

1.

B) With only a single hidden unit, the neural network produced a result which is logistic regression with a multiclassifier. As seen in previous assignments, this ultimately fits a straight line through the data.

F) The function with 4 or less hidden units in one layer lacks the complexity to fit the data to receive the results of Gaussian Bayes Discriminant Analysis. For similar results we would need to increase the complexity by adding more hidden layers and/or more hidden units.

G) The various different boundaries are due to the random initialization in weights and the shape of the data. Additionally, the decision boundaries are slightly different, but nonetheless they all classify about the same points in the region.

2.

A) I don't know

B)

$$a) \frac{\partial C}{\partial w_j} = \sum_i \frac{\partial C}{\partial z_i} \frac{\partial z_i}{\partial w_j}$$

Note: $\frac{\partial z}{\partial w} = \frac{\partial}{\partial w} (hw + w_0) = h$

$$\begin{aligned} &= \sum_i \frac{\partial C}{\partial z_i} h \\ &= \left[\frac{\partial C}{\partial z} \right]_{k,i} [H]_{k,i} \\ &= \left[\frac{\partial C}{\partial z} \right]_{k,i} [H^T]_{i,k} \\ \frac{\partial C}{\partial w} &= H^T \frac{\partial C}{\partial z} \end{aligned}$$

C)

$$a) \frac{\partial C}{\partial w_0} = \sum_i \frac{\partial C}{\partial z_i} \frac{\partial z_i}{\partial w_0}$$

Note: $\frac{\partial z}{\partial w_0} = \frac{\partial}{\partial w_0} (hw + w_0) = 1$

$$\begin{aligned} &= \sum_i \frac{\partial C}{\partial z_i} 1 \\ &= \left[\frac{\partial C}{\partial z} \right]_{k,i} 1_k \\ \frac{\partial C}{\partial w_0} &= \vec{1} \left[\frac{\partial C}{\partial z} \right] \end{aligned}$$

must be changed for matrix multiplication (reorganization)

D)

$$ad) \frac{\partial C}{\partial h_k} = \sum_i \frac{\partial C}{\partial z_i} \frac{\partial z_i}{\partial h_k}$$

$$\text{Note: } \frac{\partial z_i}{\partial h} = \frac{\partial}{\partial h} (hW + w_0) = W$$

$$= \sum_i \frac{\partial C}{\partial z_i} w_i$$

$$\frac{\partial C}{\partial H} = \frac{\partial C}{\partial Z} W^T \quad \leftarrow \text{Transpose due to } w \text{ being (hidden units, output units)} \text{ and } \frac{\partial C}{\partial Z} \text{ being the same shape as well!}$$

E)

$$de) \frac{\partial C}{\partial u} = \sum_n \frac{\partial C}{\partial h_k} \frac{\partial h_k}{\partial u_k}$$

$$\text{Note: } \frac{\partial h_k}{\partial u} = (h_k)(1-h_k)$$

$$= \sum_n \frac{\partial C}{\partial h_k} [(h_k)(1-h_k)]$$

$$= \left[\frac{\partial C}{\partial h} \right]_{n \times k} [h_k(1-h_k)]_{n \times k}$$

$$\left[\frac{\partial C}{\partial u} \right]_{n \times k} = (H_{n \times k})(1-H_{n \times k}) \left[\frac{\partial C}{\partial H} \right]_{n \times k}$$

F)

$$F) \frac{\partial C}{\partial v_n} = \sum_m \frac{\partial C}{\partial u_x} \frac{\partial u_x}{\partial v_n}$$

$$\Rightarrow \text{Note } \frac{\partial u}{\partial v} = X$$

$$= \sum_m \frac{\partial C}{\partial u_x} x_n$$

$$= \left[\frac{\partial C}{\partial u} \right]_{n \times k} X_{n \times m}$$

$$= \left[\frac{\partial C}{\partial u} \right]_{n \times k} [X^T]_{m \times n}$$

$$\frac{\partial C}{\partial v} = [X^T]_{m \times n} \left[\frac{\partial C}{\partial u} \right]_{n \times k}$$

← [matrix transpose property]

← [reorganize]

G)

$$\begin{aligned} G) \frac{\partial C}{\partial v_{0n}} &= \sum_m \frac{\partial C}{\partial u_x} \frac{\partial u_x}{\partial v_{0n}} \\ \text{Note: } \frac{\partial u}{\partial v_0} &= 1 \\ &= \sum_m \frac{\partial C}{\partial u_x} 1 \\ &= \left[\frac{\partial C}{\partial u_x} \right]_{n \times k} \vec{1}_{1 \times m} \\ \frac{\partial C}{\partial v_0} &= \vec{1}_{1 \times n} \left[\frac{\partial C}{\partial u_x} \right]_{n \times k} \end{aligned}$$

H) `dCdw0 = np.sum(dCdZ, axis=0).reshape(1,-1)`

3.

B) This is because stochastic gradient descent does more iterations of learning in our example (i.e. calculate the gradient/fix the cost more times although same number of epochs).

D) This is due to the learning rate being standardized based on the batch size. Thus, it doesn't matter how much big your dataset is, the learning rate will adjust accordingly.