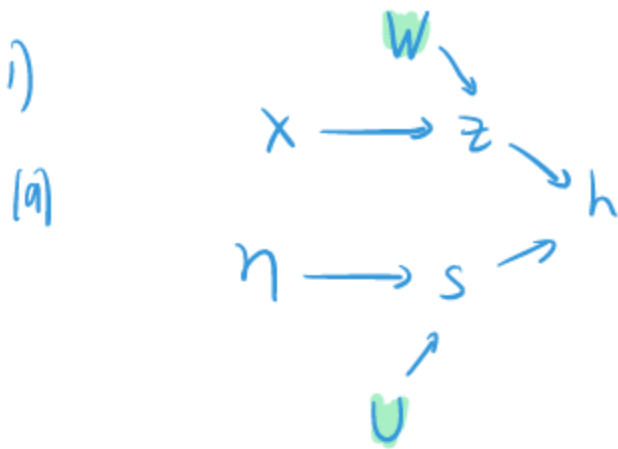


1)



$$\bar{L} = 1$$

$$\bar{h} = \bar{y} \frac{dy}{dh}$$

$$\bar{z} = \bar{h} \frac{dh}{dz}$$

$$\bar{y} = \bar{L} \frac{dL}{dy}$$

$$= \bar{y} v$$

$$= \bar{h} 100\sigma(s)$$

$$= \bar{h} \sigma(s)$$

$$= \bar{L} (y-t) \quad \bar{x} = \bar{y} \frac{dy}{dx} + \bar{z} \frac{dz}{dx}$$

$$= (y-t)$$

$$= \bar{y} r + \bar{z} W$$

$$\bar{s} = \bar{h} \frac{dh}{ds}$$

$$= \bar{h} \sigma'(s) \circ z$$

$$\bar{x} = \bar{z} \frac{dz}{dx} \quad n = \bar{s} \frac{ds}{dn}$$

$$= \bar{z} W$$

$$= \bar{s} U$$

2)

2)

(a)

$$L(\theta, \pi) = p(x, c | \theta, \pi) \\ = p(c | \pi) \prod_{j=1}^{784} p(x_j | c, \theta_{jc})$$

$$= p(c | \pi) \prod_{j=1}^{784} p(x_j | c, \theta_{jc})$$

$$= \prod_{c=0}^1 \pi_c^{t_c} \prod_{j=1}^{784} \theta_{jc}^{x_j^{(n)}} (1 - \theta_{jc})^{(1-x_j)}$$

Assuming that there are k samples in the dataset, we have:

$$L(\theta, \pi | \{t^{(i)}, x^{(i)}\}_{i=1}^k) = \prod_{i=1}^k \prod_{c=0}^1 \pi_c^{t_c^{(i)}} \prod_{j=1}^{784} \theta_{jc}^{x_j^{(i)}} (1 - \theta_{jc})^{(1-x_j^{(i)})}$$

as the likelihood function.

Taking the log-likelihood, we have:

$$\begin{aligned}
 \ell(\theta, \pi | \{t^{(i)}, x^{(i)}\}_{i=1}^k) &= \log \left[\prod_{i=1}^k \prod_{c=0}^q \pi_c^{t_c^{(i)}} \prod_{j=1}^M \theta_{jc}^{x_j^{(i)}} (1-\theta_{jc})^{(1-x_j^{(i)})} \right] \\
 &= \sum_{i=1}^k \sum_{c=0}^q \left(t_c^{(i)} \log \pi_c + \left(\sum_{j=1}^M x_j^{(i)} \log \theta_{jc} + \sum_{j=1}^M (1-x_j^{(i)}) \log (1-\theta_{jc}) \right) \right) \\
 &= \left[\sum_{i=1}^k \sum_{c=0}^q t_c^{(i)} \log \pi_c \right] + \left[\sum_{i=1}^k \sum_{c=0}^q \left(\sum_{j=1}^M x_j^{(i)} \log \theta_{jc} + (1-x_j^{(i)}) \log (1-\theta_{jc}) \right) \right]
 \end{aligned}$$

taking the derivative w.r.t. θ_{jc} , we have

$$\begin{aligned}
 \rightarrow \frac{\partial}{\partial \theta_{jc}} \sum_{i=1}^k \sum_{c=0}^q \left(\sum_{j=1}^M x_j^{(i)} \log \theta_{jc} + (1-x_j^{(i)}) \log (1-\theta_{jc}) \right) \\
 = \sum_{i=1}^k \mathbb{I}(c^{(i)}=c) \left(\frac{x_j^{(i)}}{\theta_{jc}} - \frac{(1-x_j^{(i)})}{(1-\theta_{jc})} \right) \\
 = \sum_{i=1}^k t_c^{(i)} \left(\frac{x_j^{(i)}}{\theta_{jc}} - \frac{(1-x_j^{(i)})}{(1-\theta_{jc})} \right)
 \end{aligned}$$

$$= \sum_{i=1}^k t_c^{(i)} \left(\frac{x_j^{(i)}(1-\theta_{jc}) - (1-x_j^{(i)})\theta_{jc}}{\theta_{jc}(1-\theta_{jc})} \right)$$

$$= \sum_{i=1}^k t_c^{(i)} (x_j^{(i)}(1-\theta_{jc}) - (1-x_j^{(i)})\theta_{jc})$$

$$= \sum_{i=1}^k t_c^{(i)} (x_j^{(i)} - \cancel{x_j^{(i)}\theta_{jc}} - \theta_{jc} + \cancel{x_j^{(i)}\theta_{jc}})$$

$$= \sum_{i=1}^k t_c^{(i)} (x_j^{(i)} - \theta_{jc})$$

$$\therefore \sum_{i=1}^k t_c^{(i)} x_j^{(i)} - t_c^{(i)} \theta_{jc} = 0$$

\therefore MLE for θ_{jc} is:

$$\hat{\theta}_{jc} = \frac{\sum_{i=1}^k t_c^{(i)} x_j^{(i)}}{\sum_{i=1}^k t_c^{(i)}}$$

The likelihood function for π is as follows

$$l(\theta, \pi | \{t^{(i)}, x^{(i)}\}_{i=1}^k)$$

$$= \sum_{i=1}^k \sum_{c=0}^9 \left(t_c^{(i)} \log \pi_c + \left(\sum_{j=1}^{79} x_j^{(i)} \log \theta_{jc} + (1 - x_j^{(i)}) \log (1 - \theta_{jc}) \right) \right)$$

$$= \sum_{i=1}^k \left(\sum_{c=0}^9 t_c^{(i)} \log \pi_c + t_9^{(i)} \log \left(1 - \sum_{c=0}^8 \pi_c \right) + \right.$$

$$\left. \left(\sum_{c=0}^9 \sum_{j=1}^{79} x_j^{(i)} \log \theta_{jc} + (1 - x_j^{(i)}) \log (1 - \theta_{jc}) \right) \right) \# \text{ using hint 2}$$

$$\therefore \frac{\partial l}{\partial \pi_c} = \sum_{i=1}^k \left(\sum_{c=0}^9 \frac{t_c^{(i)}}{\pi_c} - \frac{t_9^{(i)}}{1 - \sum_{c=0}^8 \pi_c} \right)$$

$$= \sum_{i=1}^k \left(\frac{t_c^{(i)}}{\pi_c} - \frac{t_9^{(i)}}{\pi_9} \right)$$

\therefore setting $\frac{\partial \mathcal{L}}{\partial \pi_c} = 0$, we have

$$\sum_{i=1}^k \left(\frac{t_c^{(i)}}{\pi_c} - \frac{t_q^{(i)}}{\pi_q} \right) = 0$$

$$\sum_{i=1}^k \frac{t_c^{(i)}}{\pi_c} - \sum_{i=1}^k \frac{t_q^{(i)}}{\pi_q} = 0$$

$$\sum_{i=1}^k \frac{t_c^{(i)}}{\pi_c} = \sum_{i=1}^k \frac{t_q^{(i)}}{\pi_q}$$

$$\therefore \frac{\hat{\pi}_c}{\hat{\pi}_q} = \frac{\sum_{i=1}^k t_c^{(i)}}{\sum_{i=1}^k t_q^{(i)}}$$

Since $\sum_{i=0}^q \pi_i = 1$, we have $\frac{\pi_1}{\pi_q} + \frac{\pi_2}{\pi_q} + \dots + \frac{\pi_q}{\pi_q} = \frac{1}{\pi_q}$

$$\therefore \frac{1}{\hat{\pi}_q} = \frac{\sum_{i=1}^k t_c^{(i)}}{\sum_{i=1}^k t_q^{(i)}}$$

$$\hat{\pi}_q = \frac{\sum_{i=1}^k t_q^{(i)}}{\sum_{i=1}^k t_c^{(i)}}$$

Thus, if we were to use the same argument for each $\hat{\pi}_1, \hat{\pi}_2, \dots, \hat{\pi}_8$, we will obtain:

$$\hat{\pi}_j = \frac{\sum_{i=1}^k t_j^{(i)}}{k}$$

(b)

Using Bayes rule we have:

$$p(c|x) = \frac{p(c)p(x|c)}{\sum_{c'} p(c')p(x|c')}$$

$$\therefore p(t|x, \theta, \pi) = \frac{p(t|\pi) p(x|t, \theta)}{\sum_{t'} p(t'|\pi) p(x|t', \theta_{jt'})}$$

$$\log(p(t|x, \theta, \pi)) = \log\left(\frac{p(t|\pi) p(x|t, \theta)}{\sum_{t'} p(t'|\pi) p(x|t', \theta_{jt'})}\right)$$

$$= \log(p(t|\pi) \prod_{j=1}^{784} p(x_j|t, \theta)) - \log\left(\sum_{t'} p(t'|\pi) p(x|t', \theta_{jt'})\right)$$

$$= \log p(t|\pi) + \sum_{j=1}^{784} \log p(x_j|t, \theta) - \log \sum_{t'=0}^9 \pi_{t'} \prod_{j=1}^{784} \theta_{jt'}^{x_j} (1-\theta_{jt'})^{(1-x_j)}$$

$$= \log \pi_{\epsilon} + \sum_{j=1}^{784} \log \theta_{j\epsilon}^{x_j} (1-\theta_{j\epsilon})^{(1-x_j)} - \log \sum_{\epsilon'=0}^9 \pi_{\epsilon'} \prod_{j=1}^{784} \theta_{j\epsilon'}^{x_j} (1-\theta_{j\epsilon'})^{(1-x_j)}$$

$$= \log \pi_{\epsilon} + \sum_{j=1}^{784} (x_j \log \theta_{j\epsilon} + (1-x_j) \log (1-\theta_{j\epsilon}))$$

$$- \left(\log \sum_{\epsilon'=0}^9 \pi_{\epsilon'} + \sum_{j=1}^{784} (x_j \log \theta_{j\epsilon'} + (1-x_j) \log (1-\theta_{j\epsilon'})) \right)$$

$$= \log \pi_{\epsilon} + \sum_{j=1}^{784} (x_j \log \theta_{j\epsilon} + (1-x_j) \log (1-\theta_{j\epsilon}))$$

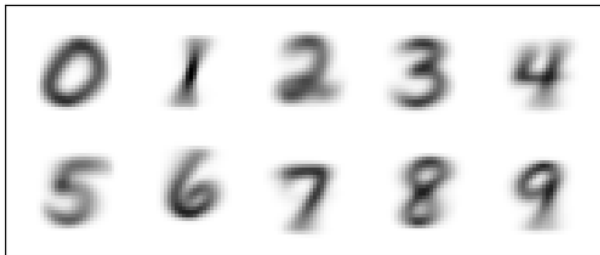
$$- \log \left[\sum_{\epsilon'=0}^9 \exp \left(\log \pi_{\epsilon'} + \sum_{j=1}^{784} (x_j \log \theta_{j\epsilon'} + (1-x_j) \log (1-\theta_{j\epsilon'})) \right) \right]$$

(c)

```
Average log-likelihood for MLE is nan
```

When computing the log-likelihood for MLE, we could encounter a divide by zero error if the count of a particular class is equal to 0 in the data.

(d)



$$(e) \hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | D)$$

$$= \arg \max_{\theta} \log p(\theta) + \log p(D | \theta)$$

$$= \arg \max_{\theta} \log p(\theta; 3, 3) + \log p(D | \theta; 3, 3)$$

$$= \arg \max_{\theta} \log \theta_{jc}^{2-1} (1-\theta_{jc})^{2-1} + \log \prod_{i=1}^k p(c | \pi) \prod_{j=1}^{784} p(x_j | c; \theta_{jc})$$

$$= \arg \max_{\theta} \log \theta_{jc}^2 (1-\theta_{jc})^2 + \log \left[\prod_{i=1}^k \frac{1}{\prod_{c=0}^9 \pi_c} t_c^{(i)} \prod_{j=1}^{784} \theta_{jc}^{x_j^{(i)}} (1-\theta_{jc})^{(1-x_j^{(i)})} \right]$$

$$= \arg \max_{\theta} 2 \log \theta_{jc} + 2 \log (1-\theta_{jc}) + \left[\sum_{i=1}^k \sum_{c=0}^9 \left(t_c^{(i)} \log \pi_c + \left(\sum_{j=1}^{784} x_j^{(i)} \log \theta_{jc} + (1-x_j^{(i)}) \log (1-\theta_{jc}) \right) \right) \right]$$

$$= \arg \max_{\theta} 2 \log \theta_{jc} + 2 \log (1-\theta_{jc}) + \left[\sum_{i=1}^k \sum_{c=0}^9 t_c^{(i)} \log \pi_c \right] + \left[\sum_{i=1}^k \sum_{c=0}^9 \left(\sum_{j=1}^{784} x_j^{(i)} \log \theta_{jc} + (1-x_j^{(i)}) \log (1-\theta_{jc}) \right) \right]$$

Taking the derivative w.r.t. θ_{jc} we have:

$$\frac{d}{d\theta_{jc}} \left(2 \log \theta_{jc} + 2 \log (1-\theta_{jc}) + \left[\sum_{i=1}^k \sum_{c=0}^9 t_c^{(i)} \log \pi_c \right] + \left[\sum_{i=1}^k \sum_{c=0}^9 \left(\sum_{j=1}^{784} x_j^{(i)} \log \theta_{jc} + (1-x_j^{(i)}) \log (1-\theta_{jc}) \right) \right] \right)$$

$$= \frac{2}{\theta_{jc}} - \frac{2}{1-\theta_{jc}} + \sum_{i=1}^k t_c^{(i)} \left(\frac{x_j^{(i)}}{\theta_{jc}} - \frac{(1-x_j^{(i)})}{(1-\theta_{jc})} \right) \quad \# \text{ from (a)}$$

having the derivative equal to 0 to solve for $\hat{\theta}_{MAP}$:

$$\frac{2}{\theta_{jc}} - \frac{2}{1-\theta_{jc}} + \sum_{i=1}^k t_c^{(i)} \left(\frac{x_j^{(i)}}{\theta_{jc}} - \frac{(1-x_j^{(i)})}{(1-\theta_{jc})} \right) = 0$$

$$\frac{2-2\theta_{jc}-2\theta_{jc}}{\theta_{jc}(1-\theta_{jc})} + \sum_{i=1}^k t_c^{(i)} \left(\frac{x_j^{(i)}(1-\theta_{jc}) - (1-x_j^{(i)})\theta_{jc}}{\theta_{jc}(1-\theta_{jc})} \right) = 0$$

$$2-2\theta_{jc}-2\theta_{jc} + \sum_{i=1}^k t_c^{(i)} x_j^{(i)} - t_c^{(i)} \theta_{jc} = 0 \quad \# \text{ from 2(a)}$$

$$2 + \sum_{i=1}^k t_c^{(i)} x_j^{(i)} = \sum_{i=1}^k t_c^{(i)} \theta_{jc} + 4\theta_{jc}$$

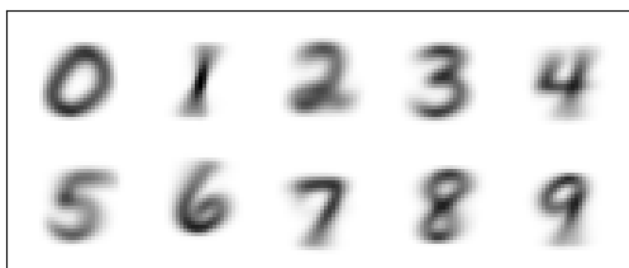
$$2 + \sum_{i=1}^k t_c^{(i)} x_j^{(i)} = \theta_{jc} \left(\sum_{i=1}^k t_c^{(i)} + 4 \right)$$

$$\hat{\theta}_{MAP} = \frac{2 + \sum_{i=1}^k t_c^{(i)} x_j^{(i)}}{4 + \sum_{i=1}^k t_c^{(i)}}$$

(f)

```
Average log-likelihood for MAP is -3.357063137860285  
Training accuracy for MAP is 0.8352166666666667  
Test accuracy for MAP is 0.816
```

(g)



3)

3)

(a)

To determine the posterior distribution, we have:

$$p(\theta|D) = \frac{p(\theta)p(D|\theta)}{\int p(\theta')p(D|\theta')d\theta'}$$

$$\therefore p(\theta|D) \propto p(\theta)p(D|\theta)$$

We know that $p(\theta) \propto \theta_1^{\alpha_1-1} \dots \theta_K^{\alpha_K-1}$

$$p(\theta) \propto \prod_{k=1}^K \theta_k^{\alpha_k-1}$$

Determining $p(D|\theta)$:

$$\begin{aligned} p(D|\theta) &= \prod_{i=1}^N p(x^{(i)}|\theta) \quad \# \text{ independently distributed} \\ &= \prod_{i=1}^N \prod_{k=1}^K \theta_k^{x_k^{(i)}} \end{aligned}$$

$$\begin{aligned}
\therefore p(\theta|D) &\propto p(\theta)p(D|\theta) \\
&= \prod_{k=1}^K \theta_k^{\alpha_k-1} \prod_{j=1}^N \prod_{k=1}^K \theta_k^{x_k^{(j)}} \\
&= \prod_{k=1}^K \theta_k^{\alpha_k-1} \prod_{k=1}^K \theta_k^{\sum_{j=1}^N x_k^{(j)}} \neq \prod_{i=1}^I \exp x_i = \exp\left(\sum_{i=1}^I x_i\right) \\
&= \prod_{k=1}^K \theta_k^{\alpha_k-1} \prod_{k=1}^K \theta_k^{N_k} \\
&= \prod_{k=1}^K \theta_k^{\alpha_k-1+N_k}
\end{aligned}$$

Yes. This is because that the prior $p(\theta)$ & the likelihood $p(D|\theta)$ have the same functional form.

(b)

$$\begin{aligned}\hat{\theta}_{\text{MAP}} &= \underset{\theta}{\operatorname{argmax}} p(\theta | D) \\ &= \underset{\theta}{\operatorname{argmax}} \prod_{k=1}^K \theta_k^{\alpha_k - 1 + N_k} \\ &= \underset{\theta}{\operatorname{argmax}} \log \left(\prod_{k=1}^K \theta_k^{\alpha_k - 1 + N_k} \right) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{k=1}^K \log(\theta_k^{\alpha_k - 1 + N_k}) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{k=1}^K (\alpha_k - 1 + N_k) \log(\theta_k)\end{aligned}$$

Setting up the Lagrangian, we have:

$$L = \sum_{k=1}^K (\alpha_k - 1 + N_k) \log(\theta_k) + \lambda \left(1 - \sum_k \theta_k \right)$$

$$\frac{dL}{d\theta} = \frac{\alpha_k - 1 + N_k}{\theta_k} - \lambda = 0$$

$$\frac{\alpha_k - 1 + N_k}{\theta_k} = \lambda$$

$$\theta_k = \frac{\alpha_k - 1 + N_k}{\lambda} \rightarrow \textcircled{1}$$

$$\frac{dL}{d\lambda} = 1 - \sum_k \theta_k \rightarrow \textcircled{2}$$

subbing $\textcircled{1}$ into $\textcircled{2}$ we have:

$$1 - \sum_k \theta_k = 1 - \sum_k \frac{\alpha_k - 1 + N_k}{\lambda} = 0$$

$$\Leftrightarrow \sum_k \frac{\alpha_k - 1 + N_k}{\lambda} = 1$$

$$\lambda = \sum_k \alpha_k - 1 + N_k$$

$$\lambda = \sum_k \alpha_k - K + N$$

$$\therefore \hat{\theta}_{k\text{MAP}} = \frac{\alpha_k - 1 + N_k}{\sum_k \alpha_k - K + N}$$

(c)

$$p(x^{(N+1)} | \mathcal{D}) = \int p(x^{(N+1)} = k | \theta) p(\theta | \mathcal{D}) d\theta$$

$$= \int \theta_k^{x_k^{(N+1)}} p(\theta_k | \mathcal{D}) d\theta_k$$

$$= \int \theta_k p(\theta_k | \mathcal{D}) d\theta_k$$

since probability of each category is θ_k

$$= E[\theta_k | \mathcal{D}]$$

$$E(x) = \int x f(x) dx$$

$$= \frac{a_k + N_k}{\sum_{k'} a_{k'} + N_{k'}}$$

4)

(a)

```
Average Conditional Log Likelihood for Training Data: -0.12462443666863034
Average Conditional Log Likelihood for Testing Data: -0.19667320325525578
```

(b)

```
Training Accuracy: 0.9814285714285714
Testing Accuracy: 0.97275
```

(c)

Top Left to Right: 0, 1, 2, 3, 4

Bottom Left to Right: 5, 6, 7, 8, 9

