

FML HW 2

Yucong Lei

Collaborated with Mengjian Hua(mh5113), Haiyang
Wang(hw1927)

1 Problem A

1.1 A.1

Proof

$$\begin{aligned}\hat{\mathfrak{R}}_S(H) &= \mathbb{E}_\sigma[\sup_{h \in H} \frac{1}{m} \sum_{1 \leq i \leq m} \sigma_i h(z_i)] \\ &\geq \mathbb{E}_\sigma[\frac{1}{m} \sum_{1 \leq i \leq m} \sigma_i h(z_i)], \forall h \in H \\ &= \frac{1}{m} \sum_{1 \leq i \leq m} \mathbb{E}_\sigma[\sigma_i h(z_i)] = 0\end{aligned}$$

Hence it is nonnegative.

1.2 A.2

Proof Let Φ_i denote a continuous function such that:

$$\begin{aligned}\Phi_i(x) &= 0, x \in (-\infty, 1] \\ \Phi_i(x) &= x - 1, x \in (1, \infty)\end{aligned}$$

It is easy to see that Φ_i is 1-Lipschitz, since:

$$\begin{aligned}|\Phi_i(x) - \Phi_i(y)| &= 0, x, y \in (-\infty, 1] \\ |\Phi_i(x) - \Phi_i(y)| &= |x - y|, x, y \in (1, \infty) \\ |\Phi_i(x) - \Phi_i(y)| &= |x - 1| \leq |x - 0| \leq |x - y|, x \in (1, \infty), y \in (-\infty, 1]\end{aligned}$$

In addition, we note that Φ_i has the following properties:

$$\begin{aligned}\Phi_i(0) &= 0 \\ \Phi_i(2) &= 1 \\ \Phi_i(1) &= 0\end{aligned}$$

Thus, $\forall h_1, h_2$, two binary classifiers, with values taken from 0,1, we have:

$$\Phi_i(h_1(z) + h_2(z)) = h_1(z)h_2(z)$$

We can now apply Talagrand's contraction lemma:

$$\begin{aligned} \hat{\mathfrak{R}}_S(\Phi_i \circ (H_1 + H_2)) &\leq 1 * \hat{\mathfrak{R}}_S((H_1 + H_2)) \\ &\leq \hat{\mathfrak{R}}_S(H_1) + \hat{\mathfrak{R}}_S(H_2) \\ \rightarrow \hat{\mathfrak{R}}_S(H) &= \hat{\mathfrak{R}}_S(H_1 * H_2) = \hat{\mathfrak{R}}_S(\Phi_i \circ (H_1 + H_2)) \leq \hat{\mathfrak{R}}_S(H_1) + \hat{\mathfrak{R}}_S(H_2) \end{aligned}$$

2 Problem B

2.1 B.1

Consider the following neural network diagram (Figure 1). Each input data $\alpha_i \in \mathbb{R}^n$, where $1 \leq i \leq m$, we obtain a binary vector of dimension k in the intermediate layer of a given neural network, denoted as $\beta_i \in \{0, 1\}^k$, $1 \leq i \leq m$.

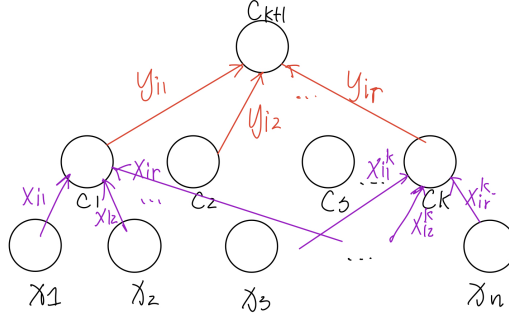


Figure 1:

Let us fix the input values, $\alpha_i \in \mathbb{R}^n$, where $1 \leq i \leq m$, and range all possible neural networks. Then the maximal distinct values are produced by ranging from different intermediate layer structures, that is, of which r binary values out of the k nodes are chosen as input value of the concept function, and ranging

from different choices of concept values. Hence the following inequality hold:

$$\begin{aligned}
& |\{(h(\alpha_1), \dots, (h(\alpha_m))), h \in H\}| \\
& \leq |\{(c \circ d(\beta_1), \dots, (c \circ d(\beta_m))), c \in C, d \in \{A \in \text{Mat}(n, \mathbb{R}); A_{i,j} = 0, 1; \sum_j A_{i_p,j} = 1, 1 \leq p \leq r\}\}| \\
& \leq \Pi_C(m)^{\binom{k}{r}} \\
& \rightarrow \Pi_H(m) \leq \Pi_C(m)^{\binom{k}{r}}
\end{aligned}$$

2.2 B.2

Let $a = 2^l \geq \binom{k}{r} \geq 1$. By the corollary of Sauer's lemma in the textbook, we have:

$$\begin{aligned}
\Pi_H(m) & \leq \Pi_C(m)^{\binom{k}{r}} \\
& \leq \left(\frac{em}{d}\right)^{da} \\
& \leq \left(\frac{8m}{d}\right)^{da}
\end{aligned}$$

Let $x = ad$, $y = \frac{8}{d}$, and $m = 2x \log_2(xy) = 2ad \log_2(8a)$. We have:

$$\begin{aligned}
x * y & > 4 \\
m & \geq 1 \\
m & \in \mathbb{Z}
\end{aligned}$$

Hence by the hint, we have: $m > x \log_2(y m)$.

$$\begin{aligned}
\Pi_H(m) & \leq \left(\frac{8m}{d}\right)^{da} \\
& \leq 2^{ad * \log_2 \frac{8m}{d}} \\
& < 2^m
\end{aligned}$$

Hence, $VCdim(H) < m = 2ad \log_2(8a)$.

2.3 B.3

We first try to compute d in terms of k, r . Then we bound $VCdim(H)$ by the inequality obtained in the last question. We will make use of the fact that the VC dimension of hyperplanes in \mathbb{R}^r is $(r+1)$.

First we prove that the VC dimension of hyperplanes through the origin in \mathbb{R}^r is no less than r , by showing that there exists r samples, $\alpha_1, \dots, \alpha_r \in \mathbb{R}^r$, such that:

$$|\{(c(\alpha_1), \dots, c(\alpha_r)), c \in C\}| = 2^r$$

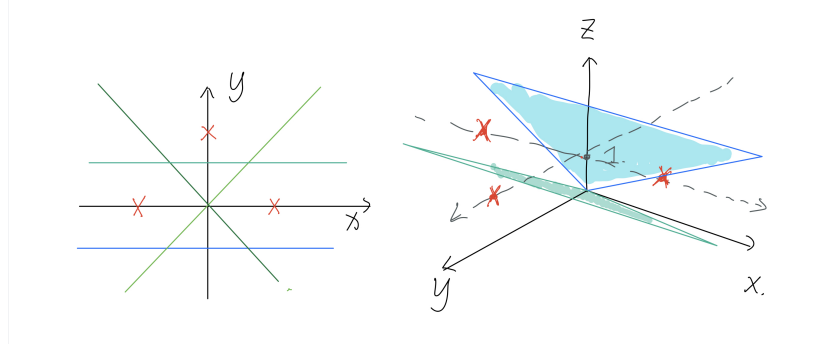


Figure 2:

Because the VC dimension of hyperplanes in \mathbb{R}^{r-1} is r , there exists r samples, $\beta_1, \dots, \beta_r \in \mathbb{R}^{r-1}$, such that:

$$|\{(d(\beta_1), \dots, d(\beta_r)), d \in \{sgn(w * \beta + b), w \in \mathbb{R}^{r-1}, b \in \mathbb{R}\}\}| = 2^r$$

We map all β to $(\beta, 1) \in \mathbb{R}^r$, and note that this mapping is injective. Thus:

$$|\{(c(\beta_1, 1), \dots, c(\beta_r, 1)), c \in \{sgn((w, b) * \alpha), w \in \mathbb{R}^{r-1}, b \in \mathbb{R}\}\}| = 2^r$$

We next prove that the VC dimension of hyperplanes through the origin in \mathbb{R}^r has to be less than $(r + 1)$, by showing that any set of $r + 1$ points in \mathbb{R}^r can only be completely classified by planes including a plane that does not pass through the origin.

We prove by contradiction. Suppose there are such $r + 1$ samples in \mathbb{R}^{r+1} , $\alpha_1, \dots, \alpha_{r+1} \in \mathbb{R}^r$, completely classified by hyperplanes through the origin. Then, we pick such a hyperplane through the origin that makes the label of every sample be the same, that is, they all lie in one side of a hyperplane through the origin. Assume without loss of generality, that this hyperplane is $x_r = 0$, and each α_i 's r -coordinate is greater than 0.

We then note that no two sample points determine a line through the origin. To see why this is the case, we assume two samples are on a line passing through the origin that is:

$$(x_1, \dots, x_r) = a(y_1, \dots, y_r), \exists a \neq 1, 0, x_r, y_r > 0.$$

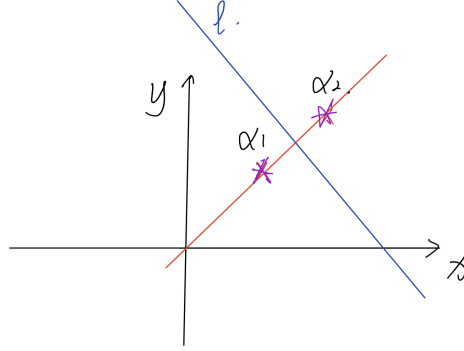


Figure 3:

If there exists a hyperplane that separates them, that is:

$$\begin{aligned}
& \exists w \in \mathbb{R}^r, b \in \mathbb{R} : \\
& w * x + b > 0 \\
& w * y + b < 0 \\
& \rightarrow \exists \lambda \in (0, 1) \\
& w * (\lambda x + (1 - \lambda)y) + b = 0 \\
& \rightarrow b = -w * (\lambda x + (1 - \lambda)y) \\
& = -w * \frac{(\|\lambda x + (1 - \lambda)y\|)x}{\|x\|} \neq 0
\end{aligned}$$

Hence we have proved that no two sample points in the $r + 1$ sample can give a line which passes through the origin. What is significant about this can be seen later on.

We define a mapping:

$$\begin{aligned}
\phi : (x_1, \dots, x_r), x_r > 0 &\mapsto \left(\frac{x_1}{x_r}, \dots, \frac{x_{r-1}}{x_r}\right) \\
|\{(c(\alpha_1), \dots, c(\alpha_{r+1})), c \in \{sgn(w * \alpha), w \in \mathbb{R}^r\}\}| &= 2^{r+1} \\
\rightarrow |\{(d(\phi(\alpha_1)), \dots, d(\phi(\alpha_{r+1}))), d \in \{sgn(w * (\phi(\alpha), 1)), w \in \mathbb{R}^r\}\}| &= 2^{r+1}
\end{aligned}$$

That is, there are $r + 1$ sample points in \mathbb{R}^{r-1} , i.e., $\phi(\alpha_1), \dots, \phi(\alpha_{r+1})$ such that they can be completely classified by hyperplanes in \mathbb{R}^{r-1} . Thus a contradiction to the fact that $VCdim(H^{r-1}) = r$. Hence we have proved that $VCdim(C) = r$.

$$\begin{aligned}
VCdim(H) &< m = 2adlog_2(8a) \\
&< 2arlog_2(8a)
\end{aligned}$$

3 Problem C

3.1 C.1

3.2 C.2,3

3.3 C.4

See abalone.data.processing.FML.HW2.ipynb

3.4 C.5

$C^* = 512, d^* = 9$. See abalone.data.processing.FML.HW2.ipynb

3.5 C.6

3.5.1

Consider the following map:

$$\hat{x} = (y_1 K(x, x_1), y_2 K(x, x_2), \dots, y_m K(x, x_m))$$

Which takes a sample $x \in \mathbb{R}^m$ to another point in the same Euclidean space. Then problem (1) can be rewritten as:

$$\begin{aligned} \min_{\alpha, b, \xi} \quad & \frac{1}{2} \|\alpha\|_2^2 + C \sum_{1 \leq i \leq m} \xi_i \\ & y_i(\alpha * \hat{x}_i + b) \geq 1 - \xi_i, 1 \leq i \leq m \\ & \xi_i, \alpha_i \geq 0 \end{aligned}$$

That is, by thinking of \hat{x} as a result of transformation of the original sample point x under the kernel function, we reformulated the problem as the primal optimization problem of SVMs for the transformed vector \hat{x} .

3.5.2

No, the positive-definiteness is not necessary, since the target function is convex, and the domain defined by the constraints are a convex region, because the domain function is linear with respect to α, ξ, b , and hence convex. To see

this:

$$\begin{aligned}
& \forall(\alpha, b, \xi), (\tilde{\alpha}, \tilde{b}, \tilde{\xi}), \forall \lambda \in [0, 1] \\
& y_i \left(\sum_{1 \leq j \leq m} \alpha_j y_j K(x_i, x_j) + b \right) + \xi_i - 1 \geq 0 \\
& y_i \left(\sum_{1 \leq j \leq m} \tilde{\alpha}_j y_j K(x_i, x_j) + \tilde{b} \right) + \tilde{\xi}_i - 1 \geq 0 \\
& \rightarrow y_i \left(\sum_{1 \leq j \leq m} (\lambda \alpha_j + (1 - \lambda) \tilde{\alpha}_j) y_j K(x_i, x_j) + (\lambda b + (1 - \lambda) \tilde{b}) \right) + (\lambda \xi + (1 - \lambda) \tilde{\xi}_i) - 1 \\
& = \lambda (y_i \left(\sum_{1 \leq j \leq m} \alpha_j y_j K(x_i, x_j) + b \right) + \xi_i - 1 \geq 0) + (1 - \lambda) (y_i \left(\sum_{1 \leq j \leq m} \tilde{\alpha}_j y_j K(x_i, x_j) + \tilde{b} \right) + \tilde{\xi}_i - 1) \\
& \geq 0
\end{aligned}$$

Since the intersection of convex sets are convex, and that $\xi_i, \alpha_i \geq 0$ are convex, our domain of optimization problem is convex.

3.5.3

First construct the Lagrange multiplier function:

$$\begin{aligned}
L &= \left(\frac{1}{2} \|\alpha\|_2^2 + C \sum_{1 \leq i \leq m} \xi_i \right) - \sum_{1 \leq i \leq m} \lambda_i (y_i (\alpha * \hat{x}_i + b) - 1 + \xi_i) - \mu * \alpha - \sigma * \xi \\
& \lambda_i, \mu_i, \sigma_i \geq 0, \forall i
\end{aligned}$$

Differentiating the function with respect to α, b, ξ , we get:

$$\begin{aligned}
\nabla_{\alpha} L &= \alpha - \sum_{1 \leq i \leq m} \lambda_{1i} \hat{x}_i - \mu = 0 \\
\nabla_b L &= - \sum_{1 \leq i \leq m} \lambda_i y_i = 0 \\
\nabla_{\xi_i} L &= C - \lambda_i - \sigma_i = 0
\end{aligned}$$

which then gives the following:

$$\begin{aligned}
\alpha &= \sum_{1 \leq i \leq m} \lambda_{1i} \hat{x}_i + \mu \\
\sum_{1 \leq i \leq m} \lambda_i y_i &= 0 \\
C &= \lambda_i + \sigma_i
\end{aligned}$$

Plug them back in the Lagrange function, we get:

$$\begin{aligned}
L &= \frac{1}{2} \left\| \sum_{1 \leq i \leq m} \lambda_{1i} \hat{x}_i + \mu \right\|^2 + C \sum_{1 \leq i \leq m} \xi_i \\
&\quad - \sum_{1 \leq i \leq m} \lambda_i [y_i (\mu * \hat{x}_i + \sum_{1 \leq j \leq m} \lambda_j y_j \hat{x}_j * \hat{x}_i + b) - 1 + \xi_i] \\
&\quad - (\mu * \mu + \mu * (\sum_{1 \leq j \leq m} \lambda_j y_j \hat{x}_j)) - \sigma * \xi \\
&= \left(\frac{1}{2} (\|\mu\|^2 + \left\| \sum_{1 \leq i \leq m} \lambda_{1i} \hat{x}_i \right\|^2) + \sum_{1 \leq i \leq m} \lambda_i y_i (\mu * \hat{x}_i) \right) \\
&\quad + C \sum_{1 \leq i \leq m} \xi_i - \sum_{1 \leq i \leq m} (\lambda_i + \sigma_i) \xi_i + \sum_{1 \leq i \leq m} \lambda_i \\
&\quad - (2 \sum_{1 \leq i \leq m} \lambda_i y_i (\mu * \hat{x}_i) + \sum_{1 \leq i, j \leq m} \lambda_i \lambda_j y_i y_j (\hat{x}_i * \hat{x}_j) + \|\mu\|^2) \\
&= -\frac{1}{2} \|\mu + \sum_{1 \leq i \leq m} \lambda_i y_i \hat{x}_i\|^2 + \sum_{1 \leq i \leq m} \lambda_i
\end{aligned}$$

Hence the dual problem is:

$$\begin{aligned}
&\max_{\mu, \lambda \in \mathbb{R}^m} -\frac{1}{2} \left\| \mu + \sum_{1 \leq i \leq m} \lambda_i y_i \hat{x}_i \right\|^2 + \sum_{1 \leq i \leq m} \lambda_i \\
&\mu_i \geq 0, 0 \leq \lambda_i \leq C, \forall i, \sum_{1 \leq j \leq m} \lambda_j y_j = 0
\end{aligned}$$