

## Submission Assignment 4

*Instructor:* Jakub Tomczak*Name:* Richie Lee, *Netid:* [ KLE680 ]

## Introduction

One of the main factors that distinguishes humans from A.I. are our emotions. In this report we will challenge this consensus, by exploring Sentiment Classification. The approach for this emotion identification will be by applying various RNN's for movie-reviews. This task could be considered developing networks with EQ (Emotional Intelligence) and as a result potentially closing the gap between humans and A.I.

## 1 Problem statement

**Data** *Introducing the IMDb dataset*

IMDb (Internet Movie Database) is known as the world's most popular source for movie reviews. However, it is also a commonly-used sequence-learning database as it is a large, challenging classification dataset with plenty of examples. The utilised dataset consists of 50.000 preprocessed and tokenized reviews.

During model development and tuning we will work with a validation split of 20.000 training instances and 5.000 validation instances. For the final test, we will apply a canonical test/train split of 25.000 each.

**Objective** *Sentiment Classification*

The objective is to accurately identify the opinions in text and to assign corresponding labels. This binary sentiment classification will assign every review to either its positive or negative class.

## 2 Methodology

First we will apply *variable batch sizes*<sup>1</sup>. In other words, we will construct (approximately) equal-sized batches in terms of tokens, while allowing the amount of sequences per batch to differ. Implementing this improves our training time by significantly reducing the required padding and as a result the amount of computations. Furthermore we can now expect roughly the same amount of computations per batch.

This report will explore RNN's<sup>2</sup> (Recurrent Neural Networks). These networks are characterized by their ability to "remember" previous inputs, allowing temporal dynamic behaviour. This is beneficial in environments such as this classification where input dependency is expected. In addition, RNN's also provide better compatibility with arbitrary-length sequences. In our experimentation, we compare 2 RNN variations, alongside a MLP. Furthermore, we apply an Embedding layer and ADAM optimizer for all variations.

1. **MLP (Multilayer Perceptron)**, often referred to as "vanilla" neural networks, these basic feed-forward models can be considered the bread and butter of deep learning. Our model will contain a single hidden layer.
2. **Elman**, this model applies an Elman layer. In addition to regular hidden layer properties, it also contains context units which hold information of previous iterations. Moreover, we also implemented a residual connection<sup>3</sup> to counter exploding and vanishing gradients<sup>4</sup>.
3. **LSTM (long-short term memory)**, this model is a RNN variation. In short, LSTM's main benefit is the "conveyor belt" consisting of only linear operations, avoiding gradient problems in the process. In addition, we also feed selected memory of the past to this belt through various "gates".

<sup>1</sup>Code snippets: Section 6.1, figure 4

<sup>2</sup>Architecture details provided in Appendix

<sup>3</sup>Code snippets: Section 6.2, figure 6

<sup>4</sup>Gradient problem details provided in Appendix

### 3 Results

**Performance after first epoch** *compares the three models in terms of average accuracy and loss*

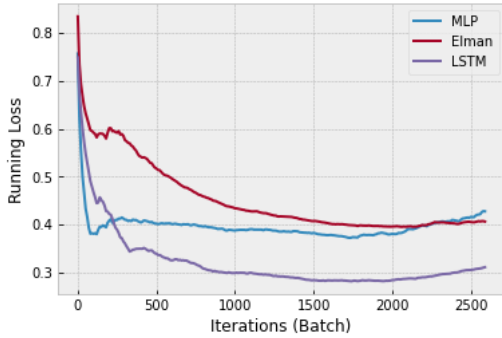
Table 1: Performance after first epoch

	Training		Test		Runtime
	Loss	Accuracy	Loss	Accuracy	
MLP ( $lr = 0.001$ )	0.4257	0.8250	0.4015	0.8455	0:34
Elman ( $lr = 0.001$ )	0.4050	0.7970	0.4056	0.8276	0:46
LSTM ( $lr = 0.001$ )	0.3102	0.8687	0.5182	0.7533	1:16

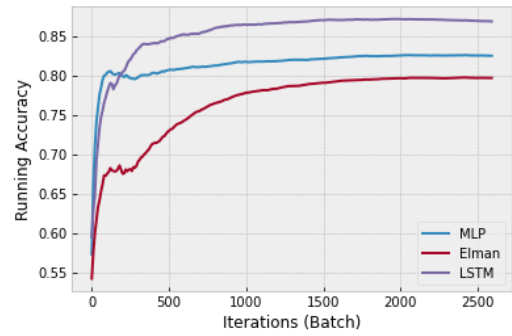
Generally we expect the LSTM to perform best, followed by the Elman and MLP respectively. However, from our models we observe the best performance using the MLP, followed by the Elman and LSTM. The reason for this unexpected result could be the difference in optimality in hyperparameters<sup>5</sup> (learning rate) among the models. Another potential driver could be insufficient training epochs, as we've decided to train single epochs (due to time constraints).

In addition we observe significant difference in training time for every model variation. These differences are not surprising as RNN's rely on higher amounts of computations. Finally, we observe significant performance difference between the LSTM's test and training performance. We believe this to be a result of overfitting<sup>6</sup>.

**Performance development during first epoch** *compares the three models in terms of average running accuracy and loss during first epoch*



(a) Loss Development



(b) Accuracy Development

Figure 1: Running performance during first (Training) epoch

Figure 1 shows how the loss and accuracy developed throughout the first epoch, by showing the average loss/accuracy up to its corresponding iteration on the x-axis. These plots give clear insight on the training speed, which in turn contributes to parameter tuning.

### 4 Conclusion

This report was intended to demonstrate the advantages and disadvantages of RNN's. However our findings contradicted the expected results to some degree. We believe this to be the result of sub-optimal hyperparameter tuning. Furthermore we also observed exploding and vanishing gradients for the Elman network, but we conclude this drawback to be manageable, as we have managed to resolve it using residual connections.

In terms of preferred models, we are unable to draw a final conclusion. This is because we are potentially overlooking RNN advantages due to the applied hyperparameters. From our findings, we can however conclude all three models to be reasonable selections with performance expectancy exceeding 75% accuracy.

<sup>5</sup>Hyperparameter tuning is discussed in Section 5, Discussion

<sup>6</sup>The occurrence of optimizing "too well", up to a point from which we start learning the training data, rather than improving the general task

## 5 Discussion

- **Hyperparameter Tuning (Learning Rates):** Due to long training times and time constraints, we were just able to test a small selection of learning rates for single epochs. Applying the models on the validation/training split yielded the following findings:

Table 2: Learning Rate Tuning (Accuracies)

Learning Rate	Training			Validation		
	0.0001	0.001	0.01	0.0001	0.001	0.01
MLP	0.7571	<b>0.8261</b>	0.7691	0.7728	<b>0.8080</b>	0.8068
Elman	0.5728	<b>0.7721</b>	0.4889	0.5948	<b>0.8113</b>	0.4952
LSTM	0.6980	<b>0.8113</b>	0.8072	0.7358	0.8028	<b>0.8238</b>

Table 2 displays significantly different performance among learning rates. Especially the RNN variations (the Elman in particular) showed sensitivity to hyperparameter tuning. For model optimization, more tuning will be required.

The MLP and Elman show clear preferences for Learning rates, however this does not hold for the LSTM. We decided to select 0.001, because of more stable training loss development (figure 2).

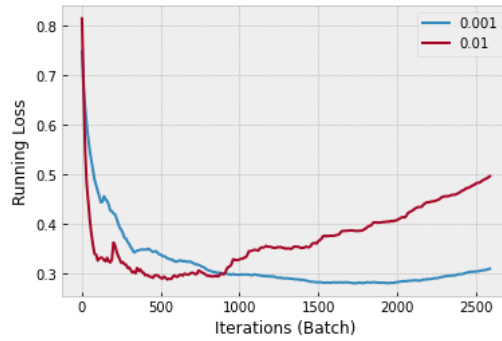
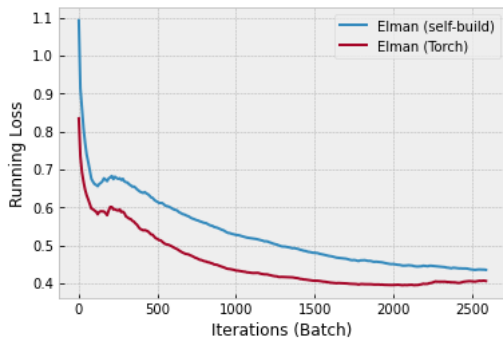
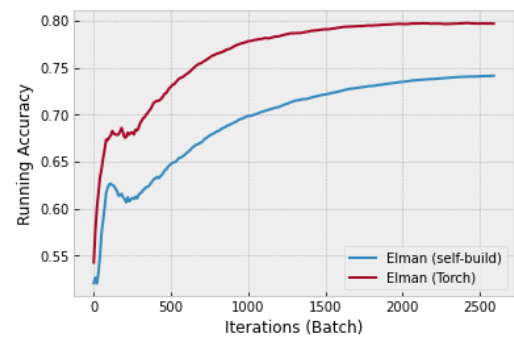


Figure 2: Loss Development LSTM

- **Comparing the self-build and Torch built-in Elman implementations:** Perhaps the most interesting finding, the performance difference among our Elman variations. (Though not closely related to the main objective of this report)



(a) Loss Development



(b) Accuracy Development

Figure 3: Comparing Training development of the "Self-build Elman" and the "Torch Built-in Elman"

Table 3: Performance after first epoch

	Training		Test		Runtime
	Loss	Accuracy	Loss	Accuracy	
Elman					
Torch ( $lr = 0.001$ )	0.4050	<b>0.7970</b>	0.4056	0.8276	0:46
Self-Build ( $lr = 0.001$ )	0.4343	0.7414	0.3323	<b>0.8562</b>	4:54

Interestingly enough, we find seemingly contradicting results in figure 3 and table 3. When evaluating the training performance, we observe better results for the Torch implementation. However, the Test-run tells us the opposite, with significantly better results for the Self-Build version.

The only difference between the models is the difference in non-linearity for the Elman layer. The self-build utilizes a *Sigmoid* activation in its Elman layer and the Torch a *Tanh*. This difference could affect performance through data compatibility.

Another aspect to consider are the different initialisations. These differences would fade away on the long term, but our testing range of 1 epoch would not catch it. Therefore we expect some degree of volatility as well.

## 6 Code Snippets

### 6.1 Variable Batch Size Implementation

```

91 # Variable Batch Size
92 """
93 Slices data into batches with approximately equal amounts of tokens
94 Pads accordingly
95 """
96
97 # Padding function
98 def pad(l, size, padding):
99     return l + [padding] * abs((len(l)-size))
100
101 # Variable batch slicing
102 def batch_slicing(x, y):
103
104     # Initialize variables (e.g. x_B = x_batched)
105     batch, x_B, y_B = [], [], []
106     current_batch, current_size = 0, 0
107     i = 1
108
109     # Maximum # Tokens in batch = # Tokens in largest sequence (x_train = sorted)
110     max_tokens = len(x[-1])
111
112     # Iterates over dataset and slices into appropriate sizes (token-based)
113     while i != len(x) + 1:
114         batch.append(x[i - 1])
115         current_size += len(x[i - 1])
116
117         # Final batch/element (longest sequence)
118         if i == len(x):
119             x_B.append(batch)
120             y_B.append(y[-1:])
121
122         # Slice batch if max tokens condition violated
123         elif current_size + len(x[i]) >= max_tokens:
124
125             # Padding for batch
126             for j in range(len(batch)):
127                 batch[j] = pad(batch[j], size = len(batch[-1]), padding = 0)
128
129             # store current batch
130             x_B.append(batch)
131             y_B.append(y[i - len(batch): i])
132
133             # reset variables
134             current_size = 0
135             batch = []
136             current_batch += 1
137
138         i += 1
139
140     # convert to Torch Tensors
141     for i in range(len(x_B)):
142         x_B[i] = torch.tensor(x_B[i], dtype = torch.long)
143         y_B[i] = torch.tensor(y_B[i], dtype = torch.long)
144
145     return x_B, y_B
146
147 # Obtained batched training data
148 x_train_B, y_train_B = batch_slicing(x_train, y_train)
149 x_val_B, y_val_B = batch_slicing(x_val, y_val)

```

Figure 4: Variable Batch Size and Padding functions

## 6.2 Model Implementations

```

153 # MLP
154 """
155 MLP network Class
156 utilizes PyTorch built-in functions
157 """
158
159 class MLP(nn.Module):
160     def __init__(self):
161         super(MLP, self).__init__()
162         self.emb = nn.Embedding(num_embeddings = len(i2w), embedding_dim = 300)
163         self.fc1 = nn.Linear(300, 300)
164         self.fc2 = nn.Linear(300, 2)
165
166     def forward(self, x):
167         x = self.emb(x)
168         x = self.fc1(x)
169         x = F.relu(x)
170         x, y = torch.max(x, dim = 1) # global max pool
171         x = self.fc2(x)
172
173         return x
174
175 mlp_net = MLP()

```

Figure 5: MLP

```

179 # Elman (without Torch built-in function)
180 """
181 Elman Network Class, and Elman (selfbuild) Layer class
182 doesn't utilize PyTorch built-in RNN function, uses "Elman_selfbuild" class instead
183 """
184
185 class Elman(nn.Module):
186     def __init__(self, insize=300, outsize=300, hsize=300):
187         super().__init__()
188         self.fc1 = nn.Linear(insize + hsize, hsize)
189         self.fc2 = nn.Linear(hsize, outsize)
190
191     def forward(self, x, hidden=None):
192         b, t, e = x.size()
193
194         # first iteration
195         if hidden is None:
196             hidden = torch.zeros(b, e, dtype = torch.float)
197
198         outputs_list = []
199         for i in range(t):
200             inputs = torch.cat([x[:, i, :], hidden], dim = 1)
201             hidden = torch.sigmoid(self.fc1(inputs))
202             outputs = self.fc2(hidden)
203             outputs_list.append(outputs[:, None, :])
204
205         return torch.cat(outputs_list, dim = 1), hidden
206

```

Figure 6: Elman (Selfbuild)

```

226 # Elman (with Torch)
227 """
228 Elman Network Class
229 utilizes PyTorch built-in function
230 """
231
232 class Elman_Torch(nn.Module):
233     def __init__(self):
234         super(Elman_Torch, self).__init__()
235         self.emb = nn.Embedding(num_embeddings = len(i2w), embedding_dim = 300)
236         self.elman = torch.nn.RNN(input_size=300, hidden_size=300) # tanh non-linear
237         self.fc2 = nn.Linear(300, 2)
238
239     def forward(self, x):
240         x = self.emb(x)
241         x_residual, hidden = self.elman(x) # residual connection
242         x = x + x_residual
243         x, y = torch.max(x, dim = 1) # global max pool
244         x = self.fc2(x)
245
246         return x
247
248 elman_torch_net = Elman_Torch()

```

Figure 7: Elman (Torch)

```

252 # LSTM (with Torch)
253 """
254 LSTM Network Class
255 utilizes PyTorch built-in function
256 """
257 class LSTM_Torch(nn.Module):
258     def __init__(self):
259         super(LSTM_Torch, self).__init__()
260         self.emb = nn.Embedding(num_embeddings = len(i2w), embedding_dim = 300)
261         self.lstm = torch.nn.LSTM(input_size=300, hidden_size=300)
262         self.fc2 = nn.Linear(300, 2)
263
264     def forward(self, x):
265         x = self.emb(x)
266         x, hidden = self.lstm(x)
267         x, y = torch.max(x, dim = 1) # global max pool
268         x = self.fc2(x)
269
270         return x
271
272 lstm_torch_net = LSTM_Torch()

```

Figure 8: LSTM (Torch)

## Appendix

- **Model Architecture specifications:**

Input → Embedding Layer (300) → [Model Dependent] → Global Maxpool  
 → Fully Connected layer (300,2) → Softmax → Cross Entropy Loss

$$[ModelDependent] = \begin{cases} MLP : & ReLu \\ Elman : & Elman(300, 300, 300), \text{ non-linearity} = \tanh \\ LSTM : & LSTM(300, 300) \end{cases}$$

- **Vanishing/Exploding gradients problem:** In our first version of the Elman (self-build) implementation, we kept running into vanishing gradient problems after roughly 1000 batches (Originating from the Sigmoid activation). For this reason we decided to implement residual connections (Code snippets, figure 6, 7) The original undesired output<sup>7</sup> looked like this:

```

[1, 1090] loss: 0.375
[1, 1100] loss: 0.411
[1, 1110] loss: 268345.881
[1, 1120] loss: nan
[1, 1130] loss: nan
[1, 1140] loss: nan

```

Figure 9: Output Gradient problems

- **Residual Connection implementation problem:** At our first attempt of implementing residual connection we ran into an interesting issue, we managed to resolve. However we still don't understand how it works. It was related to in place variable modification, disturbing the auto-gradient computations.

```

203 class Elman_Torch(nn.Module):
204     def __init__(self):
205         super(Elman_Torch, self).__init__()
206         self.emb = nn.Embedding(num_embeddings = len(i2w), embedding_dim = 300)
207         self.elman = torch.nn.RNN(input_size=300, hidden_size=300) # tanh non-
208         self.fc2 = nn.Linear(300, 2)
209
210     def forward(self, x):
211         x = self.emb(x)
212         x_residual, hidden = self.elman(x)
213
214         # Fails
215         x += x_residual
216
217         # Works
218         x = x + x_residual
219
220         x, y = torch.max(x, dim = 1) # global max pool
221         x = self.fc2(x)
222
223         return x
224
225 elman_torch_net = Elman_Torch()

```

Figure 10: Residual Connection Coding Issue (Line 214-218)

<sup>7</sup>Numbers between brackets can be interpreted as: [epoch, iteration/batch]