

## Machine Learning Exercise Sheet 05

### Linear Classification

Exercise sheets consist of two parts: In-class exercises and homework. The in-class exercises will be solved and discussed during the tutorial. The homework is for you to solve at home and further engage with the lecture content. There is no grade bonus and you do not have to upload any solutions. Note that the order of some exercises might have changed compared to last year's recordings.

## 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  $p(\mathbf{x} | y = c, \boldsymbol{\theta}_c)$  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  $c$ .

### 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

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

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

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

**Problem 1:** Consider a generative classification model for  $C$  classes defined by class probabilities  $p(y = c) = \pi_c$  and general class-conditional densities  $p(\mathbf{x} | y = c, \boldsymbol{\theta}_c)$  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  $c$ .

$$p(y^{(n)}=1 | \mathbf{x}) = \chi_c \quad p(\mathbf{x} | y=c, \boldsymbol{\theta}_c)$$

$$p(\mathbf{x}, y | \mathcal{X}, \boldsymbol{\theta}) = \prod_{n=1}^N p(\mathbf{x}^{(n)} | y^{(n)}, \boldsymbol{\theta}) \cdot p(y^{(n)} | \mathcal{X}) \quad \text{one-hot}$$

$$= \prod_{n=1}^N \prod_{c=1}^C \left[ p(\mathbf{x}^{(n)} | y_c^{(n)}, \boldsymbol{\theta}_c) p(y_c^{(n)} | \mathcal{X}_c) \right]^{y_c^{(n)}}$$

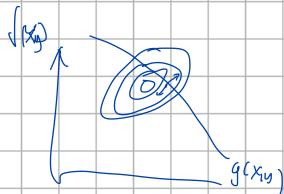
$$= \prod_{n=1}^N \prod_{c=1}^C \frac{p(\mathbf{x}^{(n)} | y_c^{(n)}, \boldsymbol{\theta}_c) p(y_c^{(n)} | \mathcal{X}_c)}{1 + \sum_{c' \neq c} p(\mathbf{x}^{(n)} | y_c^{(n)}, \boldsymbol{\theta}_{c'}) p(y_c^{(n)} = 1 | \mathcal{X}_c)}$$

log-likelihood

$$\log p(\mathbf{x}, y | \mathcal{X}, \boldsymbol{\theta}) = \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \log \chi_c + \text{const w.r.t. } \chi_c$$

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

$$\max \underbrace{\log p(\mathbf{x}, y | \mathcal{X}, \boldsymbol{\theta})}_{f} - \underbrace{\left[ \sum_{c=1}^C \chi_c - 1 \right]}_{g} = 0.$$



$$\begin{aligned} f(x, y) &= \lambda \circ g(x, y) \\ \Rightarrow f(x, y) - \lambda g(x, y) &= 0 \end{aligned}$$

$$\sum_{n=1}^N \frac{y_c^{(n)}}{\chi_c} - 1 = 0.$$

$$\sum_{c=1}^C \chi_c - 1 = 0.$$

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

$$\sum_{c=1}^C \frac{N_c}{x} = 1.$$

$$\chi_c = \frac{N_c}{x} = \frac{N_c}{N}$$

$$x = \sum_{c=1}^C N_c = N.$$

**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

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

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

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

$$\begin{aligned} p(y|x) &= p(x|y) \cdot p(y) & p(y=c) &= \lambda_c \\ p(x|y=c) &= \mathcal{N}(x | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}) \\ &= \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp\left(-\frac{1}{2} (x - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (x - \boldsymbol{\mu}_c)\right) \end{aligned}$$

$$\log p(y|x) = \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \cdot (\log \lambda_c + p(x, \boldsymbol{\theta}_c))$$

$$\begin{aligned} &= \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left[ \log \lambda_c + \left(-\frac{1}{2}\right) \left[ \log |\boldsymbol{\Sigma}| + (x - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (x - \boldsymbol{\mu}_c) \right] \right] \\ &= \sum_{n=1}^N \sum_{c=1}^C -\frac{1}{2} y_c^{(n)} \cdot \left[ -2 \log \lambda_c + D \cdot \log |\boldsymbol{\Sigma}| + (x - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (x - \boldsymbol{\mu}_c) \right] \\ &\quad \left[ (x^T \boldsymbol{\Sigma}^{-1} - \boldsymbol{\mu}_c^T \boldsymbol{\Sigma}^{-1}) (x - \boldsymbol{\mu}_c) \right]. \end{aligned}$$

$$\frac{\partial}{\partial \boldsymbol{\mu}_c} \sum_{n=1}^N -\frac{1}{2} y_c^{(n)} \left( -2 \mathbf{x}^T \boldsymbol{\Sigma}^{-1} + 2 \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_c \right) = 0.$$

$$\sum_{n=1}^N y_c^{(n)} \cdot 2 \boldsymbol{\Sigma}^{-1} \cdot \boldsymbol{\mu}_c = \sum_{n=1}^N y_c^{(n)} \cdot 2 \mathbf{x}^T \boldsymbol{\Sigma}^{-1}$$

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

$$\textcircled{2} \quad \text{Tr}(ABC) = \text{Tr}(B \cdot A) \quad \text{Tr}(a) = a.$$

$$(x - \mu_c)^T \Sigma^{-1} (x - \mu_c) = \text{Tr} \left[ \Sigma^{-1} (x - \mu_c) (x - \mu_c)^T \right].$$

$$\log |A| = -\log |A^{-1}|$$

$$\frac{d \log p}{d \bar{\varepsilon}} = \frac{d \log p}{d \bar{\varepsilon}^{-1}} \frac{d \bar{\varepsilon}^{-1}}{d \bar{\varepsilon}}$$

$$\frac{d \cdot \text{Tr}(AB)}{d A} = B^T$$

$$\bar{\Sigma}^{-1} = \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} \left[ -\log |A^{-1}| + \text{Tr} [\Sigma^{-1} (x - \mu_c) (x - \mu_c)^T] \right]$$

$$= \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} (-\bar{\Sigma}^{-1} + (x - \mu_c) (x - \mu_c)^T) = 0.$$

$$= N \bar{\Sigma}^{-1} + \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} (x - \mu_c) (x - \mu_c)^T = 0.$$

$$\bar{\Sigma}^{-1} = \frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_c^{(n)} (x - \mu_c) (x - \mu_c)^T$$

$$\bar{\Sigma}^{-1} = \frac{1}{N} \sum_{\substack{n=1 \\ y_c^{(n)}=1}}^N (x - \mu_c)^T (x - \mu_c)$$

**Problem 1:** Consider a generative classification model for  $C$  classes defined by class probabilities  $p(y = c) = \pi_c$  and general class-conditional densities  $p(\mathbf{x} | y = c, \boldsymbol{\theta}_c)$  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  $c$ .

Data likelihood

$$\begin{aligned} p(\mathbf{x} | \boldsymbol{\theta}) &= P(D | \{\pi_c, \boldsymbol{\theta}_c\}_{c=1}^C) = \prod_{n=1}^N p(\mathbf{x}^n, y^n) \\ &= \prod_{n=1}^N p(\mathbf{x}^n | y^n, \boldsymbol{\theta}) p(y^n | \boldsymbol{\pi}) \\ &= \prod_{n=1}^N \prod_{c=1}^C [p(\mathbf{x}^n | y^n=c, \boldsymbol{\theta}_c) \cdot p(y^n=c | \boldsymbol{\pi})]^{y_c^n} \end{aligned}$$

$$R^* = \underset{\boldsymbol{\pi}}{\operatorname{argmax}} \log p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{n=1}^N \sum_{c=1}^C (y_c^n \cdot \log \pi_c + C)$$

Assume  $\sum_c \pi_c = 1$

Lagrange

$$\mathcal{L} : \sum_{n=1}^N \sum_{c=1}^C (y_c^n \cdot \log \pi_c) + C - \lambda (\sum_{c=1}^C \pi_c - 1)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_c} = \frac{1}{\pi_c} \cdot y_c^n - \lambda = 0 \Rightarrow \pi_c = \frac{1}{\lambda} \sum_{n=1}^N y_c^n$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \sum_{c=1}^C \pi_c = 1$$

$$\sum_{c=1}^C \frac{N_c}{\lambda} = 1$$

$$\sum_{c=1}^C N_c = \lambda = N$$

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

$$\mathcal{L} = A \cdot \sum_{n=1}^N B_n$$

$$\frac{\partial \mathcal{L}}{\partial B_n} = \underbrace{\frac{\partial A}{\partial B_n} (B_1 + B_2 + B_3 + \dots + B_n)}_{\partial B_n}$$

$$n=1 \text{ 时. } \frac{\partial \mathcal{L}}{\partial B_1} = A \quad ; \text{ 与梯度无关}$$

$$n=2 \text{ 时. } \frac{\partial \mathcal{L}}{\partial B_2} = A.$$

$$n=N \text{ 时. } \frac{\partial \mathcal{L}}{\partial B_N} = A$$

对于  $N=1 \Rightarrow$   
 $y^n = (0 0 0 \dots 1 0)$   
 $\therefore [ \dots ]^{y^n} = \begin{cases} 1 & \text{if } y^n \text{ 有效} \\ 0 & \text{otherwise} \end{cases}$

**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

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

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.

$$\begin{aligned}
 & \text{data likelihood} \\
 \log p(D | \{\bar{x}_c, \bar{\mu}_c\}_{c=1}^C) &= \sum_{n=1}^N \sum_{c=1}^C \left( y_c^n \cdot \log \lambda_c + y_c^n \cdot \log p(x | y=c, \boldsymbol{\theta}) \right) \\
 &= \sum_{n=1}^N \sum_{c=1}^C y_c^n \cdot \left( -\frac{D}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) + \log \lambda_c \right) \\
 &= \sum_{n=1}^N \sum_{c=1}^C -\frac{1}{2} y_c^n \cdot (D \log 2\pi + \log |\boldsymbol{\Sigma}| + (\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c) - 2 \log \lambda_c) \\
 \frac{\partial}{\partial \boldsymbol{\mu}_c} &= \sum_{n=1}^N y_c^n \cdot (\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)) = 0 \\
 \underbrace{\frac{N}{N_c} \sum_{n=1}^N y_c^n \boldsymbol{\Sigma}^{-1} \mathbf{x}}_{\sum_n y_c^n \neq 0} &= \boxed{\boldsymbol{\mu}_c} = \frac{1}{N_c} \sum_{n=1}^N y_c^n \mathbf{x} \\
 \text{Tr}[(\mathbf{x} - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)] &= \text{Tr}[\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)^T (\mathbf{x} - \boldsymbol{\mu}_c)] \quad \frac{\partial \text{Tr}(AB)}{\partial A} = B^T \quad A = -\log |\boldsymbol{\Sigma}|^{-1} \\
 &= \sum_{n=1}^N \sum_{c=1}^C \frac{1}{2} y_c^n \cdot (D \log 2\pi + \log |\boldsymbol{\Sigma}| + \text{Tr}[\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)^T (\mathbf{x} - \boldsymbol{\mu}_c)] - 2 \log \lambda_c) \\
 \frac{\partial}{\partial \boldsymbol{\Sigma}} &= \frac{N}{2} \sum_{c=1}^C \frac{1}{2} y_c^n \left( -\log |\boldsymbol{\Sigma}| + \text{Tr}[\boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_c)^T (\mathbf{x} - \boldsymbol{\mu}_c)] \right) \quad \frac{\partial \log |\boldsymbol{\Sigma}|}{\partial \boldsymbol{\Sigma}} = (\boldsymbol{\Sigma}^{-1})^T \\
 &= -\frac{1}{2} \sum_{n=1}^N \sum_{c=1}^C y_c^n \left( -\boldsymbol{\Sigma}^T + (\mathbf{x} - \boldsymbol{\mu}_c)^T (\mathbf{x} - \boldsymbol{\mu}_c) \right) = 0 \\
 \boldsymbol{\Sigma}^T &= \frac{1}{\sum_{c=1}^C N_c} \sum_{n=1}^N \sum_{c=1}^C (y_c^n + (\mathbf{x} - \boldsymbol{\mu}_c)^T (\mathbf{x} - \boldsymbol{\mu}_c))
 \end{aligned}$$

Thus  $\Sigma$  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.

---

## Homework

### Linear classification

**Problem 3:** We want to create a generative binary classification model for classifying *non-negative* one-dimensional data. This means, that the labels are binary ( $y \in \{0, 1\}$ ) and the samples are  $x \in [0, \infty)$ .

We assume uniform class probabilities

$$p(y = 0) = p(y = 1) = \frac{1}{2}.$$

As our samples  $x$  are non-negative, we use exponential distributions (and not Gaussians) as class conditionals:

$$p(x | y = 0) = \text{Expo}(x | \lambda_0) \quad \text{and} \quad p(x | y = 1) = \text{Expo}(x | \lambda_1),$$

where  $\lambda_0 \neq \lambda_1$ . Assume, that the parameters  $\lambda_0$  and  $\lambda_1$  are known and fixed.

- a) Suppose you are given an observation  $x$ . What is the name of the posterior distribution  $p(y | x)$ ? You only need to provide the name of the distribution (e.g., “normal”, “gamma”, etc.), not estimate its parameters.
- b) What values of  $x$  are classified as class 1? (As usual, we assume that the classification decision is  $\hat{y} = \arg \max_k p(y = k | x)$ )

**Problem 4:** Let  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}$  be a linearly separable dataset for 2-class classification, i.e. there exists a vector  $\mathbf{w}$  such that  $\text{sign}(\mathbf{w}^T \mathbf{x})$  separates the classes. Show that the maximum likelihood parameter  $\mathbf{w}$  of a logistic regression model has  $\|\mathbf{w}\| \rightarrow \infty$ . Assume that  $\mathbf{w}$  contains the bias term.

How can we modify the training process to prefer a  $\mathbf{w}$  of finite magnitude?

**Problem 5:** Show that the softmax function is equivalent to a sigmoid in the 2-class case.

**Problem 6:** Show that the derivative of the sigmoid function  $\sigma(a) = (1 + e^{-a})^{-1}$  can be written as

$$\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a)).$$

**Problem 7:** Give a basis function  $\phi(x_1, x_2)$  that makes the data in the example below linearly separable (crosses in one class, circles in the other).

$$3$$

(1)  $p(y) \sim \text{Ber}(.)$   
 $p(y|x) \sim \text{Ber}$

(2)  $p(y=1|x) > p(y=0|x)$

$$\Rightarrow \frac{p(y=1|x)}{p(y=0|x)} > 1$$

$$\log \frac{p(y=1|x)}{p(y=0|x)} > 0$$

$$= \log \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)}$$

$$= \log \frac{p(x|y=1)}{p(x|y=0)} = \lambda_1 e^{-\lambda_1 x}$$

$$= \log \frac{p(x|y=0)}{p(x|y=1)} = \lambda_0 e^{-\lambda_0 x}$$

$$= \log \frac{\lambda_1}{\lambda_0} + \lambda_0 x - \lambda_1 x > 0 = \log \frac{\lambda_1}{\lambda_0}$$

$$(\lambda_0 - \lambda_1)x > \log \lambda_0 - \log \lambda_1$$

$$x > \frac{\log \lambda_0 - \log \lambda_1}{\lambda_0 - \lambda_1} \Rightarrow \left\{ \begin{array}{l} x \in \left( \frac{\log \lambda_0 - \log \lambda_1}{\lambda_0 - \lambda_1}, \infty \right) \\ x \in [0, \infty) \end{array} \right. \lambda_0 < \lambda_1$$

$$4$$

$E(w) = - \sum_{i=1}^N (y_i \log \sigma(w^T x_i) + (1-y_i) \log (1 - \sigma(w^T x_i)))$

$\sigma(a) = \frac{1}{1+e^{-a}}$

$w^* = \underset{w}{\operatorname{argmin}} E(w) = \frac{\partial E}{\partial w} = - \sum_{i=1}^N (y_i \cdot \frac{1}{\sigma(w^T x_i)} \cdot \frac{\partial \sigma(w^T x_i)}{\partial w} + (1-y_i) (1 - \frac{\partial \sigma(w^T x_i)}{\partial w})$

$$= - \sum_{i=1}^N (y_i (1 + e^{-w^T x_i})) \cdot (1 + e^{-w^T x_i})^{-2} \cdot e^{-w^T x_i}$$

$$+ (1-y_i) (1 + (1 + e^{-w^T x_i})^{-2} \cdot e^{-w^T x_i})$$

$$= - \sum_{i=1}^N (1 + e^{-w^T x_i})$$

$$E(w) = - \sum_{i=1}^N (y_i \log \sigma(w^T x_i) + (1-y_i) \log (1 - \sigma(w^T x_i)))$$

set  $w^T x_i > 0 \quad y_i = 1 \quad w^T x_i < 0, y_i = 0$

use scaling  $\lambda \gg 0$  make the negative log-likelihood smaller

$$\lim_{\lambda \rightarrow \infty} E(\lambda w) = - \sum_{i=1}^N \left( \lim_{\lambda \rightarrow \infty} \log \sigma(\lambda w^T x_i) + \lim_{\lambda \rightarrow \infty} (1 - \sigma(\lambda w^T x_i)) \right)$$

$\rightarrow E(w) \in [0, \infty)$

smallest  $\rightarrow 0$

我们可以看到  $E$  是一个凸函数，因为  $\log$  是凹的， $\sigma$  在  $a < 0$  时是凸的，在  $a > 0$  时是凹的。所以  $\log \sigma(a)$  在  $a > 0$  时是凹的， $\log (1 - \sigma(a))$  在  $a < 0$  时是凹的。

如果一个凸函数达到了它的最小值，那么它就有一个唯一的最小值。我们知道，当  $\lambda \rightarrow \infty$  时， $E$  趋于它的最小值，所以  $E$  不可能有一个有限的最小值，它的所有最小值只在极限中实现。由此可见，任何损失最小化问题的解决方案都具有无限的规范。

由于  $E$  是凸的，并且在某些方向上趋向于 0 的极限，我们可以通过添加任何实现其最小值的凸项，如  $w^T w$  或类似形式的权重正则化，将最小值移动到 finite 向量空间。

5

softmax function:  $\hat{f}(x)_i = \frac{\exp(x_i)}{\sum_{k=1}^C \exp(x_k)}$  LDA:

Sigmoid  $G(a) = \frac{1}{1 + \exp(-a)}$

$$= \frac{1}{1 + \exp(-w^T x - w_0)}$$

$$\hat{f}(a) = \frac{1}{1 + \exp(w^T x + w_0)}$$

$$f(a) = \frac{\exp(w^T x + w_0)}{1 + \exp(w^T x + w_0)} = \frac{\exp(w_0 + w_1 x + \dots)}{1 + \exp(w_0 + w_1 x + \dots)}$$

softmax  $\hat{f}(x_i) = \frac{\exp(w_i^T x)}{\exp(w_1^T x) + \exp(w_2^T x)}$

$$= \frac{1}{1 + \exp(w_2^T x - w_1^T x)}$$

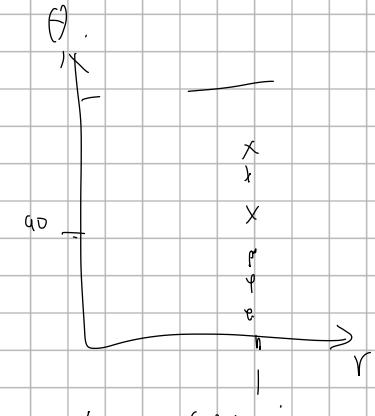
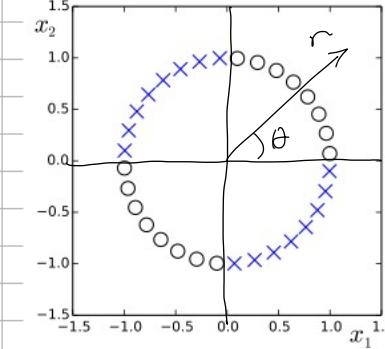
$$= \frac{1}{1 + \exp(-(w_1 - w_2)^T x)}$$

$$= G(\tilde{w}^T x)$$

$$\tilde{w} = w_1 - w_2$$

我们可以从中得出的一个结论是，如果我们有C类的参数向量  $w_c$ ，逻辑回归模型是不可识别的。这意味着在每个向量  $w_c := w_c + \tau$  中加入一个常数  $\tau \in \mathbb{R}$  将导致相同的逻辑回归模型。我们可以通过增加一个约束条件  $w_1 = 0$  来解决这个问题，这就是我们在二元分类中使用sigmoid（而不是2类softmax）时隐含的做法。

$$\begin{aligned}
 6. \frac{\partial g(u)}{\partial u} &= -\frac{1}{(1+e^{-u})^2} \cdot (-e^{-u}) \\
 &= \frac{1}{g^2(u)} \cdot \left(\frac{1}{g(u)} - 1\right) \\
 &= \frac{1}{(1+e^{-u})} \cdot \frac{e^{-u}}{(1+e^{-u})} \\
 &= g(u) \cdot \frac{1+e^{-u}-1}{(1+e^{-u})} \\
 &= g(u)(1-g(u))
 \end{aligned}$$



$$\phi(x) = (\theta)$$

$$= a_{1y} e^{(\theta)}$$

$$(x - \theta) \in (0, \pi)$$

一个例子是  $\phi(x) = x_1 x_2$ , 这使得数据可以通过超平面  $w = (1)$  来分离, 因为圆圈将被映射到正实数, 而十字架则去了负数, 也就是说, 如果  $x$  是一个圆圈,  $w^T x > 0$ , 否则  $w^T x < 0$ .

8

## Naive Bayes

**Problem 8:** In 2-class classification the decision boundary  $\Gamma$  is the set of points where both classes are assigned equal probability,

$$\Gamma = \{x \mid p(y=1|x) = p(y=0|x)\}.$$

Show that Naive Bayes with Gaussian class likelihoods produces a quadratic decision boundary in the 2-class case, i.e. that  $\Gamma$  can be written with a quadratic equation of  $x$ ,

$$\Gamma = \{x \mid x^T A x + b^T x + c = 0\},$$

for some  $A$ ,  $b$  and  $c$ .

As a reminder, in Naive Bayes we assume class prior probabilities

$$p(y=0) = \pi_0 \quad \text{and} \quad p(y=1) = \pi_1$$

and class likelihoods

$$p(x|y=c) = \mathcal{N}(x|\mu_c, \Sigma_c)$$

with per-class means  $\mu_c$  and *diagonal* (because of the feature independence) covariances  $\Sigma_c$ .

$$\frac{p(y=1|x)}{p(y=0|x)} = 1$$

$$\frac{p(x|y=1) \cdot p(y=1)}{p(x|y=0) \cdot p(y=0)} = 1$$

$$\log p(x|y=1) + \log \pi_1 - \log p(x|y=0) - \log \pi_0 = \log 1 = 0$$

$$\log N(x|\mu_1, \Sigma_1) - \log N(x|\mu_0, \Sigma_0) + \log \frac{\pi_1}{\pi_0} = 0$$

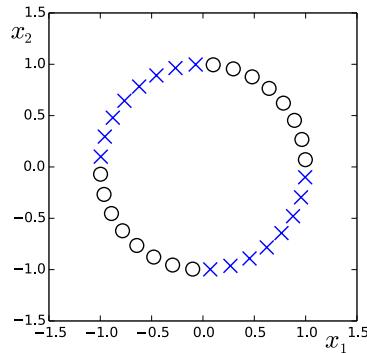
$$\begin{aligned}
 &= -\frac{1}{2} \log (2\pi)^D |\Sigma_1| - \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) + \frac{1}{2} \log (2\pi)^D |\Sigma_0| + \frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) + \log \frac{\pi_1}{\pi_0} = 0 \\
 &= (x^T \Sigma_0^{-1} - \mu_0^T \Sigma_0^{-1}) (x - \mu_0) - (x^T \Sigma_1^{-1} - \mu_1^T \Sigma_1^{-1}) (x - \mu_1) + 2 \log \frac{\pi_1}{\pi_0} + \frac{\log (2\pi)^D |\Sigma_0| - \log (2\pi)^D |\Sigma_1|}{=} \\
 &= \underline{x^T \Sigma_0^{-1} x} - \underline{2x^T \Sigma_0^{-1} \cdot \mu_0} + \underline{\mu_0^T \cdot \Sigma_0^{-1} \cdot \mu_0} - \underline{x^T \Sigma_1^{-1} x} + \underline{2x^T \Sigma_1^{-1} \cdot \mu_1} - \underline{\mu_1^T \cdot \Sigma_1^{-1} \cdot \mu_1} + 2 \log \frac{\pi_1}{\pi_0} + \dots - - - \\
 &= \underline{x^T (\Sigma_0^{-1} - \Sigma_1^{-1}) x} + \underline{x^T (2\Sigma_1^{-1} \mu_1 - 2\Sigma_0^{-1} \mu_0)} + \left( \underline{\mu_0^T \Sigma_0^{-1} \mu_0} - \underline{\mu_1^T \Sigma_1^{-1} \mu_1} + 2 \log \frac{\pi_1}{\pi_0} + \log (2\pi)^D |\Sigma_0| - \log (2\pi)^D |\Sigma_1| \right)
 \end{aligned}$$

A

b

c

for  $\Sigma_0 = \Sigma_1$ ,  $A \approx 0$ . Linear boundary



## Naive Bayes

**Problem 8:** In 2-class classification the decision boundary  $\Gamma$  is the set of points where both classes are assigned equal probability,

$$\Gamma = \{\mathbf{x} \mid p(y=1 \mid \mathbf{x}) = p(y=0 \mid \mathbf{x})\}.$$

Show that Naive Bayes with Gaussian class likelihoods produces a quadratic decision boundary in the 2-class case, i.e. that  $\Gamma$  can be written with a quadratic equation of  $\mathbf{x}$ ,

$$\Gamma = \{\mathbf{x} \mid \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c = 0\},$$

for some  $\mathbf{A}$ ,  $\mathbf{b}$  and  $c$ .

As a reminder, in Naive Bayes we assume class prior probabilities

$$p(y=0) = \pi_0 \quad \text{and} \quad p(y=1) = \pi_1$$

and class likelihoods

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

with per-class means  $\boldsymbol{\mu}_c$  and *diagonal* (because of the feature independence) covariances  $\boldsymbol{\Sigma}_c$ .