Mehryar Mohri
Foundations of Machine Learning 2022
Courant Institute of Mathematical Sciences
Homework assignment 2
March 3, 2022
Due: March 22, 2022 11:55 PM EDT

# A   VC dimension

1. We denote by $\mathcal{B}_n$ the set of all closed balls in $\mathbb{R}^n$. That is, $\mathcal{B}_n$ is the class of all subsets of the form $\left\{ x \in \mathbb{R}^n : \|x - x_0\|^2 \le r^2 \right\}$ for some $x_0 \in \mathbb{R}^n$ and $r \ge 0$.

   (a) Show that there exists a set of $n + 1$ points in $\mathbb{R}^n$ that can be shattered by $\mathcal{B}_n$. Conclude that $\mathrm{VCdim}(\mathcal{B}_n) \ge n + 1$.

   (b) Show that the VC dimension of $\mathcal{B}_n$ is at most equal to the VC dimension of hyperplanes in $\mathbb{R}^{n+1}$. Conclude that $\mathrm{VCdim}(\mathcal{B}_n) \le n + 2$.

   (c) Show that $\mathrm{VCdim}(\mathcal{B}_2) = 3$.

**Solution:**

(a) Let $x_0$ be the origin and define $x_i$, for $i \in \{1, \ldots, n\}$, as the point whose $i$th coordinate is 1 and all others are 0. Let $y_0, y_1, \ldots, y_n \in \{-1, +1\}$ be an arbitrary set of labels for $x_0, x_1 \ldots, x_n$. Let $x_*$ be the vector whose $i$th coordinate is $y_i$ and $r = \sqrt{n + \frac{y_0 - 1}{2}}$. Then the classifier defined by the closed ball of form $\|x - x_*\|^2 \le r^2$ shatters $x_0, \ldots, x_n$. Indeed, for $i = 0$,

$$\mathrm{sgn}\left( r^2 - \|x_0 - x_*\|^2 \right) = \mathrm{sign}\left( n + \frac{y_0 - 1}{2} - n \right) = \mathrm{sgn}\left( \frac{y_0 - 1}{2} \right) = y_0;$$

for $i \in \{1, \ldots, n\}$,

$$\mathrm{sgn}\left( r^2 - \|x_i - x_*\|^2 \right) = \mathrm{sgn}\left( n + \frac{y_0 - 1}{2} - \left( n - 1 + (1 - y_i)^2 \right) \right) = \mathrm{sgn}\left( \frac{y_0 + 1}{2} - (1 - y_i)^2 \right) = y_i,$$

which proves the claim. By definition, the VC dimension of $\mathcal{B}_n$ is lower bounded by $n + 1$.

(b) For any $x$, $x_0 \in \mathbb{R}^n$ and $r \ge 0$, we have

$$\|x - x_0\|^2 - r^2 \le 0 \iff w \cdot \Phi(x) + b \le 0$$

where $w = \begin{bmatrix} 1 \\ -2x_0 \end{bmatrix} \in \mathbb{R}^{n+1}$, $\Phi(x) = \begin{bmatrix} \|x\|^2 \\ x \end{bmatrix} \in \mathbb{R}^{n+1}$ and $b = \|x_0\|^2 - r^2$. Therefore, for any $m$ such that there exist $x_1, \ldots, x_m$ shattered by $\mathcal{B}_n$, $\Phi(x_1), \ldots, \Phi(x_m)$ are shattered by hyperplanes in $\mathbb{R}^{n+1}$. By definition, the VC dimension of $\mathcal{B}_n$ is at most equal to the VC dimension of hyperplanes in $\mathbb{R}^{n+1}$, that is $n + 2$.

(c) By (a), we have $\mathrm{VCdim}(\mathcal{B}_2) \ge 3$. To obtain the upper bound $\mathrm{VCdim}(\mathcal{B}_2) \le 3$, it suffices to show that no set of 4 points can be shattered by $\mathcal{B}_2$. By Radon's theorem, any set of 4 points $\mathcal{X}$ in $\mathbb{R}^2$ can be partitioned into two sets $\mathcal{X}_1$ and $\mathcal{X}_2$ such that their convex hulls intersect. Observe that when two sets of points $\mathcal{X}_1$ and $\mathcal{X}_2$ are separated by a closed ball in $\mathcal{B}_2$, their convex hulls are also separated by that closed ball. Thus, $\mathcal{X}_1$ and $\mathcal{X}_2$ cannot be separated by a closed ball in $\mathcal{B}_2$ and $\mathcal{X}$ is not shattered. Therefore, $\mathrm{VCdim}(\mathcal{B}_2) = 3$.

# B  Maximum Margin Multiple Kernel

1. Let $\mathcal{X}$ denote the input space and $\mathcal{Y} = \{1, \ldots, c\}$ a set of $c \geq 2$ classes. Let $S = ((x_1, y_1), \ldots, (x_m, y_m)) \in (\mathcal{X} \times \mathcal{Y})^m$ be a sample of size $m$. Assume that $p \geq 1$ positive semi-definite (PSD) base kernels over $\mathcal{X} \times \mathcal{X}$ are given. Consider a hypothesise set based on a kernel $K_{\boldsymbol{\mu}}$ of the form $K_{\boldsymbol{\mu}} = \sum_{k=1}^{p} \mu_k K_k$ where $\boldsymbol{\mu} = (\mu_1, \cdots, \mu_p)^{\top}$ is chosen from $\Delta_q = \{\boldsymbol{\mu} : \boldsymbol{\mu} \geq 0, \|\boldsymbol{\mu}\|_q = 1\}$ with $q \geq 1$. The multi-class maximum margin multiple kernel (M$^3$K) algorithm [CMR13] is based on the following optimization:

$$
\min_{\boldsymbol{\mu} \in \widehat{M}_q, \boldsymbol{w}, \boldsymbol{\xi}} \frac{1}{2} \sum_{y=1}^{c} \sum_{k=1}^{p} \frac{\|\boldsymbol{w}_{y,k}\|^2}{\mu_k} + C \sum_{i=1}^{m} \xi_i
$$
$$
\text{subject to: } \forall i \in [1, m], \xi_i \geq 0, \forall y \neq y_i,
$$
$$
\xi_i \geq 1 - (\boldsymbol{w}_{y_i} \cdot \Phi(x_i) - \boldsymbol{w}_y \cdot \Phi(x_i)), \tag{1}
$$

where $\cdot$ is defined as $\mathcal{A}_{m \times n} \cdot \mathcal{B}_{m \times n} = \sum_{i,j} \mathcal{A}(i,j)\mathcal{B}(i,j)$ for any two matrices $\mathcal{A}$ and $\mathcal{B}$ with the same dimension $m \times n$, $\widehat{M}_q \subset \Delta_q$ is a data-dependent set, $C \geq 0$ is a regularization parameter, $\boldsymbol{w}_y = (\boldsymbol{w}_{y,1}, \ldots, \boldsymbol{w}_{y,p})^{\top}$ is the associated hypothesis for any class $y \in \mathcal{Y}$, $\Phi(x) = (\Phi_{K_1}(x), \ldots, \Phi_{K_p}(x))^{\top}$ and $\Phi_K$ denotes a feature mapping associated to the kernel $K$.

(a) Read the Chapter 9.1 - 9.3.1 in the textbook and the paper [CMR13] to understand better the multi-class maximum margin multiple kernel (M$^3$K) algorithm. Write down the explicit expression of $\widehat{M}_q$ and briefly explain each term appearing in that expression.

(b) Show how to derive the dual optimization of M$^3$K (1):

$$
\min_{\boldsymbol{\mu} \in \widehat{M}_q} \max_{\boldsymbol{\alpha} \in \mathbb{R}^{m \times c}} \sum_{i=1}^{m} \boldsymbol{\alpha}_i \cdot \mathbf{e}_{y_i} - \frac{C}{2} \sum_{i,j=1}^{m} (\boldsymbol{\alpha}_i \cdot \boldsymbol{\alpha}_j) \sum_{k=1}^{p} \mu_k K_k(x_i, x_j)
$$
$$
\text{subject to: } \forall i \in [1, m], \boldsymbol{\alpha}_i \leq \mathbf{e}_{y_i} \wedge \boldsymbol{\alpha}_i \cdot \mathbf{1} = 0, \tag{2}
$$

where $\boldsymbol{\alpha} \in \mathbb{R}^{m \times c}$ is a matrix, $\boldsymbol{\alpha}_i$ is its $i$th row, and $\mathbf{e}_l$ is the $l$th unit vector in $\mathbb{R}^c$, $l \in [1, c]$. Prove the equivalence of primal (1) and dual (2).

(Note: you should write down every necessary step and rigorously say why the theorems apply.)

**Solution:**

(a) $\widehat{M}_q = \{\boldsymbol{\mu} : \boldsymbol{\mu} \in \Delta_q, \widehat{\gamma}_{K_{\boldsymbol{\mu}}} \geq \gamma_0\}$ where $\widehat{\gamma}_{K_{\boldsymbol{\mu}}} = \frac{1}{m} \sum_{i=1}^{m} \min_{y \neq y_i} \boldsymbol{\mu} \cdot \boldsymbol{\eta}(x_i, y_i, y)$ is the empirical multi-class kernel margin defined in the paper [CMR13] and $\gamma_0$ is a chosen constant. The additional condition $\widehat{\gamma}_{K_{\boldsymbol{\mu}}} \geq \gamma_0$ in $\widehat{M}_q$ ensures that $\boldsymbol{\mu}$ is selected such that the average empirical kernel margin is at least $\gamma_0$.

(b) We first derive the dual optimization of (1) without the outer $\min_{\boldsymbol{\mu} \in \widehat{M}_q}$. Note the constraints in (1) can be equivalently written as

$$
\forall i \in [1, m], \forall y \in [1, c], \xi_i \geq 1 - (\boldsymbol{w}_{y_i} \cdot \Phi(x_i) - \boldsymbol{w}_y \cdot \Phi(x_i)) - \delta_{yy_i},
$$

where $\delta_{yy_i} = \begin{cases} 1 & \text{if } y = y_i \\ 0 & \text{otherwise} \end{cases}$. We introduce Lagrange variables $\boldsymbol{\alpha}'_i = (\alpha'_{i1}, \alpha'_{i2}, \ldots \alpha'_{ic})^{\top} \geq 0, i \in [1, m]$ associated to these $m \times c$ constraints. Thus the Lagrangian of problem of (1) can then be defined for all $\boldsymbol{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}'$, by

$$
\mathcal{L}(\boldsymbol{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}') = \frac{1}{2} \sum_{y=1}^{c} \sum_{k=1}^{p} \frac{\|\boldsymbol{w}_{y,k}\|^2}{\mu_k} + C \sum_{i=1}^{m} \xi_i - \sum_{i=1}^{m} \sum_{y=1}^{c} \alpha'_{iy}[\boldsymbol{w}_{y_i} \cdot \Phi(x_i) - \boldsymbol{w}_y \cdot \Phi(x_i) - 1 + \xi_i + \delta_{yy_i}]. \tag{3}
$$

The KKT conditions are obtained by setting the gradient of the Lagrangian with respect to the primal variables $\boldsymbol{w}$ and $\boldsymbol{\xi}$ to zero and by writing the complementarity conditions:

$$\forall i, \, \nabla_{\xi_i} \mathcal{L} = C - \sum_{y=1}^{c} \alpha'_{iy} = 0 \Rightarrow \sum_{y=1}^{c} \alpha'_{iy} = C \tag{4}$$

$$\forall y, \, k, \, \nabla_{\boldsymbol{w}_{y,k}} \mathcal{L} = \frac{\boldsymbol{w}_{y,k}}{\mu_k} + \sum_{i=1}^{m} \alpha'_{iy} \Phi_{K_k}(x_i) - \sum_{i=1}^{m} \delta_{yy_i} \sum_{y'=1}^{c} \alpha'_{iy'} \Phi_{K_k}(x_i) = 0 \tag{5}$$

$$\Rightarrow \boldsymbol{w}_{y,k} = \mu_k \sum_{i=1}^{m} \left( \alpha'_{iy} - \delta_{yy_i} \sum_{y'=1}^{c} \alpha'_{iy'} \right) \Phi_{K_k}(x_i) \Rightarrow \boldsymbol{w}_{y,k} = \mu_k \sum_{i=1}^{m} \left( \alpha'_{iy} - C\delta_{yy_i} \right) \Phi_{K_k}(x_i) \tag{6}$$

$$\forall i, \, y, \, \alpha'_{iy}[\boldsymbol{w}_{y_i} \cdot \Phi(x_i) - \boldsymbol{w}_y \cdot \Phi(x_i) - 1 + \xi_i + \delta_{yy_i}] = 0 \Rightarrow \alpha'_{iy} = 0 \vee (\boldsymbol{w}_{y_i} - \boldsymbol{w}_y) \cdot \Phi(x_i) + \delta_{yy_i} = 1 - \xi_i \tag{7}$$

To derive the dual optimization of problem of (1) without the outer $\min_{\boldsymbol{\mu} \in \widehat{M}_q}$, we plug into the Lagrangian the definition of $\boldsymbol{w}$ in terms of the dual variables Eq. (6) and apply the constraints Eq. (4). This yields

$$\mathcal{L} = -\sum_{i=1}^{m} \boldsymbol{\alpha}'_i \cdot \mathbf{e}_{y_i} - \frac{1}{2} \sum_{i,j=1}^{m} (C\mathbf{e}_{y_i} - \boldsymbol{\alpha}'_i) \cdot (C\mathbf{e}_{y_j} - \boldsymbol{\alpha}'_j) \sum_{k=1}^{p} \mu_k K_k(x_i, x_j).$$

Note here, in addition to $\boldsymbol{\alpha}'_i \geq 0$, we must impose the constraint $\sum_{y=1}^{c} \alpha'_{iy} = C$ in Eq. (4). This leads to the following dual optimization problem of (1) without the outer $\min_{\boldsymbol{\mu} \in \widehat{M}_q}$:

$$\max_{\boldsymbol{\alpha}' \in \mathbb{R}^{m \times c}} \, -\sum_{i=1}^{m} \boldsymbol{\alpha}'_i \cdot \mathbf{e}_{y_i} - \frac{1}{2} \sum_{i,j=1}^{m} (C\mathbf{e}_{y_i} - \boldsymbol{\alpha}'_i) \cdot (C\mathbf{e}_{y_j} - \boldsymbol{\alpha}'_j) \sum_{k=1}^{p} \mu_k K_k(x_i, x_j)$$
$$\text{subject to: } \forall i \in [1, m], \, \boldsymbol{\alpha}'_i \geq 0 \wedge \boldsymbol{\alpha}'_i \cdot \mathbf{1} = C, \tag{8}$$

Then by (8), using the change of variable $\boldsymbol{\alpha}_i = e_{y_i} - \frac{\boldsymbol{\alpha}'_i}{C}$, adding back the outer $\min_{\boldsymbol{\mu} \in \widehat{M}_q}$ and omitting any additive and positive multiplicative constants, we obtain the equivalent dual optimization problem of (1) as follows:

$$\min_{\boldsymbol{\mu} \in \widehat{M}_q} \max_{\boldsymbol{\alpha} \in \mathbb{R}^{m \times c}} \sum_{i=1}^{m} \boldsymbol{\alpha}_i \cdot \mathbf{e}_{y_i} - \frac{C}{2} \sum_{i,j=1}^{m} (\boldsymbol{\alpha}_i \cdot \boldsymbol{\alpha}_j) \sum_{k=1}^{p} \mu_k K_k(x_i, x_j)$$
$$\text{subject to: } \forall i \in [1, m], \, \boldsymbol{\alpha}_i \leq \mathbf{e}_{y_i} \wedge \boldsymbol{\alpha}_i \cdot \mathbf{1} = 0, \tag{9}$$

The equivalence is straightforward by using the facts that the objective function is convex, the constrains are affine and thus qualified, the objective and constraint functions are differentiable, and the KKT conditions hold at the optimum.

# C   SVMs hand-on

(Note: please share a `GitHub` link to your open source code in the submission. Any submissions that do not have the code link will obtain a zero point. The graders will check the main lines to ensure that what was done was conceptually correct. The grade will be based on both the code and the answer.)

1. Download and install the `libsvm` software library from:

    https://www.csie.ntu.edu.tw/~cjlin/libsvm

    and briefly consult the documentation to become more familiar with the tools.

2. Download the `Abalone` data set:

    http://archive.ics.uci.edu/ml/datasets/Abalone

Use the `libsvm` scaling tool to scale the features of all the data. Use the first 3133 examples for training, the last 1044 for testing. The scaling parameters should be computed only on the training data and then applied to the test data.

3. Consider the binary classification that consists of distinguishing classes 1 through 9 from the rest. Use SVMs combined with polynomial kernels to tackle this binary classification problem.

   To do that, randomly split the training data into five equal-sized disjoint sets. For each value of the polynomial degree, $d = 1, 2, 3, 4, 5$, plot the average cross-validation error plus or minus one standard deviation as a function of $C$ (let other parameters of polynomial kernels in `libsvm` be equal to their default values), varying $C$ in powers of 3, starting from a small value $C = 3^{-k}$ to $C = 3^k$, for some value of $k$. $k$ should be chosen so that you see a significant variation in training error, starting from a very high training error to a low training error. Expect longer training times with `libsvm` as the value of $C$ increases.

4. Let $(C^*, d^*)$ be the best pair found previously. Fix $C$ to be $C^*$. Plot the five-fold cross-validation error and the test errors for the hypotheses obtained as a function of $d$. Plot the average number of support vectors obtained as a function of $d$. How many of the support vectors lie on the margin hyperplanes?

5. Fix $(C, d)$ to be $(C^*, d^*)$. Plot the training and test errors as a function of the training sample.

6. Sparse SVM. One can give two types of arguments in favor of the SVM algorithm: one based on the sparsity of the support vectors, another based on the notion of margin. Suppose that instead of maximizing the margin, we choose instead to maximize sparsity by minimizing the $L_1$ norm of the vector $\boldsymbol{\alpha}$ that defines the weight vector $\boldsymbol{w}$. This gives the following optimization problem for a kernel function $K$:

$$\min_{\boldsymbol{\alpha}, b, \boldsymbol{\xi}} \frac{1}{2} \sum_{i=1}^{m} |\alpha_i| + C \sum_{i=1}^{m} \xi_i$$
$$\text{subject to } y_i \left( \sum_{j=1}^{m} \alpha_j y_j K(\boldsymbol{x}_i, \boldsymbol{x}_j) + b \right) \geq 1 - \xi_i, i \in [1, m] \tag{10}$$
$$\xi_i, \alpha_i \geq 0, i \in [1, m].$$

   (a) Derive the equivalent dual optimization problem of (10) in terms of the feature mapping $\Phi$ associated to the kernel $K$ and write the proof clearly.

   (b) Derive the equivalent hinge loss minimization problem of (10) and write the proof clearly. Compare it with an instance of the equivalent hinge loss minimization problem of SVM shown in class.

   (c) Apply Stochastic Gradient Descent to solve the optimization problem. Plot the five-fold cross-validation training and test errors for the hypotheses obtained based on the solution $\boldsymbol{\alpha}$ as a function of $d$, for the best value of $C$ measured on the validation set.

**Solution:**

3. Take $k = 8$. Figure 1 shows the average cross-validation error plus or minus one standard deviation, with different $d$ and $C$.
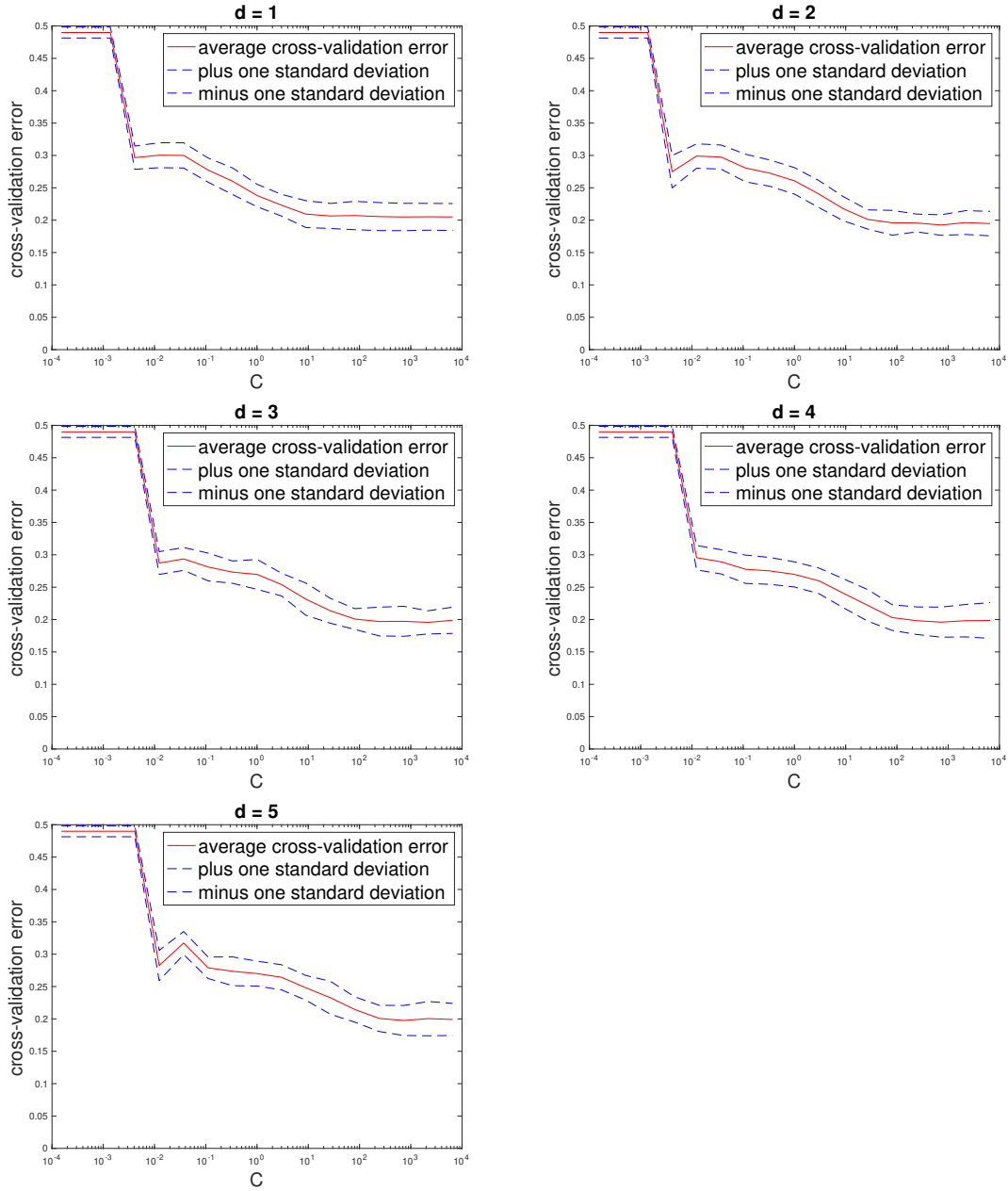
Figure 1: Average cross-validation error plus or minus one standard deviation, with different $d$ and $C$.

4. We choose $(C^*, d^*) = (3^6, 2)$.

   Fix $C$ to be $3^6$, Figure 2 shows the five-fold cross-validation error and the test errors against $d$.
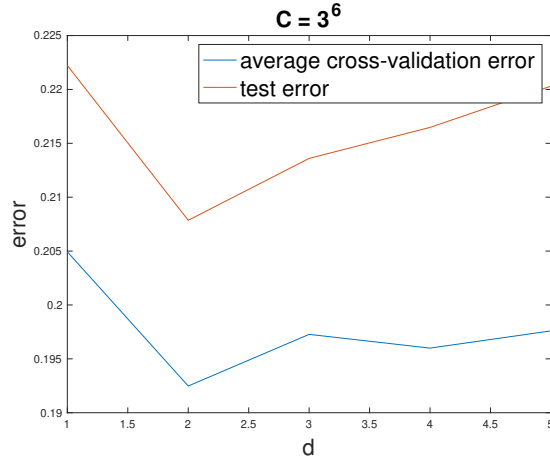
Figure 2: Average cross-validation error and test error against $d$, fixing $C$ to be $3^6$.

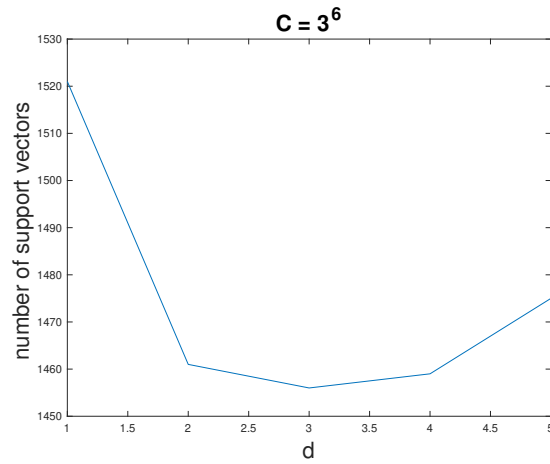Figure 3 shows the number of support vectors against $d$.



Figure 3: The number of support vectors against $d$.

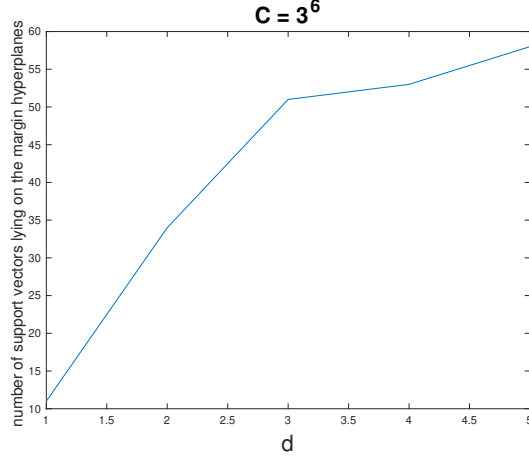Figure 4 shows the number of support vectors lying on the margin hyperplanes against $d$.

Figure 4: The number of support vectors lying on the margin hyperplanes against $d$.
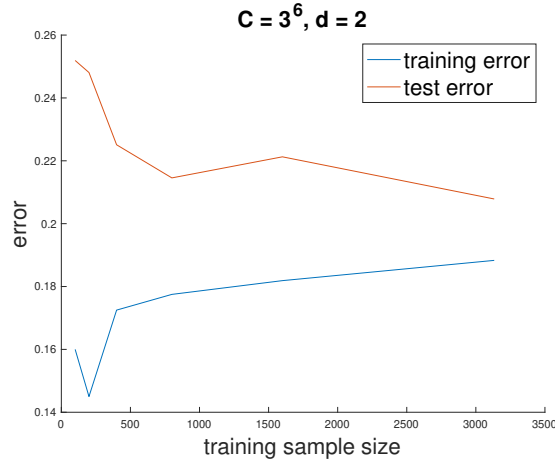
5. See Figure 5.



Figure 5: The training and test errors as a function of the training sample.

6. (a) Let $\mathbf{x}'_i = (y_1 K(x_i, x_1), \ldots, y_m K(x_i, x_m))^\top \in \mathbb{R}^m$, $i = 1, \ldots, m$. We introduce Lagrange variables $\alpha'_i \geq 0$, $i \in [1, m]$ associated to the first $m$ constraints, $\beta_i \geq 0$, $i \in [1, m]$ associated to the non-negativity constraints of the slack variables and $\gamma_i \geq 0$, $i \in [1, m]$ associated to the non-negativity constraints of the $\alpha$. Thus the Lagrangian of problem of (10) can then be defined for all $b \in \mathbb{R}$, and $\alpha, \xi, \alpha', \beta, \gamma \in \mathbb{R}^m_+$, by

$$\mathcal{L}(\alpha, b, \xi, \alpha', \beta, \gamma) = \frac{1}{2}\sum_{i=1}^m \alpha_i + C\sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha'_i[y_i(\alpha \cdot \mathbf{x}'_i + b) - 1 + \xi_i] - \sum_{i=1}^m \beta_i \xi_i - \sum_{i=1}^m \gamma_i \alpha_i. \quad (11)$$

The KKT conditions are obtained by setting the gradient of the Lagrangian with respect to the

primal variables $\alpha$, $b$, and $\xi_i$s to zero and by writing the complementarity conditions:

$$\nabla_\alpha \mathcal{L} = \frac{1}{2}\mathbf{1} - \sum_{i=1}^m \alpha_i' y_i \mathbf{x}_i' - \gamma = 0 \Rightarrow \frac{1}{2}\mathbf{1} = \sum_{i=1}^m \alpha_i' y_i \mathbf{x}_i' + \gamma \tag{12}$$

$$\nabla_b \mathcal{L} = -\sum_{i=1}^m \alpha_i' y_i = 0 \Rightarrow \sum_{i=1}^m \alpha_i' y_i = 0 \tag{13}$$

$$\nabla_{\xi_i} \mathcal{L} = C - \alpha_i' - \beta_i = 0 \Rightarrow \alpha_i' + \beta_i = C \tag{14}$$

$$\forall i,\ \alpha_i'[y_i(\alpha \cdot \mathbf{x}_i' + b) - 1 + \xi_i] = 0 \Rightarrow \alpha_i' = 0 \vee y_i(\alpha \cdot \mathbf{x}_i' + b) = 1 - \xi_i \tag{15}$$

$$\forall i,\ \beta_i \xi_i = 0 \Rightarrow \beta_i = o \vee \xi_i = 0 \tag{16}$$

$$\forall i,\ \gamma_i \alpha_i = 0 \Rightarrow \gamma_i = o \vee \alpha_i = 0 \tag{17}$$

To derive the dual optimization of problem of (10), we plug into the Lagrangian the definition of $\alpha$ in terms of the dual variables Eq. (12) and apply the constraints Eq. (13) and Eq. (14). This yields

$$\begin{aligned}
\mathcal{L} &= \frac{1}{2}\sum_{i=1}^m \alpha_i - \alpha \cdot \sum_{i=1}^m \alpha_i' y_i \mathbf{x}_i' - \sum_{i=1}^m \alpha_i' y_i b + \sum_{i=1}^m \alpha_i' - \sum_{i=1}^m \gamma_i \alpha_i && (\alpha_i' + \beta_i = C) \\
&= \frac{1}{2}\sum_{i=1}^m \alpha_i - \alpha \cdot \sum_{i=1}^m \alpha_i' y_i \mathbf{x}_i' + \sum_{i=1}^m \alpha_i' - \sum_{i=1}^m \gamma_i \alpha_i && \left(\sum_{i=1}^m \alpha_i' y_i = 0\right) \\
&= \sum_{i=1}^m \alpha_i' && \left(\frac{1}{2}\mathbf{1} = \sum_{i=1}^m \alpha_i' y_i \mathbf{x}_i' + \gamma\right)
\end{aligned}$$

Note here, in addition to $\alpha_i' \geq 0$, we must impose the constraint on the Lagrangian variable $\beta_i \geq 0$ and $\gamma_i \geq 0$. In view of Eq. (14), we obtain $0 \leq \alpha_i' \leq C$ and $\gamma_i \geq 0$. This leads to the following dual optimization problem of (10):

$$\begin{aligned}
\max_{\alpha', \gamma} \quad & \sum_{i=1}^m \alpha_i' \\
\text{subject to} \quad & \sum_{i=1}^m \alpha_i' y_i \boldsymbol{x}_i' + \gamma = \frac{1}{2}\mathbf{1} \\
& \sum_{i=1}^m \alpha_i' y_i = 0 \\
& 0 \leq \alpha_i' \leq C,\ \gamma_i \geq 0, i \in [1, m].
\end{aligned} \tag{18}$$

where $\mathbf{x}_i' = (y_1 K(x_i, x_1), \ldots, y_m K(x_i, x_m))^\top = (y_1 \Phi(x_i)\Phi(x_1), \ldots, y_m \Phi(x_i)\Phi(x_m))^\top \in \mathbb{R}^m$, $i = 1, \ldots, m$.

(b) The equivalent hinge loss minimization problem of (10) can be written as

$$\min_{\boldsymbol{\alpha}, b} \frac{1}{2}\sum_{i=1}^m |\alpha_i| + C \sum_{i=1}^m \left(1 - y_i\left(\sum_{j=1}^m \alpha_j y_j K(\boldsymbol{x}_i, \boldsymbol{x}_j) + b\right)\right)_+ \tag{19}$$

subject to $\alpha_i \geq 0, i \in [1, m]$.

Different from the equivalent hinge loss minimization problem of SVM shown in class, the problem (19) has the constraint $\alpha_i \geq 0$ and the $L_1$ norm of the vector $\boldsymbol{\alpha}$ instead of the $L_2$ norm in the objection function.
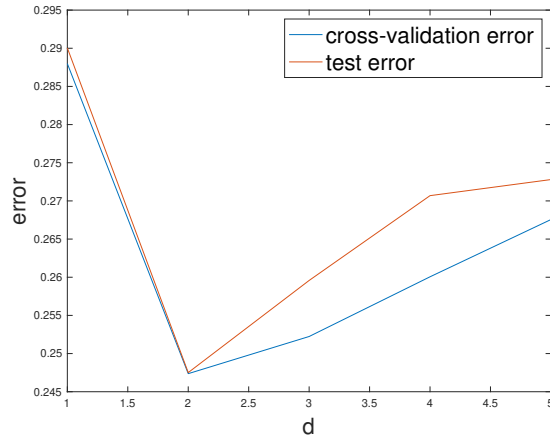
(c) See Figure 6.

Figure 6: The five-fold cross-validation training and test errors.

# References

[CMR13]  Corinna Cortes, Mehryar Mohri, and Afshin Rostamizadeh. "Multi-class classification with maximum margin multiple kernel". In: *International Conference on Machine Learning*. PMLR. 2013, pp. 46–54.