

CS-E4830 Kernel methods in machine learning, exam date 03.06.2021 / Examiner: Rohit Babbar

Instructions: You have 3 hours to complete exam. Scan (or take a picture) and send your answer sheets by 12:30pm today (03.06.2021). This is an open book exam but no additional material apart from the lecture videos/slides can be used. Consulting others to write your answers is not allowed. There are 10 questions for a total maximum of 50 points.

Questions

Q.1 (8 points) Give short (a few sentences) definitions or appropriate description of the following concepts.

- (a) Kernel functions
- (b) Empirical and expected error
- (c) Bias-variance tradeoff
- (d) Multiple Kernel Learning

Q.2 (4 points) Explain the computational advantages of using a polynomial kernel of degree two as compared to using bigram features. Under what conditions using the features directly might be more beneficial?

Q.3 (6 points in total) Assume we have the kernels $k_m(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi_m(\mathbf{x}_i), \phi_m(\mathbf{x}_j) \rangle$, $m = 1, \dots, P$ at our disposal, where $\phi_m(\mathbf{x}) = (\phi_{1m}(\mathbf{x}), \dots, \phi_{Dm}(\mathbf{x}))^T \in \mathbb{R}^D$ is the feature vector underlying the kernel k_m .

For each kernel below, write down the equation for the underlying feature vector $\tilde{\phi}_s(\mathbf{x})$, as a function of the feature vectors ϕ_m , $m = 1, \dots, P$, so that $\tilde{k}_s(\mathbf{x}_i, \mathbf{x}_j) = \langle \tilde{\phi}_s(\mathbf{x}_i), \tilde{\phi}_s(\mathbf{x}_j) \rangle$ is satisfied for each $s \in \{a, b\}$.

- (a) (2 points) $\tilde{k}_a(\mathbf{x}_i, \mathbf{x}_j) = \sum_{m=1}^P k_m(\mathbf{x}_i, \mathbf{x}_j)$
- (b) (4 points) $\tilde{k}_b(\mathbf{x}_i, \mathbf{x}_j) = (k_1(\mathbf{x}_i, \mathbf{x}_j) + 1)^2$

Q.4 (4 points) Check if $K(x, x') = \max(x, x')$ such that $x, x' \in \mathbb{R}^+$ is a valid kernel or not. If yes, prove it; give a counter-example otherwise.

Q.5 (5 points) State Representer theorem and discuss its implications for computing the prediction function values at training points $f(x_i)$ and regularizer $\|f\|_{\mathcal{H}}^2$ for solving ERM problems such as Kernel SVM and Kernel logistic regression.

Q.6 (3 points) In the statistical learning theory framework, which of the following is same for all functions in a function class, and which one changes across different functions in that class? Explain your answer in detail.

- (a) Estimation error
- (b) Approximation error

Q.7 (5 points) Recall the formulation for Kernel Logistic Regression

$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^N} \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(-y_i [K\boldsymbol{\alpha}]_i)) + \frac{\lambda}{2} \boldsymbol{\alpha}^T K \boldsymbol{\alpha}$$

Show that the objective function is convex in $\boldsymbol{\alpha}$.

Q.8 (6 points) The primal optimization problem for linear SVM formulation with squared Hinge loss $(\max(0, 1 - y\mathbf{w}^T \mathbf{x}))^2$ as the loss function is given by

$$\begin{aligned} \min_{\mathbf{w}, \xi} \quad & \frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^N \xi_i^2 \\ \text{s.t.} \quad & y_i(\mathbf{w}^T \mathbf{x}_i) \geq 1 - \xi_i, \quad i = 1, \dots, N \end{aligned}$$

Using the method of Lagrange multipliers, derive the dual of the above problem.

Q.9 (4 points) Write the formulation of Principal Component Analysis and show how it is related to eigen value problem involving co-variance matrix. Is the optimization problem convex. Explain your answer.

Q.10 (5 points) State Bochner theorem and explain how it can be used for addressing machine learning problems with large number of training samples in the context of kernel methods.