

Statistical Machine Learning 2018

Exercises, week 5

5 October 2018

TUTORIAL

Exercise 1

Consider a discrete variable x that can take K values, $x \in \{1, \dots, K\}$. If we denote the probability of $x = k$ by the parameter θ_k , then the distribution of x is given by

$$P(x = k|\boldsymbol{\theta}) = \theta_k \quad (1)$$

where $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^T$ and the parameters are constrained to satisfy

$$\theta_k \geq 0 \quad \text{and} \quad \sum_{k=1}^K \theta_k = 1 \quad (2)$$

1. Explain why the parameters should satisfy these constraints.

Now consider a dataset χ of N independent observations, $\chi = \{x_1, \dots, x_N\}$.

2. Show that the log-likelihood $\ln P(\chi|\boldsymbol{\theta})$ is of the form

$$\ln P(\chi|\boldsymbol{\theta}) = \sum_{k=1}^K m_k \ln \theta_k \quad (3)$$

What are the m_k 's (in terms of the x_i 's, k 's etc.)?

3. Show that the maximum likelihood solution $\boldsymbol{\theta}^*$ is given by

$$\theta_k^* = \frac{m_k}{N} \quad (4)$$

Hint: Use a Lagrange multiplier for the constraint $\sum_{k=1}^K \theta_k - 1 = 0$.

Exercise 2

Suppose we have two coins, A and B, and we do not know whether these coins are fair.

1. Let μ be the probability the coin comes up H(eads). Give an expression for the likelihood of a data set \mathcal{D} of N observations of independent tosses of the coin.

Suppose we have observed the following results of two series of coin tosses:

coin:	data \mathcal{D} :
A	H,T,T,H,T,T,T
B	H

2. What is the maximum likelihood estimate for μ_A , the probability that a toss with coin A results in H(eads)? And for μ_B ? Based on these maximum likelihood estimates, what is the probability that the next toss of coin A will result in H(eads)? And the next toss with coin B? Do these results make sense?
3. Let us now take a Bayesian approach. Find an expression for $p(\mu|\mathcal{D})$ using Bayes' rule and show that a prior proportional to powers of μ and $(1 - \mu)$ will lead to a posterior that is also proportional to powers of μ and $(1 - \mu)$. Are you free to choose whatever prior you like?

Such a prior exists and is called the Beta distribution with hyperparameters a and b :

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad 0 \leq \mu \leq 1 \quad (5)$$

in which $\Gamma(x)$ is the gamma function with property $\Gamma(x+1) = x\Gamma(x)$.

4. Give combinations (a, b) for a prior that expresses: a) total ignorance, b) high confidence in a reasonably fair coin. For each prior and each coin, calculate the posterior probability density of μ given the observed coin tosses \mathcal{D} and plot the results (for example by using the `betapdf` command in `MatLab`). Do these results make more sense than the ML estimates?

Exercise 3

(Bishop 2.20) A symmetric $d \times d$ real-valued matrix $\mathbf{\Sigma}$ is positive definite if the quadratic form $\mathbf{a}^T \mathbf{\Sigma} \mathbf{a}$ is strictly positive for any non-zero real value of the vector \mathbf{a} , i.e. if:

$$\mathbf{a}^T \mathbf{\Sigma} \mathbf{a} > 0, \forall \mathbf{a} \in \mathbb{R}^d \setminus \{\mathbf{0}_d\}.$$

Using the above definition, show that a necessary and sufficient condition for $\mathbf{\Sigma}$ to be positive definite is that all of the eigenvalues $\lambda_i, i \in \{1, 2, \dots, d\}$, of $\mathbf{\Sigma}$ are strictly positive.

Exercise 4

(Exercise 2.34 in Bishop) Find the maximum likelihood solution for the covariance matrix of a multivariate Gaussian by maximizing the log likelihood function

$$\ln p(\mathbf{X}|\boldsymbol{\mu}, \mathbf{\Sigma}) = -\frac{ND}{2} \ln(2\pi) - \frac{N}{2} \ln |\mathbf{\Sigma}| - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu})$$

with respect to $\mathbf{\Sigma}$. In order to perform a straightforward maximization, ignore the constraints of symmetry and positive definiteness on $\mathbf{\Sigma}$, i.e. treat $\mathbf{\Sigma}$ as if it contained D^2 free parameters instead of just $\frac{D(D+1)}{2}$.

Hint: Use the results from Appendix C in Bishop to compute the matrix derivatives.

BONUS PRACTICE

Exercise 5

The beta distribution is

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1} \quad 0 \leq \mu \leq 1 \quad (6)$$

in which $\Gamma(x)$ is the gamma function (a well defined mathematical function, see book, www exercise 1.17). The gamma function is a generalization of the factorial function $(n-1)!$ as it satisfies

$$\Gamma(x+1) = x\Gamma(x) \quad (7)$$

We are looking for an expression for the expectation value in terms of a and b

$$\langle \mu \rangle = \int_0^1 \mu \text{Beta}(\mu|a, b) d\mu \quad (8)$$

Since the beta distribution is normalised, we can start from the relation

$$\int_0^1 \mu^{a-1} (1-\mu)^{b-1} d\mu = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad (9)$$

1. Show that the expectation value is given by

$$\langle \mu \rangle = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} \quad (10)$$

Hint: you do not actually have to compute any integrals.

2. Use this result and the property $\Gamma(x+1) = x\Gamma(x)$ to show that

$$\langle \mu \rangle = \frac{a}{a+b} \quad (11)$$

Exercise 6

Find the eigenvalues and a set of mutually orthogonal eigenvectors of the symmetric matrix:

$$\begin{bmatrix} 3 & 2 & 4 \\ 2 & 0 & 2 \\ 4 & 2 & 3 \end{bmatrix}$$

Exercise 7

In kernel methods, (symmetric) positive definite matrices play an important role. As mentioned on page 295 of the book, a positive definite matrix is not the same as a matrix in which all elements are positive. In Appendix C of the book, an example of a matrix is shown that has positive elements but that has a negative eigenvalue and hence that is not positive definite. Here we will look at the converse situation.

Question: Find a 2×2 symmetric matrix that is positive definite (in other words, has two positive eigenvalues), but with at least one **negative** element. Check that the eigenvalues indeed are positive!