

**STATISTICAL MACHINE LEARNING (WS2017)**  
**SEMINAR 5**

**Assignment 1.** Consider the following parameter estimation task. You are given i.i.d. training data  $\mathcal{T}^m = \{x_i \in \mathbb{R} \mid i = 1, 2, \dots, m\}$  generated from the normal distribution

$$p_\mu(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

and the task is to estimate its unknown mean  $\mu$ .

**a)** Show that the maximum likelihood estimator is given by the arithmetic mean of the training data, i.e.

$$\mu^* = e_{ML}(\mathcal{T}^m) = \frac{1}{m} \sum_{i=1}^m x_i.$$

Prove that this estimator is unbiased.

**b)** Compute the variance of the maximum likelihood estimator, i.e.

$$\mathbb{E}_\mu [(\mu - e_{ML}(\mathcal{T}^m))^2].$$

How does it depend on  $\sigma$ ,  $\mu$  and  $m$ ?

**c)** Someone (let us call him Mr. Y) proposes an even simpler estimator - he suggests to estimate the unknown mean  $\mu$  of the distribution just by setting  $\mu^* = x_1$  and dropping the rest of the training data. Moreover, he claims that his estimator is also unbiased. Show that this is true. Why would you nevertheless consider his estimator inferior to the maximum likelihood estimator?

**Assignment 2.** Consider the exponential family

$$p_{\mathbf{u}}(x, y) = \frac{1}{Z(\mathbf{u})} \exp \langle \phi(x, y), \mathbf{u} \rangle$$

where  $\mathbf{u} \in \mathbb{R}^k$  is a parameter vector,  $\phi(x, y) \in \mathbb{R}^k$  is a generalised feature map,  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$  and  $Z(\mathbf{u}) = \sum_{x,y} \exp \langle \phi(x, y), \mathbf{u} \rangle$  is a normalising factor.

**a)** Prove that each model in this class is identifiable, provided that the affine hull of the set of vectors  $\{\phi(x, y) \mid x \in \mathcal{X}, y \in \mathcal{Y}\}$  is the entire space  $\mathbb{R}^k$  (or equivalently, there is no hyperplane containing all vectors).

**b)** Suppose that the parameter vector is bounded by  $\|\mathbf{u}\| \leq R$  and assume that the components of the vectors  $\phi(x, y)$  are bounded in some interval  $[a, b]$ . Prove the Uniform Law of Large Numbers for the corresponding Maximum Likelihood Estimator by performing the following steps

- (1) Denote the training data by  $\mathcal{T}^m = \{(x^i, y^i) \mid i = 1, 2, \dots, m\}$ , the log-likelihood of  $\mathcal{T}^m$  by  $L(\mathbf{u}, \mathcal{T}^m)$  and the expected log-likelihood by  $L(\mathbf{u}) = \mathbb{E}_v L(\mathbf{u}, \mathcal{T}^m)$ , where  $v \in \mathbb{R}^k$  is the true but unknown model.
- (2) Deduce that

$$L(\mathbf{u}, \mathcal{T}^m) - L(\mathbf{u}) = \langle \mathbb{E}_{\mathcal{T}^m} \phi - \mathbb{E}_v(\phi), \mathbf{u} \rangle$$

holds, where  $\mathbb{E}_{\mathcal{T}^m} \phi$  denotes the arithmetic mean of the vectors  $\phi(x, y)$  on the training data and  $\mathbb{E}_v(\phi)$  denotes their expectation w.r.t. the true model.

- (3) Prove that

$$\max_{\|\mathbf{u}\| \leq R} |L(\mathbf{u}, \mathcal{T}^m) - L(\mathbf{u})| = \|\mathbb{E}_{\mathcal{T}^m} \phi - \mathbb{E}_v(\phi)\| R$$

holds.

- (4) Conclude the ULLN for MLE-s in this model class.

c) Prove that the logarithm of the probability

$$\log p_{\mathbf{u}}(x, y) = \langle \phi(x, y), \mathbf{u} \rangle - \log Z(\mathbf{u})$$

is a concave function of  $\mathbf{u}$  by verifying the following steps.

- (1) prove the formulas  $\nabla_{\mathbf{u}} \langle \mathbf{a}, \mathbf{u} \rangle = \mathbf{a}$  and  $\nabla_{\mathbf{u}} f(g(\mathbf{u})) = f'(g(\mathbf{u})) \cdot \nabla_{\mathbf{u}} g(\mathbf{u})$  for the gradient by considering partial derivatives. Here,  $f$  is a function of a real variable,  $f'$  denotes its derivative and  $\mathbf{a}$  denotes a constant vector.
- (2) Prove that the gradient of  $\log Z(\mathbf{u})$  is

$$\nabla_{\mathbf{u}} \log Z(\mathbf{u}) = \sum_{x,y} p_{\mathbf{u}}(x, y) \phi(x, y) = \mathbb{E}_{\mathbf{u}}(\phi).$$

- (3) Prove that the second derivative of  $\log Z(\mathbf{u})$  is

$$\begin{aligned} \nabla_{\mathbf{u}}^2 \log Z(\mathbf{u}) &= \sum_{x,y} p_{\mathbf{u}}(x, y) \phi(x, y) \otimes \phi(x, y) - \mathbb{E}_{\mathbf{u}}(\phi) \otimes \mathbb{E}_{\mathbf{u}}(\phi) = \\ &= \mathbb{E}_{\mathbf{u}}[(\phi - \mathbb{E}_{\mathbf{u}}(\phi)) \otimes (\phi - \mathbb{E}_{\mathbf{u}}(\phi))] \end{aligned}$$

- (4) Deduce that the second derivative is positive semi-definite matrix and conclude that  $\log Z(\mathbf{u})$  is convex.

**Assignment 3.** Consider the block-coordinate ascent for Minka's lower bound of the log-likelihood (in the EM algorithm)

$$L_B(\theta, \alpha, \mathcal{T}^m) = \frac{1}{m} \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \alpha_i(y) \log p_{\theta}(x^i, y) - \frac{1}{m} \sum_{i=1}^m \sum_{y \in \mathcal{Y}} \alpha_i(y) \log \alpha_i(y).$$

The E-step requires to maximise it w.r.t.  $\alpha$ -s for fixed  $\theta$ .

**a)** Show that the maximisation decomposes into independent maximisation tasks of the type

$$\sum_{y \in \mathcal{Y}} \alpha(y) \log p_{\theta}(x^i, y) - \sum_{y \in \mathcal{Y}} \alpha(y) \log \alpha(y) \rightarrow \max_{\alpha}$$

s.t.  $\sum_{y \in \mathcal{Y}} \alpha(y) = 1$  and  $\alpha(y) \geq 0 \quad \forall y \in \mathcal{Y}$

**b)** Prove that the function

$$g(x) = \begin{cases} x \log x & \text{if } x > 0, \\ 0 & \text{if } x = 0 \\ \infty & \text{otherwise} \end{cases}$$

is convex for  $x > 0$  and implicitly accounts for the constraint  $x \geq 0$ . Conclude that the optimisation task in a) is convex.

**c)** Analyse the Lagrange dual of this task and deduce the solution  $\alpha^*(y) = p_{\theta}(y \mid x^i)$ .