

# Modeling Hotel Quality Belief in Natural Language Reviews

**Evan Shieh**  
Department of Computer Science  
*eshieh@stanford.edu*

**Alex Zamoshchin**  
Department of Computer Science  
*alexzam@stanford.edu*

## 1 Abstract

For the average traveler, finding trustworthy hotel reviews online is a tricky task - individual reviews inevitably suffer from bias, and the highest rated review isn't necessarily the most informative. Our research project focuses on modeling how exactly natural language reviews translate to beliefs about hotel quality as indicated on a 5-point ranking scale. In doing so, we draw inspiration from sentiment analysis, which tackles a tangential problem. Additionally, we leverage differences in the hotel review setting (such as skewed rankings, taglines and multiple rankings) to build more accurate classifiers. In all, we analyze the performance of several approaches (Recursive Neural Networks, Recurrent Neural Networks, Long Short-Term Memory, and Multi-task Deep Learning). Among these, our LSTM model performs the best, with a classification accuracy of 66.1% (across 5 possible outputs) and a RMSE of only 0.7313. We surmise that the LSTM achieves this performance above other models by leveraging mean pooling across all tokens in the sentence, as opposed to weighting importance of tokens by time (Recurrent) or structure (Recursive).

## 2 Introduction

We are interested in the problem of classifying hotel reviews based on reviews from online users. In this setting, ratings on a scale of 1-5 are assigned to each of five categories ("Overall", "Value", "Location", "Service" and "Cleanliness") and paired with a corresponding user review in natural language. At first blush, this problem domain appears to be quite similar to sentiment analysis, since verbal sentiment shares several parallels with expressing reviews (where low reviews are negative and high reviews are positive). Hence, our initial approach centers around using proven models from sentiment analysis (like Recursive Neural Networks). However, analyzing this approach reveals several key differences in our problem space - namely, that reviews are heavily skewed towards strongly positive, and reviews are often composed as unstructured paragraphs with grammatical abbreviations and vernacular English. Additionally, our problem space lends itself uniquely to multi-task learning, as the categories above are often interdependent and can be learned jointly. In all, we construct the data processing pipeline necessary for preprocessing reviews (such as data sanitization/filtering, and converting reviews to binary trees using the Stanford Parser) and then implement and analyze several models, including Recursive Neural Networks, Recurrent Neural Networks, Long Short-Term Memory, and multi-task learning models using both LSTMs and RNNs.

## 3 Dataset & Baseline

Our dataset is a *TripAdvisor* review dataset provided in JSON format by UIUC, containing plaintext reviews with corresponding 5-point labels for each of several categories (Location, Sleep Quality, Rooms, Service, Value, Cleanliness and Overall). Reviews are each usually several sentences long. We sanitize reviews by filtering out reviews that are missing categories, or ones that are over 800 characters long. In sum, there are 1,250,059 reviews for 12,000 hotels across the world, before 2011.

Ratings for all categories are heavily skewed towards larger ratings. An example of this for the "Overall" category is displayed in *Fig 1* below:

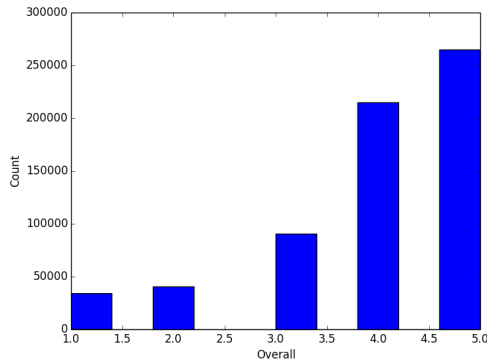


Fig 1. Histogram of Overall Ratings

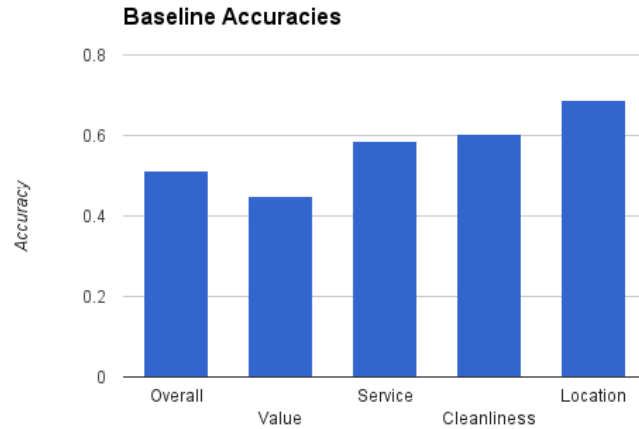


Fig 2. Baseline Accuracies for All Categories

40  
41

42 Given the large skew of classes, a strong baseline in this investigation is the *argmax* of the priors on classes.  
43 Meaning, predicting the most-common class for all sequences achieves the strong baselines for each category  
44 seen in Fig 2 above.

#### 45 **4 Related Work**

46 Our investigation builds on a number of ground-breaking models proposed in the past: Recursive Neural  
47 Networks share weights recursively over sequences modeled as binary trees, allowing them to produce  
48 predictions using information inherent in the structure of sequences. Recurrent Neural Networks form  
49 directed sequences and develop internal states which allow them to model temporal relationships, useful in  
50 language modeling and other tasks. Finally, Long Short Term Memory (LSTM) networks build on the success  
51 of Recurrent Neural Nets with the addition of memory cells, which will be discussed below. Already having  
52 been applied to model sentiment analysis, all three of these models are potentially well-motivated solutions to  
53 the task at hand - their advantages and disadvantages in practice, however, will be discussed below.

#### 54 **5 Models**

##### 55 A. Recursive Neural Network

56 Since the problem of predicting ratings shares many parallels with sentiment analysis, we first explore proven  
57 methods in the latter domain. Recursive Neural Networks, when applied to sentiment analysis, operate under  
58 the assumption that the sentiment of an overall sentence can be thought of as a recursive composition of the  
59 sentiments of its sub-phrases. A similar (and reasonable) argument can be applied to analyzing natural  
60 language reviews - logical conjunctions such as “but” will often change the entire meaning of a review as  
61 online travelers transition from making qualifying statements to expressing how they feel in reality.

62 Our model is an extension of the simple one-layer Recursive Neural Network explored in class, but with a  
63 key modification. The standard architecture relies on labeling both sentences and sub-phrases, which is not  
64 applicable to the review prediction scenario, since users express one ranking (per category) for each review.  
65 Hence, during back-propagation, error propagates from only one source (the root node). We confirm (as  
66 shown in the results section below) that taking this approach performs better than applying the same label to  
67 each sub-phrase.

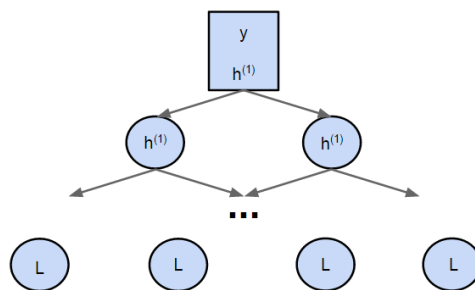


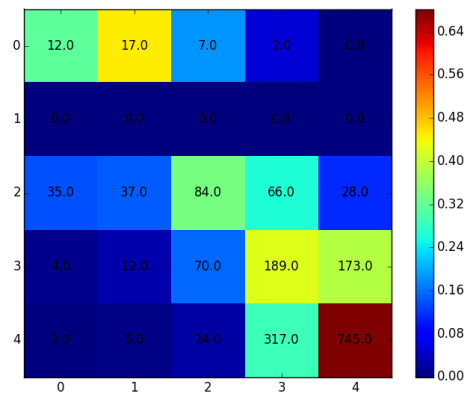
Fig 3. Depiction of Modified Recursive Neural Net

68  
69

70 Our general model-processing pipeline can be split into three stages:

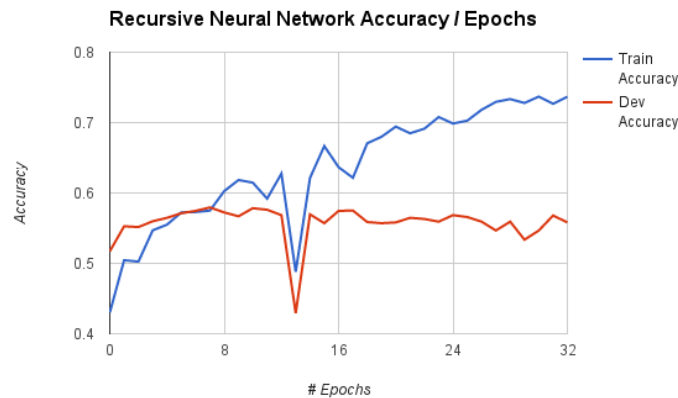
- 71 1. *Filter out invalid reviews.* Our dataset required extensive cleaning, as many data points were either  
72 invalid (i.e. containing blank reviews, missing labels, expressed in a different language) or too  
73 unwieldy to express as a single binary tree. We drop invalid reviews as well as reviews that are over  
74 800 characters long.
- 75 2. *Construct binary trees based on part-of-speech.* We leverage the Stanford Parser to construct  
76 standard parse trees based on the English PCFG, and then binarize them for purposes of our analysis.
- 77 3. *Perform classification using our Recursive Neural Network architecture.* The modified architecture  
78 above is built on top of our assignment in Python.

79 Directly applying the one-layer architecture as used for sentiment analysis yields a maximum accuracy of  
80 50.9% and a root mean squared error (RMSE) of 1.176 on our dev set, after parameter tuning. By contrast,  
81 our modified architecture achieves an improved accuracy of 56.3% with a RMSE of 0.821 on the same dev  
82 set.



83  
84 **Fig 4. Confusion Matrix for Recursive Neural Net Predictions**

85 Additionally, as the confusion matrix in *Fig 4* above shows, our modified architecture predicts a mixture of  
86 classifications from 1-star labels to 5-star labels, whereas the original Recursive Neural Network exclusively  
87 predicts labels in the 4-star to 5-star range. We surmise that this weakness of the original model is a product  
88 of over-amplifying the dataset skew - most labels are already 4 or 5-star reviews, and applying an identical  
89 label recursively forces most model parameters towards predicting the mode. Although our model suffers less  
90 from this phenomenon, to some extent artifacts exist in our model as well (as evidenced by the fact that our  
91 model never predicts 2-star ratings). Observing the training process confirms our hypothesis, since at fewer  
92 iterations both Recursive Neural Networks initially predict a mixture of all class labels.



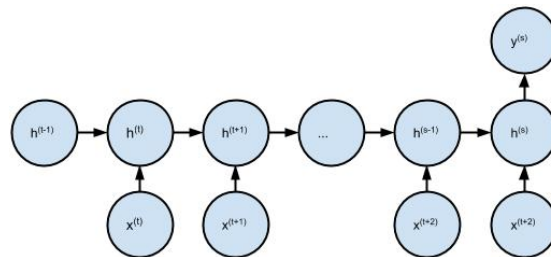
93  
94 **Fig 5. Recursive Neural Net Classification Accuracy over Time**

95 Furthermore, plotting classification accuracy for training and dev sets (*Fig 5* above) reveals that the  
96 Recursive Neural Network architecture, when applied to this scenario, does not appear to generalize well.  
97 Training accuracy increases steadily with more iterations, but dev accuracy tops out at around 56% after just  
98 a few training epochs. While this could be a result of the size of the dataset (each set contains roughly 10,000

99 examples), another plausible explanation is that reviews are not best modeled as trees. The setting of  
 100 sentiment analysis typically contains well-labeled, structured data in the form of single sentences. By  
 101 contrast, a qualitative examination of our dataset reviews reveals that reviews generally consist of multiple  
 102 sentences, often expressed in broken, vernacular English. As a result, the Recursive Neural Net architecture  
 103 expresses these reviews as long binary trees that might not be appropriate for training purposes, as child  
 104 nodes may not often be logically related to parent nodes in the casual review setting. Due to this problem  
 105 (and the amplification of the dataset skew mentioned above), we turn our attention to families of models that  
 106 make fewer assumptions about the structure of natural language reviews. One of these models is the  
 107 Recurrent Neural Network.

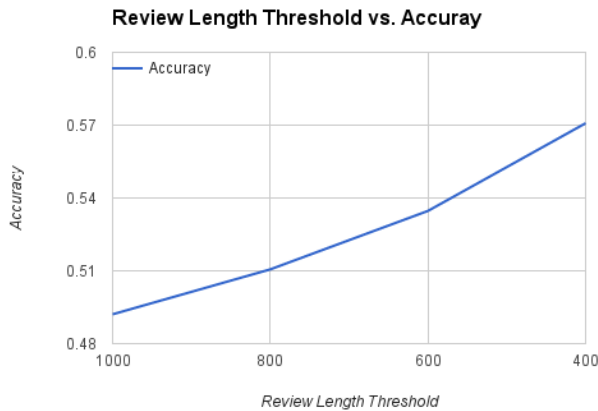
108 **B. Recurrent Neural Network**

109 Recurrent Neural Networks have grown in popularity for use in language modeling and other tasks. In this  
 110 paper we examine their use for hotel review modeling. We make a modification (similar to the modification  
 111 we made for Recursive Neural Nets) to the network structure to account for the lack of fine-grained labeling,  
 112 outputting only a single rating prediction for each sequence. Similarly, using only the rating of each review as  
 113 a whole, we back-propagate error across the entire sequence:



114  
 115 **Fig 6. Depiction of Modified Recurrent Neural Net**

116 Note that such a model is particularly prone to long-range interactions. This, combined with the fact that we  
 117 do not split sentences (in order to preserve the entirety of information contained in each review during  
 118 prediction), means we are left with extremely long sequences with only one label each. The model hence  
 119 proves to be extremely difficult to converge, achieving only a 49.20% training set accuracy and 49.74% dev  
 120 set accuracy (with 0.001 learning rate and no weight decay). Moreover, no amount of tuning is able to  
 121 achieve convergence; we attribute this problem to vanishing or exploding gradients and perform analysis to  
 122 prove this point:



123  
 124 **Fig 7. Results of Recurrent Net on Simplified Train Sets**

125 Performing thresholding on the lengths of reviews (in characters), we eliminate all reviews below a certain  
 126 threshold and train on the simplified datasets. We filter out reviews of lengths greater than 200, 400, 600, and  
 127 800 characters and analyze our hypothesis. In Fig 7 above we see the model is able to achieve improved  
 128 results of 57.08% accuracy on the training and 55.70% accuracy on the dev set, with reviews thresholded at  
 129 400 characters. This improvement is indicative of problems of vanishing or exploding gradients and long-  
 130 range interactions inherent in our model.

131 More concretely, during back-propagation the gradient is being multiplied a large number of times (for each  
 132 word in the sequence) by the associated weight matrix, causing the magnitude of the weights to have strong

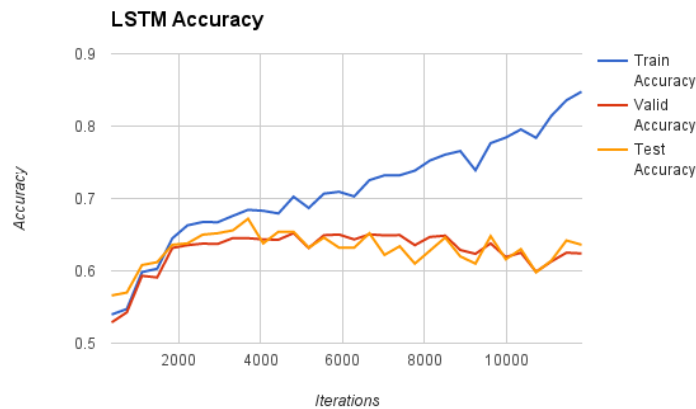
133 effects on learning. If the weights are small, this leads to vanishing gradients where the gradient signal  
 134 becomes so small such that learning stops altogether. If the weights are large, this leads to exploding  
 135 gradients where the gradient signal becomes so large such to cause learning to diverge.

136 As an initial solution to the problem of exploding gradients discussed above we examine clipping gradients at  
 137 a magnitude of 5. The hypothesis is that our huge sequences and large-range interactions will result in  
 138 gradients that are too-large, making fine-grained learning impossible. Using the same learning rate and  
 139 weight decay parameters as above, we achieve a training accuracy of 53.79% and a dev accuracy of 51.86%, a  
 140 distinct improvement over the original attempt on the unfiltered dataset.

141 C. Long Short Term Memory

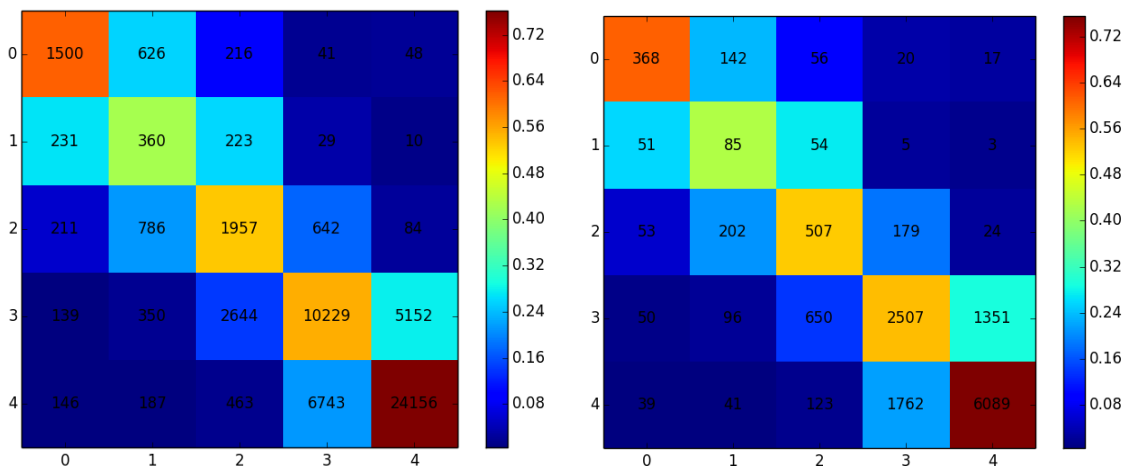
142 The analysis above and initial success of gradient clipping motivates the application of Long Short Term  
 143 Memory (LSTM) networks. LSTM networks are Recurrent Neural Networks with the addition of memory  
 144 cells - self-contained components that allow the network to learn when to remember and when to forget its  
 145 hidden states, amongst other possible outcomes. More concretely, a memory cell has four components: an  
 146 input gate, a self-recurrent connection, a forget gate, and an output gate. These four gates could potentially  
 147 allow the LSTM to better deal with the long-range interactions present in our problem. Finally, we use an  
 148 LSTM architecture with mean pooling, in which all hidden cells are directly aggregated into the Softmax  
 149 layer, allowing the network to again better model long-range interactions.

150 Training an LSTM network in Theano with 128 word embedding dimension, 0 weight decay, and adadelta  
 151 optimizer, we achieve the following results:



152  
 153 *Fig 8. LSTM Accuracy Results on Train, Dev, and Test Datasets*

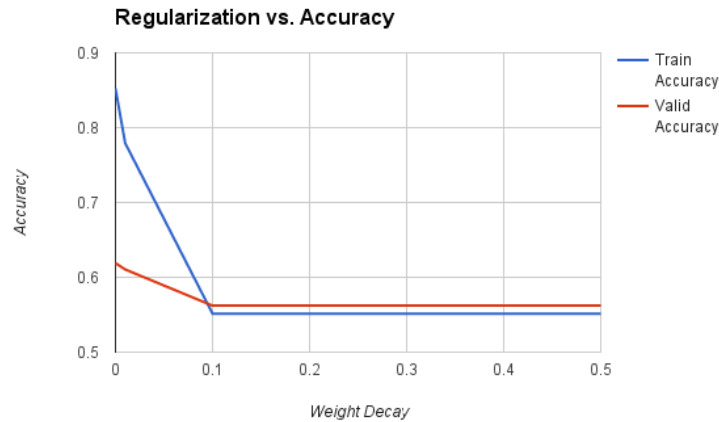
154 Hence, the greatest achieved dev accuracy was 65.21% with a corresponding test accuracy of 65.40% -  
 155 significantly improved results over those of the RNN implementation above. Furthermore, we achieve a  
 156 RMSE of 0.7313 on the dev set. The corresponding confusion matrices for these results are displayed below:



157  
 158 *Fig 9. Train Confusion Matrix (left) and Dev Confusion Matrix (right)*

159 The model performs relatively well on all labels, especially on extreme rating scores of 1 and 5. It faces some  
 160 difficulty for reviews falling in the 2-4 range (especially those scoring 2-stars) on both the training and dev  
 161 sets (perhaps indicating that our model reserves guessing 2-stars only for when it is extremely confident that  
 162 that is the case). This is similar to the weaknesses of the earlier models we examine, yet indicates that the  
 163 LSTM is still capable of learning and predicting accurately despite the dataset skew.

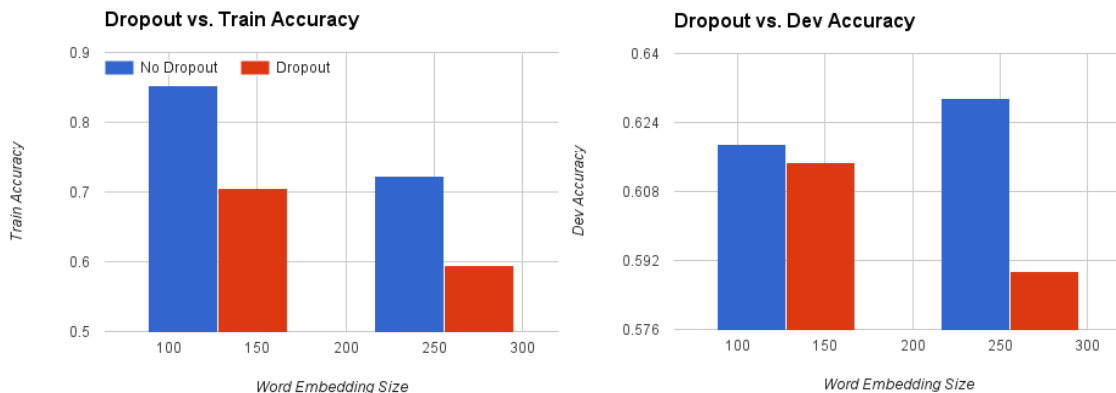
164 Clearly, the memory cells of the LSTM network are at least partially able to improve the modeling long-range  
 165 interactions inherent in the data. Yet with the improved predictive-power of the network comes a new set of  
 166 problems: the LSTM model substantially overfits the training set, as evidenced by the large disparity in train  
 167 and dev set accuracies. In an attempt to address this problem we add L2 regularization and tune the weight  
 168 decay parameter:



169  
 170

Fig 10. The Effect of Regularization on Train and Dev Accuracy

171 Surprisingly, introducing regularization does reign in the training accuracy but has no effect on improving the  
 172 dev accuracy. In fact, it seems that even a tiny amount of regularization prevents the model from establishing  
 173 enough confidence to guess any other label than 5-stars. This may be indicative of the difficulty of  
 174 classification in a skewed class scenario, and the challenge in improving the generalizability of the LSTM  
 175 model. We attempt to introduce dropout:



176  
 177

Fig 11. The Effect of Dropout on Train Accuracy (left) and Dev Accuracy (right)

178 We tune dropout at word embedding sizes of both 128 and 256, since the introduction of dropout often-times  
 179 requires larger models. Dropout proposes to address overfitting by randomly dropping units (with probability  
 180 0.5) during training. However, the theoretical improvements in performance are not realized: the model  
 181 performed better without dropout in all cases, including word embedding sizes of both 128 and 256, and on  
 182 both train and dev sets.

183 In a final attempt to address overfitting we attempt to introduce a significantly larger dataset size. Such a  
 184 dataset was not computationally feasible for most of the investigation but is a useful endeavor in addressing  
 185 overfitting:

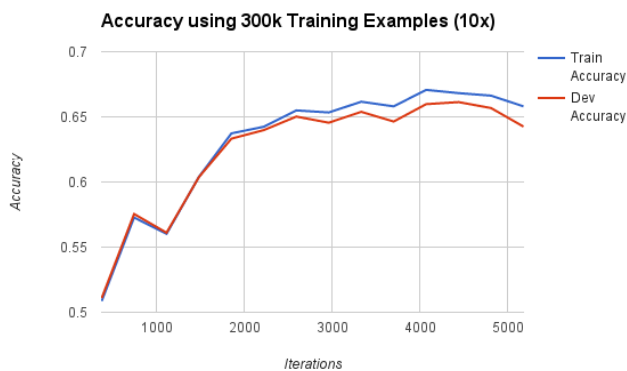


Fig 11. The Effect of Increased Dataset Size on Train and Dev Accuracies

186  
 187  
 188 The larger dataset size significantly reduces overfitting, almost completely diminishing the disparity between  
 189 train and dev accuracy. However only slight improvements in the dev and test accuracies are realized: at  
 190 66.13% dev accuracy the model achieves a corresponding test set accuracy of 61.80%. This represents the  
 191 best dev accuracy result achieved in this investigation.

192 As an additional exploration, we examine how well the model performs using review titles alone, as opposed  
 193 to both review content and title:

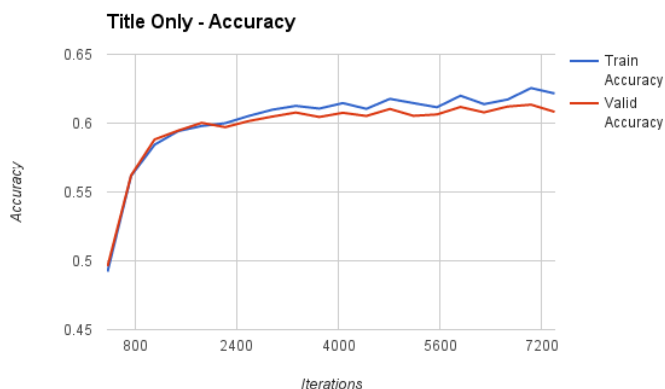


Fig 11. The Performance of the Model Using Only Title Content

194  
 195  
 196 In Fig 11 we see the optimal dev accuracy of 61.34% occurs after roughly 7000 iterations, with a  
 197 corresponding test set accuracy of 60.60%. Although the model was not able to perform as well using only  
 198 the title alone, it is evident that a substantial part of the improvement over the baseline can be attributed to  
 199 information gained from the title alone, an interesting result. Moreover, using the title alone also significantly  
 200 reduces overfitting, perhaps because the decreased quantity of words used results in a lower potential for  
 201 overfitting.

#### 202 D. Deep Multitask Learning

203 Another interesting question we explore is the hypothesis that correlated rating categories can be learned in  
 204 conjunction with one another. For instance, it is quite possible that reviews that give strongly-positive scores  
 205 for “Service” are more likely to also rate highly in terms of “Value”, and that words used to express these  
 206 sentiments are overlapping in nature. To approach this question, we turn to the theory of multi-task learning.  
 207 In broad strokes, multi-task learning acknowledges relatedness among labels by using a shared representation  
 208 of parameters to learn and predict multiple labels simultaneously (instead of training different classifiers  
 209 separately). This is especially the case in natural language processing, since we expect many of the earlier  
 210 layers to be recognizers for broad patterns present throughout the sequences. Hidden layers and word vectors  
 211 will be shared and updated by a sum of errors from multiple labels as opposed to just one.

212 More concretely, suppose there are  $C$  different categories of labels that we would like to learn simultaneously  
 213 ( $C=5$  in our example: Overall, Value, Location, Service, Cleanliness). Then, using the Recursive Neural  
 214 Network framework as an example, we implement multi-task deep learning as follows.

215 Forward propagation is extended as follows:

$$\hat{y} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_C \end{pmatrix}$$

216 where each of  $C$  probability vectors is appended together to create one output vector. Each probability vector  
 217 is calculated as follows:

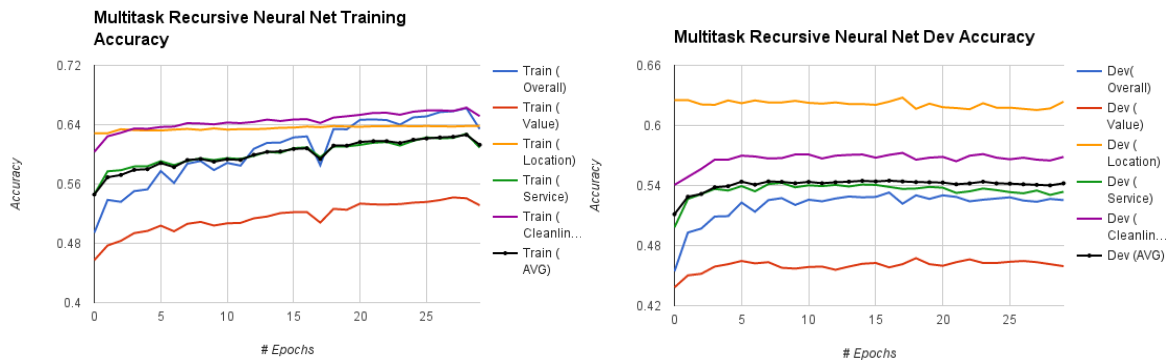
$$\hat{y}_i = \text{softmax}(W_i^{(s)}h + b_i^{(s)})$$

218 As seen from the equation above, the weight matrix  $W$  and bias vector  $b$  vary for each of  $C$  label categories,  
 219 but the hidden layer (and the parameters used to calculate the hidden layer) are shared as parameters across  
 220 all label categories. In our implementation, we append these separate per-category weight matrices and bias  
 221 vectors together in order to calculate the output probabilities (of dimension  $5 \cdot C$ ) directly.

222 Back-propagation is simply modified to include the sum of all errors across the new  $C$ -hot truth vector. Our  
 223 analysis treats errors from all categories with equal importance, but it would not be difficult to modify this  
 224 error sum to be a weighted sum of categories (treating “Overall” with greater importance, for instance).

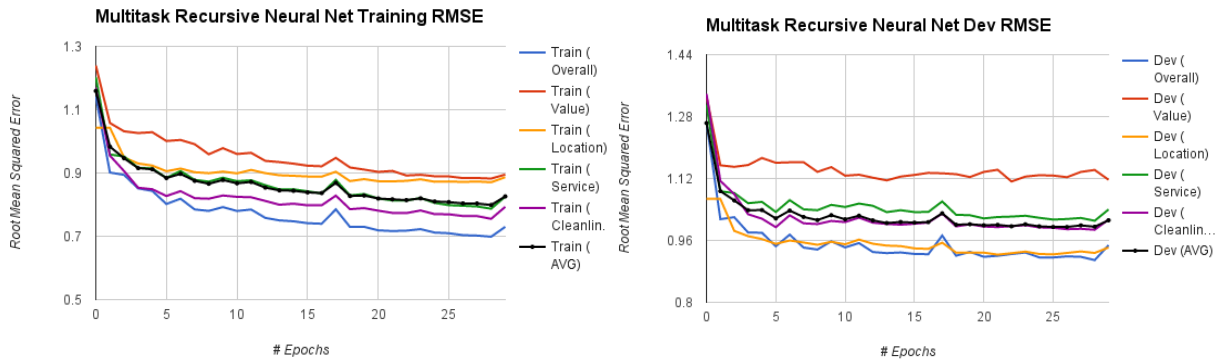
$$\delta = \sum_{i=1}^C (\hat{y}_i - y_i)$$

225 We apply multitask learning to extend our Recursive Neural Network model across the five categories  
 226 mentioned above, with varying accuracies as depicted below. Performance on the “Overall” category is  
 227 comparable to our single-task model, predicting a maximum accuracy of 53% on the dev set.



228  
 229 **Fig 12. Multi-task RNN Accuracy for Training (left) and Dev (right)**

230 As Fig 12 above indicates, the multi-task classification approach appears to benefit the model learning for the  
 231 “Overall” category (in blue) the most strongly, as classification accuracy increases by the most across epochs  
 232 for both the training and dev set (the relative differences between lines are due to varying dataset skew). This  
 233 would indicate that the “Overall” category is most highly interrelated to the other rankings categories.  
 234 Examining the RMSE (which is arguably more informative a metric and less biased by the skew of rating  
 235 labels) provides further support for this hypothesis:



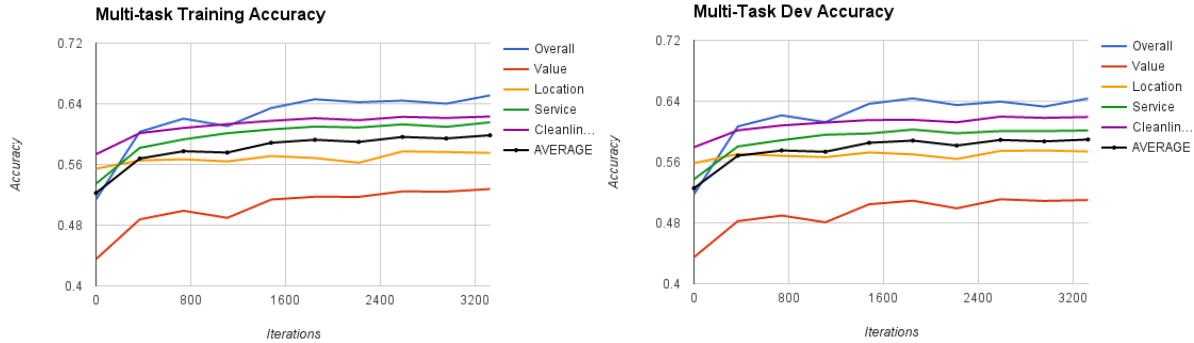
236  
 237 **Fig 13. Multi-task RNN RMSE for Training (left) and Dev (right)**

238 The above graphs show that the RMSE for the “Overall” category (in blue) experiences the greatest decrease  
 239 compared across all categories, and reliably remains the lowest as the number of training epochs increases.



240 Interestingly, the RMSE for the “Location” category (in yellow) also experiences a comparatively large  
 241 decrease on the dev set. This could indicate that “Location” is perhaps the ranking most highly correlated to  
 242 other rankings. If true, this would be an interesting trend to note, as other categories in theory should not  
 243 depend on the location of the hotel.

244 Similarly, we also hope to extend the results achieved by the LSTM network to the multi-task scenario. We  
 245 build on the LSTM network to address classification of all  $C$  tasks in the manner discussed above, the results  
 246 for which are seen below:



247  
 248 **Fig 14. Mutli-task LSTM Accuracy for Training (left) and Dev (right)**

249 Note that the best average validation accuracy is achieved after about 2600 iterations:

250

	Train	Dev	Test
Average Accuracy	59.63%	68.91%	57.20%

251 **Table 1. Mutli-task LSTM Average Accuracy Results**

252 On a more granular basis, we see the following test accuracies for each category:

253

	Overall	Value	Location	Service	Cleanliness
LSTM Accuracy	62.40%	47.40%	54.00%	60.00%	62.20%
Baseline Accuracy	51.40%	44.84%	68.62%	58.47%	60.33%

254 **Table 2. Mutli-task LSTM Test Accuracy Results for All Categories In Comparison to Baseline**

255 Comparison of the multi-task model’s test set accuracy on the “Overall” rating of 62.40% with LSTM’s prior  
 256 result of 65.40% indicates that no substantial improvement was achieved via modeling as a multi-task LSTM  
 257 task. Moreover, the model performed better on less-heavily skewed categories (Overall and Value) than on  
 258 more-heavily skewed categories (Location), indicating that perhaps the model experienced increased  
 259 performance on harder tasks at the expense of easier ones. Therefore, such a method may still be beneficial if  
 260 the objective function of the investigation is to optimize the average accuracy of *all* categories.

## 261 6 Conclusions and Future Work

262 Our work outlines a general framework for analyzing natural language, multi-label online review data. Of the  
 263 various models we examine - Recurrent Neural Nets, Recursive Neural Nets, and Long Short-Term Memory -  
 264 we find that LSTMs are best suited for the classification task, achieving a maximum classification accuracy  
 265 of 66.1% (across a label space of 5 discrete outputs) and an RMSE of 0.7313 (with the maximum error of 4).  
 266 Parts of the other two models still perform well, but are less suited for our classification task for a variety of  
 267 reasons. Recursive Neural Nets make the strong assumption that language semantics are structured  
 268 recursively, which we found not entirely true for casual, abbreviated online reviews that span multiple  
 269 sentences. Meanwhile, Recurrent Neural Nets suffer from long-range interactions and exploding gradients for  
 270 long sequences of reviews. We also analyzed both LSTM and the Recursive Neural Net models with multi-  
 271 task learning and find that strong correlations across labels do exist. With respect to multi-task learning, we  
 272 find that while overall accuracy is not improved, accuracy on harder tasks (i.e. less skewed categories) is

273 improved at the cost of accuracy on easier tasks. This actually posits multi-task learning as an effective  
274 strategy for learning against skewed datasets. Lastly, we learn on simply review titles alone using our best  
275 model (the LSTM) and find that these titles also hold surprising power in predicting rating.

276

277 Further work for our project will involve applying the (relatively reliable) models we have implemented to  
278 answer pertinent questions about online reviews. For instance, we would like to consider how reviews vary  
279 across geography, time and price range. Does a model trained on reviews from the United States have  
280 predictive power in classifying reviews from the UK? How does the price range of the hotel influence how  
281 predictable reviews for the “Value” category are? Given our work in choosing and tuning models to  
282 accurately predict hotel reviews, we are interested in turning our attention towards queries about natural  
283 language that can arguably only be tackled using statistical machine learning.

## 284 **7 References**

285

286 Bastien, Frédéric, Lamblin, Pascal, Pascanu, Razvan, Bergstra, James, Goodfellow, Ian, Bergeron, Arnaud,  
287 Bouchard, Nicolas, and Bengio, Yoshua. Theano: new features and speed improvements. NIPS  
288 Workshop on Deep Learning and Unsupervised Feature Learning, 2012.

289

290 Bergstra, James, Breuleux, Olivier, Bastien, Frédéric, Lamblin, Pascal, Pascanu, Razvan, Desjardins,  
291 Guillaume, Turian, Joseph, Warde-Farley, David, and Bengio, Yoshua. Theano: a CPU and GPU  
292 math expression compiler. In Proceedings of the Python for Scientific Computing Conference  
293 (SciPy), June 2010.

294

295 Graves, Alex. Supervised sequence labelling with recurrent neural networks. Vol. 385. Springer, 2012.

296

297 Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8), 1735-1780.

298 Gers, F. A., Schmidhuber, J., & Cummins, F. (2000). Learning to forget: Continual prediction with LSTM.  
299 *Neural computation*, 12(10), 2451-2471.

300

301 Socher, Richard; Lin, Cliff; Ng, Andrew Y.; Manning, Christopher D. "Parsing Natural Scenes and Natural  
302 Language with Recursive Neural Networks". The 28th International Conference on Machine  
303 Learning (ICML 2011).