

ECE421 - Winter 2022

Assignment 2: Neural Networks

Faranak Dayyani, Student Number: 1002373674

Part 1: Neural Networks using Numpy

1.1: Helper Functions

1. relu():

```
54  
55 def relu(x):  
56     re_lu = np.maximum(x,0)  
57     return re_lu  
58
```

Figure 1: relu() function code snippet

2. softmax():

```
60 def softmax(x):  
61  
62     x = x - x.max(axis=1, keepdims=True)  
63     z = np.exp(x) / np.exp(x).sum(axis=1, keepdims=True)  
64  
65     return z  
66
```

Figure 2: softmax() function code snippet

3. compute_layer():

```
68 def compute_layer(x, w, b):  
69  
70     product = x @ w + b  
71  
72     return product  
73
```

Figure 3: compute_layer() function code snippet

4. average_ce():

```
74 def average_ce(target, prediction):  
75  
76     avg_ce = -np.mean(target*np.log(prediction))  
77  
78     return avg_ce  
79
```

Figure 4: average_ce() function code snippet

5. grad_ce():

Analytical expression:

$$L = -\sum_{k=1}^k y_k \log(p_k) \quad , \quad \frac{\partial L}{\partial O} = ? \quad , \quad \frac{\partial L}{\partial O} = \frac{\partial L}{\partial P_i} \times \frac{\partial P_i}{\partial O_i}$$

$$p_i = \text{softmax}(O_i) = \frac{e^{O_i}}{\sum_{k=1}^k e^{O_k}} \quad , i = 1, \dots, k \text{ for } k \text{ classes}$$

$$\frac{\partial L}{\partial P_i} = -\sum_{k=1}^K y_i \frac{1}{p_i}$$

$$\text{for } i = j \Rightarrow \frac{\partial P_i}{\partial O_i} = \frac{e^{O_i} (\sum_{k=1}^k e^{O_k}) - e^{O_i} e^{O_i}}{(\sum_{k=1}^k e^{O_k})^2}$$

$$\frac{\partial P_i}{\partial O_i} = \frac{e^{O_i}}{(\sum_{k=1}^k e^{O_k})} \frac{(\sum_{k=1}^k e^{O_k})}{(\sum_{k=1}^k e^{O_k})} - \frac{e^{O_i} e^{O_i}}{(\sum_{k=1}^k e^{O_k})^2}$$

$$\frac{\partial P_i}{\partial O_i} = p_i - p_i^2 = p_i(1 - p_i)$$

$$\text{for } i \neq j \Rightarrow \frac{\partial P_i}{\partial O_j} = \frac{0 - e^{O_i} e^{O_j}}{(\sum_{k=1}^k e^{O_k})^2}$$

$$\frac{\partial P_i}{\partial O_j} = -\frac{e^{O_i}}{(\sum_{k=1}^k e^{O_k})} \frac{e^{O_j}}{(\sum_{k=1}^k e^{O_k})}$$

$$\frac{\partial P_i}{\partial O_j} = -p_i p_j$$

$$\frac{\partial L}{\partial O} = \frac{\partial L}{\partial P_i} \times \frac{\partial P_i}{\partial O_i} + \frac{\partial L}{\partial P_i} \times \frac{\partial P_i}{\partial O_j}$$

$$\frac{\partial L}{\partial O} = \left(-\sum_{i=j}^k y_i \frac{1}{p_i} \right) (p_i(1 - p_i)) + \left(-\sum_{i \neq j}^k y_i \frac{1}{p_i} \right) (-p_i p_j)$$

$$\frac{\partial L}{\partial O} = \sum_i -y_i + p_i y_i + p_i y_i = p_i \sum_i y_i - y_i = p_i - y_i$$

$$\frac{\partial L}{\partial O} = \frac{1}{N} (p - y)$$

```

81 def grad_ce(target, logits):
82
83     gr_ce = (softmax(logits) - target)/target.shape[0]
84
85     return gr_ce

```

Figure 5: grad_ce() function code snippet

1.2: Backpropagation Derivation

1. The gradient of the loss with respect to the output layer weights ($\frac{\partial L}{\partial w_o}$):

$$O = w_o h + b_o \Rightarrow \frac{\partial O}{\partial w_o} = h$$

$$\frac{\partial L}{\partial w_o} = \frac{\partial O}{\partial w_o} \times \frac{\partial L}{\partial O} = h^T \times (p - y)$$

```

90 def grad_wrt_w0(h, p, y):
91
92     p_y = grad_ce(y, p)
93     h_tran = np.transpose(h)
94     gradw0 = np.matmul(h_tran, p_y)
95
96     return gradw0

```

Figure 6: grad_wrt_w0 function code snippet

2. The gradient of the loss with respect to the output layer biases ($\frac{\partial L}{\partial b_o}$):

$$O = w_o h + b_o \Rightarrow \frac{\partial O}{\partial b_o} = 1$$

$$\frac{\partial L}{\partial b_o} = \frac{\partial O}{\partial b_o} \times \frac{\partial L}{\partial O} = 1^T \times (p - y)$$

```

101 def grad_wrt_b0(p, y):
102
103     p_y = grad_ce(y, p)
104     ones_shape = np.ones((1, y.shape[0]))
105     gradb0 = np.matmul(ones_shape, p_y)
106
107     return gradb0

```

Figure 7: grad_wrt_b0 function code snippet

3. The gradient of the loss with respect to the hidden layer weights ($\frac{\partial L}{\partial w_h}$):

$$O = w_o h + b_o \Rightarrow \frac{\partial O}{\partial h} = w_o$$

$$\frac{\partial L}{\partial O} = p - y$$

$$z = w_h x + b_h \Rightarrow \frac{\partial z}{\partial w_h} = x$$

$$\frac{\partial h}{\partial z} = \begin{cases} 0, & \text{if } z_i < 0 \\ 1, & \text{if } z_i > 0 \end{cases}$$

$$\frac{\partial L}{\partial w_h} = \frac{\partial L}{\partial O} \times \frac{\partial O}{\partial h} \times \frac{\partial h}{\partial z} \times \frac{\partial z}{\partial w_h} = (p - y) w_o^T \times \frac{\partial h}{\partial z} x^T$$

```

110 def grad_wrt_wh(p, y, x, h_wh, w_o):
111
112     h_wh[h_wh > 0] = 1
113     h_wh[h_wh < 0] = 0
114
115     x_tran = x.transpose()
116     p_y = grad_ce(y, p)
117
118     part1 = np.matmul(p_y, np.transpose(w_o))
119     part2 = h_wh * part1
120
121     grad_wh = np.matmul(x_tran, part2)
122
123     return grad_wh

```

Figure 8: grad_wrt_wh function code snippet

4. The gradient of the loss with respect to the hidden layer biases ($\frac{\partial L}{\partial b_h}$):

$$O = w_o h + b_o \Rightarrow \frac{\partial O}{\partial h} = w_o$$

$$\frac{\partial L}{\partial O} = p - y$$

$$z = w_h x + b_h \Rightarrow \frac{\partial z}{\partial b_h} = 1$$

$$\frac{\partial h}{\partial z} = \begin{cases} 0, & \text{if } z_i < 0 \\ 1, & \text{if } z_i > 0 \end{cases}$$

$$\frac{\partial L}{\partial b_h} = \frac{\partial L}{\partial o} \times \frac{\partial o}{\partial h} \times \frac{\partial h}{\partial z} \times \frac{\partial z}{\partial b_h} = (p - y)w_0^T \times \frac{\partial h}{\partial z} 1^T$$

```

127 def grad_wrt_bh(p, y, x, h_wh, w_0):
128
129     h_wh[h_wh > 0] = 1
130     h_wh[h_wh < 0] = 0
131
132     p_y = grad_ce(y, p)
133     ones_shape = np.ones((x.shape[0], 1))
134     ones_shape_tran = ones_shape.transpose()
135
136     part1 = np.matmul(p_y, np.transpose(w_0))
137     part2 = h_wh * part1
138
139     grad_bh = np.matmul(ones_shape_tran, part2)
140
141     return grad_bh

```

Figure 9: grad_wrt_bh function code snippet

1.3: Learning

The model was initialized with the values below:

- The number of hidden units are 1000 (H=1000).
- The hidden layer weights are initialized with zero mean gaussian distribution with variance of $\frac{2}{units\ in+units\ out} = \frac{2}{784+1000}$.
- The output layer weights are initialized with zero mean gaussian distribution with variance of $\frac{2}{units\ in+units\ out} = \frac{2}{1000+10}$.
- The biases are initialized to zero.
- The learning rate (α) is set to 0.1.
- γ is set to 0.99.
- The number of epochs that the model was run for is 200.

Table 1: The table below shows the loss and accuracy values for training, validation and testing data sets:

| | Loss | Accuracy |
|----------------|----------|----------|
| Training set | 0.054638 | 89.13 % |
| Validation set | 0.060305 | 87.88 % |
| Testing set | 0.060478 | 88.36 % |

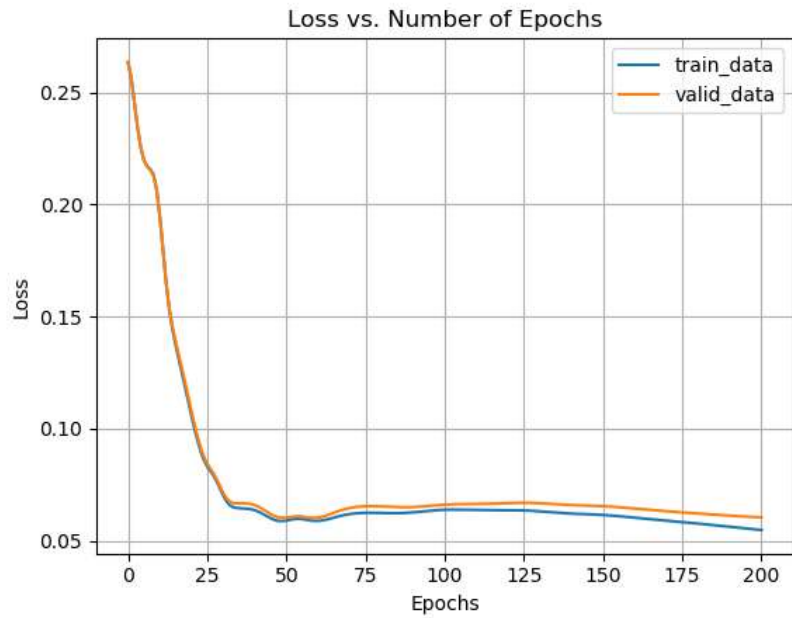


Figure 10: Loss vs. Number of Epochs for training and validation set

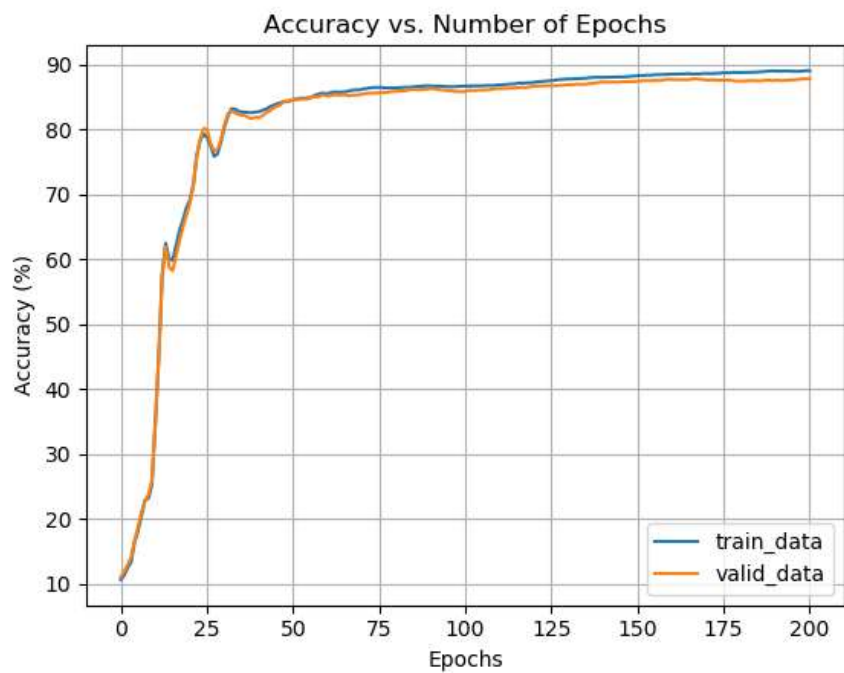


Figure 11: Accuracy vs. Number of Epochs for training and validation set

The code snippet for this part can be seen below:

```
140 def learning(train_data, train_target, valid_data, valid_target, w_o, b_o, v_o, w_h, b_h, v_h, epochs, gamma, alpha):
141
142     global i, b_o_1, b_h_1
143
144     # creating array variables for training and validation loss and accuracies
145     train_loss = []
146     train_acc = []
147     valid_loss = []
148     valid_acc = []
149
150     b_o_1 = b_o
151     b_h_1 = b_h
152
153     for i in range(epochs):
154
155         # ----- forward propagation -----
156         # training data: calculating loss and accuracy
157         # hidden layer
158         producttrain_h = compute_layer(train_data, w_h, b_h)
159         re_lu_train_h = relu(producttrain_h)
160         # output layer
161         producttrain_o = compute_layer(re_lu_train_h, w_o, b_o)
162         y_hat_train = softmax(producttrain_o)
163         avg_ce_train = average_ce(new_train, y_hat_train)
164         train_loss.append(avg_ce_train)
165         acc1 = np.argmax(y_hat_train, axis=1)
166         acc11 = np.argmax(new_train, axis=1)
167         acc111 = np.sum(acc1 == acc11)/train_data.shape[0]
168         acc111 = acc111*100
169         train_acc.append(acc111)
170
171         # validation data: calculating loss and accuracy
172         # hidden layer
173         productvalid_h = compute_layer(valid_data, w_h, b_h)
174         re_lu_valid_h = relu(productvalid_h)
175         # output layer
176         productvalid_o = compute_layer(re_lu_valid_h, w_o, b_o)
177         y_hat_valid = softmax(productvalid_o)
178         avg_ce_valid = average_ce(new_valid, y_hat_valid)
```

Figure 12: learning() function code snippet – 1

```

178     avg_ce_valid = average_ce(new_valid, y_hat_valid)
179     valid_loss.append(avg_ce_valid)
180     acc2 = np.argmax(y_hat_valid, axis=1)
181     acc22 = np.argmax(new_valid, axis=1)
182     acc222 = np.sum(acc2 == acc22)/valid_data.shape[0]
183     acc222 = acc222*100
184     valid_acc.append(acc222)
185
186
187     # ----- Backpropagation -----
188     # Output layer
189     gradw0 = grad_wrt_w0(re_lu_train_h, producttrain_o, new_train)
190     v_o = (gamma * v_o) + (alpha * gradw0)
191     w_o = w_o - v_o
192
193     gradb0 = grad_wrt_b0(producttrain_o, new_train)
194     b_o_1 = (gamma * b_o_1) + (alpha * gradb0)
195     b_o = b_o - b_o_1
196
197     # Hidden layer
198     grad_wh = grad_wrt_wh(producttrain_o, new_train, train_data, producttrain_h, w_o)
199     v_h = (gamma * v_h) + (alpha * grad_wh)
200     w_h = w_h - v_h
201
202     grad_bh = grad_wrt_bh(producttrain_o, new_train, train_data, producttrain_h, w_o)
203     b_h_1 = (gamma * b_h_1) + (alpha * grad_bh)
204     b_h = b_h - b_h_1
205
206     print("Going through the Loop: Epoch ", i+1,"/", epochs)
207
208     return w_o, b_o, w_h, b_h, train_loss, train_acc, valid_loss, valid_acc
209

```

Figure 13: learning() function code snippet – 2

```

211 if __name__ == "__main__":
212     train_data, valid_data, test_data, train_target, valid_target, test_target = load_data()
213     train_data = train_data.reshape(10000, 100)
214     valid_data = valid_data.reshape(4000, 100)
215     test_data = test_data.reshape(1000, 100)
216
217     new_train, new_valid, new_test = convert_onehot(train_target, valid_target, test_target)
218
219     # ----- section 1.1 -----
220     epochs = 1000
221     H = 1000
222     gamma = 0.99
223     alpha = 0.1
224
225     # initializing weights with Xavier initialization scheme:
226     # bias set to 0
227     # initial momentum value of 1e-5
228
229     # ----- output layer initialization -----
230     limit_o = np.sqrt(2/(H+H))
231     w_o = np.random.normal(0,0, limit_o, (H,10))
232     v_o = np.full((H,10), 1e-5)
233     b_o = np.zeros((1,10))
234
235     # ----- hidden layer initialization -----
236     limit_h = np.sqrt(2/(train_data.shape[1]+H))
237     w_h = np.random.normal(0,0, limit_h, (train_data.shape[1],H))
238     v_h = np.full((train_data.shape[1], H), 1e-5)
239     b_h = np.zeros((1, H))
240
241     w_o, b_o, w_h, b_h, train_loss, train_acc, valid_loss, valid_acc = learning(train_data, new_train, valid_data, new_valid, w_o, b_o,
242     v_o, w_h, b_h, v_h, epochs, gamma, alpha/10)
243

```

Figure 14: main() function code snippet – 1


```

245 # Training data result:
246 # ----- hidden layer -----
247 product_h_train = compute_layer(train_data, w_h, b_h)
248 relu_h_train = relu(product_h_train)
249 # ----- output layer -----
250 product_o_train = compute_layer(relu_h_train, w_o, b_o)
251 yhat_train = softmax(product_o_train)
252 avgce_train = average_ce(new_train, yhat_train)
253 train_loss.append(avgce_train)
254 acctrain1 = np.argmax(yhat_train, axis=1)
255 acctrain11 = np.argmax(new_train, axis=1)
256 acctrain111 = np.sum(acctrain1 == acctrain11)/train_data.shape[0]
257 acctrain111 = acctrain111*100
258 train_acc.append(acctrain111)
259 print(" ----- Training Data -----")
260 print("Training loss is:", train_loss[-1])
261 print("Training accuracy is:", train_acc[-1], "%")
262
263 # Validation data result:
264 # ----- hidden layer -----
265 product_h_valid = compute_layer(valid_data, w_h, b_h)
266 relu_h_valid = relu(product_h_valid)
267 # ----- output layer -----
268 product_o_valid = compute_layer(relu_h_valid, w_o, b_o)
269 yhat_valid = softmax(product_o_valid)
270 avgce_valid = average_ce(new_valid, yhat_valid)
271 valid_loss.append(avgce_valid)
272 accvalid1 = np.argmax(yhat_valid, axis=1)
273 accvalid11 = np.argmax(new_valid, axis=1)
274 accvalid111 = np.sum(accvalid1 == accvalid11)/valid_data.shape[0]
275 accvalid111 = accvalid111*100
276 valid_acc.append(accvalid111)
277 print("\n ----- Validation Data -----")
278 print("Validation loss is:", valid_loss[-1])
279 print("Validation accuracy is:", valid_acc[-1], "%")
280

```

Figure 15: main() function code snippet – 2

```

283 # Testing data result:
284 test_loss = []
285 test_acc = []
286 # ----- hidden layer -----
287 product_h_test = compute_layer(test_data, w_h, b_h)
288 relu_h_test = relu(product_h_test)
289 # ----- output layer -----
290 product_o_test = compute_layer(relu_h_test, w_o, b_o)
291 yhat_test = softmax(product_o_test)
292 avgce_test = average_ce(new_test, yhat_test)
293 test_loss.append(avgce_test)
294 acctest1 = np.argmax(yhat_test, axis=1)
295 acctest11 = np.argmax(new_test, axis=1)
296 acctest111 = np.sum(acctest1 == acctest11)/test_data.shape[0]
297 acctest111 = acctest111*100
298 test_acc.append(acctest111)
299 print("\n ----- Testing Data -----")
300 print("Testing loss is:", test_loss[-1])
301 print("Testing accuracy is:", test_acc[-1], "%")
302
303
304
305 # plotting loss for training and validation data:
306 xaxis = np.linspace(0, epochs, num=epochs)
307 plt.plot(xaxis, train_loss[0:epochs])
308 plt.plot(xaxis, valid_loss[0:epochs])
309 plt.xlabel('Epochs')
310 plt.ylabel('Loss')
311 plt.title('Loss vs. Number of Epochs')
312 plt.legend(['train_data', 'valid_data'])
313 plt.grid()
314 plt.show()
315 plt.figure()

```

Figure 16: main() function code snippet – 3

```

317 # plotting accuracy for training and validation data:
318 plt.plot(xaxis, train_acc[0:epochs])
319 plt.plot(xaxis, valid_acc[0:epochs])
320 plt.xlabel('Epochs')
321 plt.ylabel('Accuracy (%)')
322 plt.title('Accuracy vs. Number of Epochs')
323 plt.legend(['train_data', 'valid_data'])
324 plt.grid()
325 plt.show()
326

```

Figure 17: main() function code snippet - 4