

This homework is due Saturday, October 14 at 10pm.

## 1 Getting Started

You may typeset your homework in latex or submit neatly handwritten and scanned solutions. Please make sure to start each question on a new page, as grading (with Gradescope) is much easier that way! Deliverables:

1. Submit a PDF of your writeup to assignment on Gradescope, “HW[n] Write-Up”
2. Submit all code needed to reproduce your results, “HW[n] Code”.
3. Submit your test set evaluation results, “HW[n] Test Set”.

After you've submitted your homework, be sure to watch out for the self-grade form.

- (a) Before you start your homework, write down your team. Who else did you work with on this homework? List names and email addresses. In case of course events, just describe the group. How did you work on this homework? Any comments about the homework?

None. I work alone

Comments: The skeleton code is terrible! It is still buggy and not updated until last minutes. MAKE SURE IT WORKS ON YOUR

COMPUTER BEFORE RELEASING IT.

- (b) Please copy the following statement and sign next to it:

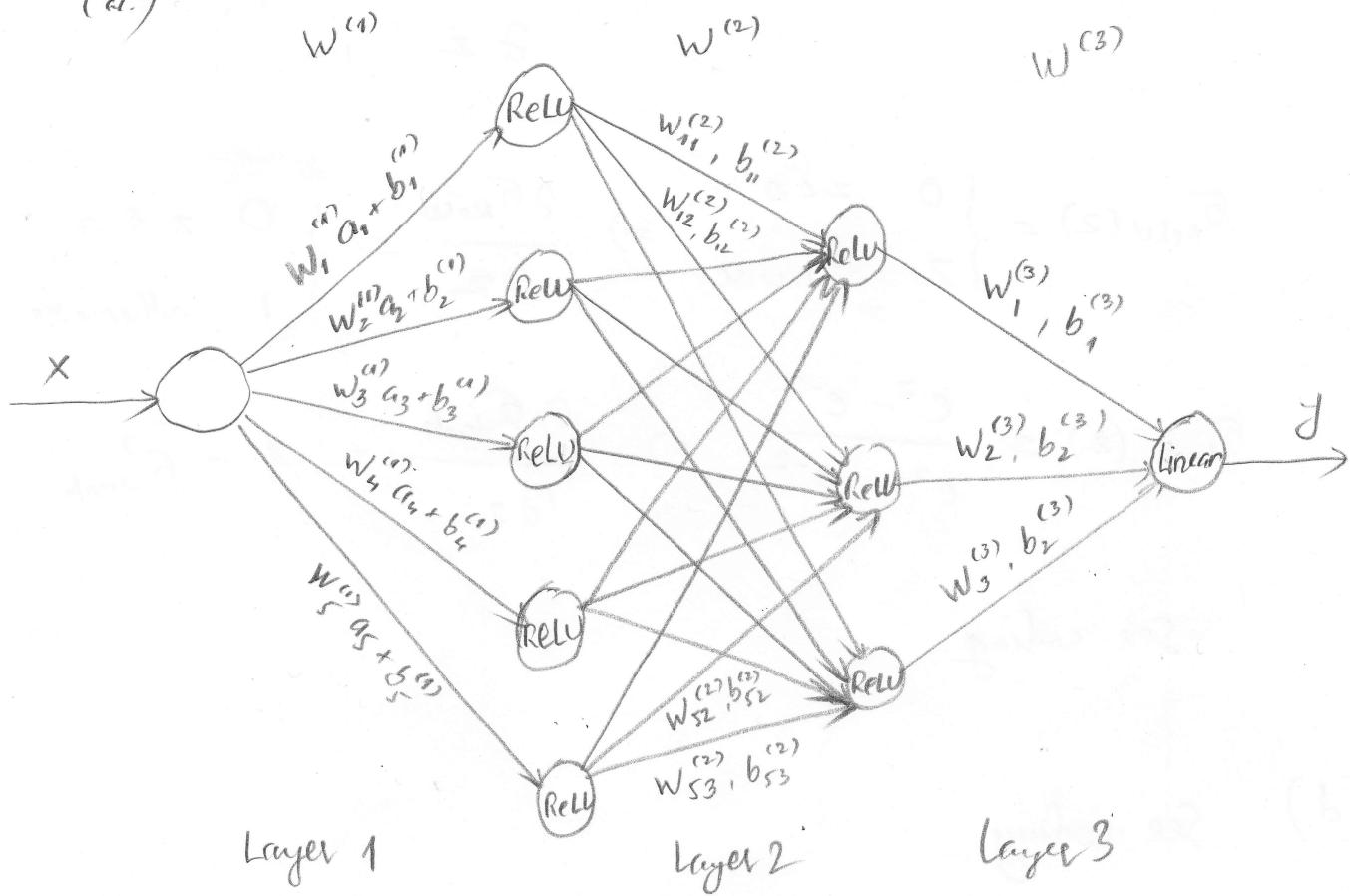
I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up.

I certify that all solutions are entirely in my words and that I have not looked at another student's solutions. I have credited all external sources in this write up.



## Problem #2

(a.)



(b.)

$$MSE(\hat{y}) = \frac{1}{2} \sum_{i=1}^n \|y_i - \hat{y}_i\|_2^2$$

$$\Rightarrow \frac{\partial MSE}{\partial \hat{y}} = \hat{y} - y$$

See coding

$$(c) \quad \tilde{\sigma}_{\text{linear}}(z) = z \Rightarrow \frac{\partial \tilde{\sigma}_{\text{linear}}}{\partial z} = 1$$

$$\tilde{\sigma}_{\text{ReLU}}(z) = \begin{cases} 0 & z < 0 \\ 1 & \text{otherwise} \end{cases} \Rightarrow \frac{\partial \tilde{\sigma}_{\text{ReLU}}}{\partial z} = \begin{cases} 0 & z < 0 \\ 1 & \text{otherwise} \end{cases}$$

$$\tilde{\sigma}_{\text{tanh}}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \Rightarrow \frac{\partial \tilde{\sigma}_{\text{tanh}}}{\partial z} = 1 - \tilde{\sigma}_{\text{tanh}}^2$$

See coding

(d) See coding

$$(e) \quad \frac{\partial \text{MSE}}{\partial \vec{a}_i} = \frac{\partial \text{MSE}}{\partial \vec{a}_{i+1}} \underbrace{\frac{\partial \vec{a}_i}{\partial \vec{a}_i}}_{\vec{z}_i = w_i \vec{a}_i + b_i}$$

$$\frac{\partial \tilde{\sigma}(z_i)}{\partial z_i} \frac{\partial z_i}{\partial a_i}$$

$$\Rightarrow \frac{\partial \text{MSE}}{\partial \vec{a}_i} = \frac{\partial \text{MSE}}{\partial \vec{a}_{i+1}} \frac{\partial \tilde{\sigma}(z_i)}{\partial z_i} w_i$$

See coding

f)

$$\frac{\partial \text{MSE}}{\partial w_i} = \frac{\partial \text{MSE}}{\partial \vec{a}_{in}} \underbrace{\frac{\partial \vec{a}_{in}}{\partial w_i}}_{\vec{a}_i} \xrightarrow{\vec{z}_i = w_i \vec{a}_i + \vec{b}_i}$$

$$\frac{\partial \delta(\vec{z}_i)}{\partial \vec{z}_i} \underbrace{\frac{\partial \vec{z}_i}{\partial w_i}}_{\vec{a}_i}$$

$$\frac{\partial \text{MSE}}{\partial w_i} = \frac{\partial \text{MSE}}{\partial \vec{a}_{in}} \frac{\partial \delta(\vec{z}_i)}{\partial \vec{z}_i} \vec{a}_i$$

Similarly

$$\frac{\partial \text{MSE}}{\partial b_i} = \frac{\partial \text{MSE}}{\partial \vec{a}_{in}} \frac{\partial \delta(\vec{z}_i)}{\partial \vec{z}_i} \frac{\partial \vec{z}_i}{\partial b_i} \xrightarrow{I}$$

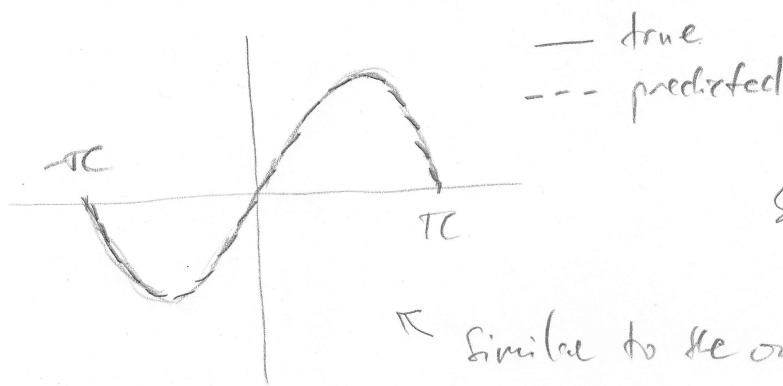
g) Output

ReLU MSE : 0.00074328

Linear MSE : 0.12079632

tanh MSE : 0.00173566

Plot :



See coding

↙ Similar to Re output

h.) not enough time to run

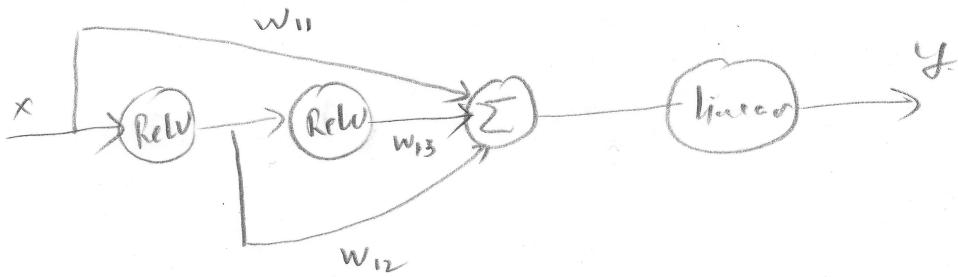
i.) not enough time to run

j.) skip

k.) skip

### Problem #3

Consider the net-work:



- 1.) Write the out-put for each layers (from  $x \rightarrow y$ )

Solution:

$$a_1 = w_{11}x + b_{11}$$

$$a_2 = w_{12} \delta_{\text{ReLU}}(x) + b_{12}$$

$$a_3 = w_{13} \delta_{\text{ReLU}}(\delta_{\text{ReLU}}(x)) + b_{13}$$

$$a_4 = w_4(a_1 + a_2 + a_3) + b_{14}$$

$$y = \delta_{\text{Linear}}(a_4) = a_4$$

where  $\delta$ 's are noise

- 2.) Define error & take the derivative of error w.r.t  
u given  $a_1, a_2, a_3$  &  $a_4$

Solution:

$$\text{MSE}(\vec{y}_p) = \frac{1}{2} \|\vec{y}_p - y\|_2^2$$

$$\frac{\partial \text{MSE}}{\partial \vec{x}} = \frac{\partial \text{MSE}}{\partial a_4} \underbrace{\frac{\partial a_4}{\partial a_1} \frac{\partial a_1}{\partial x} + \frac{\partial \text{MSE}}{\partial a_3} \frac{\partial a_3}{\partial a_2} \frac{\partial a_2}{\partial x} + \frac{\partial \text{MSE}}{\partial a_2} \frac{\partial a_2}{\partial a_1} \frac{\partial a_1}{\partial x}}_{w_{11}} + \frac{\partial \text{MSE}}{\partial a_3} \frac{\partial a_3}{\partial a_2} \frac{\partial a_2}{\partial x} + \frac{\partial \text{MSE}}{\partial a_2} \frac{\partial a_2}{\partial a_1} \frac{\partial a_1}{\partial x} + \frac{\partial \text{MSE}}{\partial a_1} \frac{\partial a_1}{\partial x}$$

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Gradient descent optimization
5 # The learning rate is specified by eta
6 class GDOptimizer(object):
7     def __init__(self, eta):
8         self.eta = eta
9
10    def initialize(self, layers):
11        pass
12
13    # This function performs one gradient descent step
14    # layers is a list of dense layers in the network
15    # g is a list of gradients going into each layer before the nonlinear activation
16    # a is a list of activations of each node in the previous layer going
17    def update(self, layers, g, a):
18        m = a[0].shape[1]
19        for layer, curGrad, curA in zip(layers, g, a):
20            # TODO
21            ######
22            ######
23            # Compute the gradients for layer.W and layer.b using the gradient for
24            # the output of the
25            # layer curA and the gradient of the output curGrad
26            # Use the gradients to update the weight and the bias for the layer
27            #
28            ######
29            # grad_W = curGrad.dot(curA.T)
30            # layer.W -= self.eta * grad_W
31            # grad_b = curGrad.dot( np.ones( (curGrad.shape[1], 1) ) )
32            # layer.b -= self.eta * grad_b
33            curZ = layer.z(curA)
34            grad_W = (curGrad * layer.activation.dx(curZ)).dot(curA.T)
35            layer.W -= self.eta * grad_W / m
36            grad_b = curGrad * layer.activation.dx(curZ)
37            layer.b = layer.b - self.eta * grad_b / m
38
39
40    # Cost function used to compute prediction errors
41    class QuadraticCost(object):
42
43        # Compute the squared error between the prediction yp and the observation y
44        # This method should compute the cost per element such that the output is the
45        # same shape as y and yp
46        @staticmethod
47        def fx(y, yp):
48            # TODO
49            ######
50            #####
51            # Implement me
52            #
53            ######
54            #####
55            return (y - yp) * (y - yp) / 2
```



```
92         #
93         #####
94         #####
95         #####
96         #####
97         #####
98         #####
99         #####
100        #####
101        #####
102        #####
103        class LinearActivation(object):
104            @staticmethod
105            def fx(z):
106                #
107                #####
108                #####
109                #####
110                #####
111                #####
112                #####
113                #####
114                #####
115                #####
116                #####
117                #####
118                #####
119                #####
120                #####
121                #####
122                #####
123                #####
124                #####
125                #####
126                #####
127                #####
128                #####
129                #####
130                #####
131                #####
132                #####
133                #####
134                ##### This class represents a single hidden or output layer in the neural network
class DenseLayer(object):
    # numNodes: number of hidden units in the layer
    # activation: the activation function to use in this layer
    def __init__(self, numNodes, activation):
        self.numNodes = numNodes
        self.activation = activation

    def getNumNodes(self):
        return self.numNodes

    # Initialize the weight matrix of this layer based on the size of the matrix W
    def initialize(self, fanIn, scale=1.0):
        s = scale * np.sqrt(6.0 / (self.numNodes + fanIn))
        self.W = np.random.normal(0, s,
                                (self.numNodes, fanIn))
```

```
135         self.b = np.random.uniform(-1,1,(self.numNodes,1))
136
137     # Apply the activation function of the layer on the input z
138     def a(self, z):
139         return self.activation.fx(z)
140
141     # Compute the linear part of the layer
142     # The input a is an n x k matrix where n is the number of samples
143     # and k is the dimension of the previous layer (or the input to the network)
144     def z(self, a):
145         return self.W.dot(a) + self.b # Note, this is implemented where we assume a
146         is k x n
147
148     # Compute the derivative of the layer's activation function with respect to z
149     # where z is the output of the above function.
150     # This derivative does not contain the derivative of the matrix multiplication
151     # in the layer. That part is computed below in the model class.
152     def dx(self, z):
153         return self.activation.dx(z)
154
155     # Update the weights of the layer by adding dw to the weights
156     def updateWeights(self, dw):
157         self.W = self.W + dw
158
159     # Update the bias of the layer by adding db to the bias
160     def updateBias(self, db):
161         self.b = self.b + db
162
163     # This class handles stacking layers together to form the completed neural network
164     class Model(object):
165
166         # inputSize: the dimension of the inputs that go into the network
167         def __init__(self, inputSize):
168             self.layers = []
169             self.inputSize = inputSize
170
171         # Add a layer to the end of the network
172         def addLayer(self, layer):
173             self.layers.append(layer)
174
175         # Get the output size of the layer at the given index
176         def getLayerSize(self, index):
177             if index >= len(self.layers):
178                 return self.layers[-1].getNumNodes()
179             elif index < 0:
180                 return self.inputSize
181             else:
182                 return self.layers[index].getNumNodes()
183
184         # Initialize the weights of all of the layers in the network and set the cost
185         # function to use for optimization
186         def initialize(self, cost, initializeLayers=True):
187             self.cost = cost
188             if initializeLayers:
189                 for i in range(0,len(self.layers)):
190                     if i == len(self.layers) - 1:
191                         self.layers[i].initialize(self.getLayerSize(i-1))
```

```
191         else:
192             self.layers[i].initialize(self.getLayerSize(i-1))
193
194     # Compute the output of the network given some input a
195     # The matrix a has shape n x k where n is the number of samples and
196     # k is the dimension
197     # This function returns
198     # yp - the output of the network
199     # a - a list of inputs for each layer of the newtork where
200     #      a[i] is the input to layer i
201     # z - a list of values for each layer after evaluating layer.z(a) but
202     #      before evaluating the nonlinear function for the layer
203 def evaluate(self, x):
204     curA = x.T
205     a = [curA]
206     z = []
207     for layer in self.layers:
208         # TODO
209         ######
210         #####
211         # Store the input to each layer in the list a
212         # Store the result of each layer before applying the nonlinear function
213         # in z
214         # Set yp equal to the output of the network
215         #
216         ######
217         ######
218         curZ = layer.z(curA)
219         z.append(curZ)
220         curA = layer.a(curZ)
221         a.append(curA)
222
223     yp = a.pop(-1)
224
225     return yp, a, z
226
227
228
229     # Compute the output of the network given some input a
230     # The matrix a has shape n x k where n is the number of samples and
231     # k is the dimension
232     def predict(self, a):
233         a, _, _ = self.evaluate(a)
234         return a.T
235
236
237
238     # Train the network given the inputs x and the corresponding observations y
239     # The network should be trained for numEpochs iterations using the supplied
240     # optimizer
241     def train(self, x, y, numEpochs, optimizer):
242
243         # Initialize some stuff
244         n = x.shape[0]
245
246
247         x = x.copy()
248         y = y.copy()
249         hist = []
250         optimizer.initialize(self.layers)
251
252         # Run for the specified number of epochs
```

```
243     for epoch in range(0, numEpochs):
244
245         # Feed forward
246         # Save the output of each layer in the list a
247         # After the network has been evaluated, a should contain the
248         # input x and the output of each layer except for the last layer
249         yp, a, z = self.evaluate(x)
250
251         # Compute the error
252         # C = self.cost.fx(yp,y.T)
253         # d = self.cost.dx(yp,y.T)
254         C = self.cost.fx(y.T, yp)
255         d = self.cost.dx(y.T, yp)
256         grad = []
257
258         # Backpropogate the error
259         idx = len(self.layers)
260         for layer, curZ in zip(reversed(self.layers), reversed(z)):
261             # TODO
262             ######
263             # Compute the gradient of the output of each layer with respect to
264             # the error
265             # grad[i] should correspond with the gradient of the output of layer i
266             #
267             ######
268             # if len(grad) == 0:
269             #     grad_i = d * layer.activation.dx(curZ)
270             #     grad.append(grad_i)
271             # else:
272             #     grad_i = (grad[-1].T.dot(last_layer_W)).T *
273             #             layer.activation.dx(curZ)
274             #     grad.append(grad_i)
275             # last_layer_W = layer.W
276             if len(grad) == 0:
277                 grad_i = layer.W.T.dot(d * layer.activation.dx(curZ))
278                 grad.append(grad_i)
279             else:
280                 grad_i = layer.W.T.dot(grad[-1] * layer.activation.dx(curZ))
281                 grad.append(grad_i)
282             grad.pop(-1)
283             grad.reverse()
284             grad.append(d)
285
286             # Update the errors
287             optimizer.update(self.layers, grad, a)
288
289             # Compute the error at the end of the epoch
290             yh = self.predict(x)
291             C = self.cost.fx(yh, y)
292             C = np.mean(C)
293             hist.append(C)
294
295         return hist
296
297 if __name__ == '__main__':
298
```

```
294     # Generate the training set
295     np.random.seed(9001)
296     x=np.random.uniform(-np.pi,np.pi,(1000,1))
297     y=np.sin(x)
298
299     activations = dict(ReLU=ReLUActivation,
300                         tanh=TanhActivation,
301                         linear=LinearActivation)
302     lr = dict(ReLU=0.02,tanh=0.02,linear=0.005)
303
304     for key in activations:
305
306         # Build the model
307         activation = activations[key]
308         model = Model(x.shape[1])
309         model.addLayer(DenseLayer(100,activation()))
310         model.addLayer(DenseLayer(100,activation()))
311         model.addLayer(DenseLayer(1,LinearActivation()))
312         model.initialize(QuadraticCost())
313
314         # Train the model and display the results
315         hist = model.train(x,y,500,GDOptimizer(eta=lr[key]))
316         yHat = model.predict(x)
317         error = np.mean(np.square(yHat - y))/2
318         print(key+' MSE: '+str(error))
319         plt.plot(hist)
320         plt.title(key+' Learning curve')
321         plt.show()
322
323         # TODO
324         ##########
325         # Plot the approximation of the sin function from all of the models
326         #
327         ##########
```