

Andrew Elysee (aelys2)

High-level discussion with Max Han (mhan36)

CS446

HW2

1.

a.

- i. Since $P(x_d = 1) = q_d$, then the probability of getting a vector x given q , is

$$P(x|q) = \prod_{d=1}^D q_d^{x_d} (1 - q_d)^{1-x_d}$$

This makes sense because if $x_d = 1$, then the probability of that happening is just q_d , but if $x_d = 0$ then the probability of that is $1 - q_d$ so we take the product of all these probabilities to get $P(x|q)$.

- ii. Since we know that each vector $x^{(i)}$ was drawn from one of our Bernoulli distributions with parameters $p^{(k)}$, the expression must be the sum of the products of the probabilities that the k th distribution was chosen and the probability of $x^{(i)}$ given $p^{(k)}$. So, we get

$$P(x^{(i)}|p, \pi) = \sum_{k=1}^K \pi_k P(x^{(i)}|p^{(k)})$$

iii.

$$P(D|\pi, p) = \prod_{i=1}^D P(x^{(i)}|p, \pi)$$

$$\log P(D|\pi, p) = \sum_{i=1}^D \log P(x^{(i)}|p, \pi)$$

b.

i.

$$P(z^{(i)}|\pi) = \prod_{k=1}^K \pi_k^{z_k^{(i)}}$$

ii.

$$P(x^{(i)}|z^{(i)}, p, \pi) = \prod_{k=1}^K P(x^{(i)}|p^{(k)})^{z_k^{(i)}}$$

iii.

$$P(Z, D|p, \pi) = \prod_{i=1}^D P(z^{(i)}, x^{(i)}|p, \pi) = \prod_{i=1}^D P(x^{(i)}|z^{(i)}, p, \pi) P(z^{(i)}|\pi)$$

$$P(Z, D|p, \pi) = \prod_{i=1}^D \left[\prod_{k=1}^K P(x^{(i)}|p^{(k)})^{z_k^{(i)}} \right] \left[\prod_{k=1}^K \pi_k^{z_k^{(i)}} \right]$$

iv.

$$\begin{aligned}\eta(z_k^{(i)}) &= E[z_k^{(i)} | x^{(i)}, p, \pi] \\ &= P(z_k^{(i)} = 1 | x^{(i)}, p, \pi)\end{aligned}$$

Using Bayes rule,

$$\begin{aligned}&= \frac{P(x^{(i)} | z_k^{(i)}, p, \pi) P(z_k^{(i)} = 1 | p, \pi)}{\sum_{j=1}^K P(x^{(i)} | z_j^{(i)}, p, \pi) P(z_j^{(i)} = 1 | p, \pi)} \\ &= \frac{\pi_k \prod_{d=1}^D (p_d^{(k)})^{x_d^{(i)}} (1 - p_d^{(k)})^{1-x_d^{(i)}}}{\sum_{j=1}^K \pi_j \prod_{d=1}^D (p_d^{(j)})^{x_d^{(i)}} (1 - p_d^{(j)})^{1-x_d^{(i)}}}\end{aligned}$$

v.

$$\begin{aligned}E[\log P(Z, D | \tilde{p}, \tilde{\pi}) | D, p, \pi] &= \log P(Z, D | \tilde{p}, \tilde{\pi}) \\ &= \sum_{i=1}^D \left[\sum_{k=1}^K z_k^{(i)} \log P(x^{(i)} | \tilde{p}^{(k)}) \right] + \left[\sum_{k=1}^K z_k^{(i)} \log \tilde{\pi}_k \right] \\ &= \sum_{i=1}^D \left[\sum_{k=1}^K z_k^{(i)} (\log P(x^{(i)} | \tilde{p}^{(k)}) + \log \tilde{\pi}_k) \right] \\ &= \sum_{i=1}^D \sum_{k=1}^K z_k^{(i)} \left[\log \prod_{d=1}^D (\tilde{p}_d^{(k)})^{x_d^{(i)}} (1 - \tilde{p}_d^{(k)})^{1-x_d^{(i)}} + \log \tilde{\pi}_k \right] \\ &= \sum_{i=1}^D \sum_{k=1}^K z_k^{(i)} \left[\sum_{d=1}^D (x_d^{(i)} \log(\tilde{p}_d^{(k)}) + (1 - x_d^{(i)}) \log(1 - \tilde{p}_d^{(k)})) + \log \tilde{\pi}_k \right]\end{aligned}$$

Since there is an imaginary expected value sign around the whole expression, $z_k^{(i)} = E[z_k^{(i)}] = \eta(z_k^{(i)})$, so we get:

$$\begin{aligned}&= \sum_{i=1}^D \sum_{k=1}^K \eta(z_k^{(i)}) \left[\sum_{d=1}^D (x_d^{(i)} \log(\tilde{p}_d^{(k)}) + (1 - x_d^{(i)}) \log(1 - \tilde{p}_d^{(k)})) \right. \\ &\quad \left. + \log \tilde{\pi}_k \right]\end{aligned}$$

c.

- i. To find the value of \tilde{p} that maximizes the E step, we must set the derivative of the E step relative to \tilde{p} to 0.

$$\begin{aligned}&\frac{\partial}{\partial \tilde{p}_d^{(k)}} E[\log P(Z, D | \tilde{p}, \tilde{\pi}) | D, p, \pi] = 0 \\ &\sum_{i=1}^D \eta(z_k^{(i)}) \left[\frac{x_d^{(i)}}{\tilde{p}_d^{(k)}} + \frac{1 - x_d^{(i)}}{1 - \tilde{p}_d^{(k)}} \right] = 0 \\ &\sum_{i=1}^D \eta(z_k^{(i)}) [x_d^{(i)}(1 - \tilde{p}_d^{(k)}) + \tilde{p}_d^{(k)}(1 - x_d^{(i)})] = 0 * \tilde{p}_d^{(k)}(1 - \tilde{p}_d^{(k)}) = 0\end{aligned}$$

$$\sum_{i=1}^D \eta(z_k^{(i)}) [x_d^{(i)} - \tilde{p}_d^{(k)}] = \sum_{i=1}^D \eta(z_k^{(i)}) x_d^{(i)} - \sum_{i=1}^D \eta(z_k^{(i)}) (\tilde{p}_d^{(k)}) = 0$$

Rearranging to solve for $\tilde{p}_d^{(k)}$ we get:

$$\tilde{p}_d^{(k)} = \frac{\sum_{i=1}^D \eta(z_k^{(i)}) x_d^{(i)}}{\sum_{i=1}^D \eta(z_k^{(i)})} = \frac{\sum_{i=1}^D \eta(z_k^{(i)}) x_d^{(i)}}{N_k}$$

- ii. First, we will set up a Lagrangian, with λ to enforce the constraint that $\sum \pi_k = 1$. Then take the derivative, and set it equal to 0.

$$L(\tilde{\pi}, \lambda) = - \sum_{i=1}^D \sum_{k=1}^K \eta(z_k^{(i)}) \log \tilde{\pi}_k + \lambda \left(\sum_{k=1}^K \tilde{\pi}_k - 1 \right)$$

$$\frac{\partial}{\partial \tilde{\pi}_k} L(\tilde{\pi}, \lambda) = - \sum_{i=1}^D \frac{\eta(z_k^{(i)})}{\tilde{\pi}_k} + \lambda = 0$$

$$\sum_{i=1}^D \frac{\eta(z_k^{(i)})}{\tilde{\pi}_k} = \lambda$$

$$\tilde{\pi}_k = \frac{\sum_{i=1}^D \eta(z_k^{(i)})}{\lambda} = \frac{N_k}{\lambda}$$

First we see that our rearranged constraint is

$$\sum_{k=1}^K \tilde{\pi}_k - 1 = \sum_{k=1}^K \frac{N_k}{\lambda} - 1 = 0 \rightarrow \frac{1}{\lambda} \sum_{k=1}^K N_k = 1 \rightarrow \sum_{k=1}^K N_k = \lambda \text{ so,}$$

solving for λ and taking the derivative, we that

$$L(\lambda) = - \sum_{i=1}^D \sum_{k=1}^K \eta(z_k^{(i)}) \log \frac{N_k}{\lambda} + \lambda \left(\sum_{k=1}^K \frac{N_k}{\lambda} - 1 \right)$$

$$= - \sum_{i=1}^D \sum_{k=1}^K \eta(z_k^{(i)}) (\log N_k - \log \lambda) + \lambda \left(\sum_{k=1}^K \frac{N_k}{\lambda} - 1 \right)$$

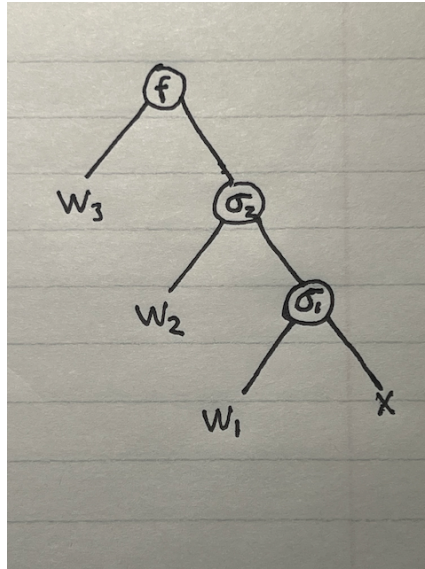
$$\frac{\partial}{\partial \lambda} L(\lambda) = \frac{1}{\lambda} \sum_{i=1}^D \sum_{k=1}^K \eta(z_k^{(i)}) - 1 = 0 \rightarrow \frac{1}{\lambda} \sum_{i=1}^D \sum_{k=1}^K \eta(z_k^{(i)}) = 1$$

$$\lambda = \sum_{i=1}^D \sum_{k=1}^K \eta(z_k^{(i)}) = \sum_{k=1}^K N_k$$

This gives us the final piece for our update for $\tilde{\pi}_k$:

$$\tilde{\pi}_k = \frac{N_k}{\lambda} = \frac{N_k}{\sum_{k'=1}^K N_{k'}}$$

a.



b. $\sigma_1(u) = \frac{1}{1+\exp(-u)} \rightarrow \frac{\partial \sigma_1}{\partial u} = -\frac{1}{(1+\exp(-u))^2} * -\exp(-u) = \frac{\exp(-u)}{(1+\exp(-u))^2}$

This is the partial derivative using only u , but with some manipulation we get:

$$\frac{\partial \sigma_1}{\partial u} = \frac{\exp(-u)}{(1+\exp(-u))^2} = \frac{1}{1+\exp(-u)} * \frac{\exp(-u)}{1+\exp(-u)} = \frac{1}{1+\exp(-u)} * \left(1 - \frac{1}{1+\exp(-u)}\right)$$

$$\frac{\partial \sigma_1}{\partial u} = \sigma_1(u)(1 - \sigma_1(u))$$

c. A forward pass refers to the calculation of the outputs from the input layers, while a backwards pass is the calculation of the change in weights using some gradient descent method.

d. $\frac{\partial f}{\partial w_3} = \sigma_2(w_2 \sigma_1(w_1 x)) = \sigma_2(x_2)$

To make the computation of this derivative easy, we should save the output of $\sigma_2(x_2)$.

e. $\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial \sigma_2} \frac{\partial \sigma_2}{\partial w_2} = \sigma_1(w_1 x) w_3 \sigma_2(w_2 \sigma_1(w_1 x)) (1 - \sigma_2(w_2 \sigma_1(w_1 x))) =$
 $w_3 \sigma_1(x_1) \sigma_2(w_2 \sigma_1(x_1)) (1 - \sigma_2(w_2 \sigma_1(x_1))) = w_3 \sigma_1(x_1) \sigma_2(x_2) (1 - \sigma_2(x_2))$

To make this calculation easier, we should save $\sigma_1(x_1)$ and $\sigma_2(x_2)$, assuming that w_3 is implied to be saved.

f. $\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial \sigma_2} \frac{\partial \sigma_2}{\partial \sigma_1} \frac{\partial \sigma_1}{\partial w_1} = w_3 w_2 x \sigma_1(w_1 x) (1 - \sigma_1(w_1 x)) \sigma_2(w_2 \sigma_1(w_1 x)) (1 - \sigma_2(w_2 \sigma_1(w_1 x))) =$
 $w_3 w_2 x \sigma_1(x_1) (1 - \sigma_1(x_1)) \sigma_2(x_2) (1 - \sigma_2(x_2)) = \frac{\partial f}{\partial w_2} * w_2 x (1 - \sigma_1(x_1))$

We should save $\sigma_1(x_1)$ and $\sigma_2(x_2)$ from the forward pass. To obtain the result as early as possible and reuse as many results as possible we should do $\frac{\partial f}{\partial w_3}, \frac{\partial f}{\partial w_2},$

and finally $\frac{\partial f}{\partial w_1}$, because we can reuse the results of $\frac{\partial f}{\partial w_3}$ in $\frac{\partial f}{\partial w_2}$ and the results of

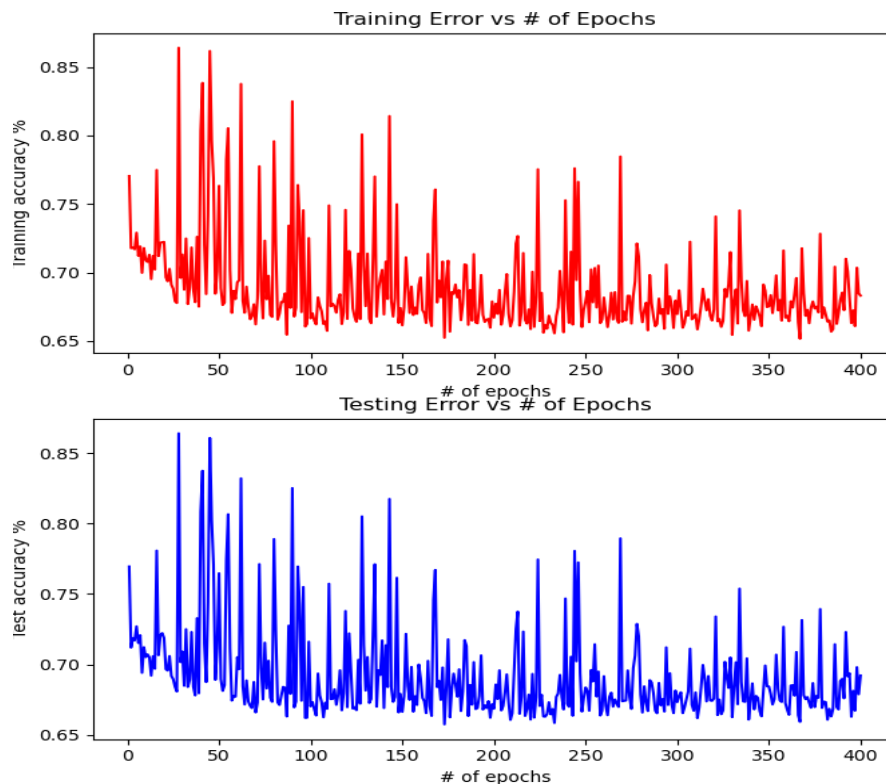
$\frac{\partial f}{\partial w_2}$ in $\frac{\partial f}{\partial w_1}$ like I showed in the equation above. This is reverse of the forward pass.

- g. After the first layer the output dimension is 20 X 24 X 24. After the max pooling layer, we get an output of dimension 20 X 12 X 12.
- h. If we undo the max pooling on the 4x4 output, we see that the output of the second convolution layer is 8x8. From this, we can derive that our filters would be of size 5x5 with a stride of 1, because $(\text{input size}) * (\text{kernel size}) + 1 = 12 - 5 + 1 = 8$. The channel dimensions is 50.
- 3.
- a. Submitted on gradescope
 - b. Submitted on gradescope
 - c. In my graphs the Testing and training errors are nearly indistinguishable but looking closely I think that as C increases, the gap between the test and training error seems to increase, this is because the low complexity of the model helps to prevent overfitting.

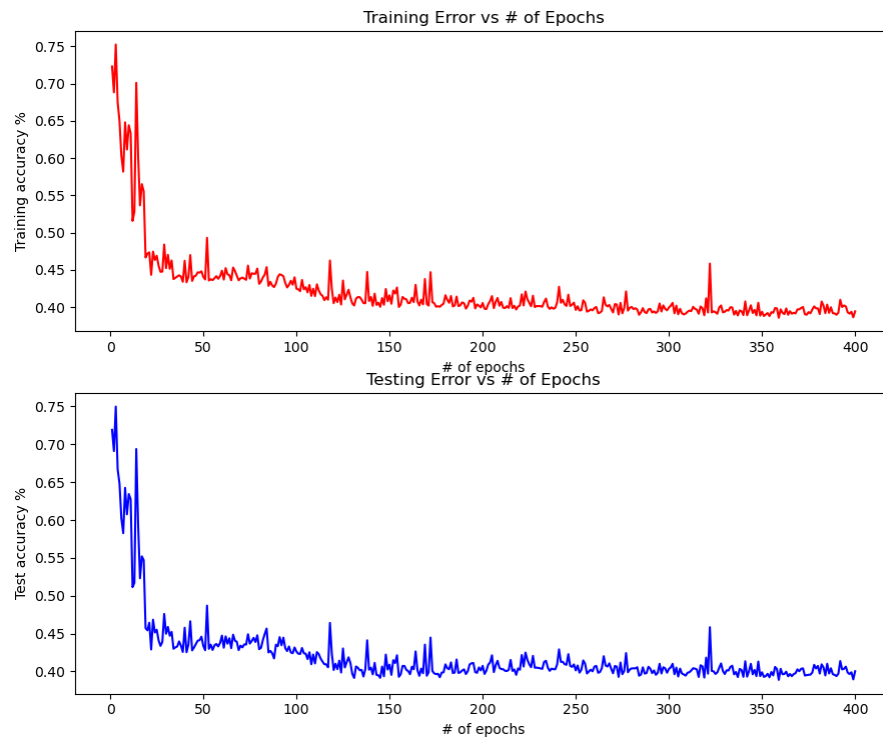
Plots are below.

Also note that the y-axis label is labeled ____ accuracy but should be ____ Error

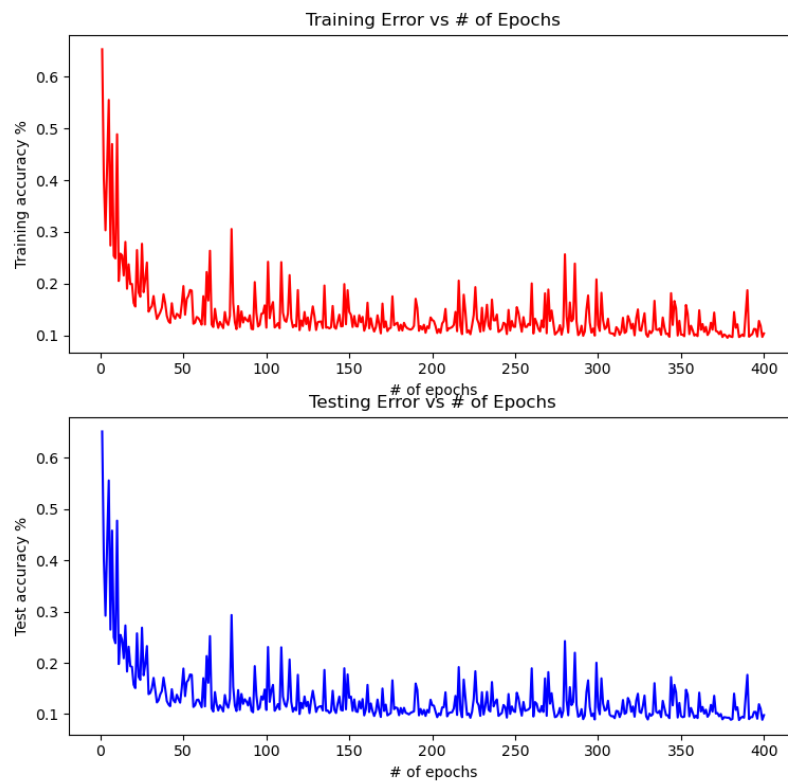
Plots for C = 1



Plots for $C = 2$



Plots for $C = 4$



- d. For $C=64$, the testing and training errors curves are significantly smoother this also start off very accurate after the first epoch. Additionally, the gap between

the testing and training sets is much more noticeable as our training set gets a lot closer to 0 than our testing error, this could just be because of the scaling of the axis though. In hindsight, I should have overlayed the graphs.

Plots for $C = 64$

