

Systems analysis and decision making

Exercises – List No. 3

Authors: A. Gonczarek, J. Tomczak

Exemplary problem

A guest of a bar drinks x litres of a beer by his first visit, and y liters of a beer by the second visit. The following quantities were observed for three random customers of the bar:

x	1	2	3
y	1	1	3

Tabela 1: Data for three clients.

We aim at predicting how many liters y a client will drink for the second visit if he drank x liters for the first time.

Linear regression: least squares

Task 1: feature extraction

Calculate design matrix $\Phi = [\phi(\mathbf{x}_1) \ \phi(\mathbf{x}_2) \ \dots \ \phi(\mathbf{x}_N)]^T$ for given definition of features:

- $\phi(x) = (1 \ x \ x^2)^T$;
- $\phi(x) = \left(1 \ e^{-(x-2)^2} \ e^{-(x-3)^2}\right)^T$;

and data given in Table 1.

Task 2: model

Linear regression model is given in the following form:

$$\bar{y}(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w},$$

where $\mathbf{w} \in \mathbb{R}^D$ is a vector of parameters. Calculate (symbolically) a vector of outputs $\bar{\mathbf{y}}$ for each feature specified in Task 1.

Task 3: learning without regularisation

For the sum-of-squares error function:

$$E(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \bar{\mathbf{y}}\|_2^2$$

calculate:

- Task 3A: gradient $\nabla_{\mathbf{w}} E(\mathbf{w})$;
- Task 3B: values of parameters minimizing the objective $E(\mathbf{w})$;
- Task 3C: values of parameters for data given in Table 1 and chosen set of features.

Task 4: learning with regularisation

For the sum-of-squares error function with ℓ_2 regularization term (Tikhonov regularization):

$$E_{\lambda}(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \bar{\mathbf{y}}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

where $\lambda > 0$ is regularization coefficient, calculate:

- Task 4A: gradient $\nabla_{\mathbf{w}} E_{\lambda}(\mathbf{w})$;
- Task 4B: values of parameters minimizing the objective $E_{\lambda}(\mathbf{w})$;
- Task 4C: values of parameters for data given in Table 1 and chosen set of features.

Task 5: Prediction

For the set of features $\phi(x) = (1 \ x \ x^2)^T$ and the vector of parameters $\mathbf{w} = (1 \ 2 \ 1)^T$, calculate the output of the model for given input $x = 4$.

HANDOUT (LINEAR REGRESSION, LISTS 3 & 4)

Bernoulli distribution:

$$B(x|\theta) = \theta^x(1-\theta)^{1-x}, \quad \text{where } x \in \{0, 1\} \text{ i } \theta \in [0, 1]$$

$$\mathbb{E}[x] = \theta$$

$$\text{Var}[x] = \theta(1-\theta)$$

Categorical (Multinoulli) distribution:

$$M(\mathbf{x}|\boldsymbol{\theta}) = \prod_{d=1}^D \theta_d^{x_d}, \quad \text{where } x_d \in \{0, 1\} \text{ i } \theta_d \in [0, 1] \text{ for all } d = 1, 2, \dots, D, \sum_{d=1}^D \theta_d = 1$$

$$\mathbb{E}[x_d] = \theta_d$$

$$\text{Var}[x_d] = \theta_d(1-\theta_d)$$

Normal distribution:

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi} \sigma} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$

$$\mathbb{E}[x] = \mu$$

$$\text{Var}[x] = \sigma^2$$

Multivariate normal distribution:

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\},$$

where \mathbf{x} is D -dimensional vector, $\boldsymbol{\mu}$ – D -dimensional vector of means, $\boldsymbol{\Sigma}$ – $D \times D$ covariance matrix

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu}$$

$$\text{Cov}[\mathbf{x}] = \boldsymbol{\Sigma}$$

Beta distribution:

$$\text{Beta}(x|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1}(1-x)^{b-1},$$

where $x \in [0, 1]$ and $a > 0$ i $b > 0$, $\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$

$$\mathbb{E}[x] = \frac{a}{a+b}$$

$$\text{Var}[x] = \frac{ab}{(a+b)^2(a+b+1)}$$

Marginal distribution:

In the continuous case:

$$p(x) = \int p(x, y) dy$$

and in the discrete case:

$$p(x) = \sum_y p(x, y)$$

Conditional distribution:

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

Marginal distribution and conditional distribution for multivariate normal distribution:

Assume $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_a & \boldsymbol{\Sigma}_c \\ \boldsymbol{\Sigma}_c^T & \boldsymbol{\Sigma}_b \end{bmatrix},$$

then we get the following dependencies:

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a|\boldsymbol{\mu}_a, \boldsymbol{\Sigma}_a),$$

$$p(\mathbf{x}_a|\mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a|\hat{\boldsymbol{\mu}}_a, \hat{\boldsymbol{\Sigma}}_a), \text{ where}$$

$$\hat{\boldsymbol{\mu}}_a = \boldsymbol{\mu}_a + \boldsymbol{\Sigma}_c \boldsymbol{\Sigma}_b^{-1}(\mathbf{x}_b - \boldsymbol{\mu}_b),$$

$$\hat{\boldsymbol{\Sigma}}_a = \boldsymbol{\Sigma}_a - \boldsymbol{\Sigma}_c \boldsymbol{\Sigma}_b^{-1} \boldsymbol{\Sigma}_c^T.$$

Bayes' rule:

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

Maximum likelihood estimator:

There are given N independent examples of \mathbf{x} from the identical distribution $p(\mathbf{x}|\theta)$, $\mathcal{D} = \{\mathbf{x}_1 \dots \mathbf{x}_N\}$.

The likelihood function is the following function:

$$p(\mathcal{D}|\theta) = \prod_{n=1}^N p(\mathbf{x}_n|\theta).$$

The logarithm of the likelihood function $p(\mathcal{D}|\theta)$ is given by the following expression:

$$\log p(\mathcal{D}|\theta) = \sum_{n=1}^N \log p(\mathbf{x}_n|\theta).$$

Maximum likelihood estimator of the parameters θ_{ML} minimizes the likelihood function:

$$p(\mathcal{D}|\theta_{ML}) = \max_{\theta} p(\mathcal{D}|\theta).$$

Maximum *a posteriori* estimator:

There are given N independent examples of \mathbf{x} from the identical distribution $p(\mathbf{x}|\theta)$, $\mathcal{D} = \{\mathbf{x}_1 \dots \mathbf{x}_N\}$.

Maximum *a posteriori* (MAP) estimator of the parameters θ_{MAP} minimizes the *a posteriori* distribution:

$$p(\theta_{MAP}|\mathcal{D}) = \max_{\theta} p(\theta|\mathcal{D}).$$

Risk in the decision making:

Risk (expected loss) is defined as follows:

$$\mathcal{R}[\bar{y}] = \iint L(y, \bar{y}(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x} dy,$$

where $L(\cdot, \cdot)$ is the loss function.

Chosen properties of matrix calculus:

For given vectors \mathbf{x} , \mathbf{y} and a matrix \mathbf{A} , which is symmetric and positive definite, the following equations hold true:

- $\frac{\partial}{\partial \mathbf{y}} (\mathbf{x} - \mathbf{y})^T \mathbf{A} (\mathbf{x} - \mathbf{y}) = -2\mathbf{A} (\mathbf{x} - \mathbf{y})$
- $\frac{\partial (\mathbf{x} - \mathbf{y})^T \mathbf{A}^{-1} (\mathbf{x} - \mathbf{y})}{\partial \mathbf{A}} = -\mathbf{A}^{-1} (\mathbf{x} - \mathbf{y}) (\mathbf{x} - \mathbf{y})^T \mathbf{A}^{-1}$
- $\frac{\partial \ln \det(\mathbf{A})}{\partial \mathbf{A}} = \mathbf{A}^{-1}$