

## Machine Learning Exercise Sheet 05

### Linear Classification

---

Exercise sheets consist of two parts: homework and in-class exercises. You solve the homework exercises on your own or with your registered group and upload it to Moodle for a possible grade bonus. The in-class exercises will be solved and explained during the tutorial. You do not have to upload any solutions of the in-class exercises.

---

### In-class Exercises

#### Multi-Class Classification

**Problem 1:** Consider a generative classification model for  $C$  classes defined by class probabilities  $p(y = c) = \pi_c$  and general class-conditional densities  $p(\mathbf{x} \mid y = c, \boldsymbol{\theta}_c)$  where  $\mathbf{x} \in \mathbb{R}^D$  is the input feature vector and  $\boldsymbol{\theta} = \{\boldsymbol{\theta}_c\}_{c=1}^C$  are further model parameters. Suppose we are given a training set  $\mathcal{D} = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$  where  $y^{(n)}$  is a binary target vector of length  $C$  that uses the 1-of- $C$  (one-hot) encoding scheme, so that it has components  $y_c^{(n)} = \delta_{ck}$  if pattern  $n$  is from class  $y = k$ . Assuming that the data points are i.i.d., show that the maximum-likelihood solution for the class probabilities  $\boldsymbol{\pi}$  is given by

$$\pi_c = \frac{N_c}{N}$$

where  $N_c$  is the number of data points assigned to class  $c$ .

The data likelihood given the parameters  $\{\pi_c, \boldsymbol{\theta}_c\}_{c=1}^C$  is

$$p(\mathcal{D} | \{\pi_c, \boldsymbol{\theta}_c\}_{c=1}^C) = \prod_{n=1}^N \prod_{c=1}^C (p(\mathbf{x}^{(n)} | \boldsymbol{\theta}_c) \pi_c)^{y_c^{(n)}}$$

and so the data log-likelihood is given by

$$\log p(\mathcal{D} | \{\pi_c, \boldsymbol{\theta}_c\}_{c=1}^C) = \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c + \text{const w.r.t. } \pi_c.$$

In order to maximize the log likelihood with respect to  $\pi_c$  we need to preserve the constraint  $\sum_c \pi_c = 1$ . For this we use the method of Lagrange multipliers where we introduce  $\lambda$  as an unconstrained additional parameter and find a local extremum of the unconstrained function

$$\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c - \lambda \left( \sum_{c=1}^C \pi_c - 1 \right).$$

instead. See wikipedia article on Lagrange multipliers for an intuition of why this works. This function is a sum of concave terms in  $\pi_c$  as well as  $\lambda$  and is therefore itself concave in these variables.

---

We can find the extremum by finding the root of the derivatives. Setting the derivative with respect to  $\pi_c$  equal to zero, we obtain

$$\pi_c = \frac{1}{\lambda} \sum_{n=1}^N y_c^{(n)} = \frac{N_c}{\lambda}.$$

Setting the derivative with respect to  $\lambda$  equal to zero, we obtain the original constraint

$$\sum_{c=1}^C \pi_c = 1$$

where we can now plug in the previous result  $\pi_c = \frac{N_c}{\lambda}$  and obtain  $\lambda = \sum_c N_c = N$ . Plugging this in turn into the expression for  $\pi_c$  we obtain

$$\pi_c = \frac{N_c}{N}$$

which we wanted to show.

## Linear Discriminant Analysis

**Problem 2:** Using the same classification model as in the previous question, now suppose that the class-conditional densities are given by Gaussian distributions with a *shared* covariance matrix, so that

$$p(\mathbf{x} \mid y = c, \boldsymbol{\theta}) = p(\mathbf{x} \mid \boldsymbol{\theta}_c) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_c, \boldsymbol{\Sigma}).$$

Show that the maximum likelihood estimate for the mean of the Gaussian distribution for class  $c$  is given by

$$\boldsymbol{\mu}_c = \frac{1}{N_c} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N \mathbf{x}^{(n)}$$

which represents the mean of the observations assigned to class  $c$ .

Similarly, show that the maximum likelihood estimate for the shared covariance matrix is given by

$$\boldsymbol{\Sigma} = \sum_{c=1}^C \frac{N_c}{N} \mathbf{S}_c \quad \text{where} \quad \mathbf{S}_c = \frac{1}{N_c} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)(\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T.$$

Thus  $\boldsymbol{\Sigma}$  is given by a weighted average of the sample covariances of the data associated with each class, in which the weighting coefficients  $N_c/N$  are the prior probabilities of the classes.

We begin by writing out the data log-likelihood.

$$\begin{aligned} & \log p(\mathcal{D} \mid \{\pi_c, \boldsymbol{\theta}_c\}_{c=1}^C) \\ &= \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \pi_c \cdot p(\mathbf{x}^{(n)} \mid \boldsymbol{\mu}_c, \boldsymbol{\Sigma}) \end{aligned}$$

Then we plug in the definition of the multivariate Gaussian

$$= \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \left( (2\pi)^{-\frac{D}{2}} \det(\mathbf{\Sigma})^{-\frac{1}{2}} \exp \left( -\frac{1}{2} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T \mathbf{\Sigma}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) \right) \right) + y^{(n)} \log \pi_c$$

and simplify.

$$= -\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left( D \log 2\pi + \log \det(\mathbf{\Sigma}) + (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T \mathbf{\Sigma}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) - 2 \log \pi_c \right)$$

This expression is concave in  $\boldsymbol{\mu}_c$ , so we can obtain the maximizer by finding the root of the derivative. With the help of the matrix cookbook, we identify the derivative with respect to  $\boldsymbol{\mu}_c$  as

$$\sum_{n=1}^N y_c^{(n)} \mathbf{\Sigma}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)$$

which we can set to 0 and solve for  $\boldsymbol{\mu}_c$  to obtain

$$\boldsymbol{\mu}_c = \frac{1}{\sum_{n=1}^N y_c^{(n)}} \sum_{n=1}^N y_c^{(n)} \mathbf{x}^{(n)} = \frac{1}{N_c} \sum_{\substack{n=1 \\ y^{(n)}=c}}^N \mathbf{x}^{(n)}.$$

To find the optimal  $\mathbf{\Sigma}$ , we need the trace trick

$$a = \text{Tr}(a) \text{ for all } a \in \mathbb{R} \quad \text{and} \quad \text{Tr}(\mathbf{ABC}) = \text{Tr}(\mathbf{BCA}).$$

With this we can rewrite

$$(\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T \mathbf{\Sigma}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) = \text{Tr} \left( \mathbf{\Sigma}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T \right)$$

and use the matrix-trace derivative rule  $\frac{\partial}{\partial \mathbf{A}} \text{Tr}(\mathbf{AB}) = \mathbf{B}^T$  to find the derivative of the data log-likelihood with respect to  $\mathbf{\Sigma}$ . Because the log-likelihood contains both  $\mathbf{\Sigma}$  and  $\mathbf{\Sigma}^{-1}$ , we convert one into the other with  $\log \det \mathbf{A} = -\log \det \mathbf{A}^{-1}$  to obtain

$$-\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left( -\log \det \mathbf{\Sigma}^{-1} + \text{Tr} \left( \mathbf{\Sigma}^{-1} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T \right) \right) + \text{const w.r.t. } \mathbf{\Sigma}.$$

Finally, we use rule (57) from the matrix cookbook  $\frac{\partial \log |\det \mathbf{X}|}{\partial \mathbf{X}} = (\mathbf{X}^{-1})^T$  and compute the derivative of the log-likelihood with respect to  $\mathbf{\Sigma}^{-1}$  as

$$-\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left( -\mathbf{\Sigma}^T + (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T \right).$$

We find the root with respect to  $\mathbf{\Sigma}$  and find

$$\mathbf{\Sigma} = \frac{1}{\sum_{n=1}^N \sum_{c=1}^C y_c^{(n)}} \left( \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T \right)^T = \frac{1}{N} \sum_{c=1}^C \sum_{\substack{n=1 \\ y^{(n)}=c}}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c) (\mathbf{x}^{(n)} - \boldsymbol{\mu}_c)^T$$

which we can immediately break apart into the representation in the instructions.