

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

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

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

P1 $p(y=c) = \pi_c$
 $p(x|y=c, \theta_c) \quad x \in D \quad \theta = \{\theta_c\}_{c=1}^C$
 $D = \{x^n, y^n\}_{n=1}^N$ $\propto p(x|y=c, \pi_c, \theta_c) \cdot p(y=c)$
 $p(D|\{\pi_c, \theta_c\}_{c=1}^C) = \prod_{n=1}^N \prod_{c=1}^C p(x^n|\theta_c) \pi_c^{y_c^n}$

$$\log p(D|\{\pi_c, \theta_c\}_{c=1}^C) = \sum_{n=1}^N \sum_{c=1}^C \left[y_c^n \log \pi_c + \frac{y_c^n \cdot \log p(x^n|\theta_c)}{N_c \pi_c} \right]$$

$\sum \pi_c = 1$ constraint

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

f.o.c $\frac{\partial}{\partial \pi_c} \frac{y_c^n}{\pi_c} - \lambda = 0 \Rightarrow \pi_c = \frac{1}{\lambda} \sum_{n=1}^N y_c^n = \frac{N_c}{\lambda}$

$$\sum_{c=1}^C \pi_c = 1 \Rightarrow \sum_{c=1}^C \frac{N_c}{\lambda} = 1 \quad \lambda = \frac{N_c}{N}$$

P2 $p(D|\{\pi_c, \theta_c\}) = \prod_{i=1}^M \prod_{c=1}^C (p(x^n|\theta_c) \cdot \pi_c)^{y_c^n} \sim N(x|\mu, \Sigma)$

$$\log p(D|\{\pi_c, \theta_c\}) = \sum_{i=1}^M \sum_{c=1}^C \left[y_c^n \ln \pi_c + y_c^n \ln p(x^n|\theta_c) \right]$$

$$= \sum_{i=1}^M \sum_{c=1}^C y_c^n \cdot \left[-\ln(2\pi)^{\frac{D}{2}} - \ln(|\Sigma|^{\frac{1}{2}}) - \frac{1}{2} (x^n - \mu_c)^T \Sigma^{-1} (x^n - \mu_c) \right]$$

$$\frac{d \log \cdot}{d \mu_c} = \sum_{i=1}^M y_c^n (x^n \Sigma^{-1} - \Sigma^{-1} \mu_c) = 0$$

$$\sum_{i=1}^M y_c^n \Sigma^{-1} (x^n - \mu_c) = 0$$

$$\sum_{i=1}^M y_c^n \Sigma^{-1} x^n = \sum_{i=1}^M y_c^n \Sigma^{-1} \mu_c$$

$$\mu_c = \frac{1}{\sum_{i=1}^M y_c^n} \sum_{i=1}^M y_c^n x^n = \frac{1}{N_c} \sum_{i=1}^M y_c^n x^n$$

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

$$\text{Trace}(ABC) = \text{Trace}(BCA)$$

$$\sum_{i=1}^M \sum_{c=1}^C y_c^n \cdot \left[-\ln(2\pi)^{\frac{D}{2}} - \ln(|\Sigma|^{\frac{1}{2}}) - \frac{1}{2} \text{Tr}[\Sigma^{-1} (x^n - \mu_c) (x^n - \mu_c)^T] \right]$$

$$\frac{d \cdot}{d \Sigma} = \sum_{i=1}^M \sum_{c=1}^C y_c^n \left[\frac{1}{2} \ln(|\Sigma|) - \frac{1}{2} (x^n - \mu_c) (x^n - \mu_c)^T \right] \quad \frac{\partial \log |\det x|}{\partial x} = (x^{-1})^T$$

$$= -\frac{1}{2} \sum_{i=1}^M \sum_{c=1}^C y_c^n \left[-\Sigma^{-1} + (x^n - \mu_c) (x^n - \mu_c)^T \right] = 0$$

$$\Sigma = \left(\frac{1}{\sum_{i=1}^M \sum_{c=1}^C y_c^n} \cdot \sum_{i=1}^M \sum_{c=1}^C y_c^n (x^n - \mu_c) (x^n - \mu_c)^T \right)^T$$

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

Thus Σ 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:

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

where $\lambda_0 \neq \lambda_1$. Assume, that the parameters λ_0 and λ_1 are known and fixed.

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

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).

P3 $p(y=c|x, \lambda) \propto p(x|y=c) p(y=c)$

Expo. Bernoulli Expo. uniform Σ'

classified as $c=1$.

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

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

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

$$\log \frac{\lambda_1 e^{-\lambda_1 x}}{\lambda_0 e^{-\lambda_0 x}} > 0$$

$$\ln \frac{\lambda_1}{\lambda_0} + (-\lambda_1 x) - (-\lambda_0 x) > 0$$

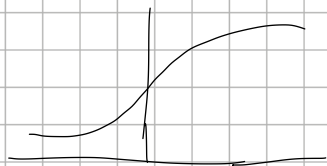
$$\ln \frac{\lambda_1}{\lambda_0} + (\lambda_0 - \lambda_1) x > \ln 1$$

$$(\lambda_0 - \lambda_1) x > \ln \frac{\lambda_0}{\lambda_1}$$

$$x > \frac{1}{\lambda_0 - \lambda_1} \ln \frac{\lambda_0}{\lambda_1}$$

$$\begin{cases} x \in \left(\frac{1}{\lambda_0 - \lambda_1} \ln \frac{\lambda_0}{\lambda_1}, +\infty \right), & \lambda_0 > \lambda_1 \\ x \in \left[0, \frac{1}{\lambda_0 - \lambda_1} \ln \frac{\lambda_0}{\lambda_1} \right), & \text{otherwise} \end{cases}$$

P4



it

P5 $p(y=1|x, w) = \frac{\exp(w_1^T x)}{\exp(w_1^T x) + \exp(w_2^T x)}$

softmax in 2-class

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

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

$$a = w_1^T x - w_2^T x = (w_1 - w_2)^T x$$

$$p(y=1|x, w) = \frac{1}{1 + \exp(-a)}$$

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

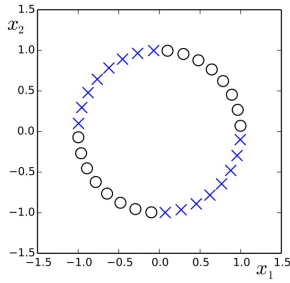
let $w_1 = 0$

p6 $G(a) = \frac{1}{1+e^{-a}} = (1-e^{-a})^{-1}$

$$\frac{\partial G(a)}{\partial a} = - (1+e^{-a})^{-2} \cdot e^{-a} \cdot (-1)$$

$$= \frac{e^{-a}}{(1+e^{-a})^2} = \frac{1}{1+e^{-a}} \cdot \frac{e^{-a}}{(1+e^{-a})} = \underline{\underline{G(a) \cdot (1-G(a))}}$$

p7



$$y = x_1 \cdot x_2$$

p8 $\frac{p(y=0|x)}{p(y=1|x)} = 1$

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

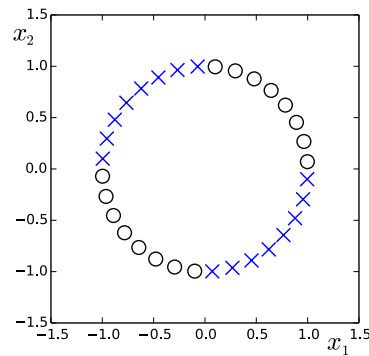
$$\log \frac{p(x|y=0) p(y=0)}{p(x|y=1) p(y=1)} = \log 1 = 0$$

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

~~$$-\ln \left(\frac{1}{(2\pi)^2} \right) - \ln \left(\frac{1}{|\Sigma_0|} \right) - \frac{1}{2} (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) + \ln \left(\frac{1}{(2\pi)^2} \right) + \ln \left(\frac{1}{|\Sigma_1|} \right) + \frac{1}{2} (x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) = \ln \frac{\pi_1}{\pi_0}$$~~

$$\frac{1}{2} \ln \left(\frac{|\Sigma_1|}{|\Sigma_0|} \right) - \frac{1}{2} \left[(x^T \Sigma_0^{-1} - \mu_0^T \Sigma_0^{-1}) (x - \mu_0) - (x^T \Sigma_1^{-1} - \mu_1^T \Sigma_1^{-1}) (x - \mu_1) \right] = \ln \frac{\pi_1}{\pi_0}$$

$$-\ln \left(\frac{1}{|\Sigma_0|} \right) - \frac{1}{2} \left[x^T \Sigma_0^{-1} x - 2x^T \Sigma_0^{-1} \mu_0 + \mu_0^T \Sigma_0^{-1} \mu_0 - x^T \Sigma_1^{-1} x + 2x^T \Sigma_1^{-1} \mu_1 - \mu_1^T \Sigma_1^{-1} \mu_1 \right]$$



Naive Bayes

Problem 8: In 2-class classification the decision boundary Γ 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 Γ 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 $\boldsymbol{\Sigma}_c$.