

COMP 4107 Assignment Two:

① Find  $\frac{\partial f}{\partial x}$  if:

a)  $f = x_1 + x_2 + x_3$ ,  $x = (x_1, x_2, x_3)$ :

Since  $f$  is a scalar &  $x$  is a vector, we use Scalar-by-Vector Conversion:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f}{\partial x_1} & \frac{\partial f}{\partial x_2} & \frac{\partial f}{\partial x_3} \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1}(x_1 + x_2 + x_3) & \frac{\partial}{\partial x_2}(x_1 + x_2 + x_3) & \frac{\partial}{\partial x_3}(x_1 + x_2 + x_3) \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$$

b)  $f = (x_1 + x_2, x_3 x_2, x_1 - x_3^2, x_2 - 8)$ ,  $x = (x_1, x_2, x_3)$ :

Since  $f$  is a vector &  $x$  is a vector, use Vector-by-Vector Conversion:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \frac{\partial f_1}{\partial x_3} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \frac{\partial f_2}{\partial x_3} \\ \frac{\partial f_3}{\partial x_1} & \frac{\partial f_3}{\partial x_2} & \frac{\partial f_3}{\partial x_3} \\ \frac{\partial f_4}{\partial x_1} & \frac{\partial f_4}{\partial x_2} & \frac{\partial f_4}{\partial x_3} \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x_1}(x_1 + x_2) & \frac{\partial}{\partial x_2}(x_1 + x_2) & \frac{\partial}{\partial x_3}(x_1 + x_2) \\ \frac{\partial}{\partial x_1}(x_3 x_2) & \frac{\partial}{\partial x_2}(x_3 x_2) & \frac{\partial}{\partial x_3}(x_3 x_2) \\ \frac{\partial}{\partial x_1}(x_1 - x_3^2) & \frac{\partial}{\partial x_2}(x_1 - x_3^2) & \frac{\partial}{\partial x_3}(x_1 - x_3^2) \\ \frac{\partial}{\partial x_1}(x_2 - 8) & \frac{\partial}{\partial x_2}(x_2 - 8) & \frac{\partial}{\partial x_3}(x_2 - 8) \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 1 & 0 \\ 0 & x_3 & x_2 \\ 1 & 0 & -2x_3 \\ 0 & 1 & 0 \end{bmatrix}$$

c)  $f = (x, 2x^3, -\cos 2x + 14)$ :

Since  $f$  is a vector &  $x$  is a scalar, use Vector-by-Scalar:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \\ \frac{\partial f_3}{\partial x} \end{bmatrix} = \begin{bmatrix} \frac{\partial}{\partial x}(x) \\ \frac{\partial}{\partial x}(2x^3) \\ \frac{\partial}{\partial x}(-\cos 2x + 14) \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \cdot 3 \cdot x^2 \\ -(-2 \sin 2x) \end{bmatrix} = \begin{bmatrix} 1 \\ 6x^2 \\ 2 \sin(2x) \end{bmatrix}$$

d)  $f = \begin{bmatrix} 9x & 3x \\ 2 & \sin(xy) \end{bmatrix}$ ,  $\frac{\partial f}{\partial x} = C$ , Hint:  $\frac{\partial xy}{\partial x} = 2cx + b$ :

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial}{\partial x}(9x) & \frac{\partial}{\partial x}(3x) \\ \frac{\partial}{\partial x}(2) & \frac{\partial}{\partial x}(\sin(xy)) \end{bmatrix} = \begin{bmatrix} 9 & 3 \\ 0 & \cos(xy) \cdot \frac{\partial(xy)}{\partial x} \end{bmatrix}$$

↳ AS  $\frac{\partial f}{\partial x} = C$  &  $\frac{\partial xy}{\partial x} = 2cx + b$  we Note:

$$= \begin{bmatrix} 9 & 3C \\ 0 & \cos(xy)(2cx + b) \end{bmatrix} = \begin{bmatrix} 9 & 3C \\ 0 & (2cx + b)\cos(xy) \end{bmatrix}$$

thus:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} 9 & 3C \\ 0 & (2cx + b)\cos(xy) \end{bmatrix}$$

$$e) f = \begin{bmatrix} 3p & 2q \\ p & 7 \end{bmatrix}, p = x_1 - 2x_2, q = x_2x_3, x = (x_1, x_2, x_3)$$

$$f = \begin{bmatrix} 3(x_1 - 2x_2) & 2(x_2x_3) \\ (x_1 - 2x_2) - (x_2x_3) & 7 \end{bmatrix} = \begin{bmatrix} 3x_1 - 6x_2 & 2x_2x_3 \\ x_1 - 2x_2 - x_2x_3 & 7 \end{bmatrix}$$

Since  $f$  is a matrix &  $x$  is a vector, use a matrix-by-vector Conversion:

thus, derive  $f$  by each  $x_i$  to get three matrices for a tensor

$$A_1 = \frac{\partial}{\partial x_1} \begin{bmatrix} 3x_1 - 6x_2 & 2x_2x_3 \\ x_1 - 2x_2 - x_2x_3 & 7 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 1 & 0 \end{bmatrix}$$

$$A_2 = \frac{\partial}{\partial x_2} \begin{bmatrix} 3x_1 - 6x_2 & 2x_2x_3 \\ x_1 - 2x_2 - x_2x_3 & 7 \end{bmatrix} = \begin{bmatrix} -6 & 2x_3 \\ -2 - x_3 & 0 \end{bmatrix}$$

$$A_3 = \frac{\partial}{\partial x_3} \begin{bmatrix} 3x_1 - 6x_2 & 2x_2x_3 \\ x_1 - 2x_2 - x_2x_3 & 7 \end{bmatrix} = \begin{bmatrix} 0 & 2x_2 \\ -x_2 & 0 \end{bmatrix}$$

So, the resulting tensor is:

$$\frac{\partial f}{\partial x} = [A_1, A_2, A_3]$$



- ② For a Perceptron model with the current weight vector  $W = (-0.3, 2.1, 1)$  & the current bias term  $W_0 = -0.2$ , the input/output pair  $(X, y) = ((1, 1, 0), -1)$  is given. If our learning constant  $\alpha = 0.2$  &  $C = 0$ , find the updated weights & the bias of the neuron after training for 3 iterations on this data point:

We know the training Rule:

$$W_i \leftarrow W_i + \Delta W_i \quad \text{where} \quad \Delta W_i = \eta(t - o)X_i$$

Since  $C = 0$ : (use bipolar mode)

$$o = \text{output} = \begin{cases} 1 & \text{if } \sum W_i * X_i > 0 \\ -1 & \text{otherwise} \end{cases}$$

As  $\alpha = 0.2$  &  $\eta = \alpha$ ,  $\eta = 0.2$

We are targeting a value of  $-1$ , thus  $t = -1$

thus, we have:

$$t = -1, o = \text{output}, \eta = 0.2 \quad \& \quad W_i * X_i = W_1 X_1 + W_2 X_2 + W_3 X_3 + W_0$$

iteration One:

$$W_1 = (-0.3)(1) + (2.1)(1) + (1)(0) - 0.2 = 1.6 > 0$$

$$\hookrightarrow o = 1$$

Our output is 1, yet  $y = -1$ , So reweight

$$\Delta W_i = 0.2(-1-1)(1, 1, 0) = -0.4(1, 1, 0) = (-0.4, -0.4, 0)$$

$$W_i \leftarrow (-0.3, 2.1, 1) + (-0.4, -0.4, 0) = (-0.7, 1.7, 1)$$

iteration Two:

$$t = -1, o = \text{output}, \eta = 0.2, W_1 = 1.6$$

$$W_2 = (-0.7)(1) + (1.7)(1) + (1)(0) + 1.6 = 2.6 > 0$$

$$\hookrightarrow o = 1 \quad \text{yet } y = -1, \text{ So reweight}$$

$$\Delta W_i = 0.2(-1-1)(1, 1, 0) = 0.2(-2)(1, 1, 0) = (-0.4, -0.4, 0)$$

$$W_i \leftarrow (-0.7, 1.7, 1) + (-0.4, -0.4, 0) = (-1.1, 1.3, 1)$$

iteration Three:

$$t = -1, o = \text{output}, \eta = 0.2, W_2 = 2.6$$

$$W_3 = (-1.1)(1) + (1.3)(1) + (1)(0) + 2.6 = 2.8 > 0$$

$$\hookrightarrow o = 1 \quad \text{yet } y = -1, \text{ So rewrite}$$

$$\Delta W_i = 0.2(-1-1)(1, 1, 0) = (-0.4, -0.4, 0)$$

$$W_i \leftarrow (-1.1, 1.3, 1) + (-0.4, -0.4, 0) = (-1.5, 0.9, 1)$$

thus, after three iterations we see that  $W = (-1.5, 0.9, 1)$  &  $W_3 = 2.8$

③ For each image, Pick the best number of layers for a hypothetical neural network classifier. Cite the reason:

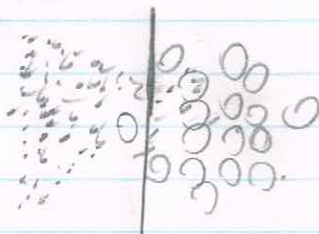
a) A single layer neural network would work best.

Reasoning:

The region in diagram A is a half-plane bounded by a Hyperplane, since it is a region separable by a single line.

For example:

Single layer



Note:

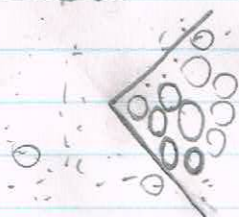
the 0's & .s are separable by a single line. Also, outliers exist, but can't easily be included using more layers as this would disrupt the main relationship.

b) A double layer neural network would work best:

The region in diagram B is a convex region since it is a region separable by a double line - V.

For example:

Two layer



Note:

the V separates the .s & 0's. Outliers are mostly scattered & can't be included by a higher-order model.

c) A single layer neural network would work, & a double can too:

The region in diagram C is a half-plane bounded by a Hyperplane, since it is a region separated by a single line.

For example:

Single or Two layer (two is more accurate)

Single-layer

double-layer

thus:

Single-layer's faster

double-layer's more accurate

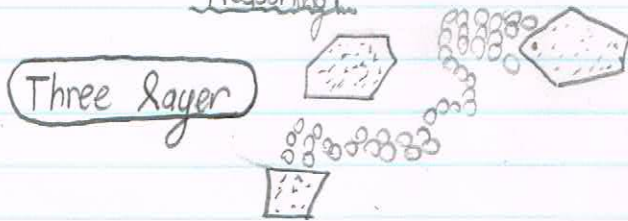
Note:

- the dashed line represents the second layer which is optional. The second layer can be used to segment out the extra-empty white space to make higher resolution results.
- The single layer model excludes the dashed line, whereas the double layer model forms a V-shape at the vertex & excludes the top half of the solid line (it forms a convex region.)



d) A Three layer neural network would work since there are multiple disconnected clusters of points. This indicates a high-complexity model should be used:

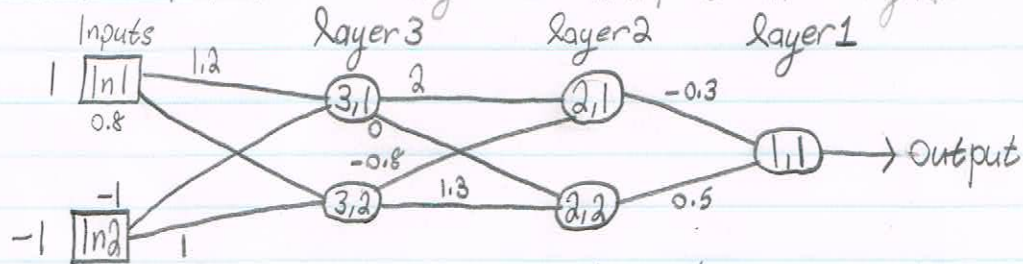
Reasoning:



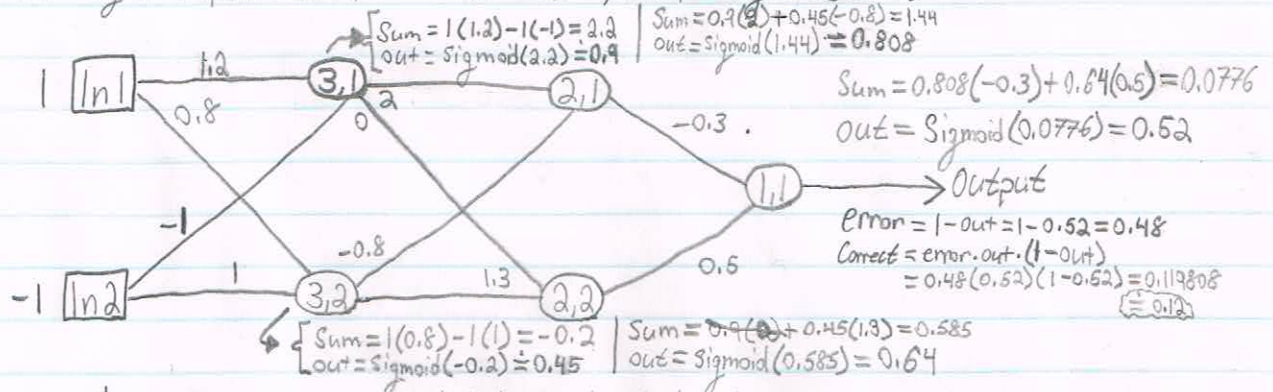
Note:

this lets us separate the 1's from the 0's, since we can skirt around the meshed regions between the 1's & 0's.

- ④ A multi-layered feed-forward network for classification is given. Suppose that no bias term exists in this particular network & all neurons have a sigmoid activation function. A data points (1,-1) with ~~label (1,-1) with~~ label 1 is given to this network to train. Calculate the weights of the neurons after 1 iteration of training for this data point with learning rate  $\alpha = 1$ :



Find Sigma outputs:  $\text{Sum}_i^k = \sum w_{ij} \cdot \text{out}_j^{k-1}$ ,  $\text{out}_i^k = \text{Sigmoid}(\text{Sum}_i^k)$

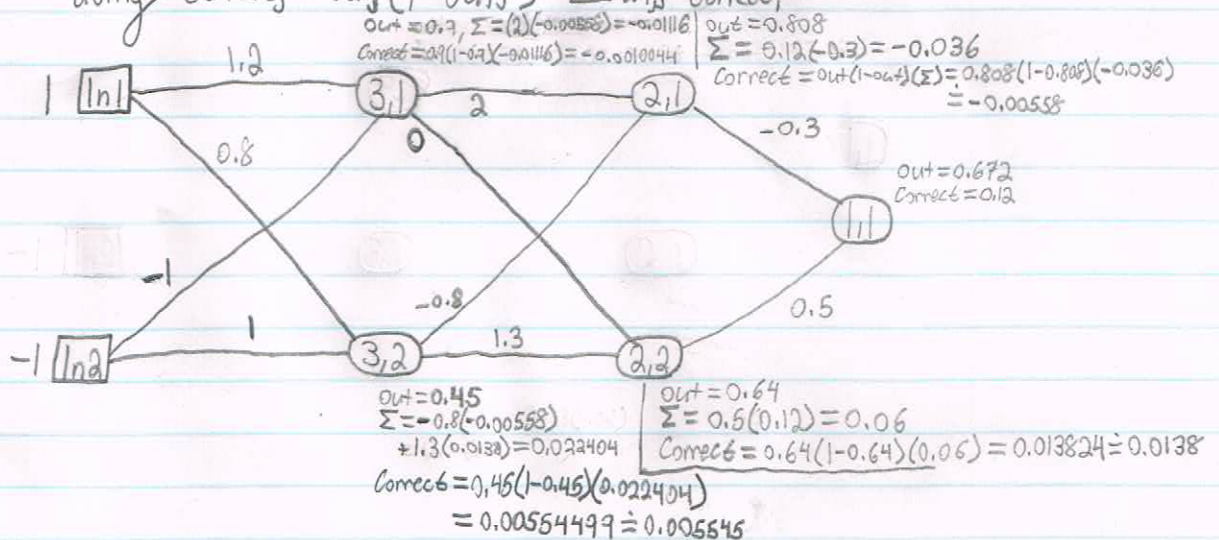


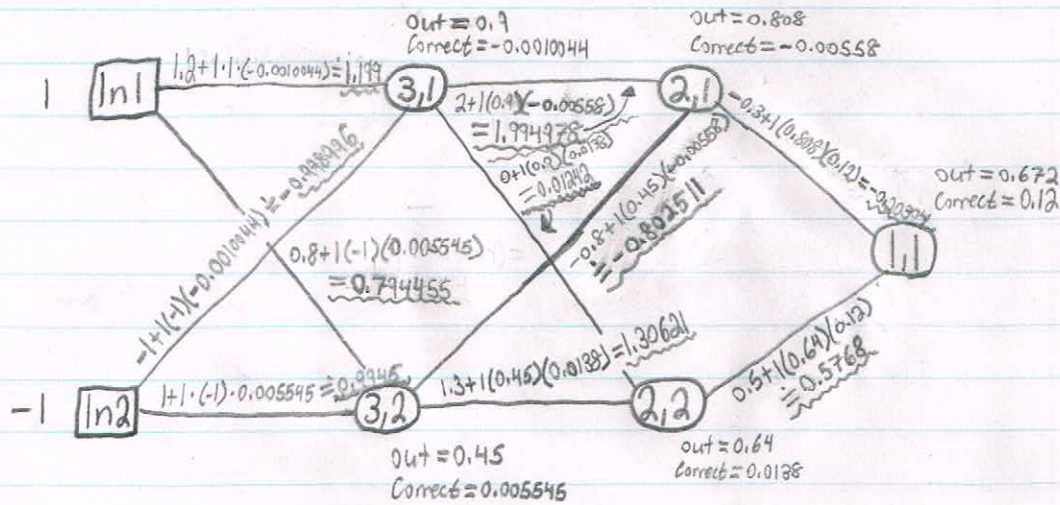
Where  $\text{LastNeuronError} = \text{desiredOutput} - \text{Output}$ , Assume  $\text{desiredOutput} = 1$

$\text{LastNeuronCorrect} = \text{LastNeuronError} \cdot \text{Output} \cdot (1 - \text{Output})$

Now, we correct the neurons based on the data from above:

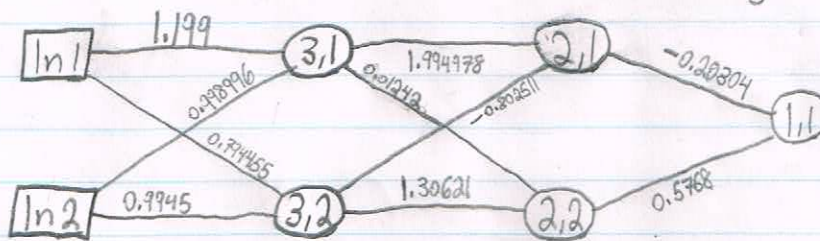
Using  $\text{Correct}_j^k = \text{out}_j^k (1 - \text{out}_j^k) \cdot \sum w_{ij} \cdot \text{Correct}_i^{k+1}$





Find the new weights using:  $W_{ij} = W_{ij} + \eta \cdot \text{Out}_j^k \cdot \text{Correct}_i^{k+1}$  given  $\eta = \alpha = 1$

So, we have the following network after 1 iteration of training:





⑤ We have trained a neural network on a dataset that consists of a train set & a test set. We decide to split the train set into 10 folds & do 10-fold Cross-validation on the training before finding the final test accuracy by retraining the best network. We get (in terms of cross validation accuracy) on the whole training set & then find the testing accuracy. After doing so, we have got these results:

Training accuracy: 81%

10-fold cross-validation accuracy: 78%

Test accuracy: 59%

Explain why the test accuracy is lower than training or validation accuracy. What do you think happened here? Which of these results is the correct result? What do you think should be done in this situation (or maybe this is not a problem in the first place)?

Why the test accuracy is lower:

The testing accuracy is lower due to overfitting. When the training accuracy is greater than the testing, the model was overfitted. This is akin to having memorized an exact pattern in the dataset without retaining the underlying trend in the dataset. Thus, we end up making an overly elaborate model that fits the training set perfectly but doesn't retain the trend (or pattern) behind the data. Therefore, when testing the data on the testing set, we find the accuracy is low since we haven't actually learned the underlying rules behind the data. This is an error due to variance.

The correct result is ultimately the testing accuracy. Although the training data is closely modeled by the program, it didn't pick up on the true trend, meaning it "memorized" the dataset but didn't "learn" anything, which defeats the point of the task.

I think that the overfitting can be resolved by reducing the variance in the sampling data, resulting in higher precision estimates in the dataset. This would reduce the rate of false treatments in the data & also prevent false variables from being in the model. We can also change the regularization parameters to increase weight decay ( $\lambda$ ), which will smooth out the model's learning rate to form a solid (rather than squiggly) curve to model the data's underlying trend with. Lastly, we can validate the dataset better. We currently are using a simple 10-fold cross-validation model (1 of 10 datasets used to validate), but we can ramp up the validation by using a 3-fold or leave-one-out model (uses 1/3 & 1/2 to validate). This would slow the program since we would be wasting  $\frac{1}{3}$  or  $\frac{1}{2}$  data for validation, but it would greatly help to prevent overfitting like we have in the current model.