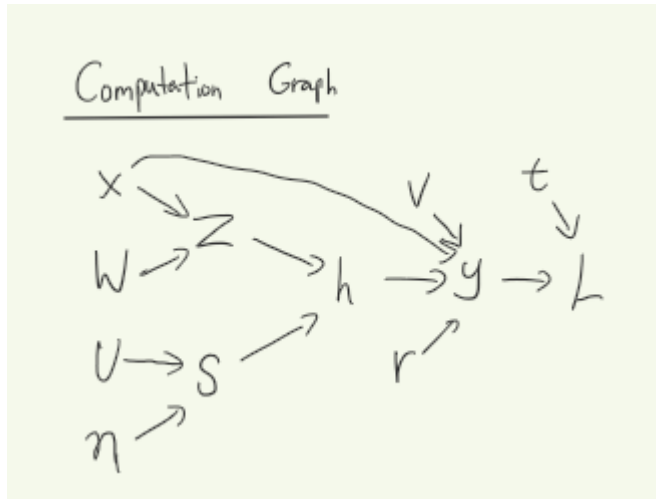


CSC311, Fall 2022, Homework 3

Huipeng Lei, 1005908154

Due Monday, November 7, 2022

Problem 1



a)

b) Given $\sigma(x) = \frac{1}{1+e^{-x}}$. We can calculate

$$\begin{aligned}
 1 - \sigma(x) &= 1 - \frac{1}{1 + e^{-x}} \\
 &= \frac{1 + e^{-x} - 1}{1 + e^{-x}} \\
 &= \frac{e^{-x}}{1 + e^{-x}}
 \end{aligned}$$

Then

$$\begin{aligned}
 \frac{d\sigma(x)}{dx} &= \frac{-(-e^{-x})}{(1 + e^{-x})^2} \\
 &= \frac{1}{1 + e^{-x}} \frac{e^{-x}}{1 + e^{-x}} \\
 &= \sigma(x)(1 - \sigma(x))
 \end{aligned}$$

c)

$$\bar{\mathcal{L}} = 1$$

$$\bar{y} = \bar{\mathcal{L}} \frac{\partial \mathcal{L}}{\partial y} = \bar{\mathcal{L}} \left(\frac{t}{y} + \frac{t-1}{1-y} \right) = \bar{\mathcal{L}} \frac{t-y}{y(1-y)}$$

$$\bar{\mathbf{v}} = \bar{y} \frac{\partial y}{\partial \mathbf{v}} = \bar{y} \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x}) (1 - \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x})) \mathbf{h}^T$$

$$\bar{\mathbf{r}} = \bar{y} \frac{\partial y}{\partial \mathbf{r}} = \bar{y} \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x}) (1 - \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x})) \mathbf{x}^T$$

$$\bar{\mathbf{h}} = \bar{y} \frac{\partial y}{\partial \mathbf{h}} = \bar{y} \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x}) (1 - \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x})) \mathbf{v}^T$$

$$\bar{\mathbf{z}} = \bar{\mathbf{h}} \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \bar{\mathbf{h}} J_{\mathbf{z}}(\mathbf{h}) = \bar{h} \begin{bmatrix} s_1 & 0 & \cdots & 0 \\ 0 & s_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_n \end{bmatrix}$$

($J_{\mathbf{z}}$ means to apply the Jacobian Matrix wrt. \mathbf{z})

$$\bar{\mathbf{s}} = \bar{h} \frac{\partial \mathbf{h}}{\partial \mathbf{s}} = \bar{h} J_{\mathbf{s}}(\mathbf{h}) = \bar{h} \begin{bmatrix} z_1 & 0 & \cdots & 0 \\ 0 & z_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & z_n \end{bmatrix}$$

($J_{\mathbf{s}}$ means to apply the Jacobian Matrix wrt. \mathbf{s})

$$\bar{\mathbf{W}} = \bar{\mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{W}} = \bar{\mathbf{z}} (\mathbf{x}^T \otimes \mathbf{I})$$

($\mathbf{x}^T \otimes \mathbf{I}$ a tensor; \otimes the Kronecker product; \mathbf{I} identity matrix)

$$\bar{\mathbf{x}} = \bar{\mathbf{z}} \frac{\partial \mathbf{z}}{\partial \mathbf{x}} + \bar{y} \frac{\partial y}{\partial \mathbf{x}} = \bar{\mathbf{z}} \mathbf{W} + \bar{y} \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x}) (1 - \sigma(\mathbf{v}^T \mathbf{h} + \mathbf{r}^T \mathbf{x})) \mathbf{r}^T$$

$$\bar{\mathbf{U}} = \bar{\mathbf{s}} \frac{\partial \mathbf{s}}{\partial \mathbf{U}} = \bar{\mathbf{s}} (\eta^T \otimes \mathbf{I})$$

($\eta^T \otimes \mathbf{I}$ a tensor; \otimes the Kronecker product; \mathbf{I} identity matrix)

$$\bar{\eta} = \bar{\mathbf{s}} \frac{\partial \mathbf{s}}{\partial \eta} = \bar{\mathbf{s}} \mathbf{U}$$

Problem 2

a) $L(\theta) = \prod_{i=1}^N (P(t^{(i)}|\pi) \prod_{j=1}^{784} P(x_j^{(i)}|c, \theta_{jc}))$. Then

$$l(\theta) = \sum_{i=1}^N \log(P(t^{(i)}|\pi)) + \sum_{i=1}^N \sum_{j=1}^{784} \log(P(x_j^{(i)}|c^{(i)}, \theta_{jc})) = \sum_{i=1}^N \log\left(\prod_{j=0}^9 \pi_j^{t_j^{(i)}}\right) + \sum_{i=1}^N \sum_{j=1}^{784} \log(P(x_j^{(i)}|c^{(i)}, \theta_{jc}))$$

Find the MLE for π

$$\begin{aligned} \frac{\partial l(\theta)}{\partial \pi} &= \begin{bmatrix} \frac{\partial l(\theta)}{\partial \pi_0} & \frac{\partial l(\theta)}{\partial \pi_1} & \dots & \frac{\partial l(\theta)}{\partial \pi_9} \end{bmatrix} \\ \frac{\partial l(\theta)}{\partial \pi_j} &= \frac{\partial}{\partial \pi_j} \sum_{i=1}^N \log\left(\prod_{j=0}^9 \pi_j^{t_j^{(i)}}\right) \\ &= \frac{\partial}{\partial \pi_j} \sum_{i=1}^N \sum_{j=0}^9 \log(\pi_j^{t_j^{(i)}}) \\ &= \frac{\partial}{\partial \pi_j} \sum_{i=1}^N (t_0^{(i)} \log(\pi_0) + \dots + t_j^{(i)} \log(\pi_j) + \dots + t_8^{(i)} \log(\pi_8) + t_9^{(i)} \log(1 - \sum_{k=0}^8 \pi_k)) \\ &= \sum_{i=1}^N \left[\frac{t_j^{(i)}}{\pi_j} - \frac{t_9^{(i)}}{1 - \sum_{k=0}^8 \pi_k} \right] \\ 0 &= \frac{1}{\pi_j} \sum_{i=1}^N t_j^{(i)} - \frac{1}{\pi_9} \sum_{i=1}^N t_9^{(i)} \\ \implies \hat{\pi}_j &= \hat{\pi}_9 \sum_{i=1}^N \frac{t_j^{(i)}}{t_9^{(i)}} = \frac{\sum_{i=1}^N t_j^{(i)}}{N}, \forall j \in \{0, \dots, 8\}, \quad (\pi_9 = \frac{\sum_{i=1}^N t_9^{(i)}}{N}) \\ &= \frac{\# \text{ class } t_j \text{ in dataset}}{\text{Total } \# \text{ of samples}} \end{aligned}$$

We can solve for

$$\begin{aligned} \hat{\pi}_9 &= 1 - \sum_{j=0}^8 \hat{\pi}_j \\ &= 1 - \frac{1}{N} \sum_{i=1}^N (t_0^{(i)} + t_1^{(i)} + \dots + t_8^{(i)}) \\ &= \frac{N - \sum_{i=1}^N (t_0^{(i)} + t_1^{(i)} + \dots + t_8^{(i)})}{N} \\ &= \frac{\sum_{i=1}^N t_9^{(i)}}{N} \end{aligned}$$

Thus,

$$\begin{aligned}\hat{\pi} &= \begin{pmatrix} \hat{\pi}_0 & \hat{\pi}_1 & \cdots & \hat{\pi}_8 & \hat{\pi}_9 \end{pmatrix} \\ &= \begin{pmatrix} \frac{\sum_{i=1}^N t_0^{(i)}}{N} & \frac{\sum_{i=1}^N t_1^{(i)}}{N} & \cdots & \frac{\sum_{i=1}^N t_8^{(i)}}{N} & \frac{\sum_{i=1}^N t_9^{(i)}}{N} \end{pmatrix}\end{aligned}$$

Find the MLE for θ

$$\begin{aligned}\frac{\partial l(\theta)}{\partial \theta_{jc}} &= \frac{\partial}{\partial \theta_{jc}} \sum_{i=1}^N \sum_{j=1}^{784} \log(P(x_j^{(i)} | c^{(i)}, \theta_{jc})) \\ &= \frac{\partial}{\partial \theta_{jc}} \sum_{i=1}^N \sum_{j=1}^{784} [x_j^{(i)} \log(\theta_{jc}) + (1 - x_j^{(i)}) \log(1 - \theta_{jc})] \\ &\quad \text{(Since } P(x_j | c, \theta_{jc}) = \theta_{jc}^{x_j} (1 - \theta_{jc})^{(1-x_j)}) \\ 0 &= \sum_{i=1}^N \mathbb{1}(c^{(i)} = c) \left(\frac{x_j^{(i)}}{\theta_{jc}} - \frac{1 - x_j^{(i)}}{1 - \theta_{jc}} \right) \\ 0 &= \frac{1}{\theta_{jc}(1 - \theta_{jc})} \sum_{i=1}^N \mathbb{1}(c^{(i)} = c) ((1 - \theta_{jc})x_j^{(i)} - \theta_{jc}(1 - x_j^{(i)})) \\ 0 &= \sum_{i=1}^N \mathbb{1}(c^{(i)} = c) (x_j^{(i)} - \theta_{jc}x_j^{(i)} - \theta_{jc} + \theta_{jc}x_j^{(i)}) \\ \sum_{i=1}^N \mathbb{1}(c^{(i)} = c) \theta_{jc} &= \sum_{i=1}^N \mathbb{1}(c^{(i)} = c) x_j^{(i)} \\ \implies \widehat{\theta}_{jc} &= \frac{\sum_{i=1}^N \mathbb{1}(c^{(i)} = c) x_j^{(i)}}{\sum_{i=1}^N \mathbb{1}(c^{(i)} = c)} \forall j = 1, 2, \dots, 784 \\ \implies \widehat{\theta}_{jc} &= \frac{N_{jc}}{N_c}\end{aligned}$$

(N_{jc} denotes # of class c images with j th pixel 1; N_c denotes # class c images)

Therefore,

$$\theta = \begin{bmatrix} \frac{N_{10}}{N_0} & \frac{N_{11}}{N_1} & \cdots & \frac{N_{19}}{N_9} \\ \frac{N_{20}}{N_0} & \frac{N_{21}}{N_1} & \cdots & \frac{N_{29}}{N_9} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{N_{7840}}{N_0} & \frac{N_{7841}}{N_1} & \cdots & \frac{N_{7849}}{N_9} \end{bmatrix}$$

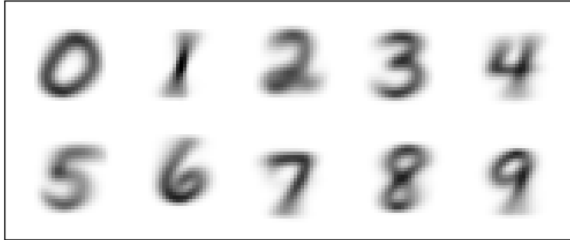
b) By Bayes' Rule, we can obtain

$$P(\mathbf{t}|\mathbf{x}, \theta, \pi) = P(c|\mathbf{x}, \theta, \pi) = \frac{P(\mathbf{x}, c|\theta, \pi)}{\sum_{c'} P(\mathbf{x}, c'|\theta, \pi)}$$

The denominator sums the conditional probability over all classes.
Then,

$$\begin{aligned}
\log P(\mathbf{t}|\mathbf{x}, \boldsymbol{\theta}, \boldsymbol{\pi}) &= \log \frac{P(\mathbf{x}, c|\boldsymbol{\theta}, \boldsymbol{\pi})}{\sum_{c'} P(\mathbf{x}, c'|\boldsymbol{\theta}, \boldsymbol{\pi})} \\
&= \log P(\mathbf{x}, c|\boldsymbol{\theta}, \boldsymbol{\pi}) - \log \sum_{c'} P(\mathbf{x}, c'|\boldsymbol{\theta}, \boldsymbol{\pi}) \\
&= \log(P(c|\boldsymbol{\pi}) \prod_{j=1}^{784} P(x_j|c, \theta_{jc})) - \log(\sum_c P(c|\boldsymbol{\pi}) \prod_{j=1}^{784} P(x_j|c, \theta_{jc})) \\
&= \log(P(c|\boldsymbol{\pi})) + \sum_{j=1}^{784} \log(P(x_j|c, \theta_{jc})) - \log(\sum_c P(c|\boldsymbol{\pi}) \prod_{j=1}^{784} P(x_j|c, \theta_{jc})) \\
&= \log(\pi_c) + \sum_{j=1}^{784} (x_j \log \theta_{jc} + (1 - x_j) \log(1 - \theta_{jc})) \\
&\quad - \log \left(\sum_{i=0}^9 \pi_i \prod_{j=1}^{784} \theta_{jc^{(i)}}^{x_j} (1 - \theta_{jc^{(i)}})^{1-x_j} \right)
\end{aligned}$$

- c) Average log-likelihood for MLE is *Nan* due to numerical errors like division by zero or log 0 error. From these errors, there exists some division by zero or log 0 error due to some θ_{jc} being 0. This is due to data sparsity of input.



d)

- e) Using a $Beta(\alpha, \beta)$ prior on each θ_{jc} , we have $P(\theta_{jc}) = \theta_{jc}^{\alpha-1} (1 - \theta_{jc})^{\beta-1}$.

Find MAP estimator for θ , $\widehat{\theta}_{MAP} = \underset{\theta}{argmax} P(\theta|x, c, \pi) = \underset{\theta}{argmax} P(\theta|c)P(x|c, \theta)$

$$P(\theta|c)P(x|c, \theta) = \theta_{jc}^{\alpha-1} (1 - \theta_{jc})^{\beta-1} \prod_{i=1}^N \prod_{j=1}^{784} \theta_{jc}^{x_j^{(i)}} (1 - \theta_{jc})^{1-x_j^{(i)}}$$

$$\begin{aligned} \log(P(\theta|c)P(x|c, \theta)) &= (\alpha - 1)\log\theta_{jc} + (\beta - 1)\log(1 - \theta_{jc}) \\ &\quad + \sum_{i=1}^N \sum_{j=1}^{784} [x_j^{(i)}\log(\theta_{jc}) + (1 - x_j^{(i)})\log(1 - \theta_{jc})] \end{aligned}$$

Take the derivative wrt. θ_{jc}

$$\begin{aligned} 0 &= \frac{\alpha - 1 - \alpha\theta_{jc} + \theta_{jc} - \beta\theta_{jc} + \theta_{jc}}{\theta_{jc}(1 - \theta_{jc})} + \frac{1}{\theta_{jc}(1 - \theta_{jc})} \sum_{i=1}^N \mathbb{1}(c^{(i)} = c)(x_j^{(i)} - \theta_{jc}) \\ &= \alpha - 1 - \theta_{jc}(\alpha + \beta - 2) + \sum_{i=1}^N \mathbb{1}(c^{(i)} = c)x_j^{(i)} - \sum_{i=1}^N \mathbb{1}(c^{(i)} = c)\theta_{jc} \end{aligned}$$

$$\begin{aligned} MAP(\widehat{\theta}_{jc}) &= \frac{\sum_{i=1}^N \mathbb{1}(c^{(i)} = c)x_j^{(i)} + \alpha - 1}{\sum_{i=1}^N \mathbb{1}(c^{(i)} = c) + \alpha + \beta - 2} \\ &= \frac{N_{jc} + \alpha - 1}{N_c + \alpha + \beta - 2} \end{aligned}$$

(N_{jc} denotes # of class c images with j th pixel 1; N_c denotes # class c images)

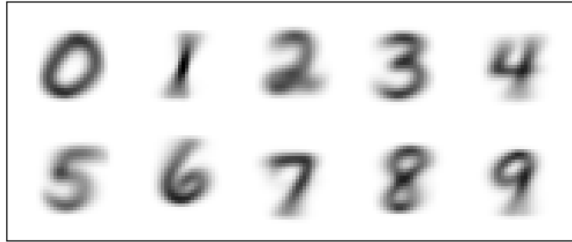
When $\alpha = 3, \beta = 3$,

$$MAP(\hat{\theta}_{jc}) = \frac{N_{jc} + 2}{N_c + 4}$$

In contrast to $\hat{\theta}_{jcMLE} = \frac{N_{jc}}{N_c}$, the numerator and denominator of the MAP estimator has more counts than the MLE, where α and β are can be think as pseudo-counts.

```
Average log-likelihood for MLE is nan
Average log-likelihood for MAP is -3.3570625208614815
Training accuracy for MAP is 0.8352166666666667
Test accuracy for MAP is 0.816
```

f)



g)

- h) The Naive Bayes assumption may be reasonable in this problem where images are independent. The Naive Bayes assumption may not be reasonable since in reality, it is almost impossible to have a set of completely independent predictors, which in this problem, the pixels may be dependent on other pixels.

Problem 3

a) From the question, we know

$$P(\mathbf{x}|\boldsymbol{\theta}) = \prod_{k=1}^K \theta_k^{x_k}$$

$$P(\boldsymbol{\theta}) \propto \theta_1^{a_1-1} \dots \theta_K^{a_K-1}$$

Then the posterior distribution can be calculated with

$$P(\boldsymbol{\theta}|\mathcal{D}) \propto P(\boldsymbol{\theta})P(\mathcal{D}|\boldsymbol{\theta})$$

$$\propto \theta_1^{a_1-1} \theta_2^{a_2-1} \dots \theta_K^{a_K-1} \prod_{i=1}^N \prod_{k=1}^K \theta_k^{x_k^{(i)}}$$

$$\propto \theta_1^{a_1-1} \theta_2^{a_2-1} \dots \theta_K^{a_K-1} \prod_{k=1}^K \theta_k^{\sum_{i=1}^N x_k^{(i)}}$$

$$\propto \theta_1^{a_1-1} \theta_2^{a_2-1} \dots \theta_K^{a_K-1} \prod_{k=1}^K \theta_k^{N_k}$$

$$\propto \theta_1^{N_1+a_1-1} \theta_2^{N_2+a_2-1} \dots \theta_K^{N_K+a_K-1}$$

We have computed the posterior distribution to be a Dirichlet distribution with

$$\boldsymbol{\theta} \sim \text{Dirichlet}(N_1 + a_1, N_2 + a_2, \dots, N_K + a_K)$$

Hence, the Dirichlet distribution is a conjugate prior for the categorical distribution.

b) The MAP estimate of $\boldsymbol{\theta}$ is

$$\hat{\boldsymbol{\theta}}_{MAP} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} P(\boldsymbol{\theta}|\mathcal{D}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \log(P(\boldsymbol{\theta}|\mathcal{D}))$$

Find $\hat{\theta}_{j_{MAP}}, \forall j = 1, 2, \dots, K-1$:

$$\log(P(\boldsymbol{\theta}|\mathcal{D})) = \log(\theta_1^{N_1+a_1-1} \theta_2^{N_2+a_2-1} \dots \theta_K^{N_K+a_K-1})$$

$$= (N_1 + a_1 - 1)\log(\theta_1) + (N_2 + a_2 - 1)\log(\theta_2) + \dots + (N_K + a_K - 1)\log(\theta_K)$$

$$= (N_1 + a_1 - 1)\log(\theta_1) + \dots + (N_K + a_K - 1)\log(1 - \sum_{k=1}^{K-1} \theta_k)$$

(Substitute $\theta_K = 1 - \sum_{k=1}^{K-1} \theta_k$)

$$\frac{\partial}{\partial \theta_j} \log(P(\boldsymbol{\theta}|\mathcal{D})) = \frac{N_j + a_j - 1}{\theta_j} - \frac{N_K + a_K - 1}{1 - \sum_{k=1}^{K-1} \theta_k}$$

$$0 = \frac{N_j + a_j - 1}{\theta_j} - \frac{N_K + a_K - 1}{\theta_K}$$

$$\frac{N_j + a_j - 1}{\theta_j} = \frac{N_K + a_K - 1}{\theta_K}$$

$$\implies \hat{\theta}_{j_{MAP}} = \hat{\theta}_{K_{MAP}} \frac{N_j + a_j - 1}{N_K + a_K - 1} \quad \forall j = 1, 2, \dots, K-1$$

We can now solve for $\hat{\theta}_{k_{MAP}}$

$$\begin{aligned}
\hat{\theta}_{k_{MAP}} &= 1 - \sum_{j=1}^{K-1} \hat{\theta}_{j_{MAP}} \\
&= 1 - \hat{\theta}_{k_{MAP}} \sum_{j=1}^{K-1} \frac{N_j + a_j - 1}{N_k + a_k - 1} \\
&= \frac{1}{1 + \sum_{j=1}^{K-1} \frac{N_j + a_j - 1}{N_k + a_k - 1}} \\
&= \frac{N_k + a_k - 1}{\sum_{j=1}^K (N_j + a_j - 1)}
\end{aligned}$$

Hence $\hat{\theta}_{k_{MAP}} = \frac{N_k + a_k - 1}{\sum_{j=1}^K (N_j + a_j - 1)}$, $\forall k = 1, 2, \dots, K$.

- c) To find the probability of \mathbf{x}^{N+1} being some class smaller than K . We first find the probability of \mathbf{x}^{N+1} being class k , i.e., $P(\mathbf{x}_k^{(N+1)} | \mathcal{D}) = \int_{\theta} P(\mathbf{x}_k^{(N+1)} | \theta) P(\theta | \mathcal{D}) d\theta$.

Since $\mathbf{x}^{(N+1)}$ is a 1-of- K encoding vector, thus we have $P(\mathbf{x}_k^{(N+1)} | \theta) = \theta_k$. Hence, the probability becomes $P(\mathbf{x}_k^{(N+1)} | \mathcal{D}) = \int_{\theta} \theta_k P(\theta | \mathcal{D}) d\theta$.

Computation for the probability of \mathbf{x}^{N+1} being some class smaller than K :

$$\begin{aligned}
P(\mathbf{x}^{(N+1)} < K) &= \sum_{k=1}^{K-1} P(\mathbf{x}_k^{(N+1)} = 1) \\
&= \sum_{k=1}^{K-1} P(\mathbf{x}_k^{(N+1)} | \mathcal{D}) \\
&= \sum_{k=1}^{K-1} \int_{\theta} \theta_k P(\theta | \mathcal{D}) d\theta \\
&= \sum_{k=1}^{K-1} \mathbb{E}[\theta_k] \quad (\text{since } \theta \sim \text{Dirichlet}(N_1 + a_1, \dots, N_K + a_K)) \\
&= \sum_{k=1}^{K-1} \frac{N_k + a_k}{\sum_{k'} (N_{k'} + a_{k'})} \quad (\theta \sim \text{Dirichlet}(N_1 + a_1, \dots, N_K + a_K))
\end{aligned}$$

Problem 4

a)

```
The average conditional log-likelihood on training set is -0.12462443666863014
The average conditional log-likelihood on test set is -0.19667320325525525
```

b)

```
The accuracy on training set is 0.9814285714285714
The accuracy on test set is 0.97275
```

c) The performance is worse compared with full-covariance matrix (lower likelihood and accuracy). Diagonal covariance matrix cannot model dependence between pixels