

EPFL | MGT-418 : Convex Optimization | Project 2

Learning the Kernel Function (graded)

(This project is due on **November 24, 2021, at 23:59**. You may form teams of up to two people. Each team should upload a single zip-file containing their report and code to Moodle. Make sure to clearly state the team members in your report.)

Description

Given training samples (\mathbf{x}_i, y_i) , $i = 1, \dots, m$, consisting of inputs $\mathbf{x}_i \in \mathbb{R}^d$ and outputs $y_i \in \{-1, 1\}$,¹ we aim to predict the outputs of test samples based on their inputs by using a linear classifier. As the training samples may not be linearly separable, we use a feature map $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}^D$ that lifts the inputs to a higher-dimensional space \mathbb{R}^D , and we solve the soft-margin support vector machine problem in the lifted space. Specifically, we solve the problem

$$\begin{aligned} & \underset{\mathbf{w} \in \mathbb{R}^D, \mathbf{s} \in \mathbb{R}^m, b \in \mathbb{R}}{\text{minimize}} && \sum_{i=1}^m s_i + \frac{\rho}{2} \|\mathbf{w}\|_2^2 \\ & \text{subject to} && y_i(\mathbf{w}^\top \Phi(\mathbf{x}_i) + b) + s_i \geq 1 \quad \forall i = 1, \dots, m \\ & && s_i \geq 0 \quad \forall i = 1, \dots, m \end{aligned} \quad (1)$$

for some $\rho > 0$. The corresponding dual problem (see lecture slide 183) can be expressed as

$$\begin{aligned} & \underset{\boldsymbol{\lambda} \in \mathbb{R}^m}{\text{maximize}} && \sum_{i=1}^m \lambda_i - \frac{1}{2\rho} \sum_{i,j=1}^m \lambda_i \lambda_j y_i y_j k(\mathbf{x}_i, \mathbf{x}_j) \\ & \text{subject to} && \sum_{i=1}^m y_i \lambda_i = 0 \\ & && 0 \leq \lambda_i \leq 1 \quad \forall i = 1, \dots, m, \end{aligned} \quad (2)$$

where $k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^\top \Phi(\mathbf{x}')$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$ denotes the kernel function. In contrast to the primal problem (1), the dual problem (2) can be solved without knowledge of the feature map Φ . Instead, it suffices to know the kernel function. Gaussian, polynomial, exponential, and tangential kernel functions are commonly used, but it is not always clear which function works best for a given problem. In practice, kernel functions are often chosen through cross-validation, that is, by solving problem (2) repeatedly for multiple kernels and regularization parameters. In this project, instead, we formulate a single optimization problem that determines both the optimal classifier and the optimal kernel function. To this end, we first choose L candidate kernel functions $\hat{k}^1, \dots, \hat{k}^L$, and then we form the candidate kernel matrices $\hat{K}^1, \dots, \hat{K}^L$ corresponding to the training samples, that is, we set $\hat{K}_{ij}^l = \hat{k}^l(\mathbf{x}_i, \mathbf{x}_j)$ for all $i, j = 1, \dots, m$ and $l = 1, \dots, L$. Next, we define the set

$$\mathcal{K} = \left\{ K \in \mathbb{S}_+^m : K = \sum_{l=1}^L \mu_l \hat{K}^l, K \succeq 0, \text{tr}(K) = c, \mu_l \geq 0 \quad \forall l = 1, \dots, L \right\}$$

of conic combinations of the candidate kernel matrices $\hat{K}^1, \dots, \hat{K}^L$, whose trace equals $c > 0$. We now express the optimal value of problem (2) in vectorized form as a function ω of the kernel matrix K ,

$$\omega(K) = \begin{cases} \max_{\boldsymbol{\lambda} \in \mathbb{R}^m} & \boldsymbol{\lambda}^\top \mathbf{1} - \frac{1}{2\rho} \boldsymbol{\lambda}^\top G(K) \boldsymbol{\lambda} \\ \text{s.t.} & \boldsymbol{\lambda}^\top \mathbf{y} = 0 \\ & 0 \leq \lambda_i \leq 1 \quad \forall i = 1, \dots, m, \end{cases} \quad (3)$$

where $G(K) \in \mathbb{S}^m$ is the matrix with entries $G_{ij}(K) = K_{ij} y_i y_j$ for all $i, j = 1, \dots, m$. We may thus determine both the optimal classifier and the optimal kernel by solving the optimization problem

$$\underset{K \in \mathcal{K}}{\text{minimize}} \quad \omega(K). \quad (4)$$

¹We assume that there is at least one sample with output -1 and at least one sample with output 1 .

Questions

- (10 points) Explain why every kernel matrix induced by any feature map Φ is positive semidefinite.
- (25 points) Show that the minimization problem in (4) is equivalent to the convex maximization problem

$$\begin{aligned} & \underset{\lambda \in \mathbb{R}^m, z \in \mathbb{R}}{\text{maximize}} && \lambda^T \mathbf{1} - cz \\ & \text{subject to} && z\hat{r}_l \geq \frac{1}{2\rho} \lambda^T G(\hat{K}^l) \lambda \quad l = 1, \dots, L \\ & && \lambda^T \mathbf{y} = 0 \\ & && 0 \leq \lambda_i \leq 1 \quad i = 1, \dots, m, \end{aligned} \quad (5)$$

where the vector $\hat{\mathbf{r}} \in \mathbb{R}^L$ has entries $\hat{r}_l = \text{tr}(\hat{K}^l)$ for all $l = 1, \dots, L$. *Hint: You may want to use the following minimax theorem. If some function $f : U \times V \rightarrow \mathbb{R}$ is convex in $u \in U$ and concave in $v \in V$, and if U and V are convex, closed, and bounded sets, then*

$$\min_{u \in U} \max_{v \in V} f(u, v) = \max_{v \in V} \min_{u \in U} f(u, v).$$

- (20 points) When solving problem (5) numerically, most solvers will not only determine the optimal primal decisions λ^* and z^* but also the optimal dual decisions. For $l = 1, \dots, L$, let μ_l^* denote the optimal dual variable of the constraint $z\hat{r}_l \geq \frac{1}{2\rho} \lambda^T G(\hat{K}^l) \lambda$. We will show that $K^* = \sum_{l=1}^L \mu_l^* \hat{K}^l$ solves problem (4) by first reformulating problem (5) as

$$\begin{aligned} & \underset{\lambda \in \mathbb{R}_+^m, \lambda^T \mathbf{y} = 0, \lambda \leq \mathbf{1}}{\text{max}} && \underset{z \in \mathbb{R}}{\text{max}} && \lambda^T \mathbf{1} - cz \\ & \text{s.t.} && && z\hat{r}_l \geq \frac{1}{2\rho} \lambda^T G(\hat{K}^l) \lambda \quad l = 1, \dots, L. \end{aligned} \quad (6)$$

- (10 points) Derive the Lagrangian dual of the inner maximization problem in (6). Explain why strong duality holds. Denote the dual variables by $\mu \in \mathbb{R}^L$.
 - (10 points) Reformulate the max-min problem derived in (a) as a min-max problem by using the minimax theorem from (2) and denote its solution by μ^* . Show that $K^* = \sum_{l=1}^L \mu_l^* \hat{K}^l$ solves (4).
4. **Radar Classification (30 points):** A lab in Goose Bay, Labrador, analyzes the ionosphere through radar signals. All signals are classified as either suitable or unsuitable for further analysis. We are given 351 classified signals with 33 features and asked to find a linear classifier for automatic signal classification.
- (5 points) Read the data file `ionosphere.data` into memory by using the scripts `read_data.py` or `read_data.m`. Use the code skeletons `main.ipynb` or `main.m` to randomly select 80% of the data for training.
 - (15 points) Define $\hat{k}^1(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^p$ as the polynomial, $\hat{k}^2(\mathbf{x}, \mathbf{x}') = \exp(-(\mathbf{x} - \mathbf{x}')^\top (\mathbf{x} - \mathbf{x}') / (2\sigma))$ as the Gaussian, and $\hat{k}^3(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$ as the linear kernel function, and construct the kernel matrices \hat{K}^l , $l = 1, 2, 3$, for all training samples. Solve the QCQP in (5) for $\rho = 2$, $p = 2$, $\sigma = 0.5$ and $c = \sum_{l=1}^3 \text{tr}(\hat{K}^l)$ with CVXPY and MOSEK in Python or with YALMIP and GUROBI in MATLAB, and record the optimal dual variables μ_1^* , μ_2^* , and μ_3^* . Use the code skeletons `kernel_learning` (in `main.ipynb`) or `kernel_learning.m`.
 - (10 points) Let $k^*(\mathbf{x}, \mathbf{x}') = \mu_1^* \hat{k}^1(\mathbf{x}, \mathbf{x}') + \mu_2^* \hat{k}^2(\mathbf{x}, \mathbf{x}') + \mu_3^* \hat{k}^3(\mathbf{x}, \mathbf{x}')$ for all $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$. As discussed in the lecture (slide 183), the label of any test sample $\mathbf{x} \in \mathbb{R}^d$ is predicted as $\text{sign}(\frac{1}{\rho} \sum_{i=1}^m \lambda_i^* y_i k^*(\mathbf{x}_i, \mathbf{x}) + b^*)$, where λ^* is an optimizer of (5) and b^* is the optimal dual variable of the constraint $\lambda^T \mathbf{y} = 0$ in (5). Use the code skeletons `SVM_predict` (in `main.ipynb`) or `SVM_predict.m`.
5. (5 points) Repeat the steps 4(a)–(c) 100 times with different seeds for the random partition of the data into training and test sets, and report the average test accuracy (correct classification rate) to Table 1.
6. (10 points) For each of the 100 training and test sets constructed in 5., solve (2) using the kernels functions \hat{k}^1 , \hat{k}^2 and \hat{k}^3 , respectively, and report the average test accuracies in Table 1. Use the code skeletons `svm_fit` (in `main.ipynb`) or `svm_predict.m`.

Kernel function	\hat{k}^1	\hat{k}^2	\hat{k}^3	$\sum_{l=1}^3 \mu_l^* \hat{k}^l$
Average accuracy				

Table 1: Average test accuracies over 100 data splits.