

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 .

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

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

Upload a single PDF file with your homework solution to Moodle by 09.12.2020, 23:59 CET. We recommend to typeset your solution (using L^AT_EX or Word), but handwritten solutions are also accepted. If your handwritten solution is illegible, it won't be graded and you waive your right to dispute that.

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.

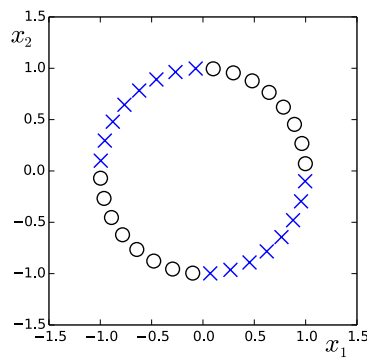
How can we modify the training process to prefer a \mathbf{w} of finite magnitude?

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).



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 $\boldsymbol{\Sigma}_c$.