

1. We are given the following:

$$E(w, v | X) = - \sum_{t=1}^N [r^t \log y^t + (1-r^t) \log (1-y^t)] + \sum_{h=1}^H \|w_h\|_2^2$$

$$y^t = \text{Sigmoid} \left(\sum_{h=1}^H v_h z_h^t + v_0 \right)$$

$$z_h^t = \text{ReLU}(w_h^T + w_{h0})$$

H1 Hint: Given $\text{ReLU}(f(x)) = \max(0, f(x))$ the derivative of $\text{ReLU}(f(x))$ is $f'(x)$ if $f(x) > 0$ and 0 otherwise.

H2 Given $y = \text{Sigmoid}(\alpha) = 1/(1+e^{-\alpha})$, $\frac{\partial y}{\partial \alpha} = y(1-y)$

$v_h = v_h + \Delta v_h$ where $\Delta v_h = -\eta \frac{\partial E}{\partial v_h}$

We find $\frac{\partial E}{\partial v_h}$

$$\frac{\partial E}{\partial v_h} = \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial \alpha^t} \frac{\partial \alpha^t}{\partial v_h}$$

We find each partial derivative in turn.

① $\frac{\partial E}{\partial y^t} = - \sum_{t=1}^N \left[\frac{r^t}{y^t} - \frac{1-r^t}{1-y^t} \right]$ (Note: the last term $\sum_{h=1}^H \|w_h\|_2^2$ has $\frac{\partial E}{\partial y^t} = 0$)

② $\frac{\partial y^t}{\partial \alpha^t} = y^t(1-y^t)$ by H2 above

Since $y^t = \text{sigmoid}(\sum_{h=1}^H V_h Z_h^t + V_0)$ we know
 $\alpha = \sum_{h=1}^H V_h Z_h^t + V_0$

So,

$$\textcircled{3} \quad \frac{\partial \alpha^t}{\partial V_h} = Z_h^t$$

$$\text{Thus } \Delta V_h = -\eta \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial \alpha^t} \frac{\partial \alpha^t}{\partial V_h}$$

$$= \eta \sum_{t=1}^N \frac{r^t}{y^t} - \frac{1-r^t}{1-y^t} \cdot y^t (1-y^t) Z_h^t \quad (\text{by } \textcircled{1} \textcircled{2} \text{ and } \textcircled{3})$$

$$= \eta \sum_{t=1}^N \frac{r^t (1-y^t) - y^t (1-r^t)}{y^t (1-y^t)} \cdot y^t (1-y^t) Z_h^t$$

$$= \eta \sum_{t=1}^N (r^t - r^t y^t - y^t + r^t y^t) Z_h^t$$

$$= \eta \sum_{t=1}^N (r^t - y^t) Z_h^t$$

$$\text{Thus, } V_h = V_h + \Delta V_h = V_h + -\eta \frac{\partial E}{\partial V_h}$$

$$= V_h + \eta \sum_{t=1}^N (r^t - y^t) Z_h^t$$

$$V_0 = V_0 + \Delta V_0 \text{ where } \Delta V_0 = -\eta \frac{\partial E}{\partial V_0}$$

We find $\frac{\partial E}{\partial V_0}$

$$\frac{\partial E}{\partial V_0} = \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial \alpha^t} \frac{\partial \alpha^t}{\partial V_0}$$

We note that $\frac{\partial E}{\partial y^t}$ and $\frac{\partial y^t}{\partial \alpha^t}$ have been derived in ① and ② above.

$$\textcircled{4} \quad \frac{\partial \alpha^t}{\partial V_0} = 1 \quad \text{Since } \alpha = \sum_{h=1}^H V_h Z_h^t + V_0$$

$$\text{Thus, } \Delta V_0 = -\eta \frac{\partial E}{\partial V_0} \frac{\partial y^t}{\partial \alpha^t} \frac{\partial \alpha^t}{\partial V_h}$$

$$= \eta \sum_{t=1}^N (r^t - y^t) \quad (\text{by } \textcircled{1}, \textcircled{2}, \text{ and } \textcircled{4})$$

$$\text{Thus, } V_0 = V_0 + \Delta V_0 = V_0 + -\eta \frac{\partial E}{\partial V_0}$$

$$= V_0 + \eta \sum_{t=1}^N (r^t - y^t)$$

$$\underline{W_h = W_h + \Delta W_h} \text{ where } \Delta W_h = -\eta \frac{\partial E}{\partial W_h}$$

We have $z_h^t = \text{ReLU}(\Delta)$ where $\Delta = W_h^T x^t + w_0$
and $y^t = \text{Sigmoid}(\alpha)$ where $\alpha = v_h z_h^t + v_0$

We want to find $\frac{\partial E}{\partial W_h}$ for both terms in $E(W, v | X)$.

For the first term:

$$\frac{\partial E}{\partial W_h} = \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial \alpha^t} \frac{\partial \alpha^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial \Delta} \frac{\partial \Delta}{\partial W_h}$$

We have already calculated $\frac{\partial E}{\partial y^t}$ and $\frac{\partial y^t}{\partial \alpha^t}$ in ① and ②

⑤ We have $\frac{\partial \alpha^t}{\partial z_h^t} = v_h$

⑥ $\frac{\partial z_h^t}{\partial \Delta}$ represents the conditional on the piecewise ReLU function. That is if $W_h^T x^t + w_0 > 0$ then ReLU will return $f'(x)$ otherwise it will return 0.

⑦ Also, $\frac{\partial \Delta}{\partial W_h} = x^t$

For the second term we want to find $\frac{\partial E}{\partial W_h}$.

We expand the sum.

$$\begin{aligned}\sum_{h=1}^H \|w_h\|_2^2 &= \|w_1\|_2^2 + \|w_2\|_2^2 + \dots + \|w_H\|_2^2 \\ &= w_1^2 + w_2^2 + \dots + w_H^2\end{aligned}$$

⑧ Taking the derivative $\frac{\partial E}{\partial w_h}$ we get

$$\begin{aligned}&= 2w_1 + 2w_2 + \dots + 2w_H \\ &= \sum_{h=1}^H 2w_h\end{aligned}$$

Thus, by ①②⑤⑥⑦⑧ we get

$$\Delta w_h = -\eta \frac{\partial E}{\partial w_h} = \eta \sum_{t=1}^T (r^t - y^t) v_h x^t + \sum 2w_h \quad \text{if } w_h^T x^t + w_0 > 0$$

or
0 otherwise

with complete update equation being

$w_h = w_h + \Delta w_h$ where Δw_h is as defined above.

$w_0 = w_0 + \Delta w_0$ is the same as above except

$$\frac{\partial \mathcal{L}}{\partial w_0} = 1 \quad \text{and} \quad \frac{\partial E}{\partial w_h} \text{ of the second term is } 0.$$

Thus,

$$w_0 := w_0 + \Delta w_0 = w_0 + -\eta \frac{\partial E}{\partial w_0}$$

$$= w_0 + \eta \sum_{t=1}^N (r^t - y^t) v_h x^t \text{ if } w_h^T x^t + w_0 > 0$$

otherwise it is just

$$= w_0 + 0 = w_0$$

2a. We determine the vector of coefficients for Figure 2.

The equation of the line which separates the two classes is

$$x_1 = 2 \Rightarrow x_1 - 2 = 0 \Rightarrow 0x_2 + x_1 - 2 = 0$$

Hence the options for the activation functions are

(1) $y = S(0x_2 + x_1 - 2)$

or

(2) $y = S(0x_2 - x_1 + 2)$

By the figure we want the activation function to return 0 for $\begin{pmatrix} 0 \\ 0 \end{pmatrix}_{x_1, x_2}$. Thus

$$S(0(0) + (0) - 2) = S(-2) = 0 \quad \text{Test (1)}$$

$$S(0(0) - (0) + 2) = S(2) = 1 \quad \text{Test (2)}$$

We see that (1) is the correct activation function.

so, $y = S(0x_2 + x_1 - 2)$

$\vec{w} \Rightarrow \vec{w} = [-2, 1, 0]^T$ where $\vec{w} = [w_0, w_1, w_2]^T$

2a. (cont'd) We determine the vector of coefficients for Figure 3.

The equation of the line which separates the two classes is

$$x_1 - 4x_2 = 0 \Rightarrow -4x_2 + x_1 + 0 = 0$$

Hence, the options for the activation functions are

$$(3) \quad y = S(-4x_2 + x_1 + 0)$$

or

$$(4) \quad y = S(4x_2 - x_1 + 0)$$

By the figure we want the activation function to return a 1 for $(1, 0)$. Thus,

$$S(-4(0) + 1 + 0) = S(1) = 1 \quad \text{Test (3)}$$

$$S(4(0) - 1 + 0) = S(-1) = 0 \quad \text{Test (4)}$$

We see that (3) is the correct activation function.

$$\text{So, } y = S(-4x_2 + x_1 + 0)$$

$$\Rightarrow \vec{w} = [0, 1, -4]^T \text{ where } \vec{w} = [w_0, w_1, w_2]^T$$

2b. We determine coefficients W and b in the two-layer perceptron in Figure 4.60 to recognize the shaded region in Figure 5.

We see that the shaded region in Figure 5 is equal to the intersection of the regions of figures 2 and 3.

Thus, the equations of the lines which represent the shaded areas are

$$Z_1 = 0x_2 + x_1 - 2$$

Q2d

$$Z_2 = -4X_2 + X_1 + 0$$

So $y = 1$ when

$$Z = 1$$

And

$$z_2 = 1$$

And so we choose a middle point between 1 and 2.

Since 1.5 splits the difference we have

$$z_2 + z_1 = 1.5 \Rightarrow z_2 + z_1 - 1.5 = 0$$

Q. 1. In (1) the subject (2.1) is a noun phrase.

2b. (cont'd) Thus, the options for the activation functions are

$$(5) \quad S(z_2 + z_1 - 1.5)$$

$$(6) \quad S(-z_2 - z_1 + 1.5)$$

We want the activation function to return a 1 for (1,1) and a 0 for (0,0). We see that (5) does both of these while (6) does not!

$$S(1+1-1.5) = S(.5) = 1$$

Tests for 5

$$S(0+0-1.5) = S(-1.5) = 0$$

$$S(-1-1+1.5) = S(-.5) = 0$$

Tests for 6

$$S(-0-0+1.5) = S(1.5) = 1$$

$$\text{Thus: } y = S(z_2 + z_1 - 1.5)$$

$$\Rightarrow \vec{v} = [-1.5, 1, 1]^T \text{ where } \vec{v} = [v_0, v_1, v_2]^T$$

and

$$W = \begin{bmatrix} -2 & 1 & 0 \\ 0 & 1 & -4 \end{bmatrix}^T$$

3.a) Programming Piece

b)

Output from hw3:

```
Validation accuracy for 4 hidden units is 0.826
Validation accuracy for 16 hidden units is 0.915
Validation accuracy for 20 hidden units is 0.912
Validation accuracy for 24 hidden units is 0.919
Validation accuracy for 32 hidden units is 0.923
Validation accuracy for 48 hidden units is 0.900
Test accuracy with 32 hidden units is 0.912
```

Output from
my code run which reports the validation accuracy by the
number of hidden units.

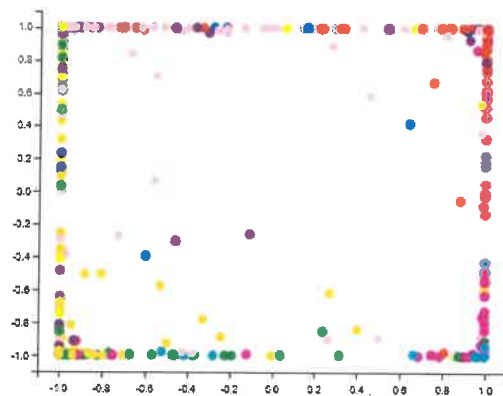
We should use 32 hidden units since this
had the highest validation accuracy.

The accuracy on the test set with 32 hidden
units is 0.912.

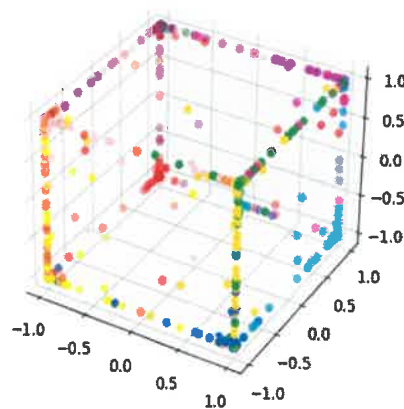
We see that as the number of hidden units increases
the validation accuracy increases, up to a certain
point. This is due to hidden units adding more
dimensions which allows for better classification. At
a certain point, though, adding more hidden units will
result in an overtrained model. Also, more hidden units
implies a more complex model so gradient descent could
find a less than optimal local minimum.

3. c)

Visualization with 2 hidden units:



Visualization with 3 hidden units:



With two hidden units we only have two dimensions to classify the data. This results in poor separation of the classes. With 3 hidden units there is another dimension which the classifier can use for separation. Thus, the separation of classes is better in the 3d plot, corresponding to 3 hidden units, than in the 2d plot, corresponding to two hidden units.