Deep Learning: Problem Formulation

1. You have a single hidden-layer neural network for a binary classification task. The input is $X \in \mathbb{R}^{n \times m}$, output $\hat{y} \in \mathbb{R}^{1 \times m}$ and true label $y \in \mathbb{R}^{1 \times m}$. The squared-error loss function is \mathcal{J} . The forward propagation equations are:

$$z = WX + B$$

$$\hat{y} = \sigma(z)$$

$$\mathcal{J} = (\hat{y} - y)^2$$

- (a) We want to compute how to change the weights W based on errors in the loss function J. Define a partial derivative D to compute this.
 Solution: We want to compute ∂J/∂W.
- (b) Express D in terms of a sequence of partial derivatives in the network, i.e., including $\partial z/\partial X$.

Solution: To compute $\partial \mathcal{J} / \partial W$, we can evaluate

$$\frac{\partial \mathcal{J}}{\partial W} = \frac{\partial \mathcal{J}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial W}$$

(c) Compute a closed-form expression for D.

Solution: We must compute each of the terms:

$$\begin{array}{rcl} \frac{\partial \mathcal{J}}{\partial \hat{y}} &=& 2(\hat{y}-y) \\ \frac{\partial \hat{y}}{\partial z} &=& \sigma(z)[1-\sigma(z)] \\ \frac{\partial z}{\partial W} &=& X \end{array}$$

This gives us $2(\hat{y} - y)\sigma(z)[1 - \sigma(z)]X$, which we can transform to

$$2(\hat{y} - y)\sigma(z)[1 - \sigma(z)]X = 2X[\sigma(z) - y]\sigma(z)[1 - \sigma(z)] \\ = 2X\sigma(WX + B)[\sigma(WX + B) - y][1 - \sigma(WX + B)]$$

- (d) Draw the computation graph corresponding to this network.
- (e) Write out the TensorFlow code for forward inference in the network.

```
1 # Python optimisation variables
2 | learning_rate = 0.5
3 \text{ epochs} = 10
4 | batch_size = 100
5
6 # declare the training data placeholders
7 | # input x - for 28 x 28 pixels = 784
8 x = tf.placeholder(tf.float32, [None, 784])
9 # declare the output data placeholder: 10 digits
   y = tf.placeholder(tf.float32, [None, 10])
11
12 |# now declare the weights connecting the input to the hidden layer
13 |W1 = tf.Variable(tf.random_normal([784, 300], stddev=0.03), name='W1')
14 b1 = tf.Variable(tf.random_normal([300]), name='b1')
15
16 # calculate the output of the hidden layer
17 | z = tf.add(tf.matmul(x, W1), b1)
18
19 |# now calculate the hidden layer output – a sigmoid activated
20 # output layer
21 | y_- = tf.nn.sigmoid(z)
22
23 MSE = tf.reduce_sum((y_- -y) **2))
24
25
26 # add an optimiser
27
   optimiser = tf.train.GradientDescentOptimizer(learning_rate=
       learning_rate).minimize(MSE)
28
29
   # finally setup the initialisation operator
   init_op = tf.global_variables_initializer()
30
31
32 # define an accuracy assessment operation
33 |correct_prediction = tf.equal(tf.argmax(y, 1), tf.argmax(y_, 1))
34 | accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float32))
```

$$\begin{aligned} \frac{d}{dx}\sigma(x) &= \frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right] \\ &= \frac{d}{dx} (1+e^x)^{-1} \\ &= (1+e^{-x})^{-2} (e^{-x}) \\ &= \frac{e^{-x}}{(1+e^{-x})^2} \\ &= \frac{1}{(1+e^{-x})} \frac{e^{-x}}{(1+e^{-x})} \\ &= \frac{1}{(1+e^{-x})} \frac{(1+e^{-x})-1}{(1+e^{-x})} \\ &= \frac{1}{(1+e^x)} \left[\frac{(1+e^x)}{(1+e^x)} - \frac{1}{(1+e^x)} \right] \\ &= \frac{1}{(1+e^x)} \left[1 - \frac{1}{(1+e^x)} \right] \\ &= \sigma(x)(1-\sigma(x)) \end{aligned}$$

2. You are solving the binary classification task of classifying images as cat vs. non-cat. You design a CNN with a single output neuron. Let the output of this neuron be z. The final output of your network, \hat{y} is given by: $\hat{y} = \sigma(ReLU(z))$.

You classify all inputs with a final value $\hat{y} \ge 0.5$ as cat images. What problem are you going to encounter?

Solution: $ReLU(z) = max(0, z) \ge 0 \ \forall z.$

 $\sigma(w) \ge 0.5 \ \forall w \ge 0.$

Hence $\hat{y} = \sigma(ReLU(z)) \ge 0.5 \ \forall z \text{ for } w = ReLU(z).$

- 3. You train a simple network in which the final output of your network, \hat{y} is given by a sigmoid activation function: $\hat{y} = \sigma(Wx + b)$.
 - (a) If you initialise the weights in W, b to be large numbers, show analytically that the network will not learn for input $x \ge 0$.
 - (b) If you initialise only the weights in b to be large numbers, will you have the same problem? Again, show why or why not.
- 4. You are given the following piece of code for forward propagation through a single hidden layer in a neural network. This layer uses the sigmoid activation. Identify and correct the error.

```
1 import numpy as np
2 def forward_prop(W, a_prev, b):
3 z = W*a_prev + b
4 a = 1/(1+np.exp(-z)) #sigmoid
5 return a
```

Solution: z = np.matmul(W, a prev) + b OR z = np.dot(W, a prev) + b