

---

# Sentiment Classification of Food Reviews

---

**Hua Feng**

Department of Electrical Engineering  
Stanford University  
Stanford, CA 94305  
fengh15@stanford.edu

**Ruixi Lin**

Department of Electrical Engineering  
Stanford University  
rlin2@stanford.edu

## Abstract

Sentiment analysis of reviews is a popular task in natural language processing. In this work, the goal is to predict the score of food reviews on a scale of 1 to 5 with two recurrent neural networks that are carefully tuned. As for baseline, we train a simple RNN for classification. Then we extend the baseline to GRU. In addition, we present two different methods to deal with highly skewed data, which is a common problem for reviews. Models are evaluated using accuracies.

## 1 Introduction

Binary classification of sentiment on reviews are an increasingly popular task in NLP. Instead of classifying positive reviews and negative reviews, we classify reviews into extremely negative, negative, neutral, positive, and extremely positive classes directly from the reviewer's score on a topic. We train a simple RNN classifier and a GRU classifier. At test time, we input a user's review as a sequence of words, and output the category of the highest softmax score as the class label. Our analysis could be a useful tool to help restaurants better understand reviewers' sentiment about food, and can be used for other tasks such as recommender systems.

## 2 Problem Statement

In order to predict review level sentiments, we label each review with a reviewer's score indicating the sentiment of the reviewer. Our task is to predict a reviewer's score on a scale of 1 to 5, where 1 indicates the reviewer extremely dislikes the food he or she mentions in the review and 5 indicates the user likes the food a lot.

### 47 **3 Related Work**

48  
49 Traditional approaches on sentiment analysis user count or word frequencies in the text  
50 which are assigned sentiment value by expert[1]. These approaches disregard the order of  
51 words. A recurrent neural network (RNN)[2] can be used for sequence labeling on sequential  
52 data of variable length, which is natural for sentiment analysis tasks where the input sentence  
53 is viewed as a sequence of tokens. Recent works explore the Gated Recurrent Units neural  
54 network(GRU)[3] on the task of sentiment classification. GRUs are a special case of the  
55 Long Short-Term(LSTM) neural network architecture. GRUs are effective in this task  
56 because of their ability to remember long time dependencies. Furthermore, GRUs are faster  
57 to train and converge than LSTM networks. For our specific task, we have not found much  
58 work on the exact problem.

### 62 **4 Dataset**

63 We work on the Amazon Fine Food Reviews dataset[4] which contains 568,454 reviews. The  
64 dataset consists of a single CSV file, which includes the id of the product, id of the reviewer,  
65 the score(rating between 1 and 5) given by the reviewer, the timestamp for the review, a brief  
66 summary of the review, and the text of the review. We extract the columns of scores and  
67 review texts as our labels and raw inputs. A sample review with score is shown below:

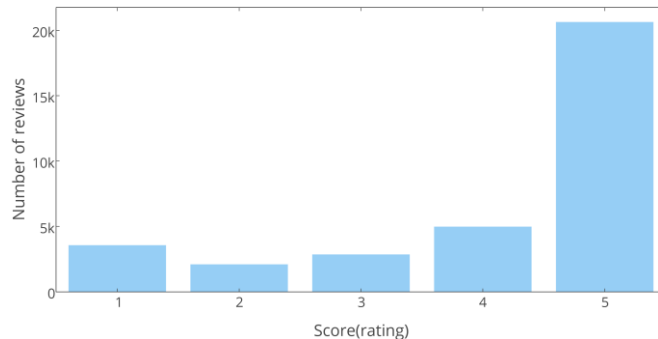
68	<b>Review</b>	<b>Score</b>
69	Product arrived labeled as Jumbo Salted Peanuts...	1
70	the peanuts were actually small sized unsalted.	
71	Not sure if this was an error or if the vendor intended	
72	to represent the product as "Jumbo".	
73		
74	I have bought several of the Vitality canned dog food	5
75	products and have found them all to be of good quality.	
76	The product looks more like a stew than a processed meat	
77	and it smells better. My Labrador is finicky and she	
78	appreciates this product better than most.	

79  
80 In order to perform mini-batch training for the neural network models, we want tokens within  
81 each slice of epoch to come from the same review. To make this happen, we need to  
82 compensate reviews with <unk>s to the maximum length of all reviews. To introduce as  
83 fewer <unk>s as possible, we do not want the reviews differ greatly in length. In this case,  
84 we would like to keep only reviews of similar lengths. We need to determine the range of  
85 lengths of reviews. In our analysis of the original dataset, we found that the average length of  
86 reviews is 80, so we choose reviews between 75 and 87 tokens and generate a dataset of  
87 34,091 reviews.

88  
89 Another problem of the dataset is that the reviews are skewed towards higher scores,  
90 especially towards the highest score, which is 5. In the 34,091 reviews, 3,550 reviews are  
91 labeled with 1, 2,085 reviews are labeled with 2, 2,844 are labeled with 3, while 4,971  
92 reviews are labeled with 4 and an even larger volume of 20,641 reviews are labeled with 5.  
93 As is shown in figure 1, score-2 class has the lowest number of reviews, which may lead to  
94 difficulty in predicting score-2. Score-5 class has the highest number of reviews as expected,

95 which is around ten times of that of score-2 class. To take care of the skewedness issue, we  
 96 introduce two resampling methods to produce a more balanced dataset. The methods will be  
 97 discussed in section 6.

98



99

100

Figure 1: Number of reviews of each score in the Amazon Food Reviews dataset.

101

102

## 103 5 Mathematical Formulations

104

### 105 5.1 Simple Recurrent Neural Network(RNN)

106

107 This baseline method is a slightly modified version of the standard RNN. Instead of  
 108 providing classification prediction at each word, we build the model to output prediction at  
 109 the end of each epoch slice. We make this modification in order to reduce the influence of  
 110 frequent words on the prediction and backpropagation.

111 Let T represents the number of steps, For each epoch slice  $x^{(t)}, \dots, x^{(t+T-1)}$ , the forward  
 112 propagation is defined as:

$$113 \quad h^{(t+k)} = \sigma(W^{(hh)}h^{(t+k-1)} + W^{(hx)}x^{(t+k)} + b_1) \quad (1)$$

$$114 \quad \hat{y}^{(t+T-1)/T} = \text{softmax}(W^{(s)}h^{(t+T-1)} + b_2) \quad (2)$$

115 Where  $k = 0, 1, \dots, T-1$ ,  $x^{(t+k)}$  is the word vector embedding for the (t+k) th word in the  
 116 review,  $h^{(t+k)}$  is the (t+k)th hidden layer and  $\hat{y}^{(t+T-1)/T}$  is the prediction output at the  
 117 (t+T-1)/T th epoch slice. Details of implementation can be seen in section 6.2.

118 Cross-entropy error is used as loss function, the expression for a corpus size of K is as  
 119 follow:

$$120 \quad J = -\frac{T}{K} \sum_{t=1}^{T/K} J^{(K_t)}(\theta) = -\frac{T}{K} \sum_{t=1}^{T/K} \sum_{c=1}^C y_{t,c} \log(y_{t,c}) \quad (3)$$

121 Where T is the number of steps, C is the total number of class and  $y_t$  is the one hot vector  
 122 representation of the label at t-th epoch slice.

123

124

$$125 \quad x^{(t)} = L_{x_t} h^{(t)} = \sigma(W^{(hh)}h^{(t-1)} + W^{(hx)}x^{(t)} + b_1^{(t)}) \hat{y} = \text{softmax}(W^{(s)}h^{t=T})$$

126

## 127 5.2 Gated Recurrent Units

128

129 The mathematical formulation of GRU at each time step is defined as follows<sup>[5]</sup>:

130

$$\begin{aligned}z^{(t)} &= \sigma(W^{(z)}x^{(t)} + U^{(z)}h^{(t-1)}) \\r^{(t)} &= \sigma(W^{(r)}x^{(t)} + U^{(r)}h^{(t-1)}) \\ \tilde{h}^{(t)} &= \tanh(r^{(t)} \circ U h^{(t-1)} + W x^{(t)}) \\ h^{(t)} &= (1 - z^{(t)}) \circ \tilde{h}^{(t)} + z^{(t)} \circ h^{(t-1)}\end{aligned}\tag{4}$$

131

132 Where  $x^{(t)}$  is the word vector embedding for input word at step t,  $z^{(t)}$  is the update gate which  
133 determines the combination of new memory and previous memory carries on to next layer,  $r^{(t)}$   
134 is the reset gate which determines the proportion of new word and previous contextual information  
135 in generating new memory,  $\tilde{h}^{(t)}$  is the new memory generated and  $h^{(t)}$  is the hidden layer at  
136 step t.

137 Since GRU has update gate to determine the importance of new memory for current state, its  
138 prediction result is less likely to be influenced by frequent word(ideally,  $z^{(t)}=1$  on frequent  
139 words without much sentiment information such as stop words ). So we output prediction at each  
140 step and use the summation of cross-entropy error at each step as loss function.

141

142

## 143 6 Experiments & Results

144 To address the skewedness problem, two different resampling methods are implemented to  
145 balance the dataset. We evaluate both resampling methods. We implement a simple RNN and  
146 a GRU with Python Tensorflow and measure the train, validation, and test accuracies of each  
147 classifier we build. We draw confusion matrices, visualize the hidden layer weights, analyze  
148 and tune hyper parameters to improve accuracies.

149

### 150 6.1 Data Pre-processing

151

#### 152 6.1.1 Sampling method 1: remove all data from the last class

153 Since the main source of data skewedness is the highest score class which has around ten  
154 times as many reviews as each of the rest of the classes, we employ a simple method to avoid  
155 the problem. We discard the data from the highest score class and redefine our task to predict  
156 the review score into one of the first 4 classes. The new dataset consisting of scores 1 to 4 is  
157 less biased towards higher scores.

158

#### 159 6.1.2 Sampling method 2: resample data from the 4- and 5-score class

160 A natural way to generate a balanced dataset is to randomly sample reviews from the skewed  
161 dataset, in which case we should sample data from the 4-score and 5-score classes.  
162 According to figure 1, we would like to obtain around 4,000 reviews for each class, so we  
163 generate 4,000 random samples from the two high score classes. Now we have a more  
164 balanced dataset.

165

### 166 6.2 Implementation of RNN

167 Word vectors are initialized as random values uniformly distributed between [-1, 1]. The  
168 number of steps is set as 8 as recommended in the course lecture. To distinguish between

169 different reviews, <EOS> is added at the end of each review. Then to ensure phrases of 8  
 170 words are from the same review within each epoch slice, we zero-pad the reviews to 88  
 171 words at the front of each review. Zero-padding is done at the beginning because if zero  
 172 padding at the end, backpropagation will come across several identical hidden layers before  
 173 propagating to an actual word, thus cause more severe vanishing gradient problem.

174  $L, W^{(hh)}, W^{(hx)}, b_1, W^{(s)}$  and  $b_2$  are updated through the training process and applied in  
 175 validation and testing. L is the embedding matrix for words.

176 The final predicted class for each review is the class with the max value in the elements of  
 177  $\hat{y}_c$ , where  $\hat{y}_c$  is the output prediction at the end of the corresponding review (identified by  
 178 EOS).

### 179 6.3 Implementation of GRU

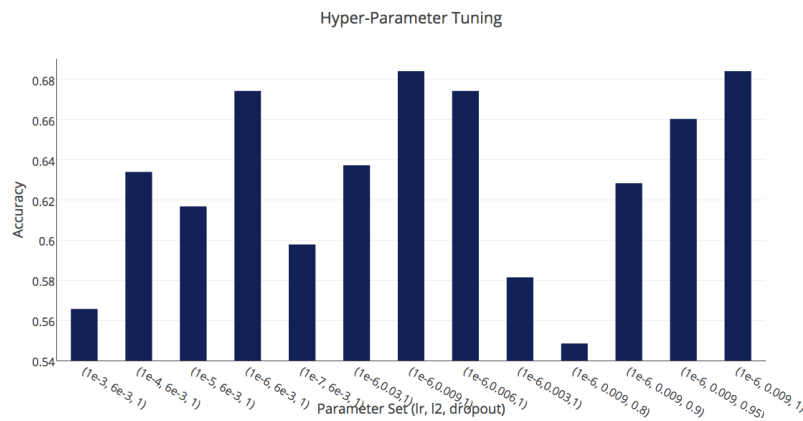
180 For GRU, we use the same dataset, number of steps and initialization strategy of word  
 181 vectors as RNN. The training is performed on dataset with/out zero-padding.

183  $L, W^{(z)}, W^{(r)}, U^{(z)}, U^{(r)}, U$  and  $W$  are updated through the training process and  
 184 applied in validation and testing. L is the embedding matrix for words.

185  
 186 The output prediction at the end of each review is used as final prediction of each class, just  
 187 like RNN, to provide a fair comparison of performance.

### 188 6.4 Hyper-Parameters Tuning

189  
 190 In order to tune and find the right hyper-parameters for our model, we divide our data into  
 191 three sets: a training set, a validation set for cross validation and a test set that will be used as  
 192 our final prediction scores. In this section, we describe how we performed our tuning and  
 193 record the accuracies depending on it. For each of the models, learning rate, L2  
 194 regularization weight and dropout value are to be tuned. Due to time and computation  
 195 resource constraints, we did not tune some parameters like hidden layer size and we were not  
 196 able to jointly optimize the parameters that would have resulted in the optimal setting.  
 197 Instead, we fix some parameters to reasonable values and tune the others. The following  
 198 figures show the tuning results.



199

200

Figure 2(a). RNN(4 classes) Hyper-Parameter Tuning

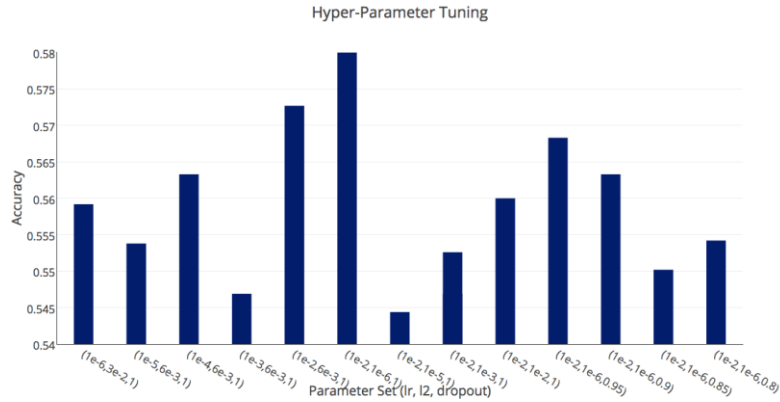


Figure 2(b). GRU(4 classes) Hyper-Parameter Tuning

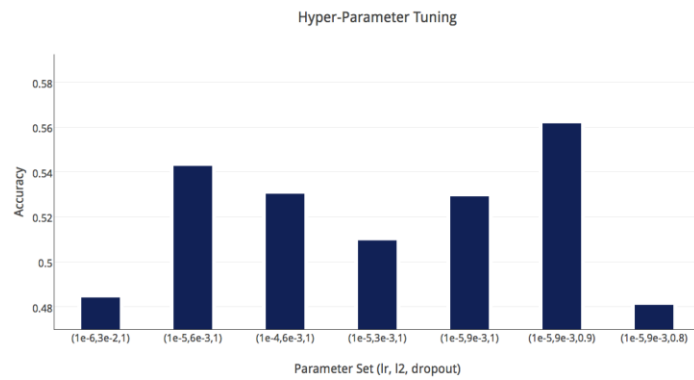


Figure 3(a). RNN(5 classes) Hyper-Parameter Tuning

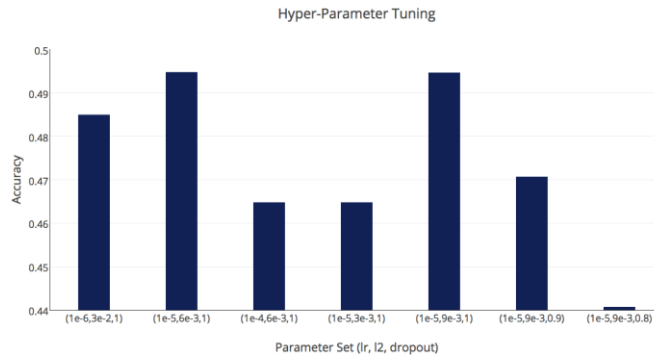


Figure 3(b). GRU(5 classes) Hyper-Parameter Tuning

The optimal set of parameters we have found for our models are as follows:  
 RNN,4 classes( $lr=10^{-6}$ ,  $l_2=0.009$ , dropout=1.0), RNN,5 classes( $lr=10^{-5}$ ,  
 $l_2=0.009$ , dropout=0.9), GRU,4 classes( $lr=0.02$ ,  $l_2=10^{-6}$ , dropout=1.0),  
 GRU,5 classes( $lr=10^{-5}$ ,  $l_2=0.006$ , dropout=1.0). With these parameters  
 obtained, we re-train our models and test the models. The test performances  
 are shown in next section.

## 6.5 Accuracies & Confusion Matrices

After tuning the hyper parameters, we use the optimal set of hyper parameters to train and test our model and evaluate the performance by accuracy. Accuracy is calculated by the

218 number of correctly labeled reviews over the total number of reviews, where the predicted  
 219 label at the end of a review is regarded as the final predicted label for that review. For  
 220 comparison purposes, we also train GRU models without zero padding.

221

Model+Resampling method	Training Accuracy	Test Accuracy
RNN(4 classes)	<b>93.35%</b>	<b>68.75%</b>
RNN(5 classes)	80.38%	51.74%
GRU(4 classes)	71.13%	55.03%
GRU(5 classes)	66.24%	44.44%
GRU(4 classes, w/o zero padding)	54.40%	42.70%
GRU(5 classes, w/o zero padding)	43.60%	35.7%

222

Table 1. Accuracies of different models

223

224 The best model is the simple RNN on the 4-class prediction task. Confusion  
 225 matrix of the train, validation and test results of this model are illustrated in  
 226 the figure below.

227

### 228 6.6 Visualization of Hidden Layer Weights

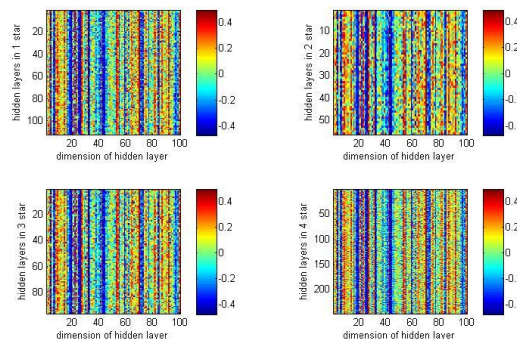
229

230 To demonstrate the effect of training under different strategies, we  
 231 present the visualization of a hidden layer at the first and last epoch in  
 232 this section.

233

234 For our modified RNN, the hidden layers for different classes looks quite  
 235 similar at epoch 0 (shown in figure 4(a)) since the word vectors are  
 236 randomly initialized. But by the last epoch of training, the hidden  
 237 layers under different labels are quite different. For instance, hidden  
 238 layers under 3 and 4 star reviews have higher values around 40th  
 239 dimension than hidden layers under 1 and 2 star.

240



241

242

Figure 4(a). Hidden Layer under RNN at Epoch 0

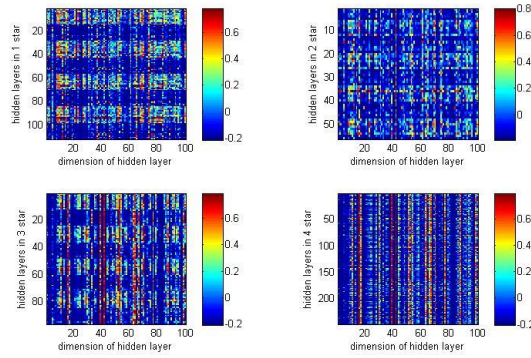


Figure 4(b). Hidden Layer under RNN at Epoch 6

For GRU, the hidden layer showed some change over the epochs, but the pattern is not as obvious as RNN.

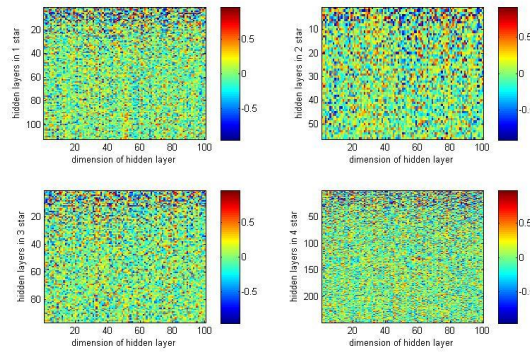


Figure 5(a). Hidden Layer under GRU at Epoch 0

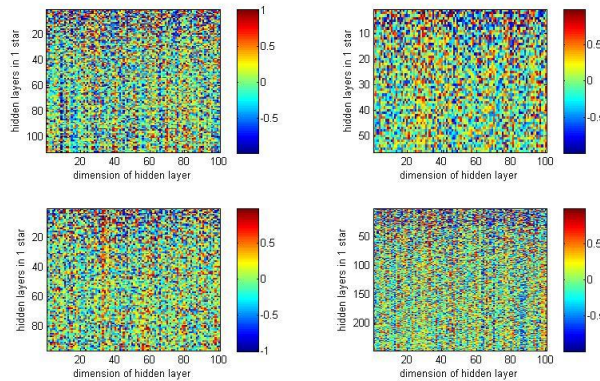


Figure 5(a). Hidden Layer under GRU at Epoch 6

## 7 Conclusion

In this paper, we present different neural network approaches including RNN and GRU for sentiment classification on Amazon Fine Food Reviews dataset and reach 68.75% test accuracy on the test set. In our experiment, we find that padding zeros to reviews proves to be useful and the zero-padded approaches outperform the approaches without zero-padding we implement. Future work might focus on trying out more RNN models, like the bidirectional RNN.



264 **References**

- 265 [1] Bo, P. (2008) Opinion Mining and Sentiment Analysis, *Foundations and trends in information*  
266 *retrieval*, 2(1–2): pp. 1–135. doi:10.1561/1500000011
- 267 [2] Mikolov, T., Karafiát, M., Burget, L., Cernocký, J., & Khudanpur, S. (2010, September). Recurrent  
268 neural network based language model. In *INTERSPEECH*, Vol. 2, pp. 3.
- 269 [3] Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2015). Gated feedback recurrent neural  
270 networks. arXiv preprint arXiv:1502.02367.
- 271 [4] McAuley, J. J., & Leskovec, J. (2013, May). From amateurs to connoisseurs: modeling the  
272 evolution of user expertise through online reviews. In *Proceedings of the 22nd international*  
273 *conference on World Wide Web* (pp. 897-908). International World Wide Web Conferences Steering  
274 Committee.
- 275 [5] Mohammadi, M., Mundra, R., Socher, R. (2015) Lecture Notes: Part IV. CS224D: Deep Learning  
276 for NLP

277  
278  
279  
280  
281  
282