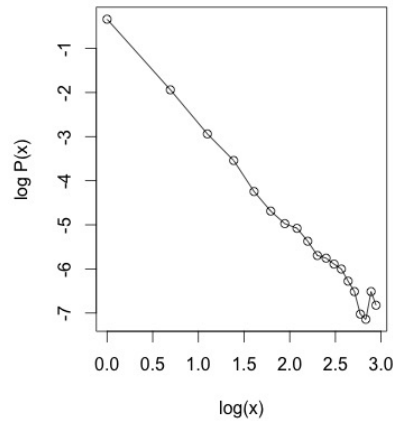


Figure 3: log-log plot of cumulative probability of reviews with upvotes = x vs x.

- Bag-of-Words (BoW) of unigrams, with max. feature dimension = 5000
- Bag-of-Words of bigrams, with max. feature dimension = 5000
- Simple RNN
- 1-Layer LSTM
- 2-Layers LSTM
- 3-Layers LSTM
- 2-Layers LSTM, with review length and rating as extra features

The neural network models are all set up with the following hyperparameters:

- Max. number words to read in review = 100
- Dropout = 0.5
- NumNeurons = 50
- Dimension of Word Embedding = 50

4.3 Glove Vectors

We compare the model performance from training our own word vectors and using pre-trained GloVe vectors. Due to our limited corpus size, our trained vectors are expected to be poor representation of the word dimensions. In fact, using GloVe vectors boosts our neural networks model performance by more than 7 percent points in accuracy. All of the following neural network models use GloVe vectors as its word vector representation.

4.4 Results and Findings

The baseline method, unigram bag-of-words with random forest, achieves a F1 score of 0.565 and an accuracy of 0.619. Bigram bag-of-words trails behind with F1 score of 0.547 and accuracy of 0.582. Both models' performances are comparable to two-layers LSTM, which achieves a F1 score of 0.549 and an accuracy of 0.65.

One-layer RNN has the worst performance: with a F1 score of 0.156 and accuracy 0.555, it borders on random guessing. It is a rather surprising result, especially considering the simplicity of the bag-of-words algorithm. Since bag-of-words capture the differences in keywords in the reviews, simple

216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269

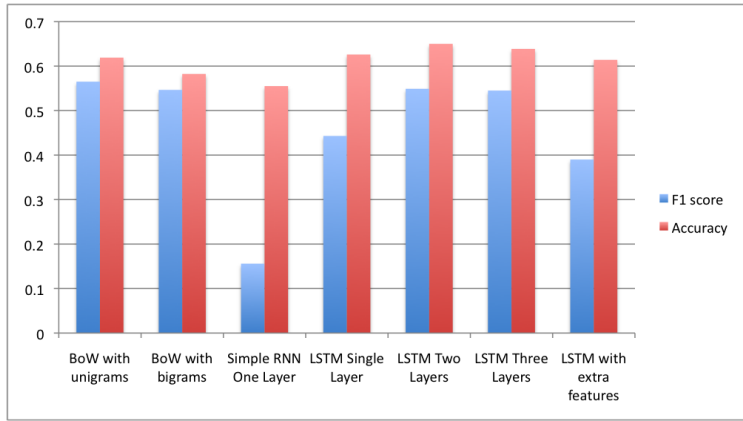


Figure 4: Experiment Results with various models

RNN the sentence structure and LSTM both of the features, it would suggest that the sentence structure and readability is of little importance when deciding whether a review is helpful or not.

Diagnosing the fitted random forest for variable importance, we get the following top predictor unigrams:

1. feature 99 : album (0.032975)
2. feature 4102 : songs (0.018699)
3. feature 2907 : music (0.015340)
4. feature 661 : cd (0.014586)
5. feature 4101 : song (0.011020)
6. feature 100 : albums (0.009514)
7. feature 3059 : one (0.008051)
8. feature 390 : best (0.007177)
9. feature 3492 : quot (0.007118)
10. feature 2533 : like (0.007077)

We repeat the same procedure and obtain the top predicting bigrams:

1. feature 1879 : hip hop (0.003301)
2. feature 2881 : one best (0.002388)
3. feature 3399 : quot quot (0.002237)
4. feature 1028 : easy use (0.001999)
5. feature 4372 : title track (0.001929)
6. feature 4854 : works well (0.001901)
7. feature 3940 : songs like (0.001882)
8. feature 1138 : every song (0.001683)
9. feature 1877 : highly recommend (0.001642)
10. feature 4847 : works great (0.001634)

The top unigram features seem to contain major keywords related to the product categories, while the top bigrams appear to convey positive sentiments. We may suggest that relevancy and sentiment reinforcement are one of the elements consumers value in reviews. Interestingly, extra features of review length and star ratings do not improve the model. This might be because the information in these two variables is already implied in the review text input.

270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323

Next, we further break down the results from the LSTM with extra features into subcategories and compare their metric scores (Figure 5). The scores for each category is similar; with the model performing slightly better on Musical Instrument. It does not appear that helpfulness evaluation is easier to identify in one product against the others.

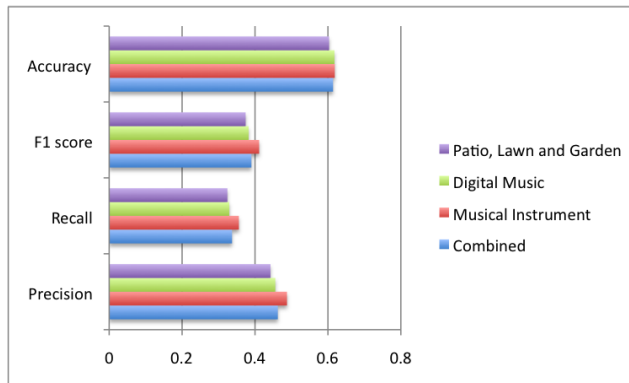


Figure 5: Performance metrics by categories

4.5 Hyperparameter Tuning

From Figure 4, while two-layers LSTM perform better than one-layer LSTM, three-layers LSTM shows no improvement over two-layers one. Very deep layer is not required in this context, as explainable variances are captured in relatively simple features such as keywords.

We also experimented limiting reviews at different lengths when being extracted from the data. Figure 6 shows the evaluation metric using maximum word lengths as 50, 100, 200 and 300. Overall it seems that various word lengths have no major impact on the model performance. As a result, word lengths of 100 are used in most of our evaluations.

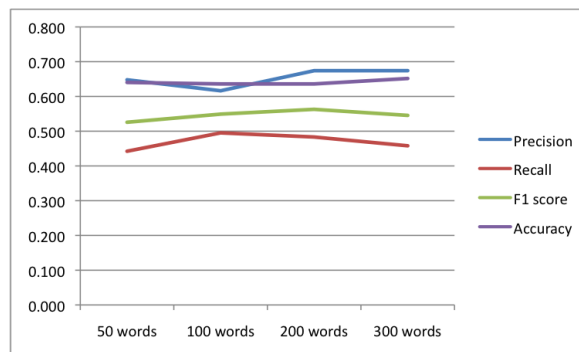


Figure 6: Performance metrics vs Review Input length

Since performance does not vary with different maximum lengths, it may suggest that users do not actually need to finish the entire review before deciding whether or not the review is helpful. We also note that in all of our models, we see significantly lower recall than precision, suggesting that many "helpful" reviews are difficult to identify from "unhelpful" ones.

5 Conclusion

From our experiment, it is suggested that helpfulness rating depends more on keywords and relevancy and surprisingly little on sentence structure and readability. In addition, it also suggests that many users may only glimpse through reviews and not reading its entirety, as various input lengths

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

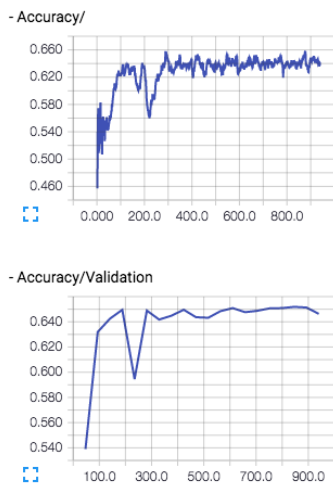


Figure 7: Accuracy vs Training Steps.



Figure 8: Loss vs Training Steps.

from the review have little effect on the performance. While preferential attachment and other biases have substantial influence, our experiment showed that review context still has impact on the perceived quality. For future works, it would be interesting to combine social network modeling technique and deep learning to understand the drivers for helpfulness ratings.

References

[1] Wan, Y., and Nakayama, M. (2012). Are Amazon.com Online Review Helpfulness Ratings Biased or Not? *Life: Web-Enabled Convergence of Commerce, Work, and Social Life* (M. J. Shaw, D. Zhang, and W. T. Yue, eds.), Springer Berlin Heidelberg, pp. 46-54.

[2] C. Danescu-Niculescu-Mizil & G. Kossinets. & J. Kleinberg & L. Lee (2009) How Opinions are Received by Online Communities: A Case Study on Amazon.com Helpfulness Votes, *Proceedings of the 18th international conference on World wide web*, pp. 141-150. New York, NY: Association for Computing Machinery.

[3] Hochreiter S. & Schmidhuber J. (1997) Long Short-Term Memory *Neural Computation Volume 9 Issue 8, November 15, 1997* pp. 1735-1780. Cambridge, MA: MIT Press

[4] Ha, Sangwook. (2015) Assessing Quality of Consumer Reviews in Mobile Application Markets: A Principal Component Analysis Approach. APA

[5] Kapoor, G. & Piramuthu, S. (2009) Sequential bias in online product reviews. *Journal of Organizational Computing and Electronic Commerce* 19 pp.85-95

[6] Wang, C. & Zhang X. & Hann I. (2010) Social bias in online product ratings: a quasi-experimental analysis *WISE 2010* St Louis, MO

[7] Rodal, J. & Xiao M. & Longoria, S. (2014) Predicting Helpfulness Rating of Amazon Reviews

[8] Bhargava S. (2015) Predicting Amazon review helpfulness

[9] Stanford University CS224D: Deep Learning for Natural Language Processing, Lecture Materials, Equations copied from: <http://cs224d.stanford.edu/lectures/CS224d-Lecture9.pdf>