Statistical Machine Learning

Exercise sheet 1

Exercise 1.1 Classification from a discrete input space. We consider a multiclass classification problem with 3 classes $\{1,2,3\}$ for data with only a single discrete descriptor in $\mathcal{X} = \{1,2,3,4\}$.

We assume that the joint probability distribution $\mathbb{P}(Y = y, X = x)$ with X taking values in \mathcal{X} and Y taking values in $\mathcal{Y} = \{1, 2, 3\}$ is specified by the following table:

	Y = 1	Y = 2	Y = 3
X = 1	0,02	0,08	0,10
X=2	0,05	$0,\!40$	$0,\!15$
X = 3	0,02	0,02	$0,\!12$
X = 4	0,02	0,01	0,01

- (a) What is the target function f^* for the 0-1 loss?
- (b) What are the values of $f^*(x)$ for x = 1, 2, 3, 4.
- (c) What is the value of the risk for the target function?

Exercise 1.2 Recap of linear models. Let $y = \mathbf{X}\beta + \varepsilon$, where $\mathbb{E}(\varepsilon) = \mathbf{0}$, $\operatorname{var}(\varepsilon) = \sigma^2 \mathbf{I}$ and \mathbf{X} is a non-random full rank matrix of size $n \times p$. This setup contains the Gauss-Markov assumptions of a linear model.

- (a) Derive the least squares estimator $\hat{\boldsymbol{\beta}} = (\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}\mathbf{X}^{\mathsf{T}}\boldsymbol{y}$.
- (b) Show that $\hat{\boldsymbol{\beta}}$ is unbiased and that the variance of $\hat{\boldsymbol{\beta}}$ is given by $\sigma^2(\mathbf{X}^{\mathsf{T}}\mathbf{X})^{-1}$.

Exercise 1.3 Linear regression for binary classification. Consider a binary classification problem with $\mathcal{X} = \mathbb{R}^n$ and $\mathcal{Y} = \mathcal{A} = \{-1, 1\}$. We model the conditional expectation of Y given X by the linear model $\mathbb{E}(Y \mid X) = X^{\mathsf{T}}\beta$.

Let $\boldsymbol{x} \in \mathbb{R}^n$ be a new input. So, we estimate $\widehat{\mathbb{E}}(Y \mid \boldsymbol{X} = \boldsymbol{x}) = \boldsymbol{x}^{\top} \widehat{\boldsymbol{\beta}}$, where $\widehat{\boldsymbol{\beta}}$ is the least-square estimate of β . We wish to estimate its class $y = f^*(\boldsymbol{x})$, where f^* is the target function corresponding to 0 - 1 loss.

- (a) Derive the linear model estimate of $\widehat{\mathbb{P}}(Y = 1 \mid \boldsymbol{X} = \boldsymbol{x})$.
- (b) Show that $\widehat{y} = \widehat{f}^*(\boldsymbol{x})$ is given by $2 \cdot 1\{\boldsymbol{x}^{\top}\widehat{\boldsymbol{\beta}} \geq 0\} 1$, where \widehat{f}^* is the estimate of f^* given by plugging-in estimated values $\widehat{\mathbb{P}}(Y = y \mid \boldsymbol{X} = \boldsymbol{x})$ of the conditional p.m.f. $\mathbb{P}(Y = y \mid \boldsymbol{X} = \boldsymbol{x})$.

Exercise 1.4 Pinball loss and quantile regression. For $\tau \in]0,1[$, the pinball function with parameter τ is the function h_{τ} given by,

$$h_{\tau}(z) = -\tau z \, 1_{\{z \le 0\}} + (1 - \tau) z \, 1_{\{z > 0\}}.$$

We consider a decision problem for which inputs, outputs and actions are all real-valued, that is $\mathcal{X} = \mathcal{Y} = \mathcal{A} = \mathbb{R}$. For $a, y \in \mathbb{R}$, we define the pinball loss by $\ell_{\tau}(a, y) = h_{\tau}(a - y)$. We assume further that

- (a) $\mathbb{E}[|Y||X=x] < \infty$ a.e. $x \in \mathbb{R}$,
- (b) and the conditional law of Y given X is absolutely continuous with respect to the Lebesgue measure. Thus the function $y \mapsto \mathbb{P}(Y \le y \mid X = x)$ is continuous, a.e. $x \in \mathbb{R}$.

Recall that for a real-valued random variable Y whose law is absolutely continuous, we define the quantile of order α or α -quantile as the unique $q_{\alpha} \in \mathbb{R}$ such that $\mathbb{P}(Y \leq q_{\alpha}) = \alpha$. Similarly, the conditional quantile of order α of Y at X = x is, under the above continuity hypothesis the unique $q_{\alpha}(x) \in \mathbb{R}$ such that

$$\mathbb{P}(Y \le q_{\alpha}(x) \mid X = x) = \alpha.$$

- (a) Plot the pinball function in R. Play around with different values of τ . Why do you think the function is called that way?
- (b) Compute the expression for the conditional risk associated with the pinball loss in terms of q_{α} .
- (c) Prove that the target function of that risk is $q_{\tau}(x)$.
- (d) We call ℓ_1 -regression or least absolute deviation regression, the regression with loss function $\ell(a,y) = |a-y|$. Deduce from the previous question what is the target function for ℓ_1 -regression.

Practical exercises

Exercise 1.5 Polynomial regression. In this exercise, we will fit a linear model to data from simreg1train.csv. In R, use the read.csv("...") function to import the data.

- (a) Using results from Exercise 1.2, compute the least squares estimates for this dataset using your statistical software and plot the fitted values. Is the model appropriate?
- (b) Calculate the empirical risk on the training set (also called *training error*) for this dataset, given by

$$\widehat{\mathcal{R}}(\widehat{f}) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, \widehat{f}(x_i)), \tag{1}$$

where $\{(x_i, y_i)\}_{i=1}^n$ is the training set, ℓ is the squared error loss and \widehat{f} is the fitted function.

- (c) For the same loss function, calculate the empirical risk on the testing set (also called *testing error*) which is also given by (1) but here $\{(x_i, y_i)\}_{i=1}^n$ is the testing set given in simregltest.csv.
- (d) We now make the model more flexible by adding features to the design matrix \mathbf{X} . Add the feature \mathbf{x}^2 into your regression model, i.e., our design matrix becomes $\mathbf{X} = (\mathbf{1} \ \mathbf{x} \ \mathbf{x}^2)$. Compute the empirical risks on the training and testing sets for this model. Discuss.
- (e) Add features up to x^k into your regression model, for k = 3, 4, ..., 10. Calculate the the empirical risks on the training and testing sets for each k = 1, ..., 10. Make a plot of the empirical risks against k. Discuss. What happens when k > 10?