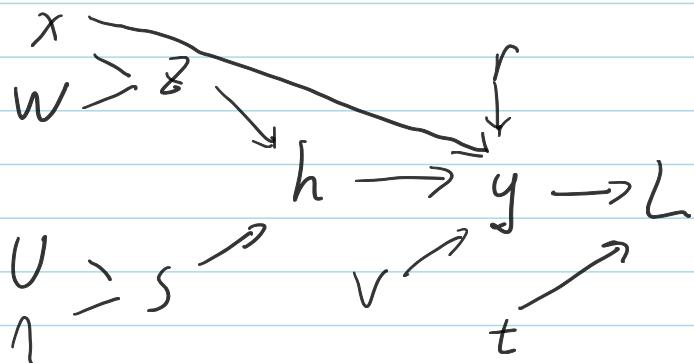


Question 1:

(a) computation graph:



b)

backprop: get the element wise forward prop.

$$z_i = \sum_j w_{ij} \cdot x_j$$

$$s_i = \sum_j v_{ij} \cdot r_j$$

$$h_i = z_i \cdot \sigma(z_i)$$

$$y = \sum_l v_{il} \cdot h_i + \sum_l r_i \cdot x_i$$

$$L = \frac{1}{2}(y - t)^2$$

get element wise back prop:

$$\bar{L} = 1$$

$$\bar{y} = \bar{L} \cdot (y - t)$$

$$\bar{r}_i = \bar{y} \cdot x_i$$

$$\bar{h}_i = \bar{y} \cdot v_i$$

$$\bar{s}_i = \bar{h}_i \cdot \bar{z}_i \cdot \sigma'(s_i)$$

$$\bar{z}_i = \bar{h}_i \cdot \bar{v}_i$$

$$\bar{w}_{ij} = \bar{z}_i \cdot \bar{x}_j$$

$$\bar{x}_i = \bar{y} \cdot r_i + \sum_j \bar{s}_j w_{ji}$$

$$\bar{v}_{il} = \sum_j \bar{s}_j U_{jl}$$

$$\bar{U}_{jl} = \bar{s}_i \cdot \eta_j$$

In vectorized form:

$$\bar{v} = \bar{y} h \quad \bar{h} = \bar{y} V$$

$$\bar{r} = \bar{y} x \quad \bar{s} = \bar{h} \circ Z \circ G'(s)$$

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

$$\bar{w} = \bar{z} \times T \quad \bar{U} = \bar{s} A^T$$

$$\bar{x} = \bar{y} r + w^T \bar{z}$$

$$\bar{\eta} = U^T \bar{s}$$

we get $\bar{v}, \bar{r}, \bar{w}, \bar{U}, \bar{x}, \bar{\eta}$ as desired

Question 2:

Using Naive Bayes.

a) Given independence,

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

from lecture, we know: log-likelihood:

$$l(\theta, \pi) = \sum_{i=1}^N \log(p(c^{(i)} | \pi)) + \sum_{i=1}^N \sum_{j=1}^{784} p(x_j | c^{(i)}, \theta_j, \pi)$$

first maximize the Bernoulli log-likelihood of labels term

Note $p(c^{(i)}) = \prod_{j=0}^q \pi_{c=j}^{t_j^{(i)}}$ ($t_j^{(i)}$) here represent the the value of j^{th} index of $t^{(i)}$

$$\begin{aligned} \frac{d}{d\pi_c} l(\theta, \pi) &= \frac{d}{d\pi_c} \sum_{i=1}^N \log(p(c^{(i)} | \pi)) \\ &= \frac{d}{d\pi_c} \sum_{i=1}^N \log\left(\prod_{j=0}^q \pi_{c=j}^{t_j^{(i)}}\right) \\ &= \frac{d}{d\pi_c} \sum_{i=1}^N \sum_{j=0}^q t_j^{(i)} \log(\pi_{c=j}) \end{aligned}$$

$$\text{replace } \pi_q = 1 - \sum_{j=0}^8 \pi_j$$

$$\begin{aligned} &= \frac{d}{d\pi_c} \sum_{i=1}^N t_q^{(i)} \log\left(1 - \sum_{j=0}^8 \pi_j\right) + \frac{d}{d\pi_c} \sum_{i=1}^N \sum_{j=0}^8 t_j^{(i)} \log(\pi_j) \\ &= \sum_{i=1}^N \frac{t_c^{(i)}}{\pi_c} - \sum_{i=1}^N \frac{t_q^{(i)}}{1 - \sum_{j=0}^8 \pi_j} \\ &= \sum_{i=1}^N \frac{t_c^{(i)}}{\pi_c} - \sum_{i=1}^N \frac{t_q^{(i)}}{\pi_q} \end{aligned}$$

Set derivative to 0:

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

$$\sum_{i=1}^N \frac{f_c(i)}{\pi_c} = \sum_{i=1}^N \frac{t_q(i)}{\pi_q}$$

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

$$\hat{\pi}_c = \hat{\pi}_q \cdot \frac{\# \text{ of sample of class } c}{\# \text{ of sample of class } q}$$

since $\sum_{i=0}^q \hat{\pi}_i = 1$

$$\hat{\pi}_q + \sum_{j=0}^8 \hat{\pi}_q \frac{\sum_{i=1}^N t_j(i)}{\sum_{i=1}^N t_q(i)} = 1$$

$$\hat{\pi}_q + \hat{\pi}_q \frac{\sum_{j=0}^8 \sum_{i=1}^N t_j(i)}{\sum_{i=1}^N t_q(i)} = 1$$

$$\hat{\pi}_q = (1 - \hat{\pi}_q) \cdot \frac{\sum_{i=1}^N t_q(i)}{\sum_{j=0}^8 \sum_{i=1}^N t_j(i)}$$

$$\hat{\pi}_q = \frac{\sum_{i=1}^N t_q(i)}{\sum_{j=0}^8 \sum_{i=1}^N t_j(i)} - \frac{\sum_{i=1}^N t_q(i)}{\sum_{j=0}^8 \sum_{i=1}^N t_j(i)} \hat{\pi}_q$$

$$\hat{\pi}_q \left(1 + \frac{\sum_{i=1}^N t_q(i)}{\sum_{j=0}^8 \sum_{i=1}^N t_j(i)} \right) = \frac{\sum_{i=1}^N t_q(i)}{\sum_{j=0}^8 \sum_{i=1}^N t_j(i)}$$

$$\begin{aligned} \hat{\pi}_q &= \frac{\sum_{i=1}^N t_q(i)}{\sum_{i=1}^N (\sum_{j=0}^8 t_j(i) + t_q(i))} \\ &= \frac{\sum_{i=1}^N t_q(i)}{N} \end{aligned}$$

(only one of index of t is 1)

$$\hat{\pi}_q = \frac{\# \text{ of sample of class } q}{\# \text{ of samples}}$$

Hence for c in $1 \sim 8$

$$\hat{\pi}_c = \hat{\pi}_q \cdot \frac{\# \text{ of sample of class } c}{\# \text{ of sample of class } q}$$

$$\frac{1}{\pi_c} = \frac{\# \text{of sample of class } c}{\# \text{of samples}}.$$

notice this form works for c in $1 \sim 9$

maximize Bernoulli log-likelihood for feature x_j

notice $p(x_j | C, \Theta_{jC}) = \Theta_{jC}^{x_j} (1 - \Theta_{jC})^{1-x_j}$ since $x_j \in \{0, 1\}$

$$\frac{d}{d\Theta_{jC}} L(\Theta, \pi) = \frac{d}{d\Theta_{jC}} \sum_{i=1}^N \sum_{j=1}^{784} \log (\Theta_{jC}^{x_{j,i}} (1 - \Theta_{jC}^{x_{j,i}})^{1-x_{j,i}})$$

$$= \frac{d}{d\Theta_{jC}} \sum_{i=1}^N \sum_{j=1}^{784} x_{j,i} \log \Theta_{jC}^{x_{j,i}}$$

$$+ \frac{d}{d\Theta_{jC}} \sum_{i=1}^N \sum_{j=1}^{784} (1 - x_{j,i}) \log (1 - \Theta_{jC}^{x_{j,i}})$$

$$= \sum_{i=1}^N \sum_{j=1}^{784} \frac{d}{d\Theta_{jC}} x_{j,i} \log \Theta_{jC}^{x_{j,i}} + \sum_{i=1}^N \sum_{j=1}^{784} \frac{d}{d\Theta_{jC}} (1 - x_{j,i}) \log (1 - \Theta_{jC}^{x_{j,i}})$$

$$= \sum_{i=1}^N \frac{d}{d\Theta_{jC}} x_{j,i} \log \Theta_{jC}^{x_{j,i}} + \sum_{i=1}^N \frac{d}{d\Theta_{jC}} (1 - x_{j,i}) \log (1 - \Theta_{jC}^{x_{j,i}})$$

$$= \sum_{i \in \{k \mid i \leq N \mid C^{(i)} = c\}} \frac{x_{j,i}}{\Theta_{jC}} - \sum_{i \notin \{k \mid i \leq N \mid C^{(i)} = c\}} \frac{(1 - x_{j,i})}{1 - \Theta_{jC}}$$

setting it to zero,

$$\sum_{i \in \{k \mid i \leq N \mid C^{(i)} = c\}} \frac{x_{j,i}}{\Theta_{jC}} - \sum_{i \notin \{k \mid i \leq N \mid C^{(i)} = c\}} \frac{(1 - x_{j,i})}{1 - \Theta_{jC}} = 0$$

$$\frac{\sum_{i \in \{k \mid i \leq N \mid C^{(i)} = c\}} x_{j,i}}{\Theta_{jC}} = \frac{\sum_{i \notin \{k \mid i \leq N \mid C^{(i)} = c\}} (1 - x_{j,i})}{1 - \Theta_{jC}}$$

$$(1 - \hat{\pi}_{jC})_{i \in \{1 \leq i \leq N | C^{(i)} = C\}} \sum_{j \in \{1 \leq j \leq J | C^{(i)} = C\}} x_{j|i} \\ = \hat{\theta}_{jC} \sum_{i \in \{1 \leq i \leq N | C^{(i)} = C\}} (1 - x_{j|i})$$

$$\hat{\pi}_{jC} = \frac{\sum_{i \in \{1 \leq i \leq N | C^{(i)} = C\}} 1}{\sum_{i \in \{1 \leq i \leq N | C^{(i)} = C\}} \sum_{j \in \{1 \leq j \leq J | C^{(i)} = C\}} x_{j|i}}$$

$$\hat{\pi}_{jC} = \frac{\sum_{i \in \{1 \leq i \leq N | C^{(i)} = C\}} x_{j|i}}{\sum_{i \in \{1 \leq i \leq N | C^{(i)} = C\}} 1}$$

$$= \frac{\# \text{ samples of class } C \text{ with } x_{j|i} = 1}{\# \text{ samples of class } C}$$

Hence we derived form of $\hat{\pi}_{jC}$ and $\hat{\theta}_{jC}$

b)

Using Bayes' Rules:

$$P(t|x,\Theta, \pi_c) = \frac{p(x|t, \Theta) \cdot P(t|\pi_c)}{\sum_{c=0}^q p(c|\pi_c) \cdot P(x|t_c=1, \Theta)}$$

we know $P(t|\pi_c) = \prod_{c=0}^q \pi_c^{t_c}$ | say A | 1 (say B) -

$$P(x|t, \Theta) = \prod_{j=1}^q p(x_j|t, \Theta)$$

$$= \prod_{j=1}^q \prod_{c=0}^q (\Theta_{jc}^{x_j} (1-\Theta_{jc})^{1-x_j})^{t_c}$$

$$P(t|x, \Theta, \pi_c) = \frac{\prod_{j=1}^q \prod_{c=0}^q (\Theta_{jc}^{x_j} (1-\Theta_{jc})^{1-x_j})^{t_c} \prod_{c=0}^q \pi_c^{t_c}}{\sum_{c=0}^q \pi_c \prod_{j=1}^q \Theta_{jc}^{x_j} \cdot (1-\Theta_{jc})^{1-x_j}}$$

for log-likelihood (hard):

$$\log(P(t|x, \Theta, \pi_c)) = \log \frac{\prod_{j=1}^q \prod_{c=0}^q (\Theta_{jc}^{x_j} (1-\Theta_{jc})^{1-x_j})^{t_c} \prod_{c=0}^q \pi_c^{t_c}}{\sum_{c=0}^q \pi_c \prod_{j=1}^q \Theta_{jc}^{x_j} \cdot (1-\Theta_{jc})^{1-x_j}}$$
$$= \log \left(\prod_{j=1}^q \prod_{c=0}^q (\Theta_{jc}^{x_j} (1-\Theta_{jc})^{1-x_j})^{t_c} \prod_{c=0}^q \pi_c^{t_c} \right) - \log \left(\sum_{c=0}^q \pi_c \prod_{j=1}^q \Theta_{jc}^{x_j} \cdot (1-\Theta_{jc})^{1-x_j} \right)$$

Separate it into 2 parts: Simplify each

$$\log \left(\prod_{j=1}^q \prod_{c=0}^q (\Theta_{jc}^{x_j} (1-\Theta_{jc})^{1-x_j})^{t_c} \prod_{c=0}^q \pi_c^{t_c} \right)$$
$$= \sum_{j=1}^q \sum_{c=0}^q t_c (x_j \log \Theta_{jc} + (1-x_j) \log (1-\Theta_{jc})) + \sum_{c=0}^q t_c \log \pi_c$$

$$\begin{aligned}
 & \log \left(\sum_{c=0}^q \pi_c \prod_{j=1}^q \theta_{jc}^{x_j} \cdot (1-\theta_{jc})^{1-x_j} \right) \\
 &= \log \sum_{c=0}^q \pi_c \cdot e^{\sum_{j=1}^q (x_j \log \theta_{jc} + (1-x_j) \log (1-\theta_{jc}))} \\
 &= \log \sum_{c=0}^q \pi_c \cdot e^{\sum_{j=1}^q (x_j \log \theta_{jc} + (1-x_j) \log (1-\theta_{jc}))}
 \end{aligned}$$

Hence the log-likelihood $\log p(t|x, \psi, \pi_c)$

$$\begin{aligned}
 &= \sum_{j=1}^q \sum_{c=0}^q t_c (x_j \log \theta_{jc} + (1-x_j) \log (1-\theta_{jc})) + \sum_{c=0}^q t_c \log \pi_c \\
 &= \log \sum_{c=0}^q \pi_c \cdot e^{\sum_{j=1}^q (x_j \log \theta_{jc} + (1-x_j) \log (1-\theta_{jc}))}
 \end{aligned}$$

c) in order to calculate the average log likelihood. with MLE
need to calculate the individual log likelihood.

from part b, we know need to calculate $\hat{\theta}_{jc}$

from part a, we know $\hat{\theta}_{jc} = \frac{\# \text{ samples of class } C \text{ with } x_j=1}{\# \text{ samples of class } C}$

we can see that around corner. for each image.

$x_j=0$ as white. Hence $\# \text{ samples of class } C \text{ with } x_j=1$ is 0 in some case

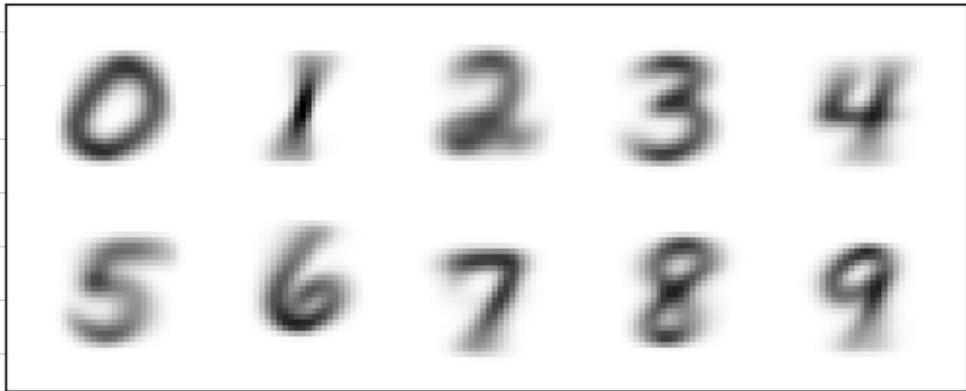
Hence $\hat{\theta}_{jc} = 0$ and $\log(0)$ is undefined.

Since we involve undefined term when calculating average.

the average is also undefined.

d)

MLG-estimator $\hat{\theta}$ for each class



e) from lecture, we know,

$$\hat{\theta} \text{ MAP} = \underset{\theta}{\operatorname{argmax}} \log p(\mathbf{t}) + \log p(\mathbf{D}|\mathbf{t})$$

$$\text{let } f = \log p(\mathbf{t}) + \log p(\mathbf{D}|\mathbf{t})$$

Since $P(\mathbf{t})$ follows Beta(3, 3),

$$P(\mathbf{t}) = \text{constant} \cdot \prod_{j=1}^N \prod_{c=0}^q \theta_{jc}^{t_c^{(j)}} (1 - \theta_{jc})^{1-t_c^{(j)}}$$

$$p(\mathbf{D}|\theta) = \prod_{i=1}^N \prod_{j=1}^q p(x_j^{(i)} | t^{(i)}, \theta)$$

$$= \prod_{i=1}^N \prod_{j=1}^q \prod_{c=0}^q t_c^{t_c^{(i)}} x_j^{t_c^{(i)}} (1 - \theta_{jc})^{t_c^{(i)}(1 - x_j^{(i)})}$$

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \sum_{j=1}^q \sum_{c=0}^q (2 \log \theta_{jc} + 2 \log (1 - \theta_{jc}))$$

$$+ \sum_{i=1}^N \sum_{j=1}^q \sum_{c=0}^q \frac{t_c^{t_c^{(i)}} x_j^{t_c^{(i)}} \cdot \log (\theta_{jc}) + t_c^{1-t_c^{(i)}} (1 - x_j^{(i)})}{\log (1 - \theta_{jc})}$$

$$= \underset{\theta}{\operatorname{argmax}} \left(\sum_{j=1}^q \sum_{c=0}^q \left(\left(2 + \sum_{i=1}^N t_c^{t_c^{(i)}} x_j^{t_c^{(i)}} \right) \log (\theta_{jc}) + \left(2 + \sum_{i=1}^N t_c^{1-t_c^{(i)}} (1 - x_j^{(i)}) \right) \log (1 - \theta_{jc}) \right) \right)$$

differentiate by θ_{jc} and set it to 0

$$\frac{d}{d\theta_{jc}} = \frac{d}{d\theta_{jc}} \sum_{j=1}^q \sum_{c=0}^q \left(\left(2 + \sum_{i=1}^N t_c^{t_c^{(i)}} x_j^{t_c^{(i)}} \right) \log (\theta_{jc}) + \left(2 + \sum_{i=1}^N t_c^{1-t_c^{(i)}} (1 - x_j^{(i)}) \right) \log (1 - \theta_{jc}) \right)$$

$$= 0$$

$$= \frac{2 + \sum_{i=1}^N t_c^{(i)} x_j^{(i)}}{\hat{\Theta}_{jC}} - \frac{2 + \sum_{i=1}^N t_c^{(i)} (1 - x_j^{(i)})}{1 - \hat{\Theta}_{jC}}$$

we get $\hat{\Theta}_{jC} = \frac{2 + \sum_{i=1}^N t_c^{(i)} x_j^{(i)}}{4 + \sum_{i=1}^N t_c^{(i)}}$

$$\hat{\Theta}_{jC} = \frac{2 + \# \text{ of samples, that } x_j = 1 \text{ and class } C}{4 + \# \text{ of samples, and class } C}$$

Hence $\hat{\Theta}_{jC} = \frac{2 + \# \text{ of samples, that } x_j = 1 \text{ and class } C}{4 + \# \text{ of samples, and class } C}$

for each i, j entry for Θ map

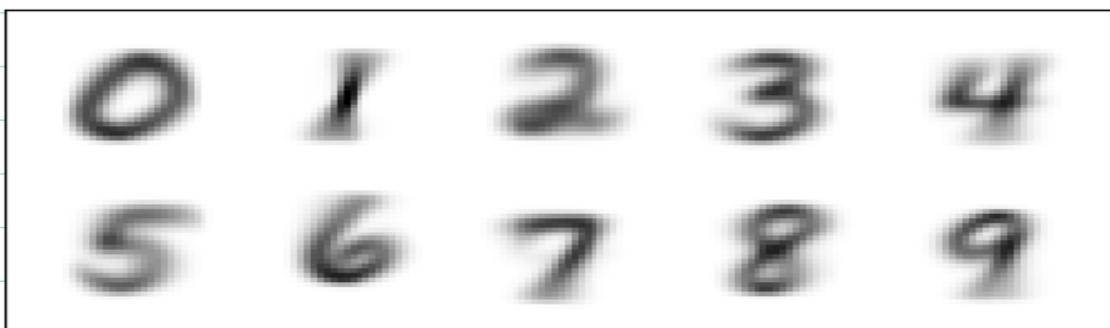
f)

Average log-likelihood for MAP is -3.357063137860287

Training accuracy for MAP is 0.8352166666666667

Test accuracy for MAP is 0.816

g) MAP estimator $\hat{\theta}$



Question 3:

a) based on given,

$$P(D|\theta) = \prod_{i=1}^N \prod_{k=1}^K \theta_k^{x_{ki}\alpha_k}, \quad P(\theta) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

$$\text{since } P(G|D) \propto P(\theta) \cdot P(D|G)$$

$$P(G|D) \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1} \prod_{i=1}^N \prod_{k=1}^K \theta_k^{x_{ki}\alpha_k} \\ \propto \prod_{k=1}^K \theta_k^{\alpha_k - 1} \prod_{k=1}^K \theta_k^{\sum_{i=1}^N x_{ki}}$$

$$\propto \prod_{k=1}^K \theta_k^{\alpha_k - 1} \prod_{k=1}^K \theta_k^{N_k}$$

$$\propto \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1}$$

$$\text{Hence } P(G|D) \propto \prod_{k=1}^K \theta_k^{N_k + \alpha_k - 1},$$

which is a conjugate prior for the categorical distribution

b)

$$\hat{\theta}_{\text{MAP}} = \operatorname{argmax}_{\theta} P(\theta | D)$$

$$= \operatorname{argmax}_{\theta}$$

$$= \operatorname{argmax}_{\theta} \log \left(\prod_{k=1}^K \frac{1}{\theta_k^{N_k + \alpha_k - 1}} \right)$$

$$= \operatorname{argmax} \sum_{k=1}^K (N_k + \alpha_k - 1) \log \theta_k$$

$$\text{let } f = \sum_{k=1}^K (N_k + \alpha_k - 1) \log \theta_k,$$

differentiate on θ_j ($j \neq k$)

$$\frac{d}{d\theta_j} f = \frac{d}{d\theta_j} \sum_{k=1}^{K-1} (\alpha_k - 1 + N_k) \log (\theta_k)$$

$$+ \frac{d}{d\theta_j} (\alpha_k - 1 + N_k) \log (\theta_k).$$

$$= \frac{d}{d\theta_j} \sum_{k=1}^{K-1} (\alpha_k - 1 + N_k) \log (\theta_k) + \\ \frac{d}{d\theta_j} (\alpha_k - 1 + N_k) \log \left(1 - \sum_{k=1}^{K-1} \theta_k \right)$$

$$(\text{by } \sum_k \theta_k = 1)$$

$$= \frac{\alpha_j - 1 + N_j}{\theta_j} - \frac{\alpha_k - 1 + N_k}{1 - \sum_{k=1}^{K-1} \theta_k}$$

$$= \frac{\alpha_j - 1 + N_j}{\theta_j} - \frac{\alpha_k - 1 + N_k}{\theta_k}$$

Set derivative to 0

$$\frac{\alpha_j - 1 + V_j}{\hat{\theta}_j} = \frac{\alpha_k - 1 + V_k}{\hat{\theta}_k}$$

$$\hat{\theta}_j = \frac{\hat{\theta}_k (\alpha_j - 1 + V_j)}{\alpha_k - 1 + V_k}$$

since $\sum_{k=1}^K \hat{\theta}_k = 1$

$$\hat{\theta}_k + \sum_{j=1}^{k-1} \hat{\theta}_j = 1$$

$$\hat{\theta}_k + \sum_{j=1}^{k-1} \frac{\hat{\theta}_k (\alpha_j - 1 + V_j)}{\alpha_k - 1 + V_k} = 1$$

$$\sum_{j=1}^K \frac{\hat{\theta}_k (\alpha_j - 1 + V_j)}{\alpha_k - 1 + V_k} = 1$$

$$\hat{\theta}_k = \frac{\alpha_k - 1 + V_k}{\sum_{j=1}^K \alpha_j - 1 + V_j}$$

Hence $\hat{\theta}_j = \frac{\alpha_k - 1 + V_k}{\sum_{j=1}^K \alpha_j - 1 + V_j} \cdot \frac{\alpha_j - 1 + V_j}{\alpha_k - 1 + V_k}$

$$\hat{\theta}_j = \frac{\alpha_j - 1 + V_j}{\sum_{j=1}^K \alpha_j - 1 + V_j}$$

Hence for all $j \in \{j \in \mathbb{Z} \mid 1 \leq j \leq K\}$

$$\hat{\theta}_j = \frac{\alpha_j - 1 + V_j}{\sum_{j=1}^K \alpha_j - 1 + V_j}$$

c) want to know

$$P(X_K^{D+1} = 1 | D)$$

by definition,

$$P(X_K^{D+1} = 1 | D) = \int P(X_K^{D+1} = 1 | \Theta) P(\Theta | D)$$

$$= \int \Theta_K \prod_{j=1}^K \Theta_j^{\alpha_j - 1 + N_j}$$

In addition, if $\Theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$

$$E(\Theta_K) = \int \Theta_K \cdot P(\Theta) \quad \text{since } \Theta \text{ is continuous,}$$

$$= \int \Theta_K \prod_{j=1}^K \Theta_j^{\alpha_j - 1}$$

$$= \frac{\alpha_k}{\sum_{k'} \alpha_{k'}}$$

if $\Theta \sim \text{Dirichlet}(\alpha_1 + N_1, \dots, \alpha_k + N_k),$

$$E(\Theta_K) = \int \Theta_K \cdot P(\Theta)$$

$$= \int \Theta_K \prod_{j=1}^K \Theta_j^{\alpha_j + N_j - 1} = P(X_K^{D+1} = 1 | D) \quad (\text{has same form})$$

Since if $\Theta \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k)$

$$E(\Theta_K) = \frac{\alpha_k}{\sum_{k'} \alpha_{k'}}$$

then if $\theta \sim \text{Dirichlet}(\alpha_1 + V_1, \dots, \alpha_k + V_k)$

$$E(\theta_k) = \frac{\alpha_k + V_k}{\sum_{k'} \alpha_{k'} + V_{k'}} = P(X_k^{D+1} = 1 | D)$$

Hence probability of the $N+1$ outcome was k ,

$$P(X_k^{D+1} = 1 | D) = \frac{\alpha_k + V_k}{\sum_{k'} \alpha_{k'} + V_{k'}}$$

Question 4

a)

average log-likelihood for training set is -0.12462443666863016

average log-likelihood for training set is -0.19667320325525606

b)

training set accuracy is 0.9814285714285714

test set accuracy is 0.97275

c)

