

## CSC311 HOMEWORK 4

Q1.

(1)

By Bayes' rule

$$\begin{aligned} p(y = k | \mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\sigma}) &= \frac{p(\mathbf{x} | y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})p(y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})}{p(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\sigma})} \\ &= \frac{p(\mathbf{x} | y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})p(y = k | \boldsymbol{\mu}, \boldsymbol{\sigma})p(\boldsymbol{\mu}, \boldsymbol{\sigma})}{p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\sigma})p(\boldsymbol{\mu}, \boldsymbol{\sigma})} \\ &= \frac{p(\mathbf{x} | y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})p(y = k | \boldsymbol{\mu}, \boldsymbol{\sigma})}{p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\sigma})} \end{aligned}$$

By the law of total probability,

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\sigma}) = \sum_{k'} p(\mathbf{x} | y = k', \boldsymbol{\mu}, \boldsymbol{\sigma}) \cdot p(y = k') = \sum_{k'} \alpha_{k'} p(\mathbf{x} | y = k', \boldsymbol{\mu}, \boldsymbol{\sigma})$$

And since  $p(y = k) = \alpha_k$  which doesn't depend on  $\boldsymbol{\mu}, \boldsymbol{\sigma}$ ,  $y = k$  and  $\boldsymbol{\mu}, \boldsymbol{\sigma}$  independent.

So  $p(y = k, \boldsymbol{\mu}, \boldsymbol{\sigma}) = p(y = k) \cdot p(\boldsymbol{\mu}, \boldsymbol{\sigma})$ .

$$p(y = k | \boldsymbol{\mu}, \boldsymbol{\sigma}) = \frac{p(y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})}{p(\boldsymbol{\mu}, \boldsymbol{\sigma})} = \frac{p(y = k) \cdot p(\boldsymbol{\mu}, \boldsymbol{\sigma})}{p(\boldsymbol{\mu}, \boldsymbol{\sigma})} = p(y = k)$$

So

$$\begin{aligned} p(y = k | \mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\sigma}) &= \frac{p(\mathbf{x} | y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})p(y = k | \boldsymbol{\mu}, \boldsymbol{\sigma})}{p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\sigma})} = \frac{p(\mathbf{x} | y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})p(y = k)}{\sum_{k'} \alpha_{k'} p(\mathbf{x} | y = k', \boldsymbol{\mu}, \boldsymbol{\sigma})} \\ &= \frac{\alpha_k p(\mathbf{x} | y = k, \boldsymbol{\mu}, \boldsymbol{\sigma})}{\sum_{k'} \alpha_{k'} p(\mathbf{x} | y = k', \boldsymbol{\mu}, \boldsymbol{\sigma})} \end{aligned}$$

(2)

$$\ell(\boldsymbol{\theta}; D) = -\log p(y^{(1)}, \mathbf{x}^{(1)}, \dots, y^{(N)}, \mathbf{x}^{(N)} | \boldsymbol{\theta})$$

since the data are iid,

$$\ell(\boldsymbol{\theta}; D) = -\log \prod_{i=1}^N p(y^{(i)}, \mathbf{x}^{(i)} | \boldsymbol{\theta})$$

By Bayes' rule

$$\begin{aligned} &= -\log \prod_{i=1}^N \frac{p(\mathbf{x}^{(i)} | y^{(i)}, \boldsymbol{\theta})p(y^{(i)}, \boldsymbol{\theta})}{p(\boldsymbol{\theta})} = -\log \prod_{i=1}^N \frac{p(\mathbf{x}^{(i)} | y^{(i)}, \boldsymbol{\theta})p(y^{(i)} | \boldsymbol{\theta})p(\boldsymbol{\theta})}{p(\boldsymbol{\theta})} \\ &= -\log \prod_{i=1}^N p(\mathbf{x}^{(i)} | y^{(i)}, \boldsymbol{\theta})p(y^{(i)} | \boldsymbol{\theta}) = -\sum_{i=1}^N \log p(\mathbf{x}^{(i)} | y^{(i)}, \boldsymbol{\theta}) + \log p(y^{(i)} | \boldsymbol{\theta}) \end{aligned}$$

Substituting terms given in question, we get

$$\begin{aligned}
\ell(\boldsymbol{\theta}; \mathbf{D}) &= - \sum_{i=1}^N \log \left( \left( \prod_{j=1}^D 2\pi\sigma_j^2 \right)^{-\frac{1}{2}} \exp \left\{ - \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(i)} - \mu_{y^{(i)}j})^2 \right\} \right) + \log \alpha_{y^{(i)}} \\
&= - \sum_{i=1}^N -\frac{1}{2} \sum_{j=1}^D \log 2\pi\sigma_j^2 - \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(i)} - \mu_{y^{(i)}j})^2 + \log \alpha_{y^{(i)}} \\
&= \frac{N}{2} \sum_{j=1}^D \log 2\pi\sigma_j^2 + \sum_{i=1}^N \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(i)} - \mu_{y^{(i)}j})^2 - \sum_{i=1}^N \log \alpha_{y^{(i)}} \\
&= \frac{N}{2} \sum_{j=1}^D \log 2\pi + \frac{N}{2} \sum_{j=1}^D \log \sigma_j^2 + \sum_{i=1}^N \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(i)} - \mu_{y^{(i)}j})^2 - \sum_{i=1}^N \log \alpha_{y^{(i)}}
\end{aligned}$$

(3)

From (2) we know that

$$\ell(\boldsymbol{\theta}; \mathbf{D}) = \frac{N}{2} \sum_{j=1}^D \log 2\pi + \frac{N}{2} \sum_{j=1}^D \log \sigma_j^2 + \sum_{i=1}^N \sum_{j=1}^D \frac{1}{2\sigma_j^2} (x_j^{(i)} - \mu_{y^{(i)}j})^2 - \sum_{i=1}^N \log \alpha_{y^{(i)}}$$

Therefore,

$$\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \mu_{ki}} = \frac{\partial}{\partial \mu_{ki}} \sum_{j=1}^N \sum_{i'=1}^D \frac{1}{2\sigma_{i'}^2} (x_{i'}^{(j)} - \mu_{y^{(j)}i'})^2$$

(here I change the i and j used in the previous summations to j and i' to avoid name confusion)

Since only terms that satisfy  $y^{(j)} = k$  and  $i' = i$  contribute to the derivative,

$$\begin{aligned}
&= \sum_{j=1}^N 1(y^{(j)} = k) \frac{-2}{2\sigma_i^2} (x_i^{(j)} - \mu_{ki}) \text{ where } 1(y^{(i)} = k) \text{ is an indicator function} \\
&= - \sum_{j=1}^N 1(y^{(j)} = k) \frac{1}{\sigma_i^2} (x_i^{(j)} - \mu_{ki})
\end{aligned}$$

Similarly,

$$\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \sigma_i^2} = \frac{\partial}{\partial \sigma_i^2} \left( \frac{N}{2} \sum_{i'=1}^D \log \sigma_{i'}^2 + \sum_{j=1}^N \sum_{i'=1}^D \frac{1}{2\sigma_{i'}^2} (x_{i'}^{(j)} - \mu_{y^{(j)}i'})^2 \right)$$

(here I change the i and j used in the previous summations to j and i' to avoid name confusion)

Since only terms that satisfy  $i' = i$  contribute to the derivative,

$$= \frac{N}{2\sigma_i^2} - \sum_{j=1}^N \frac{1}{2\sigma_i^4} (x_i^{(j)} - \mu_{y^{(j)}i})^2$$

Hence,

$$\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \mu_{ki}} = - \sum_{j=1}^N 1(y^{(j)} = k) \frac{1}{\sigma_i^2} (x_i^{(j)} - \mu_{ki})$$

$$\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \sigma_i^2} = \frac{N}{2\sigma_i^2} - \sum_{j=1}^N \frac{1}{2\sigma_i^4} (x_i^{(j)} - \mu_{y^{(j)}i})^2$$

(4)

To find the maximum likelihood estimate for  $\boldsymbol{\mu}$ , we set  $\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \mu_{ki}} = 0$  to find the maximum likelihood estimate for  $\mu_{ki}$ .

$$\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \mu_{ki}} = - \sum_{j=1}^N 1(y^{(j)} = k) \frac{1}{\sigma_i^2} (x_i^{(j)} - \widehat{\mu}_{ki}) = 0$$

$$\sum_{j=1}^N 1(y^{(j)} = k) (x_i^{(j)} - \widehat{\mu}_{ki}) = 0$$

$$\sum_{j=1}^N 1(y^{(j)} = k) x_i^{(j)} - \sum_{j=1}^N 1(y^{(j)} = k) \widehat{\mu}_{ki} = 0$$

Let  $n$  be the number of  $y^{(j)}$ s that are equal to  $k$  i.e. the number of data points whose class is assigned as  $k$ ,

$$\sum_{j=1}^N 1(y^{(j)} = k) x_i^{(j)} - n \widehat{\mu}_{ki} = 0$$

$$\widehat{\mu}_{ki} = \frac{1}{n} \sum_{j=1}^N 1(y^{(j)} = k) x_i^{(j)}$$

Therefore,

$$\widehat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{j=1}^N 1(y^{(j)} = k) \mathbf{x}^{(j)}$$

where  $k$  is each class and  $n$  is the number of data points whose class label is  $k$

To find the maximum likelihood estimate for  $\boldsymbol{\sigma}$ , we set  $\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \sigma_i^2} = 0$  to find the maximum likelihood estimate for  $\sigma_i^2$ .

$$\frac{\partial \ell(\boldsymbol{\theta}; \mathbf{D})}{\partial \sigma_i^2} = \frac{N}{2\widehat{\sigma}_i^2} - \sum_{j=1}^N \frac{1}{2\widehat{\sigma}_i^4} (x_i^{(j)} - \mu_{y^{(j)}i})^2 = 0$$

$$N\widehat{\sigma}_i^2 - \sum_{j=1}^N (x_i^{(j)} - \mu_{y^{(j)}i})^2 = 0$$

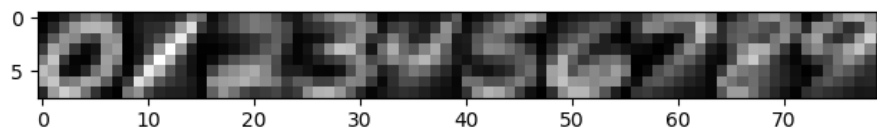
$$\widehat{\sigma}_i^2 = \frac{1}{N} \sum_{j=1}^N (x_i^{(j)} - \mu_{y^{(j)}i})^2$$

Therefore,

$$\hat{\sigma} = \frac{1}{N} \sum_{j=1}^N \left( x^{(j)} - \mu_{y^{(j)}} \right)^2$$

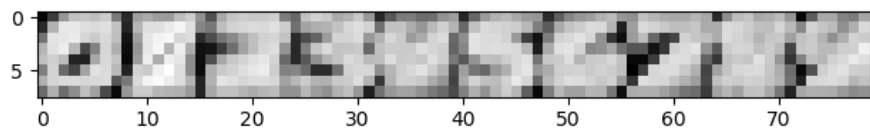
Q2.

2.0



2.1

(1)



(2)

The average conditional log-likelihood for train set:

-0.12462443666862928

The average conditional log-likelihood for test set:

-0.19667320325525448

(3)

The accuracy on the train set:

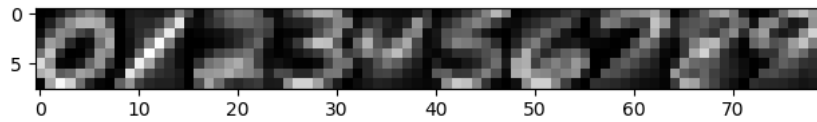
0.9814285714285714

The accuracy on the test set:

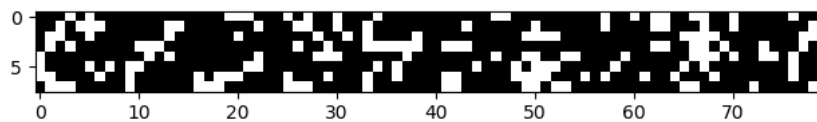
0.97275

2.2

(3)



(4)



(5)

The average conditional log-likelihood for train set:

-0.9437538618002541

The average conditional log-likelihood for train set:

-0.9872704337253588

(6)

The accuracy on the train set:

0.7741428571428571

The accuracy on the test set:

0.76425

2.3

Conditional Gaussian Classifier performs better than Naive Bayes Classifier which matches my expectation.

The accuracy of Naive Bayes Classifier highly depends on the assumption that the features are conditionally independent given the class. If the assumption is not correct, the accuracy of this classifier will be very low.

For this dataset, this assumption can't be hold. For example, for class 1, if a pixel is in the diagonal line and is set to ON, then its neighbors are also highly possible to be ON since that what a "1" looks like: pixels in the diagonal line is set to ON.

Therefore, the assumption can be matched. It is expected that Conditional Gaussian Classifier performs better than Naive Bayes Classifier.