

Solutions to Exercises

Chapter 2

2.1 Two-oracle variant of the PAC model

- Assume that \mathcal{C} is efficiently PAC-learnable using \mathcal{H} in the standard PAC model using algorithm \mathcal{A} . Consider the distribution $\mathcal{D} = \frac{1}{2}(\mathcal{D}_- + \mathcal{D}_+)$. Let $h \in \mathcal{H}$ be the hypothesis output by \mathcal{A} . Choose δ such that:

$$\mathbb{P}[R_{\mathcal{D}}(h) \leq \epsilon/2] \geq 1 - \delta.$$

From

$$\begin{aligned} R_{\mathcal{D}}(h) &= \mathbb{P}_{x \sim \mathcal{D}}[h(x) \neq c(x)] \\ &= \frac{1}{2} \left(\mathbb{P}_{x \sim \mathcal{D}_-}[h(x) \neq c(x)] + \mathbb{P}_{x \sim \mathcal{D}_+}[h(x) \neq c(x)] \right) \\ &= \frac{1}{2}(R_{\mathcal{D}_-}(h) + R_{\mathcal{D}_+}(h)), \end{aligned}$$

it follows that:

$$\mathbb{P}[R_{\mathcal{D}_-}(h) \leq \epsilon] \geq 1 - \delta \quad \text{and} \quad \mathbb{P}[R_{\mathcal{D}_+}(h) \leq \epsilon] \geq 1 - \delta.$$

This implies two-oracle PAC-learning with the same computational complexity.

- Assume now that \mathcal{C} is efficiently PAC-learnable in the two-oracle PAC model. Thus, there exists a learning algorithm \mathcal{A} such that for $c \in \mathcal{C}$, $\epsilon > 0$, and $\delta > 0$, there exist m_- and m_+ polynomial in $1/\epsilon$, $1/\delta$, and $\text{size}(c)$, such that if we draw m_- negative examples or more and m_+ positive examples or more, with confidence $1 - \delta$, the hypothesis h output by \mathcal{A} verifies:

$$\mathbb{P}[R_{\mathcal{D}_-}(h) \leq \epsilon] \geq 1 - \delta \quad \text{and} \quad \mathbb{P}[R_{\mathcal{D}_+}(h) \leq \epsilon] \geq 1 - \delta.$$

Now, let \mathcal{D} be a probability distribution over negative and positive examples. If we could draw m examples according to \mathcal{D} such that $m \geq \max\{m_-, m_+\}$, m polynomial in $1/\epsilon$, $1/\delta$, and $\text{size}(c)$, then two-oracle PAC-learning would imply standard PAC-learning:

$$\begin{aligned} \mathbb{P}[R_{\mathcal{D}}(h)] &\leq \mathbb{P}[R_{\mathcal{D}}(h)|c(x) = 0] \mathbb{P}[c(x) = 0] + \mathbb{P}[R_{\mathcal{D}}(h)|c(x) = 1] \mathbb{P}[c(x) = 1] \\ &\leq \epsilon(\mathbb{P}[c(x) = 0] + \mathbb{P}[c(x) = 1]) = \epsilon. \end{aligned}$$

If \mathcal{D} is not too biased, that is, if the probability of drawing a positive example, or that of drawing a negative example is more than ϵ , it is not hard to show, using Chernoff bounds or just Chebyshev's inequality, that drawing a polynomial number of examples in $1/\epsilon$ and $1/\delta$ suffices to guarantee that $m \geq \max\{m_-, m_+\}$ with high confidence.

Otherwise, \mathcal{D} is biased toward negative (or positive examples), in which case returning $h = h_0$ (respectively $h = h_1$) guarantees that $\mathbb{P}[R_{\mathcal{D}}(h)] \leq \epsilon$.

To show the claim about the not-too-biased case, let S_m denote the number of positive examples obtained when drawing m examples when the probability of a positive example is ϵ . By Chernoff bounds,

$$\mathbb{P}[S_m \leq (1 - \alpha)m\epsilon] \leq e^{-m\epsilon\alpha^2/2}.$$

We want to ensure that at least m_+ examples are found. With $\alpha = \frac{1}{2}$ and $m = \frac{2m_+}{\epsilon}$,

$$\mathbb{P}[S_m > m_+] \leq e^{-m_+/4}.$$

Setting the bound to be less than or equal to $\delta/2$, leads to the following condition on m :

$$m \geq \min\left\{\frac{2m_+}{\epsilon}, \frac{8}{\epsilon} \log \frac{2}{\delta}\right\}$$

A similar analysis can be done in the case of negative examples. Thus, when \mathcal{D} is not too biased, with confidence $1 - \delta$, we will find at least m_- negative and m_+ positive examples if we draw m examples, with

$$m \geq \min\left\{\frac{2m_+}{\epsilon}, \frac{2m_-}{\epsilon}, \frac{8}{\epsilon} \log \frac{2}{\delta}\right\}.$$

In both solutions, our training data is the set T and our learned concept $L(T)$ is the tightest circle (with minimal radius) which is consistent with the data.

2.2 PAC learning of hyper-rectangles

The proof in the case of hyper-rectangles is similar to the one given presented within the chapter. The algorithm selects the tightest axis-aligned hyper-rectangle containing all the sample points. For $i \in [2n]$, select a region r_i such that $\mathbb{P}_{\mathcal{D}}[r_i] = \epsilon/(2n)$ for each edge of the hyper-rectangle R . Assuming that $\mathbb{P}_{\mathcal{D}}[R - R'] > \epsilon$, argue that R' cannot meet all r_i s, so it must miss at least one. The probability that none of the m sample points falls into region r_i is $(1 - \epsilon/2n)^m$. By the union bound, this shows that

$$\mathbb{P}[R(R') > \epsilon] \leq 2n(1 - \epsilon/2n)^m \leq 2n \exp\left(-\frac{\epsilon m}{2n}\right). \quad (\text{E.35})$$

Setting δ to the right-hand side shows that for

$$m \geq \frac{2n}{\epsilon} \log \frac{2n}{\delta}, \quad (\text{E.36})$$

with probability at least $1 - \delta$, $R_{\mathcal{D}}(R') \leq \epsilon$.

2.3 Concentric circles

Suppose our target concept c is the circle around the origin with radius r . We will choose a slightly smaller radius s by

$$s := \inf\{s' : P(s' \leq \|x\| \leq r) < \epsilon\}.$$

Let A denote the annulus between radii s and r ; that is, $A := \{x : s \leq \|x\| \leq r\}$. By definition of s ,

$$P(A) \geq \epsilon. \quad (\text{E.37})$$

In addition, our generalization error, $P(c \Delta L(T))$, must be small if T intersects A . We can state this as

$$P(c \Delta L(T)) > \epsilon \implies T \cap A = \emptyset. \quad (\text{E.38})$$

Using (E.37), we know that any point in T chosen according to P will “miss” region A with probability at most $1 - \epsilon$. Defining $error := P(c \Delta L(T))$, we can combine this with (E.38) to see that

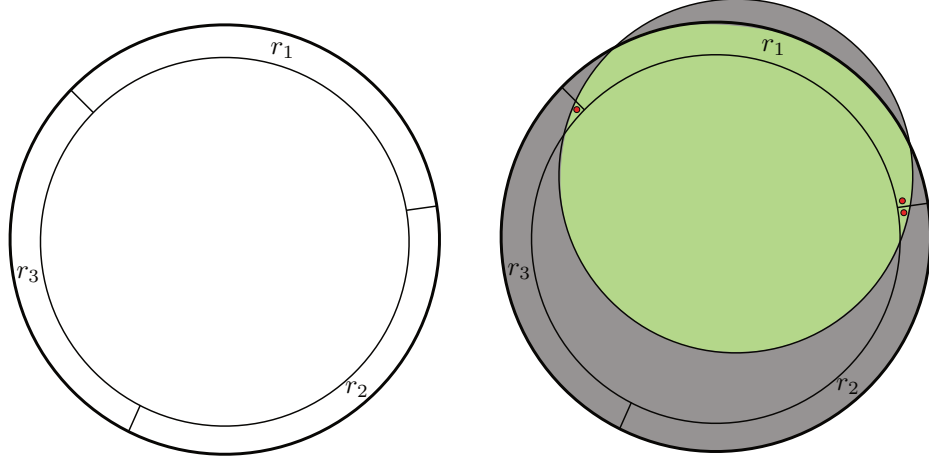
$$P(error > \epsilon) \leq P(T \cap A = \emptyset) \leq (1 - \epsilon)^m \leq e^{-m\epsilon}.$$

Setting δ to be greater than or equal to the right-hand side leads to $m \geq \frac{1}{\epsilon} \log(\frac{1}{\delta})$.

2.4 Non-concentric circles

As in the previous example, it is natural to assume the learning algorithm operates by returning the smallest circle which is consistent with the data. Gertrude is relying on the logical implication

$$error > \epsilon \implies T \cap r_i = \emptyset \text{ for some } i, \quad (\text{E.39})$$

**Figure E.5**

Counter-example shows error of tightest circle in gray.

which is not necessarily true here. Figure E.5 illustrates a counterexample. In the figure, we have one training point in each region r_i . The points in r_1 and r_2 are very close together, and the point in r_3 is very close to region r_1 . On this training data (some other points may be included outside the three regions r_i), our learned circle is the “tightest” circle including these points, and hence one diameter approximately traverses the corners of r_1 . In the figure, the gray regions are the error of this learned hypotheses versus the target circle, which has a thick border. Clearly, the error may be greater than ϵ even while $T \cap r_i \neq \emptyset$ for any i ; this contradicts (E.39) and invalidates poor Gertrude’s proof.

2.5 Triangles

As in the case of axis-aligned rectangles, consider three regions r_1, r_2, r_3 , along the sides of the target concept as indicated in figure E.6. Note that the triangle formed by the points A'', B'', C'' is similar to ABC (same angles) since $A''B''$ must be parallel to AB , and similarly for the other sides.

Assume that $\mathbb{P}[ABC] > \epsilon$, otherwise the statement would be trivial. Consider a triangle $A'B'C'$ similar to ABC and consistent with the training sample and such that it meets all three regions r_1, r_2, r_3 .

Since it meets r_1 , the line $A'B'$ must be below $A''B''$. Since it meets r_2 and r_3 , A' must be in r_2 and B' in r_3 (see figure E.6). Now, since the angle $A'B'C'$ is equal to $A''B''C''$, C' must be necessarily above C'' . This implies that triangle $A'B'C'$ contains $A''B''C''$, and thus $\text{error}(A'B'C') \leq \epsilon$.

$$\text{error}(A'B'C') > \epsilon \implies \exists i \in \{1, 2, 3\}: A'B'C' \cap r_i = \emptyset.$$

Thus, by the union bound,

$$\mathbb{P}[\text{error}(A'B'C') > \epsilon] \leq \sum_{i=1}^3 \mathbb{P}[A'B'C' \cap r_i = \emptyset] \leq 3(1 - \epsilon/3)^m \leq 3e^{-3m\epsilon}.$$

Setting δ to match the right-hand side gives the sample complexity $m \geq \frac{3}{\epsilon} \log \frac{3}{\delta}$.

2.8 Learning intervals

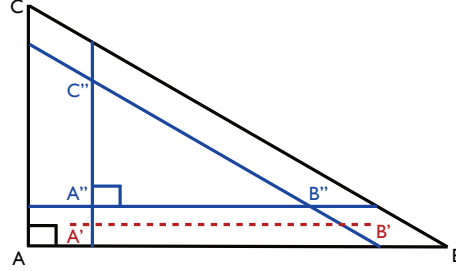


Figure E.6

Rectangle triangles.

Given a sample S , one algorithm consists of returning the tightest closed interval I_S containing positive points. Let $I = [a, b]$ be the target concept. If $\mathbb{P}[I] < \epsilon$, then clearly $R(I_S) < \epsilon$. Assume that $\mathbb{P}[I] \geq \epsilon$. Consider two intervals I_L and I_R defined as follows:

$$I_L = [a, x] \quad \text{with } x = \inf\{x: \mathbb{P}[a, x] \geq \epsilon/2\}$$

$$I_R = [x', b] \quad \text{with } x' = \sup\{x': \mathbb{P}[x', b] \geq \epsilon/2\}.$$

By the definition of x , the probability of $[a, x[$ is less than or equal to $\epsilon/2$, similarly the probability of $]x', b]$ is less than or equal to $\epsilon/2$. Thus, if I_S overlaps both with I_L and I_R , then its error region has probability at most ϵ . Thus, $R(I_S) > \epsilon$ implies that I_S does not overlap with either I_L or I_R , that is either none of the training points falls in I_L or none falls in I_R . Thus, by the union bound,

$$\begin{aligned} \mathbb{P}[R(I_S) > \epsilon] &\leq \mathbb{P}[S \cap I_L = \emptyset] + \mathbb{P}[S \cap I_R = \emptyset] \\ &\leq 2(1 - \epsilon/2)^m \leq 2e^{-m\epsilon/2}. \end{aligned}$$

Setting δ to match the right-hand side gives the sample complexity $m = \frac{2}{\epsilon} \log \frac{2}{\delta}$ and proves the PAC-learning of closed intervals. \square

2.9 Learning union of intervals

Given a sample S , our algorithm consists of the following steps:

- Sort S in ascending order.
- Loop through sorted S , marking where intervals of consecutive positively labeled points begin and end.
- Return the union of intervals found on the previous step. This union is represented by a list of tuples that indicate start and end points of the intervals.

This algorithm works both for $p = 2$ and for a general p . We will now consider the problem for \mathcal{C}_2 . To show that this is a PAC-learning algorithm we need to distinguish between two cases.

The first case is when our target concept is a disjoint union of two closed intervals: $I = [a, b] \cup [c, d]$. Note, there are two sources of error: false negatives in $[a, b]$ and $[c, d]$ and also false positives in (b, c) . False positives may occur if no sample is drawn from (b, c) . By linearity of expectation and since these two error regions are disjoint, we have that $R(h_S) = R_{FP}(h_S) + R_{FN,1}(h_S) + R_{FN,2}(h_S)$, where

$$R_{FP}(h_S) = \mathbb{P}_{x \sim \mathcal{D}}[x \in h_S, x \notin I],$$

$$R_{FN,1}(h_S) = \mathbb{P}_{x \sim \mathcal{D}}[x \notin h_S, x \in [a, b]],$$

$$R_{FN,2}(h_S) = \mathbb{P}_{x \sim \mathcal{D}}[x \notin h_S, x \in [c, d]].$$

Since we need to have that at least one of $R_{\text{FP}}(h_S)$, $R_{\text{FN},1}(h_S)$, $R_{\text{FN},2}(h_S)$ exceeds $\epsilon/3$ in order for $R(h_S) > \epsilon$, by union bound

$$\begin{aligned} \mathbb{P}(R(h_S) > \epsilon) &\leq \mathbb{P}(R_{\text{FP}}(h_S) > \epsilon/3 \text{ or } R_{\text{FN},1}(h_S) > \epsilon/3 \text{ or } R_{\text{FN},2}(h_S) > \epsilon/3) \\ &\leq \mathbb{P}(R_{\text{FP}}(h_S) > \epsilon/3) + \sum_{i=1}^2 \mathbb{P}(R_{\text{FN},i}(h_S) > \epsilon/3) \end{aligned} \quad (\text{E.40})$$

We first bound $\mathbb{P}(R_{\text{FP}}(h_S) > \epsilon/3)$. Note that if $R_{\text{FP}}(h_S) > \epsilon/3$, then $\mathbb{P}((b, c)) > \epsilon/3$ and hence

$$\mathbb{P}(R_{\text{FP}}(h_S) > \epsilon/3) \leq (1 - \epsilon/3)^m \leq e^{-m\epsilon/3}.$$

Now we can bound $\mathbb{P}(R_{\text{FN},i}(h_S) > \epsilon/3)$ by $2e^{-m\epsilon/6}$ using the same argument as in the previous question. Therefore,

$$\mathbb{P}(R(h_S) > \epsilon) \leq e^{-m\epsilon/3} + 4e^{-m\epsilon/6} \leq 5e^{-m\epsilon/6}.$$

Setting, the right-hand side to δ and solving for m yields that $m \geq \frac{6}{\epsilon} \log \frac{5}{\delta}$.

The second case that we need to consider is when $I = [a, d]$, that is, $[a, b] \cap [c, d] \neq \emptyset$. In that case, our algorithm reduces to the one from exercise 2.8 and it was already shown that only $m \geq \frac{2}{\epsilon} \log \frac{2}{\delta}$ samples is required to learn this concept. Therefore, we conclude that our algorithm is indeed a PAC-learning algorithm.

Extension of this result to the case of \mathcal{C}_p is straightforward. The only difference is that in (E.40), one has two summations for $p - 1$ regions of false positives and $2p$ regions of false negatives. In that case sample complexity is $m \geq \frac{2(2p-1)}{\epsilon} \log \frac{3p-1}{\delta}$.

Sorting step of our algorithm takes $O(m \log m)$ time and steps (b) and (c) are linear in m , which leads to overall time complexity $O(m \log m)$.

2.10 Consistent hypotheses

Since PAC-learning with L is possible for any distribution, let \mathcal{D} be the uniform distribution over \mathcal{Z} . Note that, in that case, the cost of an error of a hypothesis h on any point $z \in \mathcal{Z}$ is $\mathbb{P}_{\mathcal{D}}[z] = 1/m$. Thus, if $R_{\mathcal{D}}(h) < 1/m$, we must have $R_{\mathcal{D}}(h) = 0$ and h is consistent. Thus, choose $\epsilon = 1/(m+1)$. Then, for any $\delta > 0$, with probability at least $1 - \delta$ over samples S with $|S| \geq P((m+1), 1/\delta)$ points (where P is some fixed polynomial) the hypothesis h_S returned by L is consistent with Z since $R_{\mathcal{D}}(h_S) \leq 1/(m+1)$.

2.11 Senate laws

(a) The true error in the consistent case is bounded as follows:

$$R_{\mathcal{D}}(h) \leq \frac{1}{m} (\log |\mathcal{H}| + \log \frac{1}{\delta}). \quad (\text{E.41})$$

For $\delta = .05$, $m = 200$ and $|\mathcal{H}| = 2800$, $R_{\mathcal{D}}(h) \leq 5.5\%$.

(b) The true error in the inconsistent case is bounded as:

$$R_{\mathcal{D}}(h) \leq \hat{R}_{\mathcal{D}}(h) + \sqrt{\frac{1}{2m} (\log 2|\mathcal{H}| + \log \frac{1}{\delta})}. \quad (\text{E.42})$$

For $\delta = .05$, $\hat{R}_{\mathcal{D}}(h) = m'/m = .1$, $m = 200$ and $|\mathcal{H}| = 2800$, $R_{\mathcal{D}}(h) \leq 27.05\%$.

2.12 Bayesian bound. For any fixed $h \in \mathcal{H}$, by Hoeffding's inequality, for any $\delta > 0$,

$$\mathbb{P} \left[R(h) - \hat{R}_S(h) \geq \sqrt{\frac{\log \frac{1}{p(h)\delta}}{2m}} \right] \leq p(h)\delta. \quad (\text{E.43})$$

By the union bound,

$$\begin{aligned} \mathbb{P} \left[\exists h: R(h) - \hat{R}_S(h) \geq \sqrt{\frac{\log \frac{1}{p(h)\delta}}{2m}} \right] &\leq \sum_{h \in \mathcal{H}} \mathbb{P} \left[R(h) - \hat{R}_S(h) \geq \sqrt{\frac{\log \frac{1}{p(h)\delta}}{2m}} \right] \\ &\leq \sum_{h \in \mathcal{H}} p(h)\delta = \delta. \end{aligned}$$

In the case of a finite hypothesis set and a uniform prior $p(h) = 1/|\mathcal{H}|$, the bound coincides with the one presented in the chapter.

2.13 Learning with an unknown parameter.

(a) By definition of acceptance,

$$\begin{aligned} \mathbb{P}[h \text{ is accepted}] &= \mathbb{P}[\hat{R}_S(h) \leq 3/4\epsilon] \\ &\leq \mathbb{P}[\hat{R}_S(h) \leq 3/4 R(h)] && (R(h) \geq \epsilon) \\ &\leq \exp \left(-\frac{n}{2} R(h) (1/4)^2 \right) && (\text{Chernoff bound}) \\ &= \exp \left(-\frac{R(h)}{\epsilon} \log \frac{2^{i+1}}{\delta} \right) && (\text{def. of } n) \\ &= \exp \left(-\log \frac{2^{i+1}}{\delta} \right) = \frac{\delta}{2^{i+1}}. && (R(h) \geq \epsilon) \end{aligned}$$

(b) By definition, $\mathbb{P}[h \text{ is rejected}] = \mathbb{P}[\hat{R}_S(h) \geq \frac{3}{4}\epsilon]$. Since $R(h) \leq \epsilon/2$, $\mathbb{P}[h \text{ is rejected}] \leq \mathbb{P}[\hat{R}_S(h) \geq \frac{3}{4}\epsilon \mid R(h) = \epsilon/2]$. By the Chernoff bounds, we can thus write

$$\begin{aligned} \mathbb{P}[h \text{ is rejected}] &\leq \exp \left(-\frac{n}{3} \frac{\epsilon}{2} (1/2)^2 \right) && (\text{Chernoff bound}) \\ &= \exp \left(-\frac{4}{3} \log \frac{2^{i+1}}{\delta} \right) && (\text{def. of } n) \\ &\leq \exp \left(-\log \frac{2^{i+1}}{\delta} \right) = \frac{\delta}{2^{i+1}}. \end{aligned}$$

(c) The estimate \tilde{s} is then an upper bound on s and thus, by definition of algorithm \mathcal{B} , $\mathbb{P}[R(h_i) \leq \epsilon/2] \geq 1/2$. If a hypothesis h has error at least $\epsilon/2$ it is rejected with probability at most $\delta/2^{i+1} \leq \delta/4 \leq 1/4$, therefore, it is accepted with probability at $3/4$. Thus, for $\tilde{s} \geq s$, $\mathbb{P}[h_i \text{ is accepted}] \geq 1/2 \times 3/4 = 3/8$.

(d) By the previous question, the probability that algorithm \mathcal{B} fails to halt while $\tilde{s} \geq s$ is at most $1 - 3/8 = 5/8$. Thus, the probability that it does not halt after j iterations is at most $(5/8)^j \leq (5/8)^{\log \frac{2}{\delta} / \log \frac{8}{5}} = \exp \left(\log \frac{2}{\delta} / \log \frac{8}{5} \log \frac{5}{8} \right) = \delta/2$.

(e) By definition,

$$\begin{aligned} \tilde{s} \geq s &\iff \lfloor 2^{(i-1)/\log \frac{2}{\delta}} \rfloor \geq s \\ &\iff 2^{(i-1)/\log \frac{2}{\delta}} \geq s \\ &\iff \frac{i-1}{\log \frac{2}{\delta}} \geq \log_2 s \\ &\iff i \geq 1 + (\log_2 s) \log \frac{2}{\delta} \\ &\iff i \geq \lceil 1 + (\log_2 s) \log \frac{2}{\delta} \rceil. \end{aligned}$$

(f) In view of the two previous questions, with probability at least $1 - \delta/2$, algorithm \mathcal{B} halts after at most j' iterations. The probability that the hypothesis it returns be accepted while its error is greater than ϵ is at most $\delta/2^{j'+1} \leq \delta/2$. Thus, with probability $1 - \delta$, the algorithm halts and the hypothesis it returns has error at most ϵ . \square

Chapter 3

3.1 Growth function of intervals in \mathbb{R}

Let $\{x_1, \dots, x_m\}$ be a set of m distinct ordered real numbers and let I be an interval. If $I \cap \{x_1, \dots, x_m\}$ contains k numbers, their indices are of the form $j, j+1, \dots, j+k-1$. How many different set of indices are possible? The answer is $m-k+1$ since j can take values in $\{1, \dots, m-k+1\}$. Thus, the total number of dichotomies for a set of size m is:

$$1 + \sum_{k=1}^m (m-k+1) = 1 + \frac{m(m+1)}{2} = \binom{m}{0} + \binom{m}{1} + \binom{m}{2},$$

which matches the general bound on shattering coefficients.

3.2 Growth function and Rademacher complexity of thresholds in \mathbb{R}

Given m distinct points on the line, there are at most $m+1$ choices of the threshold θ leading to distinct classifications: between two points or before/after all points. Since there are two choices (classifying those to the right as positive or those to the left), there are $2(m+1)$ choices. Thus, $\Pi_m(\mathcal{H}) \leq 2m$ and, by the bound on the Rademacher complexity shown in the chapter, $\mathfrak{R}_m(\mathcal{H}) \leq \sqrt{\frac{2 \log(2m)}{m}}$.

3.3 Growth function of linear combinations

- (a) $\{\mathcal{X}^+ \cup \{\mathbf{x}_{m+1}\}, \mathcal{X}^-\}$ and $\{\mathcal{X}^+, \mathcal{X}^- \cup \{\mathbf{x}_{m+1}\}\}$ are linearly separable by a hyperplane going through the origin if and only if there exists $\mathbf{w}_1 \in \mathbb{R}^d$ such that

$$\forall \mathbf{x} \in \mathcal{X}^+, \mathbf{w}_1 \cdot \mathbf{x} > 0 \quad \forall \mathbf{x} \in \mathcal{X}^-, \mathbf{w}_1 \cdot \mathbf{x} < 0, \text{ and } \mathbf{w}_1 \cdot \mathbf{x}_{m+1} > 0 \quad (\text{E.44})$$

and there exists $\mathbf{w}_2 \in \mathbb{R}^d$ such that

$$\forall \mathbf{x} \in \mathcal{X}^+, \mathbf{w}_2 \cdot \mathbf{x} > 0 \quad \forall \mathbf{x} \in \mathcal{X}^-, \mathbf{w}_2 \cdot \mathbf{x} < 0, \text{ and } \mathbf{w}_2 \cdot \mathbf{x}_{m+1} < 0. \quad (\text{E.45})$$

For any $\mathbf{w}_1, \mathbf{w}_2$, the function $f: (t \mapsto t\mathbf{w}_1 + (1-t)\mathbf{w}_2) \cdot \mathbf{x}_{m+1}$ is continuous over $[0, 1]$. (E.44) and (E.45) hold iff $f(0) < 0$ and $f(1) > 0$, that is iff there exists $\mathbf{w} = t_0\mathbf{w}_1 + (1-t_0)\mathbf{w}_2$ linearly separating $\{\mathcal{X}^+, \mathcal{X}^-\}$ and such at $\mathbf{w} \cdot \mathbf{x}_{m+1} = 0$.

- (b) Repeating the formula, we obtain $C(m, d) = \sum_{k=0}^{m-1} \binom{m-1}{k} C(1, d-k)$. Since, $C(1, n) = 2$ if $n \geq 1$ and $C(1, n) = 0$ otherwise, the result follows.

- (c) This is a direct application of the result of the previous question.

3.4 Lower bound on growth function

Let X be an arbitrary set. Consider the set of all subsets of X of size less than or equal to d . The indicator functions of these sets forms a hypothesis class whose shattering coefficient is equal to the general upper bound.

3.5 Finer Rademacher upper bound

Following the proof given in the chapter and using Jensen's inequality (at the last inequality), we can write:

$$\begin{aligned} \widehat{\mathfrak{R}}_m(\mathcal{G}) &= \mathbb{E}_{S, \sigma} \left[\sup_{g \in \mathcal{G}} \frac{1}{m} \begin{bmatrix} \sigma_1 \\ \vdots \\ \sigma_m \end{bmatrix} \cdot \begin{bmatrix} g(z_1) \\ \vdots \\ g(z_m) \end{bmatrix} \right] \\ &\leq \mathbb{E}_S \left[\frac{\sqrt{m} \sqrt{2 \log |\{(g(z_1), \dots, g(z_m)) : g \in \mathcal{G}\}|}}{m} \right] \quad (\text{Massart's Lemma}) \\ &= \mathbb{E}_S \left[\frac{\sqrt{m} \sqrt{2 \log \Pi(\mathcal{G}, S)}}{m} \right] \\ &\leq \frac{\sqrt{m} \sqrt{2 \log \mathbb{E}_S [\Pi(\mathcal{G}, S)]}}{m} = \sqrt{\frac{2 \log \mathbb{E}_S [\Pi(\mathcal{G}, S)]}{m}}. \end{aligned}$$

Note that in the application of Jensen's inequality, we are using the fact that $\sqrt{\log(x)}$ is concave, which is true because it is the composition of functions that are concave functions

and the outer function is non-decreasing. It is not true in general that the composition of any two concave functions is concave.

3.6 Singleton hypothesis class

(a)

$$\begin{aligned}\mathfrak{R}_m(\mathcal{H}) &= \mathbb{E}_{S \sim \mathcal{D}^m, \sigma} \left[\sup_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \sigma_i h(z_i) \right] = \mathbb{E}_{S \sim \mathcal{D}^m, \sigma} \left[\frac{1}{m} \sum_{i=1}^m \sigma_i h_0(z_i) \right] \\ &= \mathbb{E}_{S \sim \mathcal{D}^m} \left[\frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\sigma} [\sigma_i] h_0(z_i) \right] = 0,\end{aligned}$$

where the final equality follows from the fact $\mathbb{E}[\sigma] = 0$ for a Rademacher random variable σ .

(b) Consider the set $A = \mathbf{x}_0$, where \mathbf{x}_0 is any vector in \mathbb{R}^m with $\|\mathbf{x}_0\| = r$. Then,

$$\mathbb{E}_{\sigma} \left[\frac{1}{m} \sup_{\mathbf{x} \in A} \sum_{i=1}^m \sigma_i x_i \right] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\sigma} [\sigma_i] x'_i = 0$$

and,

$$\frac{r \sqrt{2 \log |A|}}{m} = 0.$$

3.7 Two function hypothesis class

(a) $\text{VCdim}(\mathcal{H}) = 1$ since \mathcal{H} can shatter one point and clearly at most one. Observe that

$$\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) = \sup_{h \in \mathcal{H}} \left(\sum_{i=1}^m \sigma_i \right) h(x_1) = \left| \sum_{i=1}^m \sigma_i \right|. \quad (\text{E.46})$$

Thus, by Jensen's inequality,

$$\begin{aligned}\widehat{\mathfrak{R}}_S(\mathcal{H}) &= \frac{1}{m} \mathbb{E}_{\sigma} \left[\left| \sum_{i=1}^m \sigma_i \right| \right] \\ &\leq \frac{1}{m} \left[\mathbb{E}_{\sigma} \left[\left(\sum_{i=1}^m \sigma_i \right)^2 \right] \right]^{1/2} \\ &= \frac{1}{m} \left[\mathbb{E}_{\sigma} \left[\sum_{i=1}^m \sigma_i^2 \right] \right]^{1/2} \quad (\mathbb{E}[\sigma_i \sigma_j] = 0 \text{ for } i \neq j) \\ &= \frac{1}{\sqrt{m}}.\end{aligned}$$

By the Khintchine inequality, the upper bound is tight modulo the constant $1/\sqrt{2}$. The upper bound coincides with $\sqrt{d/m}$.

(b) $\text{VCdim}(\mathcal{H}) = 1$ since \mathcal{H} can shatter x_1 and clearly at most one point. By definition,

$$\begin{aligned}\widehat{\mathfrak{R}}_S(\mathcal{H}) &= \frac{1}{m} \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}} \sigma_1 h(x_1) - \sum_{i=2}^m \sigma_i \right] \\ &= \frac{1}{m} \mathbb{E}_{\sigma_1} \left[\sup_{h \in \mathcal{H}} \sigma_1 h(x_1) \right] \quad (\mathbb{E}[\sigma_i] = 0) \\ &= \frac{1}{m} \mathbb{E}_{\sigma_1} [1] = \frac{1}{m}.\end{aligned}$$

Here $\widehat{\mathfrak{R}}_S(\mathcal{H})$ is a clearly more favorable quantity than $\sqrt{d/m} = \sqrt{1/m}$.

3.8 Rademacher identities

(a) If $\alpha \geq 0$, then

$$\sup_{h \in \alpha \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) = \sup_{h \in \mathcal{H}} \sum_{i=1}^m \alpha \sigma_i h(x_i) = \alpha \sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i),$$

otherwise if $\alpha < 0$, then

$$\sup_{h \in \alpha \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) = \sup_{h \in \alpha \mathcal{H}} \sum_{i=1}^m \alpha \sigma_i h(x_i) = (-\alpha) \sup_{h \in \mathcal{H}} \sum_{i=1}^m (-\sigma_i) h(x_i).$$

Since σ_i and $-\sigma_i$ have the same distribution, this shows that $\mathfrak{R}_m(\alpha \mathcal{H}) = |\alpha| \mathfrak{R}_m(\mathcal{H})$.

(b) The second equality follows from:

$$\begin{aligned} & \mathfrak{R}_m(\mathcal{H} + \mathcal{H}') \\ &= \frac{1}{m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i (h(x_i) + h'(x_i)) \right] \\ &= \frac{1}{m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i h(x_i) + \sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i h'(x_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) \right] + \frac{1}{m} \mathbb{E}_{\sigma, S} \left[\sup_{h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i h'(x_i) \right]. \end{aligned}$$

(c) For the inequality, using the identity $\max(a, b) = \frac{1}{2}[a + b + |a - b|]$ valid for all $a, b \in \mathbb{R}$, and the sub-additivity of sup we can write:

$$\begin{aligned} & \mathfrak{R}_m(\{\max(h, h') : h \in \mathcal{H}, h' \in \mathcal{H}'\}) \\ &= \frac{1}{2m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i [h(x_i) + h'(x_i) + |h(x_i) - h'(x_i)|] \right] \\ &\leq \frac{1}{2} [\mathfrak{R}_m(\mathcal{H}) + \mathfrak{R}_m(\mathcal{H}')] + \frac{1}{2m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i |h(x_i) - h'(x_i)| \right]. \end{aligned}$$

Since the absolute value function is 1-Lipschitz, by the contraction lemma, the third term can be bounded as follows

$$\begin{aligned} & \frac{1}{2m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i |h(x_i) - h'(x_i)| \right] \\ &\leq \frac{1}{2m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i [h(x_i) - h'(x_i)] \right] \\ &\leq \frac{1}{2m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i h(x_i) + \sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m -\sigma_i h'(x_i) \right] \\ &= \frac{1}{2m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m \sigma_i h(x_i) \right] + \frac{1}{2m} \mathbb{E}_{\sigma, S} \left[\sup_{h \in \mathcal{H}, h' \in \mathcal{H}'} \sum_{i=1}^m -\sigma_i h'(x_i) \right] \\ &= \frac{1}{2} [\mathfrak{R}_m(\mathcal{H}) + \mathfrak{R}_m(\mathcal{H}')], \end{aligned}$$

using the fact that σ_i and $-\sigma_i$ follow the same distribution.

3.9 Rademacher complexity of intersection of concepts

Observe that for any $h_1 \in \mathcal{H}_1$ and $h_2 \in \mathcal{H}_2$, we can write $h_1 h_2 = (h_1 + h_2 - 1)1_{h_1 + h_2 - 1 \geq 0} = (h_1 + h_2 - 1)_+$. Since $x \mapsto (x - 1)_+$ is 1-Lipschitz over $[0, 2]$, by Talagrand's contraction lemma, the following holds: $\mathfrak{R}_S(\mathcal{H}) \leq \mathfrak{R}_S(\mathcal{H}_1 + \mathcal{H}_2) \leq \mathfrak{R}_S(\mathcal{H}_1) + \mathfrak{R}_S(\mathcal{H}_2)$.

Since concepts can be identified with indicator functions, the intersection of two concepts can be identified with the product two such indicator functions. In view of that, by the result just proven and after taking expectations, the following holds:

$$\mathfrak{R}_m(\mathcal{C}) \leq \mathfrak{R}_m(\mathcal{C}_1) + \mathfrak{R}_m(\mathcal{C}_2).$$

3.10 Rademacher complexity of prediction vector

(a) The following proves the result:

$$\begin{aligned} \widehat{\mathfrak{R}}_S(\mathcal{H}) &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{f \in \{h, -h\}} \sum_{i=1}^m \sigma_i f(x_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} [\max\{\boldsymbol{\sigma} \cdot \mathbf{u}, -\boldsymbol{\sigma} \cdot \mathbf{u}\}] \\ &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} [|\boldsymbol{\sigma} \cdot \mathbf{u}|] \\ &\leq \frac{1}{m} \sqrt{\mathbb{E}_{\boldsymbol{\sigma}} [|\boldsymbol{\sigma} \cdot \mathbf{u}|^2]} \quad (\text{by Jensen's inequality}) \\ &= \frac{1}{m} \sqrt{\mathbb{E}_{\boldsymbol{\sigma}} \left[\sum_{i,j=1}^m \sigma_i \sigma_j u_i u_j \right]} \\ &= \frac{1}{m} \sqrt{\sum_{i,j=1}^m \mathbb{E}_{\boldsymbol{\sigma}} [\sigma_i \sigma_j] u_i u_j} \\ &= \frac{1}{m} \sqrt{\sum_{i=1}^m u_i^2} \quad (\mathbb{E}_{\boldsymbol{\sigma}} [\sigma_i \sigma_j] = \mathbb{E}_{\boldsymbol{\sigma}} [\sigma_i] \mathbb{E}_{\boldsymbol{\sigma}} [\sigma_j] = 0 \text{ for } i \neq j) \\ &= \frac{\|\mathbf{u}\|}{m}. \end{aligned}$$

Thus, $\widehat{\mathfrak{R}}_S(\mathcal{H}) \leq \frac{\sqrt{n}}{m}$. For $n = 1$, $\widehat{\mathfrak{R}}_S(\mathcal{H}) \leq \frac{1}{m}$ while for $n = m$, $\widehat{\mathfrak{R}}_S(\mathcal{H}) \leq \frac{1}{\sqrt{m}}$.

(b) The empirical Rademacher complexity of $\mathcal{F} + h$ can be expressed as follows:

$$\begin{aligned} \widehat{\mathfrak{R}}_S(\mathcal{F} + h) &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(x_i) + \sigma_i h(x_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(x_i) \right] + \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sum_{i=1}^m \sigma_i h(x_i) \right] \\ &= \widehat{\mathfrak{R}}_S(\mathcal{F}) + \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\boldsymbol{\sigma}} [\sigma_i] h(x_i) = \widehat{\mathfrak{R}}_S(\mathcal{F}). \end{aligned}$$

The empirical Rademacher complexity of $\mathcal{F} \pm h$ can be upper bounded as follows using the first question:

$$\begin{aligned} \widehat{\mathfrak{R}}_S(\mathcal{F} + h) &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(x_i) + \sup_{s \in \{-1, +1\}} s \sum_{i=1}^m \sigma_i h(x_i) \right] \\ &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{f \in \mathcal{F}} \sum_{i=1}^m \sigma_i f(x_i) \right] + \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{s \in \{-1, +1\}} s \sum_{i=1}^m \sigma_i h(x_i) \right] \\ &\leq \widehat{\mathfrak{R}}_S(\mathcal{F}) + \frac{\|\mathbf{u}\|}{m}. \end{aligned}$$

3.11 Rademacher complexity of regularized neural networks

(a)

$$\begin{aligned}
\widehat{\mathfrak{R}}_S(\mathcal{H}) &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{w}\|_1 \leq \Lambda', \|\mathbf{u}_j\|_2 \leq \Lambda} \sum_{i=1}^m \sigma_i \sum_{j=1}^{n_2} w_j \sigma(\mathbf{u}_j \cdot \mathbf{x}_i) \right] \\
&= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{w}\|_1 \leq \Lambda', \|\mathbf{u}_j\|_2 \leq \Lambda} \sum_{j=1}^{n_2} w_j \sum_{i=1}^m \sigma_i \sigma(\mathbf{u}_j \cdot \mathbf{x}_i) \right] \\
&= \frac{\Lambda'}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{u}_j\|_2 \leq \Lambda} \max_{j \in [n_2]} \left| \sum_{i=1}^m \sigma_i \sigma(\mathbf{u}_j \cdot \mathbf{x}_i) \right| \right] \quad (\text{all the weight put on largest term}) \\
&= \frac{\Lambda'}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{u}_j\|_2 \leq \Lambda, j \in [n_2]} \left| \sum_{i=1}^m \sigma_i \sigma(\mathbf{u}_j \cdot \mathbf{x}_i) \right| \right] \\
&= \frac{\Lambda'}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{u}\|_2 \leq \Lambda} \left| \sum_{i=1}^m \sigma_i \sigma(\mathbf{u} \cdot \mathbf{x}_i) \right| \right].
\end{aligned}$$

(b) By Talagrand's lemma, since σ is L -Lipschitz, the following inequality holds:

$$\begin{aligned}
\widehat{\mathfrak{R}}_S(\mathcal{H}) &\leq \frac{\Lambda' L}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \left| \sum_{i=1}^m \sigma_i \mathbf{u} \cdot \mathbf{x}_i \right| \right] \\
&= \frac{\Lambda' L}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{h \in \mathcal{H}} \sup_{s \in \{-1, +1\}} s \sum_{i=1}^m \sigma_i \mathbf{u} \cdot \mathbf{x}_i \right] \quad (\text{def. of abs. value}) \\
&= \Lambda' L \widehat{\mathfrak{R}}_S(\mathcal{H}').
\end{aligned}$$

(c)

$$\begin{aligned}
\widehat{\mathfrak{R}}_S(\mathcal{H}') &= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{u}\|_2 \leq \Lambda, s \in \{-1, +1\}} \sum_{i=1}^m \sigma_i s \mathbf{u} \cdot \mathbf{x}_i \right] \\
&= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{u}\|_2 \leq \Lambda} \left| \sum_{i=1}^m \sigma_i \mathbf{u} \cdot \mathbf{x}_i \right| \right] \quad (\text{def. of abs. val.}) \\
&= \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{\|\mathbf{u}\|_2 \leq \Lambda} \left| \mathbf{u} \cdot \sum_{i=1}^m \sigma_i \mathbf{x}_i \right| \right] \\
&= \frac{\Lambda}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2 \right] \quad (\text{Cauchy-Schwarz eq. case}).
\end{aligned}$$

The last equality holds by setting $\mathbf{u} = \frac{\Lambda \sum_{i=1}^m \sigma_i \mathbf{x}_i}{\left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|}$.

(d)

$$\begin{aligned}
\widehat{\mathfrak{R}}_S(\mathcal{H}') &= \frac{\Lambda}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2 \right] \\
&\leq \frac{\Lambda}{m} \sqrt{\mathbb{E}_{\boldsymbol{\sigma}} \left[\left\| \sum_{i=1}^m \sigma_i \mathbf{x}_i \right\|_2^2 \right]} && \text{(Jensen's ineq.)} \\
&= \frac{1}{m} \sqrt{\sum_{i,j=1}^m \mathbb{E}_{\boldsymbol{\sigma}} [\sigma_i \sigma_j] (\mathbf{x}_i \cdot \mathbf{x}_j)} \\
&= \frac{\Lambda}{m} \sqrt{\sum_{i,j=1}^m 1_{i=j} (\mathbf{x}_i \cdot \mathbf{x}_j)} && \text{(independence of } \sigma_i \text{'s)} \\
&= \frac{\Lambda}{m} \sqrt{\sum_{i=1}^m \|\mathbf{x}_i\|_2^2}.
\end{aligned}$$

(e) In view of the previous questions,

$$\widehat{\mathfrak{R}}_S(\mathcal{H}) \leq \Lambda' L \widehat{\mathfrak{R}}_S(\mathcal{H}') \leq \frac{\Lambda' \Lambda L}{m} \sqrt{\sum_{i=1}^m \|\mathbf{x}_i\|_2^2} \leq \frac{\Lambda' \Lambda L}{m} \sqrt{mr^2} = \frac{\Lambda' \Lambda L r}{\sqrt{m}}.$$

3.12 Rademacher complexity

Consider the simple case where \mathcal{H} is reduced to the constant hypothesis $h_1: x \mapsto 1$ and $h_{-1}: x \mapsto -1$. Then, by definition of the empirical Rademacher complexity,

$$\widehat{\mathfrak{R}}_S(\mathcal{H}) = \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} [\max\{\sum_{i=1}^m \sigma_i, \sum_{i=1}^m -\sigma_i\}] = \frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} [\sum_{i=1}^m \sigma_i]$$

Let $X = \sum_{i=1}^m \sigma_i$. Note that $\mathbb{E}[X^2] = \mathbb{E}[\sum_{i,j=1}^m \sigma_i \sigma_j]$. For any $i \neq j$, since σ_i and σ_j are independent, $\mathbb{E}[\sigma_i \sigma_j] = \mathbb{E}[\sigma_i] \mathbb{E}[\sigma_j] = 0$. Thus,

$$\mathbb{E}[X^2] = \sum_{i=1}^m \mathbb{E}[\sigma_i^2] = \sum_{i=1}^m \mathbb{E}[\sigma_i^2] = m.$$

Now, by Hölder's inequality,

$$m = \mathbb{E}[X^2] = \mathbb{E}[|X|^{2/3} |X|^{4/3}] \leq \mathbb{E}[|X|^{2/3}] \mathbb{E}[|X|^{4/3}]^{1/3}.$$

Thus,

$$\begin{aligned}
\mathbb{E}[|X|] &\geq \frac{m^{3/2}}{\mathbb{E}[X^4]^{1/2}} = \frac{m^{3/2}}{\sqrt{\mathbb{E}[\sum_{i=1}^m \sigma_i^4 + 3 \sum_{i \neq j} \sigma_i^2 \sigma_j^2]}} = \frac{m^{3/2}}{\sqrt{m + 3m(m-1)}} \\
&= \frac{m^{3/2}}{\sqrt{m(3m-2)}} \geq \frac{m^{3/2}}{\sqrt{m(3m)}} = \frac{\sqrt{m}}{\sqrt{3}}.
\end{aligned}$$

This shows that

$$\widehat{\mathfrak{R}}_S(\mathcal{H}) \geq \frac{\sqrt{m}}{\sqrt{3}}.$$

Since $\mathfrak{R}_m(\mathcal{H}) \geq \widehat{\mathfrak{R}}_S(\mathcal{H}) + O(\frac{1}{\sqrt{m}})$, it implies $\mathfrak{R}_m(\mathcal{H}) \geq O(\frac{1}{\sqrt{m}})$, which contradicts $\mathfrak{R}_m(\mathcal{H}) \leq O(\frac{1}{m})$.

Note that for the lower bound, we could have used instead a more general result (Khinchine's inequality) which states that for any $\mathbf{a} \in \mathbb{R}^m$,

$$\mathbb{E}[|\boldsymbol{\sigma} \cdot \mathbf{a}|] \geq \frac{\|\mathbf{a}\|_2}{\sqrt{2}}.$$

3.13 VC-dimension of union of k intervals

The VC-dimension of this class is $2k$. It is not hard to see that any $2k$ distinct points on the real line can be shattered using k intervals; it suffices to shatter each of the k pairs of consecutive points with an interval. Assume now that $2k+1$ distinct points $x_1 < \dots < x_{2k+1}$ are given. For any $i \in [2k+1]$, label x_i with $(-1)^{i+1}$, that is alternatively label points with 1 or -1 . This leads to $k+1$ points labeled positively and requires $2k+1$ intervals to shatter the set, since no interval can contain two consecutive points. Thus, no set of $2k+1$ points can be shattered by k intervals, and the VC-dimension of the union of k intervals is $2k$.

3.14 VC-dimension of finite hypothesis sets

With a finite set \mathcal{H} , at most $2^{|\mathcal{H}|}$ dichotomies can be defined.

3.15 VC-dimension of subsets

The set of three points $\{0, 3/4, 3/2\}$ can be fully shattered as follows:

+++	$\alpha = -2$
++-	$\alpha = 0$
+ - +	$\alpha = -1$
+ - -	$\alpha = 3/2 - 2 + \epsilon$
- + +	$\alpha = 3/4 - 2$
- + -	$\alpha = \epsilon$
- - +	$\alpha = 3/2$
- - -	$\alpha = 3/2 + \epsilon$,

where ϵ is a small number, e.g., $\epsilon = .1$. No set of four points $x_1 < x_2 < x_3 < x_4$ can be labeled by $+ - + -$. This is because the three leftmost labels $+ - +$ imply that $\alpha + 2 \leq x_3$ and thus also $\alpha + 2 < x_4$. Thus, the VC-dimension of the set of subsets I_α is 3. Note that this does not coincide with the number of parameters used to describe the class.

3.16 VC-dimension of axis-aligned squares and triangles

- (a) It is not hard to see that the set of 3 points with coordinates $(1,0)$, $(0,1)$, and $(-1,0)$ can be shattered by axis-aligned squares: e.g., to label positively two of these points, use a square defined by the axes and with those two points as corners. Thus, the VC-dimension is at least 3. No set of 4 points can be fully shattered. To see this, let P_T be the highest point, P_B the lowest, P_L the leftmost, and P_R the rightmost, assuming for now that these can be defined in a unique way (no tie) – the cases where there are ties can be treated in a simpler fashion. Assume without loss of generality that the difference d_{BT} of y -coordinates between P_T and P_B is greater than the difference d_{LR} of x -coordinates between P_L and P_R . Then, P_T and P_B cannot be labeled positively while P_L and P_R are labeled negatively. Thus, the VC-dimension of axis-aligned squares in the plane is 3.
- (b) Check that the set of 4 points with coordinates $(1,0)$, $(0,1)$, $(-1,0)$, and $(0,-1)$ can be shattered by such triangles. This is essentially the same as the case with axis aligned rectangles. To see that no five points can be shattered, the same example or argument as for axis-aligned rectangles can be used: labeling all points positively except from the one within the interior of the convex hull is not possible (for the degenerate cases where no point is in the interior of the convex hull is simpler, this is even easier to see). Thus, the VC-dimension of this family of triangles is 4.

3.17 VC-dimension of closed balls in \mathbb{R}^n .

Let $B(a, r)$ be the ball of radius r centered at $a \in \mathbb{R}^n$. Then $x \in B(a, r)$ iff

$$\sum_{i=1}^n \|x_i\|^2 - 2 \sum_{i=1}^n a_i x_i + \sum_{i=1}^n a_i^2 - r \leq 0, \quad (\text{E.47})$$

which is equivalent to

$$\langle W, X \rangle + B \leq 0, \quad (\text{E.48})$$

with $W = \begin{bmatrix} 1 \\ -2a_1 \\ \dots \\ -2a_n \end{bmatrix}$, $X = \begin{bmatrix} \sum_{i=1}^n \|x_i\|^2 \\ x_1 \\ \dots \\ x_n \end{bmatrix}$, and $B = \sum_{i=1}^n a_i^2 - r$. The VC-dimension of

closed balls in \mathbb{R}^n is thus at most equal to the VC-dimension of hyperplanes in \mathbb{R}^{n+1} , that is, $n+2$.

3.18 VC-dimension of ellipsoids

The general equation of ellipsoids in \mathbb{R}^n is

$$(\mathbf{X} - \mathbf{X}_0)^\top \mathbf{A} (\mathbf{X} - \mathbf{X}_0) \leq 1, \quad (\text{E.49})$$

where $\mathbf{X}, \mathbf{X}_0 \in \mathbb{R}^n$ and $\mathbf{A} = (a_{ij}) \in \mathbb{S}_+^n$ is a positive semidefinite symmetric matrix. This can be rewritten as

$$\mathbf{X}^\top \mathbf{A} \mathbf{X} - 2\mathbf{X}^\top \mathbf{A} \mathbf{X}_0 + \mathbf{X}_0^\top \mathbf{A} \mathbf{X}_0 \leq 1, \quad (\text{E.50})$$

or $\sum_{i,j=1}^n a_{ij}(x_i x_j + x_j x_i) - \sum_{i=1}^n 2(\mathbf{A} \mathbf{X}_0)_i x_i + (\mathbf{X}_0^\top \mathbf{A} \mathbf{X}_0 - 1) \leq 0$ using the fact that \mathbf{A} is symmetric. Let $a_i = -2(\mathbf{A} \mathbf{X}_0)_i$ for $i \in [n]$ and let $b = \mathbf{X}_0^\top \mathbf{A} \mathbf{X}_0 - 1$. Then this can be viewed in terms of the following equations of hyperplanes in $\mathbb{R}^{n(n+1)/2+n}$

$$\mathbf{W}^\top \mathbf{Z} + b \leq 0, \quad (\text{E.51})$$

with

$$\mathbf{W} = \begin{bmatrix} a_1 \\ \dots \\ a_n \\ a_{11} \\ \dots \\ a_{ij} \\ \dots \\ a_{nn} \end{bmatrix}, \quad \mathbf{Z} = \begin{bmatrix} x_1 \\ \dots \\ x_n \\ x_1 x_1 + x_1 x_1 \\ \dots \\ x_i x_j + x_j x_i \\ \dots \\ x_n x_n + x_n x_n \end{bmatrix} \left. \vphantom{\begin{bmatrix} a_1 \\ \dots \\ a_n \\ a_{11} \\ \dots \\ a_{ij} \\ \dots \\ a_{nn} \end{bmatrix}} \right\} n(n+1)/2 + n \quad (\text{E.52})$$

Since the VC-dimension of hyperplanes in $\mathbb{R}^{n(n+1)/2+n}$ is $n(n+1)/2+n+1 = (n+1)(n/2+1)$, the VC-dimension of ellipsoids in \mathbb{R}^n is bounded by $(n+1)(n+2)/2$.

3.19 VC-dimension of a vector space of real functions

Show that no set of size $m = r+1$ can be shattered by \mathcal{H} . Let x_1, \dots, x_m be m arbitrary points. Define the linear mapping $l: F \rightarrow \mathbb{R}^m$ defined by:

$$l(f) = (f(x_1), \dots, f(x_m))$$

Since the dimension of $\dim(F) = m-1$, the rank of l is at most $m-1$, and there exists $\alpha \in \mathbb{R}^m$ orthogonal to $l(F)$:

$$\forall f \in F, \sum_{i=1}^m \alpha_i f(x_i) = 0$$

We can assume that at least one α_i is negative. Then,

$$\forall f \in F, \sum_{i: \alpha_i \geq 0} \alpha_i f(x_i) = - \sum_{i: \alpha_i < 0} \alpha_i f(x_i)$$

Now, assume that there exists a set $\{x: f(x) \geq 0\}$ selecting exactly the x_i s on the left-hand side. Then all the terms on the left-hand side are non-negative, while those on the right-hand side are negative, which cannot be. Thus, $\{x_1, \dots, x_m\}$ cannot be shattered.

3.20 VC-dimension of sine functions

- (a) Fix $x \in \mathbb{R}$ and suppose there exists an ω that realizes the labeling $--+-$. Thus $\sin(\omega x) < 0$, $\sin(2\omega x) < 0$, $\sin(3\omega x) \geq 0$ and $\sin(4\omega x) < 0$. We will show that this implies $\sin^2(\omega x) < \frac{1}{2}$ and $\sin^2(\omega x) \geq \frac{3}{4}$, a contradiction.

Using the identity $\sin(2\theta) = 2\sin(\theta)\cos(\theta)$ and the fact that $\sin(4\omega x) < 0$ we have

$$2\sin(2\omega x)\cos(2\omega x) = \sin(4\omega x) < 0.$$

Since $\sin(2\omega x) < 0$ we can divide both sides of this inequality by $2\sin(2\omega x)$ to conclude $\cos(2\omega x) > 0$. Applying the identity $\cos(2\theta) = 1 - 2\sin^2(\theta)$ yields $1 - 2\sin^2(\omega x) > 0$, or $\sin^2(\omega x) < \frac{1}{2}$.

Using the identity $\sin(3\theta) = 3\sin(\theta) - 4\sin^3(\theta)$ and the fact that $\sin(3\omega x) \geq 0$ we have

$$3\sin(\omega x) - 4\sin^3(\omega x) = \sin(3\omega x) \geq 0$$

Since $\sin(\omega x) < 0$ we can divide both sides of this inequality by $\sin(\omega x)$ to conclude $3 - 4\sin^2(\omega x) \leq 0$, or $\sin^2(\omega x) \geq \frac{3}{4}$.

- (b) For any $m > 0$, consider the set of points (x_1, \dots, x_m) with arbitrary labels $(y_1, \dots, y_m) \in \{-1, +1\}^m$. Now, choose the parameter $\omega = \pi(1 + \sum_{i=1}^m 2^i y'_i)$ where $y'_i = \frac{1 - y_i}{2}$. We show that this single parameter will always correctly classify the entire sample for any $m > 0$ and choice of labels. For any $j \in [m]$ we have,

$$\omega x_j = \omega 2^{-j} = \pi(2^{-j} + \sum_{i=1}^m 2^{i-j} y'_i) = \pi(2^{-j} + (\sum_{i=1}^{j-1} 2^{i-j} y'_i) + y'_j + (\sum_{i=j+1}^m 2^i y'_i)).$$

The last term can be dropped from the sum, since it contributes only multiples of 2π . Since $y'_i \in \{0, 1\}$ the remaining term $\pi(2^{-j} + (\sum_{i=1}^{j-1} 2^{i-j} y'_i) + y'_j) = \pi(\sum_{i=1}^{j-1} 2^{-i} y'_i + 2^{-j} + y'_j)$ can be upper and lower bounded as follows:

$$\begin{aligned} \pi(\sum_{i=1}^{j-1} 2^{-i} y'_i + 2^{-j} + y'_j) &\leq \pi(\sum_{i=1}^j 2^{-i} + y'_j) < \pi(1 + y'_j), \\ \pi(\sum_{i=1}^{j-1} 2^{-i} y'_i + 2^{-j} + y'_j) &> \pi y'_j. \end{aligned}$$

Thus, if $y_j = 1$ we have $y'_j = 0$ and $0 < \omega x_j < \pi$, which implies $\text{sgn}(\omega x_j) = 1$. Similarly, for $y_j = -1$ we have $y'_j = 1$ and $\pi < \omega x_j < 2\pi$, which implies $\text{sgn}(\omega x_j) = -1$.

3.21 VC-dimension of union of halfspaces

3.22 VC-dimension of intersection of halfspaces

Let $m \geq 0$. Note the general fact that for any concept class $\mathcal{C} = \{c_1 \cap c_2 : c_1 \in \mathcal{C}_1, c_2 \in \mathcal{C}_2\}$,

$$\Pi_{\mathcal{C}}(m) \leq \Pi_{\mathcal{C}_1}(m) \Pi_{\mathcal{C}_2}(m). \quad (\text{E.53})$$

Indeed, fix a set \mathcal{X} of m points. Let $\mathcal{Y}_1, \dots, \mathcal{Y}_k$ be the traces of \mathcal{C}_1 on \mathcal{X} . By definition of $\Pi_{\mathcal{C}_1}(\mathcal{X})$, $k \leq \Pi_{\mathcal{C}_1}(\mathcal{X}) \leq \Pi_{\mathcal{C}_1}(m)$. By definition of $\Pi_{\mathcal{C}_2}(\mathcal{Y}_i)$, the traces of \mathcal{C}_2 on a subset \mathcal{Y}_i are at most $\Pi_{\mathcal{C}_2}(\mathcal{Y}_i) \leq \Pi_{\mathcal{C}_2}(m)$. Thus, the traces of \mathcal{C} on \mathcal{X} are at most

$$k \Pi_{\mathcal{C}_2}(\mathcal{Y}_i) \leq \Pi_{\mathcal{C}_1}(m) \Pi_{\mathcal{C}_2}(m). \quad (\text{E.54})$$

For the particular case of \mathcal{C}_k , using Sauer's lemma, this implies that

$$\Pi_{\mathcal{C}_k}(m) \leq (\Pi_{\mathcal{C}_1}(m))^k \leq \left(\frac{em}{n+1}\right)^{k(n+1)}. \quad (\text{E.55})$$

If $(em/(n+1))^{k(n+1)} < 2^m$, then the VC-dimension of \mathcal{C}_k is less than m . If the VC-dimension of \mathcal{C}_k is m , then $\Pi_{\mathcal{C}_k}(m) = 2^m \leq (em/(n+1))^{k(n+1)}$. These inequalities give an upper bound and a lower bound on $\text{VCdim}(\mathcal{C}_k)$. As an example, using the inequality: $\forall x \in \mathbb{N} - \{3\}, \log_2(x) \leq x/2$, one can verify that:

$$\text{VCdim}(\mathcal{C}_k) \leq 2(n+1)k \log(3k). \quad (\text{E.56})$$

3.23 VC-dimension of intersection concepts

- (a) Fix a set \mathcal{X} of m points. Let $\mathcal{Y}_1, \dots, \mathcal{Y}_k$ be the set of intersections of the concepts of \mathcal{C}_1 with \mathcal{X} . By definition of $\Pi_{\mathcal{C}_1}(\mathcal{X})$, $k \leq \Pi_{\mathcal{C}_1}(\mathcal{X}) \leq \Pi_{\mathcal{C}_1}(m)$. By definition of $\Pi_{\mathcal{C}_2}(\mathcal{Y}_i)$, the intersection of the concepts of \mathcal{C}_2 with \mathcal{Y}_i are at most $\Pi_{\mathcal{C}_2}(\mathcal{Y}_i) \leq \Pi_{\mathcal{C}_2}(m)$. Thus, the

number of sets intersections of concepts of \mathcal{C} with \mathcal{X} is at most

$$k\Pi_{\mathcal{C}_2}(\mathcal{Y}_i) \leq \Pi_{\mathcal{C}_1}(m) \Pi_{\mathcal{C}_2}(m). \quad (\text{E.57})$$

- (b) In view of the result proved in the previous question, $\Pi_{\mathcal{C}_s}(m) \leq (\Pi_{\mathcal{C}_1}(m))^s$. By Sauer's lemma, this implies

$$\Pi_{\mathcal{C}_s}(m) \leq \left(\frac{em}{d}\right)^{sd}. \quad (\text{E.58})$$

If $\left(\frac{em}{d}\right)^{sd} < 2^m$, then the VC-dimension of \mathcal{C}_s is less than m . Thus, it suffices to show this inequality holds with $m = 2ds \log_2(3s)$. Plugging in that value for m and taking the \log_2 yield:

$$ds \log_2(2es \log_2(3s)) < 2ds \log_2(3s) \quad (\text{E.59})$$

$$\Leftrightarrow \log_2(2es \log_2(3s)) < 2 \log_2(3s) = \log_2(9s^2) \quad (\text{E.60})$$

$$\Leftrightarrow 2es \log_2(3s) < 9s^2 \quad (\text{E.61})$$

$$\Leftrightarrow \log_2(3s) < \frac{9s}{2e}. \quad (\text{E.62})$$

This last inequality holds for $s = 2$: $\log_2(6) \approx 2.6 < 9/(2e) \approx 3.3$. Since the functions corresponding to the left-hand-side grows more slowly than the one corresponding to the right-hand-side (compare derivatives for example), this implies that the inequality holds for all $s \geq 2$.

3.24 VC-dimension of union of concepts

- (a) When $\mathcal{C} = \mathcal{A} \cup \mathcal{B}$, $\Pi_{\mathcal{C}}(\mathcal{X}) \leq \Pi_{\mathcal{A}}(\mathcal{X}) + \Pi_{\mathcal{B}}(\mathcal{X})$ for any set \mathcal{X} , since dichotomies in $\Pi_{\mathcal{C}}(\mathcal{X})$ can be generated by \mathcal{A} or by \mathcal{B} . Thus, for all m , $\Pi_{\mathcal{C}}(m) \leq \Pi_{\mathcal{A}}(m) + \Pi_{\mathcal{B}}(m)$.
- (b) For $m \geq d_{\mathcal{A}} + d_{\mathcal{B}} + 2$, by Sauer's lemma,

$$\begin{aligned} \Pi_{\mathcal{C}}(m) &\leq \sum_{i=0}^{d_{\mathcal{A}}} \binom{m}{i} + \sum_{i=0}^{d_{\mathcal{B}}} \binom{m}{i} = \sum_{i=0}^{d_{\mathcal{A}}} \binom{m}{i} + \sum_{i=0}^{d_{\mathcal{B}}} \binom{m-i}{i} \\ &= \sum_{i=0}^{d_{\mathcal{A}}} \binom{m}{i} + \sum_{i=m-d_{\mathcal{B}}}^m \binom{m}{i} \end{aligned} \quad (\text{E.63})$$

$$\leq \sum_{i=0}^{d_{\mathcal{A}}} \binom{m}{i} + \sum_{i=d_{\mathcal{A}}+2}^{d_{\mathcal{B}}} \binom{m}{i} \quad (\text{E.64})$$

$$< \sum_{i=0}^m \binom{m}{i} = 2^m. \quad (\text{E.65})$$

Thus, the VC-dimension of \mathcal{C} is strictly less than $d_{\mathcal{A}} + d_{\mathcal{B}} + 2$:

$$\text{VCdim}(\mathcal{C}) \leq d_{\mathcal{A}} + d_{\mathcal{B}} + 1. \quad (\text{E.66})$$

3.25 VC-dimension of symmetric difference of concepts

Fix a set \mathcal{S} . We can show that the number of classifications of \mathcal{S} using \mathcal{H} is the same as when using $\mathcal{H}\Delta A$. The set of classifications obtained using \mathcal{H} can be identified with $\{\mathcal{S} \cap h : h \in \mathcal{H}\}$ and the set of classifications using $\mathcal{H}\Delta A$ can be identified with $\{\mathcal{S} \cap (h\Delta A) : h \in \mathcal{H}\}$. Observe that for any $h \in \mathcal{H}$,

$$\mathcal{S} \cap (h\Delta A) = (\mathcal{S} \cap h) \Delta (\mathcal{S} \cap A). \quad (\text{E.67})$$

Figure E.7 helps illustrate this equality in a special case. Now, in view of this inequality, if $\mathcal{S} \cap (h\Delta A) = \mathcal{S} \cap (h'\Delta A)$ for $h, h' \in \mathcal{H}$, then

$$(\mathcal{S} \cap h) \Delta \mathcal{B} = (\mathcal{S} \cap h') \Delta \mathcal{B}, \quad (\text{E.68})$$

with $\mathcal{B} = \mathcal{S} \cap A$. Since two sets that have the same symmetric differences with respect to a set \mathcal{B} must be equal, this implies

$$\mathcal{S} \cap h = \mathcal{S} \cap h'. \quad (\text{E.69})$$

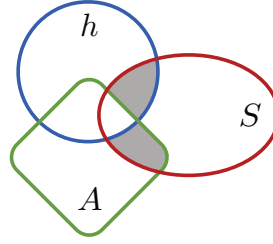
**Figure E.7**

Illustration of $(h\Delta A) \cap S = (h \cap S) \Delta (A \cap S)$ shown in gray.

This shows that ϕ defined by

$$\phi: S \cap \mathcal{H} \rightarrow S \cap (\mathcal{H}\Delta A)$$

$$S \cap h \mapsto S \cap (h\Delta A)$$

is a bijection, and thus that the sets $S \cap \mathcal{H}$ and $S \cap (\mathcal{H}\Delta A)$ have the same cardinality.

3.26 Symmetric functions

- (a) For $i = 0, \dots, n$, let $x_i \in \{0, 1\}^n$ be defined by $x_i = (\underbrace{1, \dots, 1}_{i \text{ 1's}}, 0, \dots, 0)$. Then, $\{x_0, \dots, x_n\}$

can be shattered by \mathcal{C} . Indeed, let $y_0, \dots, y_n \in \{0, 1\}$ be an arbitrary labeling of these points. Then, the function h defined by:

$$h(x) = y_i \quad (\text{E.70})$$

for all x with i 1's is symmetric and $h(x_i) = y_i$. Thus, $\text{VCdim}(\mathcal{C}) \geq n + 1$. Conversely, a set of $n + 2$ points cannot be shattered by \mathcal{C} , since at least two points would then have the same number of 1's and will not be distinguishable by \mathcal{C} . Thus,

$$\text{VCdim}(\mathcal{C}) = n + 1. \quad (\text{E.71})$$

- (b) Thus, in view of the theorems presented in class, a lower bound on the number of training examples needed to learn symmetric functions with accuracy $1 - \epsilon$ and confidence $1 - \delta$ is

$$\Omega\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{n}{\epsilon}\right), \quad (\text{E.72})$$

and an upper bound is:

$$O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{n}{\epsilon} \log \frac{1}{\epsilon}\right), \quad (\text{E.73})$$

which is only within a factor $\frac{1}{\epsilon}$ of the lower bound.

- (c) For a training data $(z_0, t_0), \dots, (z_m, t_m) \in \{0, 1\}^n \times \{0, 1\}$ define h as the symmetric function such that $h(z_i) = t_i$ for all $i = 0, \dots, m$.

3.27 VC-dimension of neural networks

- (a) Let $\Pi_u(m)$ denote the growth function at a node u in the intermediate layer. For a fixed set of values at the intermediate layer, using the concept class \mathcal{C} the output node can generate at most $\Pi_{\mathcal{C}}(m)$ distinct labelings. There are $\prod_u \Pi_u(m)$ possible sets of values at the intermediate layer since, by definition, for a sample of size m , at most $\Pi_u(m)$ distinct values are possible at each u . Thus, at most $\Pi_{\mathcal{C}}(m) \times \prod_u \Pi_u(m)$ labelings can be generated by the neural network and $\Pi_{\mathcal{H}}(m) \leq \Pi_{\mathcal{C}}(m) \prod_u \Pi_u(m)$.
- (b) For any intermediate node u , $\Pi_u(m) = \Pi_{\mathcal{C}}(m)$. Thus, for $\tilde{k} = k + 1$, $\Pi_{\mathcal{H}}(m) \leq \Pi_{\mathcal{C}}(m)^{\tilde{k}}$. By Sauer's lemma, $\Pi_{\mathcal{C}}(m) \leq \left(\frac{em}{d}\right)^d$, thus $\Pi_{\mathcal{H}}(m) \leq \left(\frac{em}{d}\right)^{d\tilde{k}}$. Let $m = 2\tilde{k}d \log_2(ek)$. In

view of the inequality given by the hint and $e\tilde{k} > 4$, this implies $m > d\tilde{k} \log_2 \left(\frac{em}{d} \right)$, that is $2^m > \left(\frac{em}{d} \right)^{d\tilde{k}}$. Thus, the VC-dimension of \mathcal{H} is less than

$$2\tilde{k}d \log_2(e\tilde{k}) = 2(k+1)d \log_2(e(k+1)).$$

- (c) For threshold functions, the VC-dimension of \mathcal{C} is r , thus, the VC-dimension of \mathcal{H} is upper bounded by

$$2(k+1)r \log_2(e(k+1)).$$

3.28 VC-dimension of convex combinations

Following the hint, we can think of this family of functions as a one hidden layer neural network, where the hidden layer is represented by the functions $h_i \in \mathcal{H}$, and the top layer is a threshold function characterized by $(\alpha_1, \dots, \alpha_T)$. Denote this class of threshold functions by Δ_T . From the solution of exercise 3.27(a) we can bound the growth function of \mathcal{F}_T by:

$$\Pi_{\mathcal{F}_T}(m) \leq \Pi_{\Delta_T}(m) (\Pi_{\mathcal{H}}(m))^T.$$

From the solution to exercise 3.27(c), the VC dimension of Δ_T is at most T , and we may further denote the VC dimension of \mathcal{H} by d . Applying Sauer's lemma to the growth function yields:

$$\Pi_{\Delta_T}(m) \leq \left(\frac{em}{T} \right)^T, \quad \Pi_{\mathcal{H}}(m) \leq \left(\frac{em}{d} \right)^d.$$

Thus, we have that

$$\Pi_{\mathcal{F}_T}(m) \leq \left(\frac{em}{T} \right)^T \left(\frac{em}{d} \right)^{Td}.$$

Finally, we may apply the hint in exercise 3.27(b) with $m = \max\{4T \log_2(2e), 2Td \log_2(eT)\}$ to see that

$$\left(\frac{em}{T} \right)^T \left(\frac{em}{d} \right)^{Td} < 2^{4T \log_2(2e) + 2Td \log_2(eT)},$$

so that the VC Dimension of \mathcal{F}_T is bounded by:

$$2T(2 \log_2(2e) + d \log_2(eT)).$$

Note that a coarser but relatively simpler bound would be to write:

$$\left(\frac{em}{T} \right)^T \left(\frac{em}{d} \right)^{Td} < (em)^{T(d+1)},$$

and to apply the hint in exercise 3.27(b) with $m = 2T(d+1) \log_2(eT(d+1))$. Notice that this is actually asymptotically optimal in T and d up to log terms.

3.29 Infinite VC-dimension

- (a) Theorem 3.20 shows that there exists a distribution that can force an error of $\Omega(\frac{d}{m})$. Thus, for an infinite VC-dimension, this lower bound requires an infinite number of points to achieve a bounded error and thus implies that PAC-learning is not possible.
- (b) Here is a description of the algorithm. Let M be the maximum value observed after drawing m points and let p be the probability that a point greater than M be drawn. The probability that all points drawn be smaller than or equal to M is

$$(1-p)^m \leq e^{-pm}. \quad (\text{E.74})$$

Setting $\delta/2$ to match the upper bound, yields $\delta/2 = e^{-pm}$, that is

$$p = \frac{1}{m} \log \frac{2}{\delta}. \quad (\text{E.75})$$

To bound p by $\epsilon/2$, we can impose the following

$$\frac{1}{m} \log \frac{2}{\delta} \leq \frac{\epsilon}{2}. \quad (\text{E.76})$$

Thus, with confidence at least $1 - \delta/2$, the probability that a point greater than M be drawn is at most $\epsilon/2$ if L draws $m \geq \frac{2}{\epsilon} \log \frac{2}{\delta}$ points.

In the second stage, the problem is reduced to a finite VC-dimension M . Since PAC-learning with $(\epsilon/2, \delta/2)$ is possible for a finite dimension, this guarantees the (ϵ, δ) -PAC-learning of the full algorithm.

3.30 VC-dimension generalization bound – realizable case

- (a) Let
- $h_0 \in \mathcal{H}_S$
- , then we have the following set of inequalities:

$$\begin{aligned}
\mathbb{P} \left[\sup_{h \in \mathcal{H}_S} |\widehat{R}_S(h) - \widehat{R}_{S'}(h)| > \frac{\epsilon}{2} \right] &\geq \mathbb{P} \left[|\widehat{R}_S(h_0) - \widehat{R}_{S'}(h_0)| > \frac{\epsilon}{2} \right] \\
&= \mathbb{P} \left[\widehat{R}_{S'}(h_0) > \frac{\epsilon}{2} \right] \\
&\geq \mathbb{P} \left[\widehat{R}(h_0) > \frac{\epsilon}{2} \mid R(h_0) > \epsilon \right] \mathbb{P}[R(h_0) > \epsilon] \\
&> \mathbb{P} \left[B(m, \epsilon) > \frac{m\epsilon}{2} \right] \mathbb{P}[R(h_0) > \epsilon].
\end{aligned}$$

The second inequality follows from the fact that for any two random events A and B , $\mathbb{P}[A] \geq \mathbb{P}[A \wedge B] = \mathbb{P}[A|B] \mathbb{P}[B]$. The final equality follows, since the event we are concerned with is the probability that we get at least a fraction of $\epsilon/2$ errors on a sample of size m when the true probability of error is at least ϵ . In the case the true error rate equals ϵ , this exactly describes the probability that $B(m, \epsilon) \geq m\epsilon/2$.

- (b) We apply Chebyshev's inequality to the binomial random variable
- $B(m, \epsilon)$
- , which has mean
- $m\epsilon$
- variance
- $m\epsilon(1 - \epsilon)$
- .

$$\begin{aligned}
\mathbb{P} \left[B(m, \epsilon) \leq \frac{m\epsilon}{2} \right] &= \mathbb{P} \left[m\epsilon - B(m, \epsilon) \geq \frac{m\epsilon}{2} \right] \\
&\leq \frac{m\epsilon(1 - \epsilon)}{(m\epsilon/2)^2} = \frac{4(1 - \epsilon)}{m\epsilon} \leq \frac{4}{m\epsilon} \leq \frac{1}{2}
\end{aligned}$$

where the last inequality uses the assumption $m\epsilon \geq 8$. Thus, this shows that $\mathbb{P}[B(m, \epsilon) > m\epsilon/2] > 1 - 1/2 = 1/2$. Plugging the bound into part (a) completes the question.

- (c) There are
- $\binom{2m}{l}$
- total ways to distribute the
- l
- error over the sample
- T
- and
- $\binom{m}{l}$
- way to distribute the error such that the only hit
- S'
- . Thus, the probability of all error falling only into
- S'
- is bounded as

$$\frac{\binom{m}{l}}{\binom{2m}{l}} = \prod_{i=0}^{l-1} \frac{m-i}{2m-i} \leq \prod_{i=0}^{l-1} \frac{m-i}{2m-2i} \leq \frac{1}{2^l}.$$

- (d) This follows from

$$\begin{aligned}
&\mathbb{P}_{\substack{T \sim \mathcal{D}^{2m}; \\ T \rightarrow (S, S')}} \left[\widehat{R}_S(h) = 0 \wedge \widehat{R}_{S'}(h) > \frac{\epsilon}{2} \right] \\
&= \mathbb{P}_{\substack{T \sim \mathcal{D}^{2m}; \\ T \rightarrow (S, S')}} \left[\widehat{R}_S(h) = 0 \wedge \widehat{R}_{S'}(h) > \frac{\epsilon}{2} \wedge \widehat{R}_T(h) > \frac{\epsilon}{2} \right] \\
&= \mathbb{P}_{\substack{T \sim \mathcal{D}^{2m}; \\ T \rightarrow (S, S')}} \left[\widehat{R}_S(h) = 0 \wedge \widehat{R}_{S'}(h) > \frac{\epsilon}{2} \mid \widehat{R}_T(h) > \frac{\epsilon}{2} \right] \mathbb{P}[\widehat{R}_T(h) > \frac{\epsilon}{2}] \\
&\leq \mathbb{P}_{\substack{T \sim \mathcal{D}^{2m}; \\ T \rightarrow (S, S')}} \left[\widehat{R}_S(h) = 0 \wedge \widehat{R}_{S'}(h) > \frac{\epsilon}{2} \mid \widehat{R}_T(h) > \frac{\epsilon}{2} \right] \\
&\leq 2^{-l} \leq 2^{-\frac{m\epsilon}{2}}.
\end{aligned}$$

- (e) Using the definition of the growth function, we can provide the following union bound that is then in turn bounded using corollary 3.18:

$$\mathbb{P}_{\substack{T \sim \mathcal{D}^{2m}; \\ T \rightarrow (S, S')}} \left[\exists h \in \mathcal{H}: \widehat{R}_S(h) = 0 \wedge \widehat{R}_{S'}(h) > \frac{\epsilon}{2} \right] \leq \Pi_{\mathcal{H}}(2m) 2^{-\frac{m\epsilon}{2}} \leq \left(\frac{2em}{d} \right)^d 2^{-m\epsilon/2}.$$

Combining part (a) through (e), we finally have,

$$\mathbb{P}[R(h) > \epsilon] \leq 2 \left(\frac{2em}{d} \right)^d 2^{-m\epsilon/2}. \quad (\text{E.77})$$

Setting the right-hand side equal to δ and solving for ϵ show that with probability at least $1 - \delta$

$$\epsilon \leq \frac{1}{m} \left(d \log \frac{2em}{d} + \log \frac{1}{\delta} + \log 2 \right) \frac{2}{\log 2}.$$

3.31 Covering number generalization bound.

(a) First split the term into two separate terms:

$$\begin{aligned} |L_S(h_1) - L_S(h_2)| &\leq |R(h_1) - R(h_2)| + |\widehat{R}_S(h_1) - \widehat{R}_S(h_2)| \\ &= \left| \mathbb{E}_{x,y} [(h_1(x) - y)^2 - (h_2(x) - y)^2] \right| + \left| \frac{1}{m} \sum_{i=1}^m (h_1(x_i) - y_i)^2 - (h_2(x_i) - y_i)^2 \right|. \end{aligned}$$

Then, expanding the term

$$\begin{aligned} (h_1(x) - y)^2 - (h_2(x) - y)^2 &= (h_1(x) - h_2(x))(h_1 + h_2 - 2y) \\ &= (h_1(x) - h_2(x))((h_1 - y) + (h_2 - y)) \leq \|h_1 - h_2\|_\infty 2M, \end{aligned}$$

allows us to bound both the empirical and true error, resulting in a total bound of $4M\|h_1 - h_2\|_\infty$.

(b) This follows by splitting the event into the union of several smaller events and then using the sum rule,

$$\begin{aligned} \mathbb{P}_S \left[\sup_{h \in \mathcal{H}} |L_S(h)| \geq \epsilon \right] \\ = \mathbb{P}_S \left[\bigvee_{i=1}^k \sup_{h \in B_i} |L_S(h)| \geq \epsilon \right] \leq \sum_{i=1}^k \mathbb{P}_S \left[\sup_{h \in B_i} |L_S(h)| \geq \epsilon \right]. \end{aligned}$$

(c) For any i let h_i be the center of ball B_i with radius $\frac{\epsilon}{8M}$. Note that for any $h \in \mathcal{H}$ we have $|L_S(h) - L_S(h_i)| \leq 4M\|h - h_i\|_\infty \leq \epsilon/2$. Thus, if for any $h \in B_i$ we have $|L_S(h)| \geq \epsilon$ it must be the case that $|L_S(h_i)| \geq \epsilon/2$, which shows the inequality.

To complete the bound, we use Hoeffding's inequality applied to the random variables $(h(x_i) - y_i)^2/m \leq M^2/m$, which guarantees

$$\mathbb{P}_S \left[|L_S(h_i)| \geq \frac{\epsilon}{2} \right] \leq 2 \exp \left(\frac{-m\epsilon^2}{2M^4} \right).$$

Chapter 5

5.1 Soft margin hyperplanes

(a) The corresponding dual problem is:

$$\max_{\alpha, \beta} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) - \sum_{i=1}^m \frac{(\alpha_i + \beta_i)^{k/(k-1)}}{(kC)^{1/(k-1)}} \left(1 - \frac{1}{k}\right)$$

subject to:

$$\sum_{i=1}^m \alpha_i y_i = 0 \quad \alpha \geq 0 \quad \beta \geq 0.$$

(b) Here we see that the objective function is more complex requiring an optimization over both α and β and there is the additional constraint $\beta \geq 0$.

For $k = 2$ the additional term of interest is $-\sum_{i=1}^m (\alpha_i + \beta_i)^2$ (to see this, note that the Hessian is negative semidefinite), which is jointly concave in α_i and β_i , which allows for convex optimization techniques.

5.2 Tighter Rademacher bound

Proceed as in the proof of theorem 5.9, but choose $\rho_k = 1/\gamma^k$. For any $\rho \in (0, 1)$, there exists $k \geq 1$ such that $\rho \in (\rho_k, \rho_{k-1}]$, with $\rho_0 = 1$. For that k , $\rho \leq \rho_{k-1} = \gamma\rho_k$, thus $1/\rho_k \leq \gamma/\rho$ and $\log k = \sqrt{\log \log_\gamma(1/\rho_k)} \leq \sqrt{\log \log_2(\gamma/\rho)}$. Furthermore, for any $h \in \mathcal{H}$, $\widehat{R}_{S, \rho_k}(h) \leq \widehat{R}_{S, \rho}(h)$. Thus,

$$\mathbb{P} \left[\exists k: R(h) - \widehat{R}_{S, \rho}(h) > \frac{2\gamma}{\rho} \mathfrak{R}_m(\mathcal{H}) + \sqrt{\frac{\log \log_\gamma(\gamma/\rho)}{m}} + \epsilon \right] \leq 2 \exp(-2m\epsilon^2),$$

which proves the statement.

5.3 Importance weighted SVM

The modified primal optimization problem can be written as

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i p_i \\ & \text{subject to} && y_i [w \cdot x_i + b] \geq 1 - \xi_i. \end{aligned}$$

The Lagrangian holding for all $w, b, \alpha_i \geq 0, \beta_i \geq 0$ is then

$$\begin{aligned} L(w, b, \alpha) &= \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i p_i \\ &\quad - \sum_{i=1}^m \alpha_i [y_i (w \cdot x_i + b) - 1 + \xi_i] - \sum_{i=1}^m \beta_i \xi_i. \end{aligned} \tag{E.78}$$

Then $\frac{\partial L}{\partial w}$ and $\frac{\partial L}{\partial b}$ are the same as for the regular non-separable SVM optimization problem. We also have $\frac{\partial L}{\partial \xi_i} = Cp_i - \alpha_i - \beta_i$. Thus, to satisfy the KKT conditions we have for all $i \in [m]$,

$$w = \sum_{i=1}^m \alpha_i y_i x_i \tag{E.79}$$

$$\sum_{i=1}^m \alpha_i y_i = 0 \tag{E.80}$$

$$\alpha_i + \beta_i = Cp_i \tag{E.81}$$

$$\alpha_i [y_i (w \cdot x_i + b) - 1 + \xi_i] = 0 \tag{E.82}$$

$$\beta_i \xi_i = 0. \tag{E.83}$$

Plugging equation E.79 into equation E.78, we get

$$\begin{aligned} L &= \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i x_i \right\|^2 + C \sum_{i=1}^m \xi_i p_i - \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \\ &\quad - \sum_{i=1}^m \alpha_i y_i b + \sum_{i=1}^m \alpha_i - \sum_{i=1}^m \alpha_i \xi_i - \sum_{i=1}^m \beta_i \xi_i. \end{aligned} \tag{E.84}$$

Using equation E.81, we can simplify:

$$L = \sum_{i=1}^m \alpha_i - \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i x_i \right\|^2,$$

meaning that the objective function is the same as in the regular SVM problem. The difference is in the constraints on the optimization. Recall that our dual form holds for $\beta_i \geq 0$. Using again equation E.81, our optimization problem is to maximize L subject to the constraints:

$$\forall i \in [m], 0 \leq \alpha_i \leq Cp_i \wedge \sum_{i=1}^m \alpha_i y_i = 0.$$

5.4 Sequential Minimal Optimization (SMO).

- (a) Starting from equation (5.33) and removing all terms that are constant with respect to α_1 and α_2 yields the desired result.
- (b) Substituting into Ψ_1 , we have:

$$\Psi_2 = \gamma - s\alpha_2 + \alpha_2 - \frac{1}{2}K_{11}(\gamma - s\alpha_2)^2 - \frac{1}{2}K_{22}\alpha_2^2 - sK_{12}(\gamma - s\alpha_2)\alpha_2 - y_1(\gamma - s\alpha_2)v_1 - y_2\alpha_2v_2.$$

We take the derivative to find the equation for the stationary point as follows:

$$\frac{d\Psi_2}{d\alpha_2} = -s + 1 + sK_{11}(\gamma - s\alpha_2) - K_{22}\alpha_2 - sK_{12}(\gamma - s\alpha_2) + s^2K_{12}\alpha_2 + y_1sv_1 - y_2v_2 = 0.$$

Noting that $s^2 = 1$ and rearranging terms yields the statement of interest.

- (c) By definition of $f(\cdot)$ we see that $v_1 = f(\mathbf{x}_1) - y_1\alpha_1^*K_{11} - y_2\alpha_2^*K_{12}$ and similarly, $v_2 = f(\mathbf{x}_2) - y_1\alpha_1^*K_{12} - y_2\alpha_2^*K_{22}$. Using these equations along with the identities $\alpha_1^* = \gamma - s\alpha_2^*$ and $y_1 = sy_2$ we have

$$\begin{aligned} v_1 - v_2 &= f(\mathbf{x}_1) - f(\mathbf{x}_2) - y_1\alpha_1^*(K_{11} - K_{12}) - y_2\alpha_2^*(K_{12} - K_{22}) \\ &= f(\mathbf{x}_1) - f(\mathbf{x}_2) - y_1(\gamma - s\alpha_2^*)(K_{11} - K_{12}) - y_2\alpha_2^*(K_{12} - K_{22}) \\ &= f(\mathbf{x}_1) - f(\mathbf{x}_2) - sy_2(\gamma - s\alpha_2^*)(K_{11} - K_{12}) - y_2\alpha_2^*(K_{12} - K_{22}) \\ &= f(\mathbf{x}_1) - f(\mathbf{x}_2) - sy_2\gamma(K_{11} - K_{12}) + y_2\alpha_2^*(K_{11} - K_{12}) - y_2\alpha_2^*(K_{12} - K_{22}) \\ &= f(\mathbf{x}_1) - f(\mathbf{x}_2) - sy_2\gamma(K_{11} - K_{12}) + y_2\alpha_2^*\eta. \end{aligned}$$

- (d) Combining the results from (b) and (c), we have

$$\begin{aligned} \eta\alpha_2 &= s(K_{11} - K_{12})\gamma + y_2[f(\mathbf{x}_1) - f(\mathbf{x}_2) - sy_2\gamma(K_{11} - K_{12}) + y_2\alpha_2^*\eta] - s + 1 \\ &= y_2[f(\mathbf{x}_1) - f(\mathbf{x}_2) + y_2\alpha_2^*\eta] - s + 1 \\ &= \alpha_2^*\eta + y_2[f(\mathbf{x}_1) - f(\mathbf{x}_2) - y_1 + y_2] \\ &= \alpha_2^*\eta + y_2[(y_2 - f(\mathbf{x}_2)) - (y_1 - f(\mathbf{x}_1))]. \end{aligned}$$

Dividing both sides by η yields the desired result.

- (e) Clipping is required to ensure that the new values of α_1 and α_2 satisfy the inequality constraints $0 \leq \alpha_1, \alpha_2 \leq C$. The lower bound of 0 follows directly from this inequality constraint, as does the upper bound of C . Moreover, when $s = +1$, the lower bound of $\gamma - C$ ensures that α_2^{clip} is large enough such that $\alpha_2^{clip} + \alpha_1 = \gamma$ while respecting the constraint $\alpha_1 \leq C$. Similarly, the upper bound γ ensures that α_2^{clip} is small enough such that $\alpha_2^{clip} + \alpha_1 = \gamma$ while respecting the constraint $\alpha_1 \geq 0$.

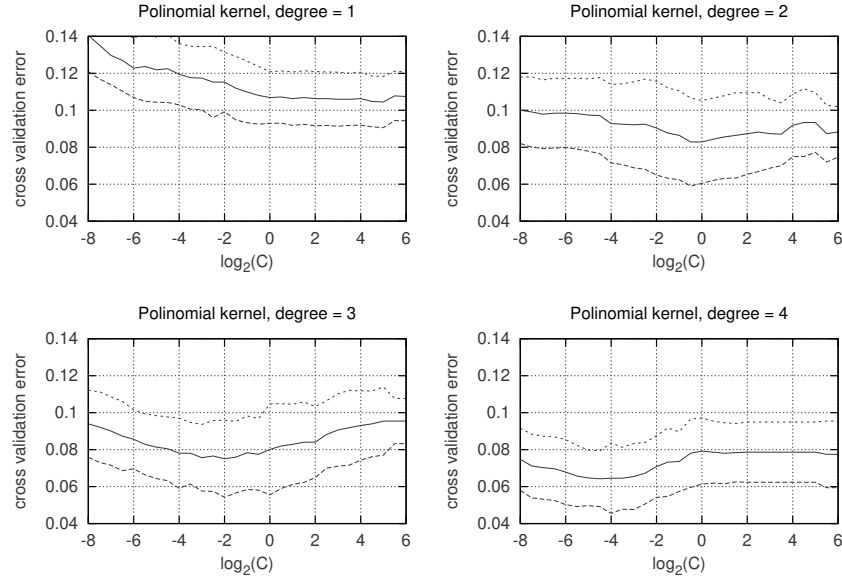
5.5 SVM hands-on

- (a) Download and install the `libsvm` software library from:

<http://www.csie.ntu.edu.tw/~cjlin/libsvm/>

- (b) Concatenate and scale data.

```
cat \
  satimage/satimage.scale.t \
  satimage/satimage.scale.val > data/train
libsvm-2.88/svm-scale \
  data/train > data/train.scaled
```

**Figure E.8**

Average cross-validation error plus or minus one standard deviation for different values of the trade-off constant C and of the degree of the polynomial kernel

- (c) Run 10-fold cross-validation, for different values of the degree d of the polynomial kernel and of the trade-off constant C . We test $d = 1, 2, 3, 4$ and $\log_2(C) = -8, -7.5, \dots, 5.5, 6$. We step the value of the trade-off constant logarithmically as suggested by:

<http://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>

Which gives the cross-validation error plots shown in figure E.8.

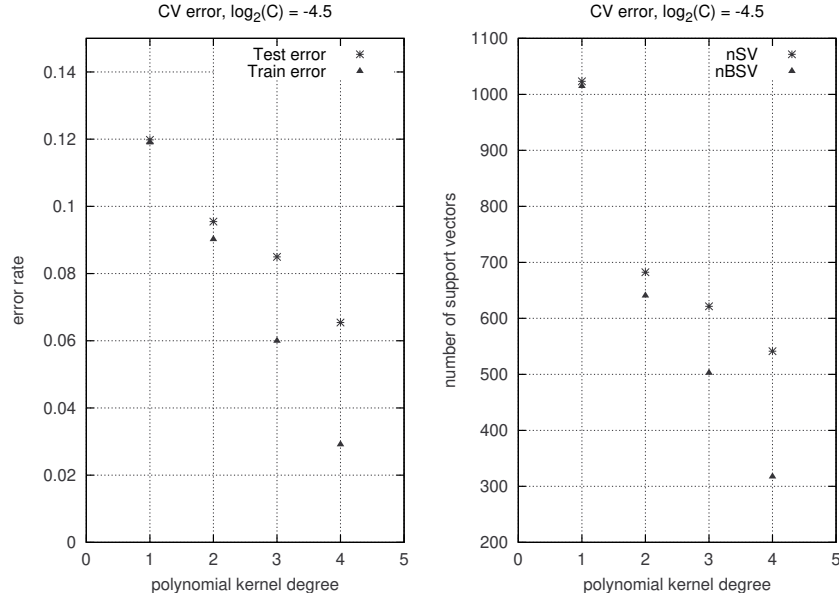
The best values of trade-off constant C are:

$d = 1 \setminus \text{colon } C^* = 2^{(+5.0)} = 32$	$\text{cv-err} = 10.4\% \pm 1.4\%$
$d = 2 \setminus \text{colon } C^* = 2^{(+0.0)} = 1$	$\text{cv-err} = 8.3\% \pm 2.2\%$
$d = 3 \setminus \text{colon } C^* = 2^{(-2.0)} = .25$	$\text{cv-err} = 7.5\% \pm 2.1\%$
$d = 4 \setminus \text{colon } C^* = 2^{(-4.5)} = .0442$	$\text{cv-err} = 6.4\% \pm 1.5\%$

The best C measured on the validation set is $C^* = 2^{-4.5}$, with degree $d^* = 4$, which gives an average error rate of $6.4\% \pm 1.5\%$.

- (d) The trade-off constant is fixed to $C^* = 2^{-4.5}$, and 10-fold cross-validation is run for degrees 1 through 4. In figure E.9 we plot the resulting cross-validation training and test errors and the average number of support vectors (nSV is the number of support vectors, nBSV is the number of bounded support vectors, i.e. whose dual variable is equal to the trade-off constant).
- (e) Support vectors always lie on the margin hyperplanes when their dual variable is smaller than C . This happens for all the support vectors (SV) that are not bounded (BSV). Our measurement gives the following averages:

$$d = 1 \setminus \text{colon } \text{nSV} - \text{nBSV} = 8.5$$

**Figure E.9**

Average cross-validation error rates and average number of support vectors (nSV) and of bounded support vectors (nBSV) as a function of the degree of the polynomial kernel.

$$\begin{aligned} d = 2 & \text{ \colon nSV - nBSV = 41.8} \\ d = 3 & \text{ \colon nSV - nBSV = 118.9} \\ d = 4 & \text{ \colon nSV - nBSV = 223.8} \end{aligned}$$

- (f) We can consider the more general problem of assigning a weight k_i to every sample that will multiply its misclassification penalty. The optimization problem becomes:

$$\begin{aligned} \text{minimize} \quad & \frac{1}{2} \|w\|^2 + \sum_{i=1}^m k_i \xi_i \\ \text{subject to} \quad & y_i(w x_i + b) \leq 1 - \xi_i, \quad \xi_i \geq 0, \quad \forall i \in 1, \dots, m. \end{aligned}$$

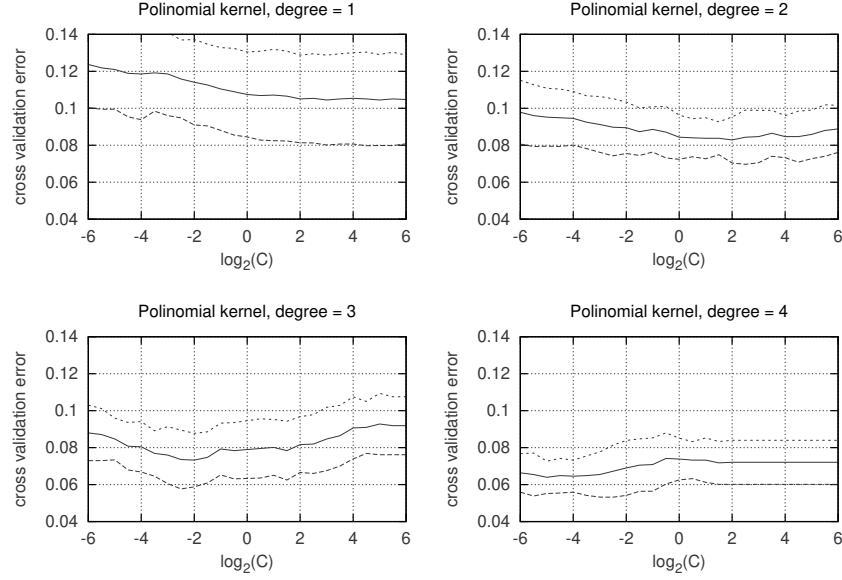
Moving to the dual exactly as shown in the chapter we obtain:

$$\begin{aligned} \text{maximize} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i \cdot x_j \\ \text{subject to} \quad & 0 \leq \alpha_i \leq C k_i, \quad \forall i \in 1, \dots, m. \end{aligned}$$

Setting $k_i = k$ for $y_i = -1$ and $k_i = 1$ for $y_i = 1$ gives the desired problem.

- (g) If k is an integer we can repeat every negative training point in the data k times. False positives will thus get penalized k times as much as false negatives.
- (h) Repeating the training for $k = 1, 2, 4, 8, 16$, we find the following results:

$$\begin{aligned} k = 1 & \text{ \colon } (d^*, C^*) = (4, 2^{(-4.5)}), \quad \text{cv-err} = 6.4\% \pm 1.5\% \\ k = 2 & \text{ \colon } (d^*, C^*) = (4, 2^{(-5.0)}), \quad \text{cv-err} = 6.4\% \pm 0.9\% \end{aligned}$$

**Figure E.10**

As in figure E.8, but false positives are penalized twice as much as false negatives.

$$\begin{aligned}
 k = 4 \quad & \text{colon} \quad (d^*, C^*) = (4, 2^{(-5.5)}), \quad \text{cv-err} = 6.1\% \pm 1.7\% \\
 k = 8 \quad & \text{colon} \quad (d^*, C^*) = (4, 2^{(-3.5)}), \quad \text{cv-err} = 6.2\% \pm 1.0\% \\
 k = 16 \quad & \text{colon} \quad (d^*, C^*) = (4, 2^{(-3.5)}), \quad \text{cv-err} = 6.2\% \pm 1.4\%
 \end{aligned}$$

Plots equivalent to the ones in figure E.8 are given in figure E.10, figure E.11, figure E.12, and figure E.13. We obtain the best average accuracy for $k = 4$.

5.6 Sparse SVM

(a) Let

$$\mathbf{x}'_i = (y_1 \mathbf{x}_i \cdot \mathbf{x}_1, \dots, y_m \mathbf{x}_i \cdot \mathbf{x}_m).$$

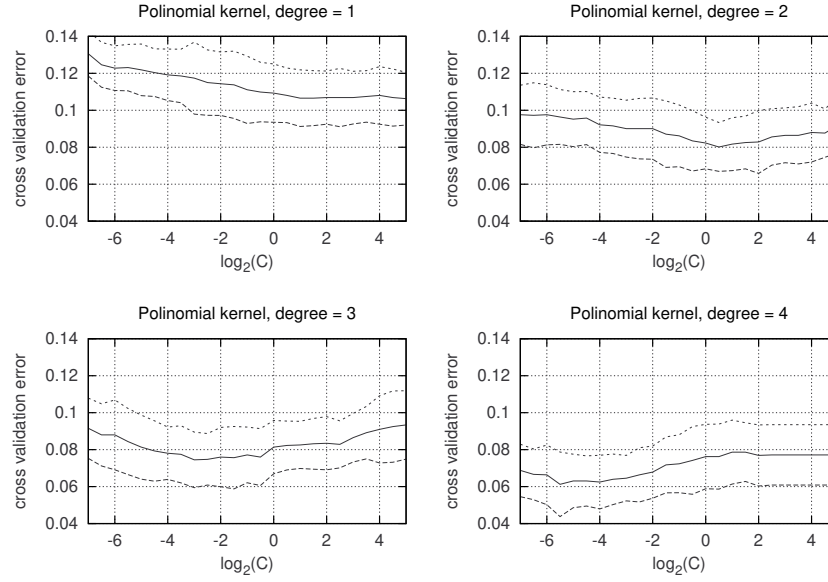
Then the optimization problem becomes

$$\begin{aligned}
 \min_{\alpha, b, \xi} \quad & \frac{1}{2} \|\alpha\|^2 + C \sum_{i=1}^m \xi_i \\
 \text{subject to} \quad & y_i (\alpha \cdot \mathbf{x}'_i + b) \geq 1 - \xi \\
 & \xi_i, \alpha_i \geq 0, i \in [m],
 \end{aligned}$$

which is the standard formulation of the primal SVM optimization problem on samples \mathbf{x}'_i , modulo the non-negativity constraints on α_i .

(b) The Lagrangian of (1) for all $\alpha_i \geq 0, \xi_i \geq 0, b, \alpha'_i \geq 0, \beta_i \geq 0, \gamma_i \geq 0, i \in [m]$ is

$$L = \frac{1}{2} \|\alpha\|^2 + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha'_i (y_i (\alpha \cdot \mathbf{x}'_i + b) - 1 + \xi_i) - \sum_{i=1}^m \beta_i \xi_i - \sum_{i=1}^m \gamma_i \alpha_i,$$

**Figure E.11**

As in figure E.8, but false positives are penalized four times as much as false negatives.

and the KKT conditions are

$$\nabla_{\alpha} L = 0 \quad \Leftrightarrow \quad \alpha = \sum_{i=1}^m \alpha'_i y_i \mathbf{x}'_i + \gamma$$

$$\nabla_b L = 0 \quad \Leftrightarrow \quad \sum_{i=1}^m \alpha'_i y_i = 0$$

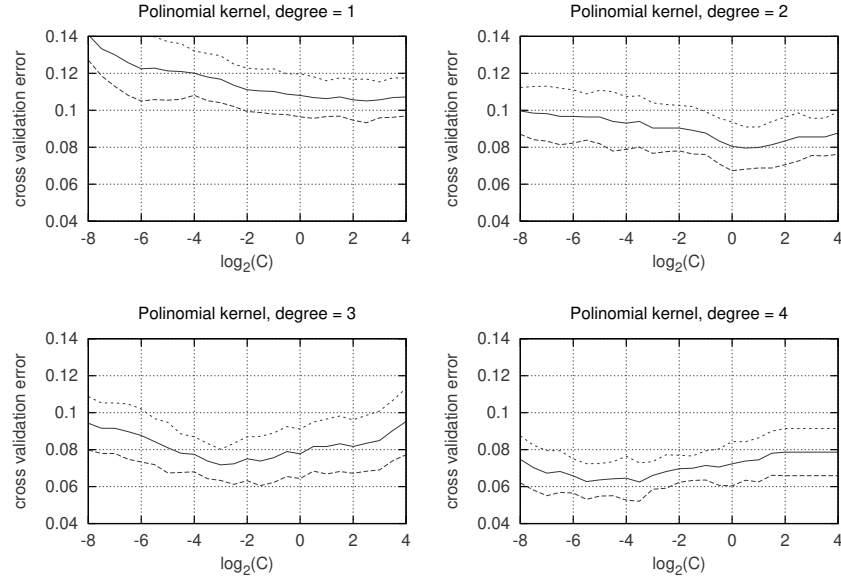
$$\nabla_{\xi_i} L = 0 \quad \Leftrightarrow \quad \alpha'_i + \beta_i = C$$

and

$$\alpha'_i (y_i (\alpha \cdot \mathbf{x}'_i + b) - 1 - \xi_i) = 0$$

$$\beta_i \xi_i = 0$$

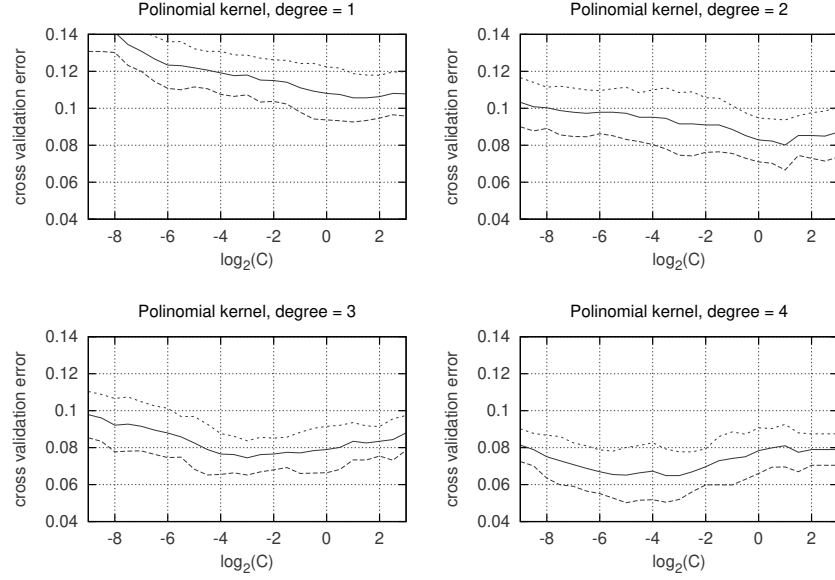
$$\gamma_i \alpha_i = 0.$$

**Figure E.12**

As in figure E.8, but false positives are penalized eight times as much as false negatives.

Using the KKT conditions on L we get

$$\begin{aligned}
 L &= \frac{1}{2} \left(\sum_{i=1}^m \alpha'_i y_i \mathbf{x}'_i + \gamma \right) \cdot \left(\sum_{j=1}^m \alpha'_j y_j \mathbf{x}'_j + \gamma \right) + C \sum_{i=1}^m \xi_i \\
 &\quad - \sum_{i=1}^m \alpha'_i \left(y_i \left(\left(\sum_{j=1}^m \alpha'_j y_j \mathbf{x}'_j + \gamma \right) \cdot \mathbf{x}'_i + b \right) - 1 + \xi_i \right) \\
 &\quad - \sum_{i=1}^m \beta_i \xi_i - \sum_{i=1}^m \gamma_i \alpha_i \\
 &= -\frac{1}{2} \sum_{i=1}^m \alpha'_i y_i \mathbf{x}'_i \cdot \left(\sum_{j=1}^m \alpha'_j y_j \mathbf{x}'_j + \gamma \right) + \frac{1}{2} \gamma \cdot \alpha \\
 &\quad + \sum_{i=1}^m C \xi_i - \alpha'_i (y_i b - 1 + \xi_i) \\
 &= \sum_{i=1}^m \alpha'_i - \frac{1}{2} \sum_{i,j=1}^m \alpha'_i \alpha'_j y_i y_j \mathbf{x}'_i{}^\top (\mathbf{x}'_j + \gamma) \\
 &\quad + \sum_{i=1}^m (C - \alpha'_i) \xi_i - \sum_{i=1}^m \alpha'_i y_i b.
 \end{aligned}$$

**Figure E.13**

As in figure E.8, but false positives are penalized sixteen times as much as false-negatives.

Thus the dual optimization problem is

$$\begin{aligned} \max_{\alpha', \gamma} \quad & \sum_{i=1}^m \alpha'_i - \frac{1}{2} \sum_{i,j=1}^m \alpha'_i \alpha'_j y_i y_j \mathbf{x}'_i \cdot (\mathbf{x}'_j + \gamma) \\ \text{subject to} \quad & \sum_{i=1}^m \alpha'_i y_i = 0 \\ & 0 \leq \alpha'_i \leq C, \gamma_i \geq 0, i \in [m]. \end{aligned}$$

5.7 VC-dimension of canonical hyperplanes

- (a) By definition of $\{\mathbf{x}_1, \dots, \mathbf{x}_d\}$, for all $\mathbf{y} = (y_1, \dots, y_d) \in \{-1, +1\}^d$, there exists \mathbf{w} such that,

$$\forall i \in [d], 1 \leq y_i (\mathbf{w} \cdot \mathbf{x}_i).$$

Summing up these inequalities yields

$$d \leq \mathbf{w} \cdot \sum_{i=1}^d y_i \mathbf{x}_i \leq \|\mathbf{w}\| \left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\| \leq \Lambda \left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\|.$$

- (b) Since this inequality holds for all $\mathbf{y} \in \{-1, +1\}^d$, it also holds on expectation over y_1, \dots, y_d drawn i.i.d. according to a uniform distribution over $\{-1, +1\}$. In view of the independence assumption, for $i \neq j$ we have $\mathbb{E}[y_i y_j] = \mathbb{E}[y_i] \mathbb{E}[y_j]$. Thus, since the

distribution is uniform, $\mathbb{E}[y_i y_j] = 0$ if $i \neq j$, $\mathbb{E}[y_i y_j] = 1$ otherwise. This gives

$$\begin{aligned}
 d &\leq \Lambda \mathbb{E}_{\mathbf{y}} \left[\left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\|^2 \right] && \text{(taking expectations)} \\
 &\leq \Lambda \left[\mathbb{E}_{\mathbf{y}} \left[\left\| \sum_{i=1}^d y_i \mathbf{x}_i \right\|^2 \right] \right]^{\frac{1}{2}} && \text{(Jensen's inequality)} \\
 &= \Lambda \left[\sum_{i,j=1}^d \mathbb{E}_{\mathbf{y}} [y_i y_j] (\mathbf{x}_i \cdot \mathbf{x}_j) \right]^{\frac{1}{2}} \\
 &= \Lambda \left[\sum_{i=1}^d (\mathbf{x}_i \cdot \mathbf{x}_i) \right]^{\frac{1}{2}}.
 \end{aligned}$$

Thus, $\sqrt{d} \leq \Lambda r$, which completes the proof.

(c) In view of the previous inequality, we can write $d \leq \Lambda [dr^2]^{\frac{1}{2}} = \Lambda r \sqrt{d}$.

Chapter 6

6.1 This follows directly the Cauchy-Schwarz inequality:

$$|\Phi(x)^\top \Phi(y)| \leq \|\Phi(x)\| \|\Phi(y)\|, \quad (\text{E.85})$$

where Φ is a feature mapping associated to K .

6.2 (a) Since $\cos(x - y) = \cos x \cos y + \sin x \sin y$, $K(x, y)$ can be written as the dot product of the vectors

$$\Phi(x) = \begin{bmatrix} \cos x \\ \sin x \end{bmatrix} \quad \text{and} \quad \Phi(y) = \begin{bmatrix} \cos y \\ \sin y \end{bmatrix}; \quad (\text{E.86})$$

thus it is PDS.

(b) This is a consequence of the fact that the kernel of the previous question is PDS since $\sum_{i,j} c_i c_j \cos(x_i^2 - x_j^2) = \sum_{i,j} c_i c_j \cos(x'_i - x'_j)$, with $x'_i = x_i^2$ for all i .

The solution for this question is similar to that of the previous question.

(c) Since the product and sum of PDS kernels is PDS, it suffices to show that $k: x \mapsto \cos(x^2 - y^2)$ is PDS over $\mathbb{R} \times \mathbb{R}$, which was proven in part (b).

(d) For any $x_1, \dots, x_m \in \mathbb{R}$, $c_1, \dots, c_m \in \mathbb{R}$, and $a \geq 0$, let $f(a)$ be defined by

$$f(a) = \sum_{i,j} c_i c_j K(x_i, x_j) a^{x_i + x_j}. \quad (\text{E.87})$$

Then, $f'(a) = \sum_{i,j} c_i c_j a^{x_i + x_j - 1} = \frac{1}{a} \|(c_i a^{x_i})_i\|^2 \geq 0$. Therefore f is monotonically increasing. Thus, $f(1) \geq f(0)$, that is $\sum_{i,j} c_i c_j K(x_i, x_j) \geq 0$.

(e) Rewriting the cosine in terms of the dot product, we have

$$K(\mathbf{x}, \mathbf{x}') = \cos \angle(\mathbf{x}, \mathbf{x}') = \frac{\mathbf{x} \cdot \mathbf{x}'}{\|\mathbf{x}\| \|\mathbf{x}'\|}.$$

Thus, the cosine kernel is just a scaling of the standard dot product, which is a PDS kernel. Hence, the cosine kernel is also PDS.

(f) For all $x, x' \in \mathbb{R}$,

$$\begin{aligned}
 [\sin(x' - x)]^2 &= 1 - [\cos(x' - x)]^2 \\
 &= 1 - [\cos x' \cos x + \sin x' \sin x]^2 \\
 &= 1 - (\mathbf{u}(x') \cdot \mathbf{u}(x))^2,
 \end{aligned}$$

where $\mathbf{u}(x) = (\cos x, \sin x)^\top$ for all $x \in \mathcal{X}$. Observe that $\|\mathbf{u}(x)\| = 1$ for all $x \in \mathcal{X}$. Thus, $[\sin(x' - x)]^2 = \frac{1}{2} \|\mathbf{u}(x') - \mathbf{u}(x)\|^2$ and $K(x, x') = e^{-\frac{\lambda}{2} \|\mathbf{u}(x') - \mathbf{u}(x)\|^2}$. Since the Gaussian kernel is known to be PDS, K is also PDS (the fact that Gaussian kernels are PDS can be shown easily by observing that they are the normalized kernels associated to the kernel obtained by applying exp to standard inner products).

(g) It suffices to show that K is the normalized kernel associated to the kernel K' defined by

$$\forall(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^N \times \mathbb{R}^N, K'(\mathbf{x}, \mathbf{y}) = e^{\phi(\mathbf{x}, \mathbf{y})}$$

where $\phi(\mathbf{x}, \mathbf{y}) = \frac{1}{\sigma} [\|\mathbf{x}\| + \|\mathbf{y}\| - \|\mathbf{x} - \mathbf{y}\|]$, and to show that K' is PDS. For the first part, observe that

$$\frac{K'(\mathbf{x}, \mathbf{y})}{\sqrt{K'(\mathbf{x}, \mathbf{x})K'(\mathbf{y}, \mathbf{y})}} = e^{\phi(\mathbf{x}, \mathbf{y}) - \frac{1}{2}\phi(\mathbf{x}, \mathbf{x}) - \frac{1}{2}\phi(\mathbf{y}, \mathbf{y})} = e^{-\frac{\|\mathbf{x} - \mathbf{y}\|}{\sigma}}.$$

To show that K' is PDS, it suffices to show that ϕ is PDS, since composition with a power series with non-negative coefficients (here exp) preserve the PDS property. Now, for any $c_1, \dots, c_n \in \mathbb{R}$, let $c_0 = -\sum_{i=1}^n c_i$, then, we can write

$$\begin{aligned} \sum_{i,j=1}^n c_i c_j \phi(\mathbf{x}_i, \mathbf{x}_j) &= \frac{1}{\sigma} \sum_{i,j=1}^n c_i c_j [\|\mathbf{x}_i\| + \|\mathbf{x}_j\| - \|\mathbf{x}_i - \mathbf{x}_j\|] \\ &= \frac{1}{\sigma} \left[-\sum_{i=1}^n c_0 c_i \|\mathbf{x}_i\| + -\sum_{i=1}^n c_0 c_j \|\mathbf{x}_j\| - \sum_{i,j=1}^n c_i c_j \|\mathbf{x}_i - \mathbf{x}_j\| \right] \\ &= -\frac{1}{\sigma} \sum_{i,j=0}^n c_i c_j \|\mathbf{x}_i - \mathbf{x}_j\|, \end{aligned}$$

with $\mathbf{x}_0 = 0$. Now, for any $z \in \mathbb{R}$, the following equality holds:

$$z^{\frac{1}{2}} = \frac{1}{2\Gamma(\frac{1}{2})} \int_0^{+\infty} \frac{1 - e^{-tz}}{t^{\frac{3}{2}}} dt.$$

Thus,

$$\begin{aligned} -\frac{1}{\sigma} \sum_{i,j=0}^n c_i c_j \|\mathbf{x}_i - \mathbf{x}_j\| &= \frac{1}{2\Gamma(\frac{1}{2})} \int_0^{+\infty} -\frac{1}{\sigma} \sum_{i,j=0}^n c_i c_j \frac{1 - e^{-t\|\mathbf{x}_i - \mathbf{x}_j\|}}{t^{\frac{3}{2}}} dt \\ &= \frac{1}{2\Gamma(\frac{1}{2})} \int_0^{+\infty} \frac{1}{\sigma} \frac{\sum_{i,j=0}^n c_i c_j e^{-t\|\mathbf{x}_i - \mathbf{x}_j\|}}{t^{\frac{3}{2}}} dt. \end{aligned}$$

Since a Gaussian kernel is PDS, the inequality $\sum_{i,j=0}^n c_i c_j e^{-t\|\mathbf{x}_i - \mathbf{x}_j\|} \geq 0$ holds and the right-hand side is non-negative. Thus, the inequality $-\frac{1}{\sigma} \sum_{i,j=0}^n c_i c_j \|\mathbf{x}_i - \mathbf{x}_j\| \geq 0$ holds, which shows that ϕ is PDS. \square

Alternatively, one can also apply the theorem on page 43 of the lecture slides on kernel methods to reduce the problem to showing that the norm $G(x, y) = \|x - y\|$ is a NDS function. This can be shown through a direct application of the definition of NDS together with the representation of the norm given in the hint.

(h) Let $\langle \cdot, \cdot \rangle$ denote the inner product over the set of all measurable functions on $[0, 1]$:

$$\langle f, g \rangle = \int_0^1 f(t)g(t)dt.$$

Note that $k_1(x, y) = \int_0^1 1_{t \in [0, x]} 1_{t \in [0, y]} dt = \langle 1_{[0, x]}, 1_{[0, y]} \rangle$, and $k_2(x, y) = \int_0^1 1_{t \in [x, 1]} 1_{t \in [y, 1]} dt = \langle 1_{[x, 1]}, 1_{[y, 1]} \rangle$. Since they are defined as inner products, both k_1 and k_2 are PDS. Note that $k_1(x, y) = \min(x, y)$ and $k_2(x, y) = 1 - \max(x, y)$. Observe that $K = k_1 k_2$, thus K is PDS as a product of two PDS kernels.

- (i) $f: x \mapsto \frac{1}{\sqrt{1-x}}$ admits the Taylor series expansion

$$f(x) = \sum_{n=0}^{\infty} \binom{1/2}{n} (-1)^n x^n$$

for $|x| < 1$, where $\binom{1/2}{n} = \frac{\frac{1}{2}(\frac{1}{2}-1)\cdots(\frac{1}{2}-n+1)}{n!}$. Observe that $\binom{1/2}{n}(-1)^n > 0$ for all $n \geq 0$, thus, the coefficients in the power series expansion are all positive. Since the radius of the convergence of the series is one and that by the Cauchy-Schwarz inequality $|\mathbf{x}' \cdot \mathbf{x}| \leq \|\mathbf{x}'\| \|\mathbf{x}\| < 1$ for $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, by the closure property theorem presented in class, $(\mathbf{x}, \mathbf{x}') \mapsto f(\mathbf{x}' \cdot \mathbf{x})$ is a PDS kernel.

Now, let $a_n = \binom{1/2}{n}(-1)^n$. Then, for $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$,

$$\begin{aligned} f(\mathbf{x}' \cdot \mathbf{x}) &= \sum_{n=0}^{\infty} a_n \left(\sum_{i=1}^N x_i x'_i \right)^n \\ &= \sum_{n=0}^{\infty} a_n \sum_{s_1+\cdots+s_N=n} \binom{n}{s_1, \dots, s_N} (x_{i_1} x'_{i_1})^{s_1} \cdots (x_{i_N} x'_{i_N})^{s_N} \\ &= \sum_{s_1, \dots, s_N \geq 0} a_{s_1+\cdots+s_N} \binom{s_1+\cdots+s_N}{s_1, \dots, s_N} (x_{i_1} x'_{i_1})^{s_1} \cdots (x_{i_N} x'_{i_N})^{s_N}, \end{aligned}$$

where the sums can be permuted since the series is absolutely summable. Thus, we can write $f(\mathbf{x}' \cdot \mathbf{x}) = \Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$ with

$$\Phi(\mathbf{x}) = \left(\sqrt{a_{s_1+\cdots+s_N} \binom{s_1+\cdots+s_N}{s_1, \dots, s_N}} x_{i_1}^{s_1} \cdots x_{i_N}^{s_N} \right)_{s_1, \dots, s_N \geq 0}.$$

Φ is a mapping to an infinite-dimensional space.

- (j) For $x \geq 0$, the integral $\int_0^{+\infty} e^{-sx} e^{-s} ds$ is well defined and

$$\int_0^{+\infty} e^{-sx} e^{-s} ds = \int_0^{+\infty} e^{-s(1+x)} ds = \left[-\frac{e^{-s(1+x)}}{1+x} \right]_0^{+\infty} = \frac{1}{1+x}. \quad (\text{E.88})$$

Since the Gaussian kernel, $(x, y) \mapsto e^{-\frac{\|x-y\|^2}{\sigma^2}}$ is PDS for any $\sigma \neq 0$, the kernel $(x, y) \mapsto \int_0^{+\infty} e^{-\frac{s\|x-y\|^2}{\sigma^2}} e^{-s} ds$ is also PDS since for any x_1, \dots, x_m in \mathbb{R}^N and c_1, \dots, c_m in \mathbb{R} ,

$$\sum_{i,j=1}^m c_i c_j e^{-\frac{s\|x_i-x_j\|^2}{\sigma^2}} \geq 0 \Rightarrow \int_0^{+\infty} \sum_{i,j=1}^m c_i c_j e^{-\frac{s\|x_i-x_j\|^2}{\sigma^2}} e^{-s} ds \geq 0. \quad (\text{E.89})$$

By (E.88), $\int_0^{+\infty} e^{-\frac{s\|x-y\|^2}{\sigma^2}} e^{-s} ds = \frac{1}{1+\frac{\|x-y\|^2}{\sigma^2}}$ for all x, y in \mathbb{R}^N , which concludes the proof.

- (k) Observe that for all $x, y \in \mathbb{R}$,

$$\int_0^{+\infty} 1_{t \in [0, |x|]} 1_{t \in [0, |y|]} dt = \min\{|x|, |y|\}. \quad (\text{E.90})$$

Thus, $(x, y) \mapsto \min(|x|, |y|)$ is a PDS kernel over $\mathbb{R} \times \mathbb{R}$ since for any x_1, \dots, x_m in \mathbb{R}^N and c_1, \dots, c_m in \mathbb{R} ,

$$\begin{aligned} & \sum_{i,j=1}^m c_i c_j \min\{|x_i|, |x_j|\} \\ &= \int_0^{+\infty} \sum_{i,j=1}^m c_i c_j \sum_{k=1}^N 1_{t \in [0, |x_i|]} 1_{t \in [0, |x_j|]} dt \\ &= \int_0^{+\infty} \left\| \begin{bmatrix} c_1 1_{t \in [0, |x_1|]} \\ \vdots \\ c_m 1_{t \in [0, |x_m|]} \end{bmatrix} \right\|^2 dt \geq 0. \end{aligned}$$

Thus, $(x, y) \mapsto \sum_{i=1}^N \min(|x_i|, |y_i|)$ is a PDS kernel over $\mathbb{R}^N \times \mathbb{R}^N$ as a sum of PDS kernels. Its composition with $x \mapsto e^{\frac{x}{\sigma^2}}$ with admits a power series with an infinite radius of convergence and non-negative coefficients is thus also PDS, which concludes the proof.

6.3 Graph kernel For any $i, j \in \mathcal{V}$, let $\mathcal{E}[i, j]$ denote the set of edges between i and j which is either reduced to one edge or is empty. For convenience, let $\mathcal{V} = \{1, \dots, n\}$ and define the matrix $\mathbf{W} = (W_{ij}) \in \mathbb{R}^{n \times n}$ by $W_{ij} = 0$ if $\mathcal{E}[i, j] = \emptyset$, $W_{ij} = w[e]$ if $e \in \mathcal{E}[i, j]$. Then, we can write

$$K(p, q) = \sum_{e \in \mathcal{E}[p, r], e' \in \mathcal{E}[r, q]} w[e] w[e'] = \sum_r W_{pr} W_{rq} = W_{pq}^2.$$

Let $\mathbf{K} = (K_{pq})$ denote the kernel matrix. Since \mathbf{W} is symmetric For any vector $\mathbf{X} \in \mathbb{R}^n$, $\mathbf{X}^\top \mathbf{K} \mathbf{X} = \mathbf{X}^\top \mathbf{W}^2 \mathbf{X} = \mathbf{X}^\top \mathbf{W}^\top \mathbf{W} \mathbf{X} = \|\mathbf{W} \mathbf{X}\|^2 \geq 0$. Thus, the eigenvalues of \mathbf{K} are non-negative. The same holds similarly for the kernel matrix restricted to any subset of \mathcal{V} , thus K is PDS.

6.4 Symmetric difference kernel

$k: (\mathcal{A}, \mathcal{B}) \mapsto |\mathcal{A} \cap \mathcal{B}|$ is a PDS kernel over $2^{\mathcal{X}}$ since $k(\mathcal{A}, \mathcal{B}) = \sum_{x \in \mathcal{X}} 1_{\mathcal{A}}(x) 1_{\mathcal{B}}(x) = \Phi(\mathcal{A}) \cdot \Phi(\mathcal{B})$, where, for any subset \mathcal{A} , $\Phi(\mathcal{A}) \in \{0, 1\}^{|\mathcal{X}|}$ is the vector whose coordinate indexed by $x \in \mathcal{X}$ is 1 if $x \in \mathcal{A}$, 0 otherwise. Since k is PDS, $K' = \exp(k)$ is also PDS (PDS property preserved by composition with power series). Since K is the result of the normalization of K' , it is also PDS.

6.5 Set kernel

Let Φ_0 be a feature mapping associated to k_0 , then $K'(\mathcal{A}, \mathcal{B}) = \sum_{x \in \mathcal{A}, x' \in \mathcal{B}} \phi_0(x) \cdot \phi_0(x') = \left(\sum_{x \in \mathcal{A}} \phi_0(x) \right) \cdot \left(\sum_{x' \in \mathcal{B}} \phi_0(x') \right) = \Psi(\mathcal{A}) \cdot \Psi(\mathcal{B})$, where for any subset \mathcal{A} , $\Psi(\mathcal{A}) = \sum_{x \in \mathcal{A}} \phi_0(x)$.

6.6 (a) Note that $K(x, y) = 1 - \cos^2(x - y)$. Now, for any $x_1, \dots, x_m \in \mathbb{R}$ and $c_1, \dots, c_m \in \mathbb{R}$ such that $\sum_{i=1}^m c_i = 0$,

$$\sum_{i,j=1}^m c_i c_j K(x_i, x_j) = \sum_{i,j=1}^m c_i c_j - \sum_{i,j=1}^m c_i c_j \cos^2(x_i - x_j) \quad (\text{E.91})$$

$$= - \sum_{i,j=1}^m c_i c_j \cos^2(x_i - x_j), \quad (\text{E.92})$$

since $\sum_{i,j=1}^m c_i c_j = (\sum_{i=1}^m c_i)^2 = 0$. It was already shown in a previous question that $K'(x, y) = \cos(x - y)$ defines a PDS kernel. Since the product of two PDS kernels is PDS, K'^2 is also PDS. This implies that $\sum_{i,j=1}^m c_i c_j \cos^2(x_i - x_j) \geq 0$, and $\sum_{i,j=1}^m c_i c_j K(x_i, x_j) \leq 0$, which can also be verified directly by expanding $\cos^2(x_i - x_j)$.

(b) As presented in this chapter, K is NDS iff $\exp(-tK)$ is PDS for all $t > 0$. Here, for any $t > 0$, $\exp(-tK) = (x + y)^{-t}$. It was already shown in a previous question that $(x + y)^{-1}$, which corresponds to the case $t = 1$, is PDS. Since a product of PDS kernels is PDS, $(x + y)^{-n}$ also defines a PDS kernel for any positive integer n . To see the general case, it suffices to show that $(x + y)^{-t}$ is PDS for any $0 < t < 1$.

6.7 Consider the Gram matrix defined as $\mathbf{K}_{ij} = K(x_i, x_j)$. It is clear that \mathbf{K} will have all zeros on the diagonal. Hence $\text{tr}(\mathbf{K}) = 0$. When $\mathbf{K} \neq 0$, this means it must have at least one negative eigenvalue. Hence K is not PDS.

6.8 The kernel K is NDS. In fact, more generally, $\|x - y\|^p$ for $0 < p \leq 2$ defines an NDS kernel. This follows a general fact: if a kernel K is NDS, then K^α is NDS for any $0 < \alpha < 1$. Indeed, if K is NDS, then $\exp(-tK)$ is PDS for all $t > 0$. Then $-\exp(-tK)$ is NDS and so is $1 - \exp(-tK)$. Now, the following formula holds for all $0 < \alpha < 1$:

$$K^\alpha = c_\alpha \int_0^{+\infty} [1 - \exp(-tK)] \frac{dt}{t^{\alpha+1}}, \quad (\text{E.93})$$

where c_α is a positive constant. Thus, K^α is NDS for all $0 < \alpha < 1$.

6.9 Let $n \geq 1$, $\{c_1, \dots, c_n\} \subseteq \mathbb{R}$ and $\{x_1, \dots, x_n\} \subseteq \mathcal{H}$ with $\sum_{i=1}^n c_i = 0$. Then,

$$\sum_{i,j=1}^n c_i c_j K(x_i, x_j) = \left(\sum_{i=1}^n c_i \right)^2 - \left\langle \sum_{i=1}^n c_i x_i, \sum_{j=1}^n c_j x_j \right\rangle = -\left\| \sum_{i=1}^n c_i x_i \right\|^2 \leq 0. \quad (\text{E.94})$$

6.10 It is straightforward to see that the kernel K defined over $\mathbb{R}_+ \times \mathbb{R}_+$ by $K(x, y) = x + y$ is negative definite symmetric (NDS) and that $K(x, x) \geq 0$ for all $x \geq 0$. Thus, for any $0 < \alpha \leq 1$, K^α is also NDS. Thus, $\exp(-tK^\alpha)$ is PDS for all $t > 0$, in particular $t = 1$. Assume now that K_p is PDS for some $p > 1$. Then, for all $\{c_1, \dots, c_n\} \subseteq \mathbb{R}$, $\{x_1, \dots, x_n\} \subseteq X$, and $t > 0$,

$$\sum_{i,j=1}^n c_i c_j e^{-t(x_i + x_j)^p} = \sum_{i,j=1}^n c_i c_j e^{-(t^{1/p} x_i + t^{1/p} x_j)^p} \geq 0. \quad (\text{E.95})$$

Thus, $\Psi(x, y) = (x + y)^p$ is NDS. This contradicts

$$\sum_{i,j=1}^2 c_i c_j (x_i + x_j)^p = 2^p - 2 > 0 \quad (\text{E.96})$$

for $x_1 = 0$, $x_2 = 1$, and $c_1 + c_2 = 0$ with $c_1 = 1$.

6.11 Explicit mappings

(a) Because \mathbf{K} is positive semidefinite, it can be diagonalized as $\mathbf{K} = \mathbf{S}\mathbf{\Lambda}\mathbf{S}^\top$ where $\mathbf{\Lambda}$ is a diagonal matrix of \mathbf{K} 's eigenvalues and \mathbf{S} is the matrix of \mathbf{K} 's eigenvectors. Further decomposing, we get $\mathbf{K} = \mathbf{S}\mathbf{\Lambda}^{1/2}\mathbf{\Lambda}^{1/2}\mathbf{S}^\top$. We then have

$$K(x_i, x_j) = \mathbf{K}_{ij} = (\mathbf{S}\mathbf{\Lambda}\mathbf{S}^\top)_{ij} = (\mathbf{\Lambda}^{1/2}\mathbf{S}_i) \cdot (\mathbf{\Lambda}^{1/2}\mathbf{S}_j)$$

where \mathbf{S}_i is the i th eigenvector of \mathbf{K} . Thus the kernel map $\Phi(x_i) = \mathbf{\Lambda}^{1/2}\mathbf{S}_i$ clearly satisfies the desired condition.

(b) For any $\alpha_1, \dots, \alpha_m \in \mathbb{R}$, we have

$$\sum_{i,j=1}^m \alpha_i \alpha_j \mathbf{K}_{ij} = \left\| \sum_{i=1}^m \alpha_i \Phi(x_i) \right\|^2 \geq 0$$

6.12 Explicit polynomial kernel mapping

It is clear that K can be written as a linear combinations of all monomials $x_1^{k_1} \dots x_N^{k_N}$ with $\sum_{j=1}^N k_j \leq d$. The dimension of the feature space is thus the number of such monomials, $f(N, d)$, that is the number of ways of adding N non-negative integers to obtain a sum of at most d . Note that any sum of $N-1$ integers less than or equal to d can be uniquely completely to be equal to d by adding one more term. This defines in fact a bijection between sums of $N-1$ integers with value at most d and sums of N integers with value equal to d . Thus, the number of sums of N integers exactly equal to d is $f(N-1, d)$.

Now, since a sum of N terms less than or equal to d is either equal to d or less than or equal to $d-1$,

$$f(N, d) = f(N-1, d) + f(N, d-1).$$

The result then follows by induction on $N + d$, using $f(1, 0) = f(0, 1) = 1$.

By the binomial identity, $K(\mathbf{x}, \mathbf{x}') = (\mathbf{x} \cdot \mathbf{x}' + c)^d = \sum_{i=0}^d \binom{d}{i} c^{d-i} (\mathbf{x} \cdot \mathbf{x}')^i$. The weight assigned to each k_i , $i \leq d$, is thus $\binom{d}{i} c^{d-i}$. Increasing c decreases the weight of k_i , particularly that of k_i s with larger i .

6.13 High-dimensional mapping

- (a) By definition of K' , for all $x, x' \in \mathcal{X}$, $\mathbb{E}_{\mathcal{D}}[K'(x, x')] = K(x, x')$. Replacing one index $i \in I$ by i' affects $K'(x, x')$ by at most $\frac{1}{n}(|[\Phi(x)]_i[\Phi(x')]_i| + |[\Phi(x)]_{i'}[\Phi(x')]_{i'}|) \leq \frac{2R^2}{n}$. The result thus follows directly McDiarmid's inequality.
- (b) Since \mathbf{K} and \mathbf{K}' are symmetric, it suffices to prove the statement for the entries of these matrices that are above the diagonal. By the union bound and the previous question, the following holds:

$$\mathbb{P}_{I \sim \mathcal{D}^n}[\exists i, j \in [m]: |\mathbf{K}'_{ij} - \mathbf{K}_{ij}| > \epsilon] \leq 2 \frac{m(m+1)}{2} e^{\frac{-n\epsilon^2}{2R^4}} = m(m+1)e^{\frac{-n\epsilon^2}{2R^4}}.$$

Setting $\delta > 0$ to match the upper bound leads directly to the inequality claimed.

6.14 Classifier based kernel.

- (a) By definition, given the distribution \mathcal{D} , h^* is defined as:

$$h^* = \underset{h: X \rightarrow \{-1, +1\}}{\operatorname{argmin}} \operatorname{error}_{\mathcal{D}}(h). \quad (\text{E.97})$$

K^* is clearly positive definite symmetric since $K^*(x, x')$ is defined as the dot product (in dimension one) of the features vectors $h^*(x)$ and $h^*(x')$.

- (b) The general expression of the solution is

$$h(x) = \operatorname{sgn}\left(\sum_{i=1}^m \alpha_i K^*(x, x_i) + b\right). \quad (\text{E.98})$$

Here, it is easy to see both in the separable and non-separable case that the solution is simply:

$$h(x) = \operatorname{sgn}(K^*(x, x_+)), \quad (\text{E.99})$$

where x_+ is such that $h^*(x_+) = +1$. One support vector is enough. The solution can be rewritten as

$$h(x) = h^*(x). \quad (\text{E.100})$$

The generalization error of the solution is thus that of the Bayes classifier (it is optimal). The data is separable iff the Bayes error is zero.

- (c) A kernel of this type is always positive definite symmetric since $K(x, x')$ is defined as a dot product of the feature vectors $h(x)$ and $h(x')$.

6.15 Image classification kernel

- (a) Observe that $\min(|u|^\alpha, |u'|^\alpha) = \int_0^{+\infty} 1_{t \in [0, |u'|^\alpha]} 1_{t \in [0, |u|^\alpha]} dt$, which shows that $(u, u') \mapsto \min(|u|^\alpha, |u'|^\alpha)$ is PDS.
- (b) Since $K_\alpha(x, x') = \sum_{k=1}^N \min(|x_k|^\alpha, |x'_k|^\alpha)$, K_α is PDS as a sum of N PDS kernels.

6.16 Fraud detection

Let $\langle \cdot, \cdot \rangle$ be the bilinear function defined over the space of all functions defined by:

$$\langle f, g \rangle = \mathbb{E}(f(x) - \mathbb{E}f(x))(g(x) - \mathbb{E}g(x)). \quad (\text{E.101})$$

Since the covariance matrix of n random variables is positive semidefinite and symmetric, this defines a dot product. Observe that $\mathbb{P}[U] = \mathbb{E}1_U$ and $\mathbb{P}[U \wedge V] = \mathbb{E}1_U 1_V$, thus

$$K(U, V) = \langle 1_U, 1_V \rangle. \quad (\text{E.102})$$

This shows that K is a positive definite symmetric kernel.

6.17 Relationship between NDS and PDS kernels

- Assume $\exp(-tK)$ is PDS. Then it is easy to see that $-\exp(-tK)$ is NDS. It is also straightforward to show that shifting by a constant and scaling by a positive constant $t > 0$ preserves the NDS property, thus, $\frac{1 - \exp(-tK)}{t}$ is NDS. Finally, we examine the one-sided limit of the expression, which is also NDS,

$$\lim_{t \rightarrow 0^+} \frac{1 - \exp(-tK)}{t} = - \left. \frac{\partial \exp(-xK)}{\partial x} \right|_{x=0} = K,$$

showing that K is NDS.

- Assume K is NDS. Theorem 6.16 then implies for some PDS function K' and x_0 :

$$\begin{aligned} -K(x, x') &= K'(x, x') - K(x, x_0) - K(x', x_0) + K(x_0, x_0) \\ \iff e^{-K(x, x')} &= e^{K'(x, x')} \underbrace{e^{-K(x, x_0)} e^{-K(x', x_0)} e^{K(x_0, x_0)}}_{\phi(x)\phi(x')}, \end{aligned}$$

where K' is PDS. Since K' is PDS $\exp(K')$ is also PDS, it is easy to show that any function of the form $K(x, x') = \phi(x)\phi(x')$ is PDS and finally the constant function $K(x, x') = \exp(K(x_0, x_0))$ is also PDS. Using the closure of PDS functions with respect to multiplication completes the proof.

6.18 Metrics and kernels

- (a) If K is an NDS kernel, then by theorem 6.16 the kernel K' defined for any $x_0 \in \mathcal{X}$ by:

$$K'(x, x') = \frac{1}{2}[K(x, x_0) + K(x', x_0) - K(x, x')]$$

is a PDS kernel ($K(x_0, x_0) = 0$). Let \mathbb{H} be the reproducing Hilbert space associated to K' . There exists a mapping $\Phi(x)$ from \mathcal{X} to \mathbb{H} such that $\forall x, x' \in \mathcal{X}, K'(x, x') = \Phi(x) \cdot \Phi(x')$. Then,

$$\begin{aligned} \|\Phi(x) - \Phi(x')\|^2 &= K'(x, x) + K'(x', x') - 2K'(x, x') \\ &= \frac{1}{2}[2K(x, x_0) - K(x, x)] + \\ &\quad \frac{1}{2}[2K(x', x_0) - K(x', x')] - \\ &\quad [K(x, x_0) + K(x', x_0) - K(x, x')] \\ &= K(x, x') \end{aligned}$$

It is then straightforward to show that \sqrt{K} is a metric.

- (b) Suppose that $K(x, x') = \exp(-|x - x'|^p)$, $x, x' \in \mathbb{R}$, is positive definite for $p > 2$. Then, for any $t > 0$, $\{x_1, \dots, x_n\} \subseteq X$, $\{c_1, \dots, c_n\} \subseteq \mathbb{R}$,

$$\sum_{i,j=1}^n c_i c_j \exp(-t|x_j - x_k|^p) = \sum_{i,j=1}^n c_i c_j \exp(-|t^{1/p}x_j - t^{1/p}x_k|^p) \geq 0$$

Thus, by theorem 6.17, $K'(x, x') = |x - x'|^p$ is an NDS kernel. But, $\sqrt{K'}$ is not a metric for $p > 2$ since it does not verify the triangle inequality (take $x = 1$, $x' = 2$, $x'' = 3$), which contradicts part (a).

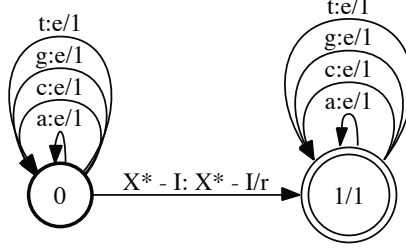
- (c) If $a < 0$ or $b < 0$, $a\|x\|^2 + b < 0$ for some non-null vectors x . For such values, $K(x, x) = \tanh(a\|x\|^2 + b) < 0$. The kernel is thus not PDS and the SVM training may not converge to an optimal value. The equivalent neural network may also converge to a local minimum.

6.19 Sequence kernels

- (a) $X^* - I$ is a regular language and can be represented by a finite automaton. K can thus be defined by

$$\forall x, y \in X^*, \quad K(x, y) = [[T \circ T^{-1}]](x, y), \quad (\text{E.103})$$

where T is the weighted transducer shown in figure E.14. Thus, K is a rational kernel and in view of the theorem 6.21, it is positive definite symmetric.

**Figure E.14**

Weighted transducer T . e represents the empty string, and $r = \rho$. $X^* - I$ stands for a finite automaton accepting $X^* - I$.

- (b) Let $M_{X^* - I}$ be the minimal automaton representing $X^* - I$. The transducer T of figure E.14 can be constructed using $M_{X^* - I}$. Then, $|T| = |M_{X^* - I}| + 8$. Using composition of weighted transducers, the running time complexity of the computation of the algorithm is:

$$O(|x||y||T \circ T^{-1}|) = O(|x||y||T|^2) = O(|x||y||M_{X^* - I}|^2). \quad (\text{E.104})$$

- (c) The set of strings Y over the alphabet X of length less than n form a regular language since they can be described by:

$$Y = \bigcup_{i=0}^{n-1} X^i. \quad (\text{E.105})$$

Thus, $Y_1 = Y \cap (X^* - I)$ and $Y_2 = (X^* - I) - Y_1$ are also regular languages. It suffices to replace in the transducer T of figure E.14 the transition labeled with $X^* - I : X^* - I/\rho$ with two transitions:

- $Y_1 : Y_1/\rho_1$, and
- $Y_2 : Y_2/\rho_2$,

with the same origin and destination states and with Y_1 and Y_2 denoting finite automata representing them. The kernel is thus still rational and PDS since it is of the form $T' \circ T'^{-1}$.

6.20 n -gram kernel

6.21 Mercer's condition

We prove the contrapositive, that if K is not PDS then Mercer's condition necessarily does not hold. Consider \mathbf{c} and a sample (x_1, \dots, x_m) such that $\sum_{i,j=1}^m c_i c_j K(x_i, x_j) < 0$, which exists by the non-PDS assumption. Then define a function c_σ that is the sum of m Gaussians $\sum_{i=1}^m c_i \mathcal{N}(x_i, \sigma)$. For sufficiently small σ and $\forall i \in [m]$, $c_\sigma(x_i) \approx c_i p_i$, where p_i is equal to the number of occurrences of x_i in the sample, and approximately zero elsewhere. Thus, for some ϵ_σ , where $|\epsilon_\sigma| \rightarrow 0$ as $\sigma \rightarrow 0$, $\int \int_{\mathcal{X} \times \mathcal{X}} c(x) c(x') K(x, x') dx dx' \leq \sum_{i,j=1}^m c_i c_j K(x_i, x_j) + \epsilon_\sigma$, which is less than 0 for sufficiently small σ .

6.22 Anomaly detection

- (a) The Lagrangian for the optimization problem is:

$$\mathcal{L}(\mathbf{c}, r, \boldsymbol{\alpha}) = r^2 + \sum_{i=1}^m \alpha_i (\|\Phi(x_i) - \mathbf{c}\|^2 - r^2).$$

By the KKT conditions, the optimal solution must satisfy

$$\begin{aligned}\nabla_r \mathcal{L} &= 2r - \sum_{i=1}^m \alpha_i 2r = 0 \iff \sum_{i=1}^m \alpha_i = 1 \\ \nabla_{\mathbf{c}} \mathcal{L} &= \sum_{i=1}^m \alpha_i (2\mathbf{c} - 2\Phi(x_i)) = 0 \iff \mathbf{c} = \sum_{i=1}^m \alpha_i \Phi(x_i) / \left(\sum_{i=1}^m \alpha_i \right).\end{aligned}$$

Combining the two statements above, we see that at the optimum the solution must satisfy $\mathbf{c} = \sum_{i=1}^m \alpha_i \Phi(x_i)$. Finally, adding the explicit constraint $\sum_{i=1}^m \alpha_i = 1$ to the optimization problem and substituting the expression for \mathbf{c} into the Lagrangian results in the dual optimization problem.

- (b) Since we know $\mathbf{c} = \sum_{i=1}^m \alpha_i \Phi(x_i)$, $\alpha_i \geq 0$ and $\sum_{i=1}^m \alpha_i = 1$, we have

$$\begin{aligned}\|\mathbf{c}\| &= \left\| \sum_{i=1}^m \alpha_i \Phi(x_i) \right\| \leq \max_i \|\Phi(x_i)\| \leq M \\ r &= \max_i \|\mathbf{c} - \Phi(x_i)\| \leq 2 \max_i \|\Phi(x_i)\| \leq 2M.\end{aligned}$$

- (c) Note that the losses over \mathcal{D} coincide with the losses in standard binary classification problems in the case when all points in the source domain are labeled with the positive class. Thus, we may directly apply theorem 5.9. We simply need to bound the empirical Rademacher complexity of \mathcal{H} :

$$\begin{aligned}m\hat{\mathcal{R}}(\mathcal{H}) &= \mathbb{E}_{\sigma} \left[\sup_{h \in \mathcal{H}} \sum_{i=1}^m \sigma_i h(x_i) \right] \\ &= \mathbb{E}_{\sigma} \left[\sup_{r, \mathbf{c}} \sum_{i=1}^m \sigma_i (r^2 - \|\Phi(x_i) - \mathbf{c}\|^2) \right] \\ &\leq \mathbb{E}_{\sigma} \left[\sup_r \sum_{i=1}^m \sigma_i r^2 \right] + \mathbb{E}_{\sigma} \left[\sum_{i=1}^m \|\Phi(x_i)\|^2 \right] \\ &\quad + \mathbb{E}_{\sigma} \left[\sup_{\mathbf{c}} \sum_{i=1}^m \|\mathbf{c}\|^2 \right] + \mathbb{E}_{\sigma} \left[\sup_{\mathbf{c}} \sum_{i=1}^m 2\mathbf{c} \cdot \Phi(x_i) \right].\end{aligned}$$

We bound each of these four terms separately. First we have

$$\mathbb{E}_{\sigma} \left[\sum_{i=1}^m \|\Phi(x_i)\|^2 \right] = \mathbb{E}_{\sigma} \left[\sum_{i=1}^m \sigma_i K(x_i, x_i) \right] = \sum_{i=1}^m \mathbb{E}_{\sigma} [\sigma_i K(x_i, x_i)] = 0.$$

Next we bound

$$\begin{aligned}\mathbb{E}_{\sigma} \left[\sup_{\mathbf{c}} \sum_{i=1}^m \sigma_i \|\mathbf{c}\|^2 \right] &= \mathbb{E}_{\sigma} \left[C^2 \sum_{i=1}^m \sigma_i 1_{\sum_{i=1}^m \sigma_i > 0} \right] \\ &\leq C^2 \mathbb{E}_{\sigma} \left[\left| \sum_{i=1}^m \sigma_i \right| \right] \\ &\leq C^2 \sqrt{\mathbb{E}_{\sigma} \left[\left| \sum_{i=1}^m \sigma_i \right|^2 \right]} \quad (\text{Jensen's ineq.}) \\ &= C^2 \sqrt{\mathbb{E}_{\sigma} \left[\sum_{i \neq j} \sigma_i \sigma_j + \sum_{i=1}^m \sigma_i^2 \right]} \\ &= C^2 \sqrt{\mathbb{E}_{\sigma} \left[\sum_{i=1}^m \sigma_i^2 \right]} = C^2 \sqrt{m}.\end{aligned}$$

The term containing r is bound similarly,

$$\left[\sup_r \sum_{i=1}^m \sigma_i r^2 \right] \leq R^2 \sqrt{m}.$$

The final term is bounded as follows

$$\begin{aligned} 2 \mathbb{E}_{\sigma} \left[\sup_{\mathbf{c}} \sum_{i=1}^m \sigma_i \mathbf{c} \cdot \Phi(x_i) \right] &\leq \mathbb{E}_{\sigma} \left[\sup_{\mathbf{c}} \|\mathbf{c}\| \left\| \sum_{i=1}^m \sigma_i \Phi(x_i) \right\| \right] \quad (\text{Cauchy-Schwarz ineq.}) \\ &= C \mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^m \sigma_i \Phi(x_i) \right\| \right] \leq C \sqrt{\text{Tr}[\mathbf{K}]}. \end{aligned}$$

Collecting terms and plugging into theorem 5.9 completes the generalization guarantee.

- (d) The solution follows very similar steps as in part (a), but using the same modifications as seen in the analysis of soft-margin SVMs.

Chapter 7

7.1 VC-dimension of the hypothesis set of AdaBoost.

7.2 Alternative objective functions.

- (a) The direction e_u taken by coordinate descent after $T-1$ rounds is the argmin_u of:

$$\begin{aligned} \left. \frac{dL(\alpha + \beta e_u)}{d\beta} \right|_{\beta=0} &= - \sum_{i=1}^m y_i h_u(x_i) \Phi'(-y_i f(x_i)) \\ &\propto - \sum_{i=1}^m y_i h_u(x_i) \frac{\Phi'(-y_i f(x_i))}{\sum_{i=1}^m \Phi'(-y_i f(x_i))} \\ &\propto - \sum_{i=1}^m y_i h_u(x_i) \mathcal{D}_{T-1}(i) \\ &= -(1 - 2\epsilon_u), \end{aligned}$$

where $\mathcal{D}_{T-1}(i) = \frac{1}{m} \frac{\Phi'(-y_i f(x_i))}{\sum_{i=1}^m \Phi'(-y_i f(x_i))}$. Thus, the base classifier h_u selected at each round is the one with the minimal error rate over the training data.

- (b) Only the logistic loss (Φ_4) verifies the hypotheses (assuming log with base 2).
(c) The step size β is the solution of:

$$\begin{aligned} \frac{dL(\alpha + \beta e_u)}{d\beta} &= - \sum_{i=1}^m y_i h_u(x_i) \Phi'_4(-y_i f(x_i) - \beta y_i h_u(x_i)) = 0 \\ \Leftrightarrow - \sum_{i=1}^m y_i h_u(x_i) \frac{1}{1 + \exp(y_i f(x_i) + \beta y_i h_u(x_i))} &= 0. \end{aligned}$$

Using β as the step function, the algorithm differs from AdaBoost, only by the choice of the distribution update at each boosting round, here:

$$\mathcal{D}_{t+1}(i) = \frac{1}{Z_t} \frac{1}{1 + \exp(y_i f(x_i))}.$$

7.3 Update guarantee

By the weak learning assumption, there exists a hypothesis $h \in \mathcal{H}$ whose \mathcal{D}_{t+1} -error is less than half. Examine the empirical error of h_t for the distribution \mathcal{D}_{t+1} . Since $Z_t = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$

and $\alpha_t = \frac{1}{2} \log \frac{1-\epsilon_t}{\epsilon_t}$,

$$\begin{aligned}\hat{R}_{\mathcal{D}_{t+1}}(h_t) &= \sum_{i=1}^m \frac{\mathcal{D}_t(i) e^{-\alpha_t y_i h_t(x_i)}}{Z_t} 1_{y_i h_t(x_i) < 0} \\ &= \sum_{y_i h_t(x_i) < 0}^m \frac{\mathcal{D}_t(i) e^{\alpha_t}}{Z_t} \\ &= \frac{e^{\alpha_t}}{Z_t} \sum_{y_i h_t(x_i) < 0}^m \mathcal{D}_t(i) \\ &= \frac{\sqrt{\frac{1-\epsilon_t}{\epsilon_t}}}{2\sqrt{\epsilon_t(1-\epsilon_t)}} \epsilon_t = \frac{1}{2}.\end{aligned}$$

This shows that h_t cannot be selected at round $t+1$.

7.4 Weighted instances

For all α , $F(\alpha) = \sum_{i=1}^m w_i e^{-y_i f(x_i)} = \sum_{i=1}^m w_i e^{-y_i \sum_{t=1}^T \alpha_t h_t(x_i)}$, with $w_i \geq 0$. F is convex as a non-negative linear combination of convex functions and is clearly differentiable. Applying coordinate descent to this function leads to the same algorithm as AdaBoost with the only difference that

$$\mathcal{D}_1(i) \leftarrow \frac{w_i}{\sum_{i=1}^m w_i}. \quad (\text{E.106})$$

7.5 By definition the unnormalized correlation is given by

$$\sum_{i=1}^m \mathcal{D}_{t+1}(i) y_i h_t(x_i) = \sum_{i=1}^m \frac{\mathcal{D}_t(i) e^{-\alpha_t y_i h_t(x_i)} y_i h_t(x_i)}{Z_t} = \frac{1}{Z_t} \frac{dZ_t}{d\alpha_t}, \quad (\text{E.107})$$

since $Z_t = \sum_{i=1}^m \mathcal{D}_t(i) e^{-\alpha_t y_i h_t(x_i)}$. Recall that α_t minimizes Z_t , thus $\frac{dZ_t}{d\alpha_t} = 0$. This shows that the distribution at round $t+1$ and the vector of margins at round t are uncorrelated.

7.6 It is not hard to see that the base hypotheses in this problem can be defined to be threshold functions based on the first or second axis, or constant functions (horizontal or vertical thresholds outside the convex hull of all the points).

The hypotheses selected by AdaBoost are therefore chosen from this set. It can be shown that the hypotheses selected in two consecutive rounds of AdaBoost are distinct (see exercise 7.3). Furthermore, h_t and $-h_t$ cannot be selected in consecutive rounds, since misclassified and correctly classified points by h_t are assigned the same distribution mass. Thus, at each round a distinct hypothesis is chosen. The points at coordinate $(-1, -1)$ are misclassified by all these base hypotheses.

The algorithm should be stopped when the best ϵ_t found is $1/2$. It can be shown, then, that the error of the final classifier returned on the training set is $\frac{1}{4}(1-\epsilon)$, since it misclassifies exactly the points at at coordinate $(+1, -1)$.

7.7 Noise-tolerant AdaBoost.

- (a) Let $G_1(x) = e^x$ and $G_2(x) = x+1$. G_1 and G_2 are continuously differentiable over \mathbb{R} and $G'_1(0) = G'_2(0)$. Thus, G is differentiable over \mathbb{R} . Note that $G' \geq 0$.

Both G_1 and G_2 are convex, thus

$$G(y) - G(x) \geq G'(x)(y-x) \quad (\text{E.108})$$

for $x, y \leq 0$ or $x, y \geq 0$. Assume now that $y \leq 0$ and $x \geq 0$, then

$$\begin{aligned}G(y) - G(x) &= e^y - (x+1) \geq (y+1) - (x+1) = G'(x)(y-x), \\ \text{since } G'(x) &= 1. \text{ Thus } G \text{ is convex.}\end{aligned} \quad (\text{E.109})$$

- (b) The direction e_u taken by coordinate descent after $T - 1$ rounds is the argmin_u of:

$$\left. \frac{dG(\alpha + \beta e_u)}{d\beta} \right|_{\beta=0} = - \sum_{i=1}^m y_i h_u(x_i) G'(-y_i f(x_i)) \quad (\text{E.110})$$

$$(\text{since } G' \geq 0) \propto - \sum_{i=1}^m y_i h_u(x_i) \frac{G'(-y_i f(x_i))}{\sum_{i=1}^m G'(-y_i f(x_i))} \quad (\text{E.111})$$

$$\propto - \sum_{i=1}^m y_i h_u(x_i) \mathcal{D}_{T-1}(i) \quad (\text{E.112})$$

$$= -(1 - 2\epsilon_u), \quad (\text{E.113})$$

with $\mathcal{D}_{T-1}(i) = \frac{1}{m} \frac{G'(-y_i f(x_i))}{\sum_{i=1}^m G'(-y_i f(x_i))}$. Thus, the base classifier h_u selected at each round is the one with the minimal error rate over the training data.

The step size β is the solution of:

$$\frac{dF(\alpha + \beta e_u)}{d\beta} = - \sum_{i=1}^m y_i h_u(x_i) G'(-y_i f(x_i) - \beta y_i h_u(x_i)) = 0, \quad (\text{E.114})$$

which can be solved numerically. A closed-form solution can be given under certain conditions, e.g., if

$$\beta \leq \rho = \min_{i \in [m]} |f(x_i)|. \quad (\text{E.115})$$

- (c) By definition of the objective function, this algorithm is less aggressively reducing the empirical error rate than AdaBoost.

7.8 Simplified AdaBoost

- (a) As in the standard case, we can show that

$$\hat{R}(h) \leq \prod_{t=1}^T Z_t, \quad (\text{E.116})$$

and that

$$Z_t = (1 - \epsilon_t)e^{-\alpha} + \epsilon_t e^{\alpha}. \quad (\text{E.117})$$

By definition of γ and the fact that $e^{\alpha} - e^{-\alpha} > 0$ for all $\alpha > 0$,

$$Z_t = \epsilon_t(e^{\alpha} - e^{-\alpha}) + e^{-\alpha} \quad (\text{E.118})$$

$$\leq (1 - \gamma)(e^{\alpha} - e^{-\alpha}) + e^{-\alpha} \quad (\text{E.119})$$

$$= \left(\frac{1}{2} - \gamma\right)e^{\alpha} + \left(\frac{1}{2} + \gamma\right)e^{-\alpha} = u(\alpha). \quad (\text{E.120})$$

$u(\alpha)$ is minimized for

$$\left(\frac{1}{2} - \gamma\right)e^{\alpha} = \left(\frac{1}{2} + \gamma\right)e^{-\alpha}, \quad (\text{E.121})$$

that is, for

$$\alpha = \frac{1}{2} \log \frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma}. \quad (\text{E.122})$$

Tighter bounds on the product of the Z_t s can lead to better values for α .

- (b) As in the standard case, at round t , the probability mass assigned to correctly classified points is $p_+ = (1 - \epsilon_t)e^{-\alpha}$ and the probability mass assigned to the misclassified points is $p_- = \epsilon_t e^{\alpha}$. Thus,

$$\frac{p_-}{p_+} = \frac{\epsilon_t}{1 - \epsilon_t} \frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma} \leq \frac{\frac{1}{2} - \gamma}{\frac{1}{2} + \gamma} \frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma} = 1. \quad (\text{E.123})$$

This contrasts with AdaBoost's property.

(c)

$$Z_t \leq \left(\frac{1}{2} - \gamma\right)e^\alpha + \left(\frac{1}{2} + \gamma\right)e^{-\alpha} \quad (\text{E.124})$$

$$= \left(\frac{1}{2} - \gamma\right)\sqrt{\frac{\frac{1}{2} + \gamma}{\frac{1}{2} - \gamma}} + \left(\frac{1}{2} + \gamma\right)\sqrt{\frac{\frac{1}{2} - \gamma}{\frac{1}{2} + \gamma}} \quad (\text{E.125})$$

$$= 2\sqrt{\left(\frac{1}{2} + \gamma\right)\left(\frac{1}{2} - \gamma\right)}. \quad (\text{E.126})$$

Thus, the empirical error can be bounded as follows:

$$\widehat{R}_S(h) \leq \prod_{t=1}^T Z_t \quad (\text{E.127})$$

$$\leq [2\sqrt{\left(\frac{1}{2} + \gamma\right)\left(\frac{1}{2} - \gamma\right)}]^T \quad (\text{E.128})$$

$$= (1 - 4\gamma^2)^{T/2} \quad (\text{E.129})$$

$$\leq e^{-2\gamma^2 T}. \quad (\text{E.130})$$

(d) If $\widehat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m 1_{y_i f(x_i) \leq 0} \leq \frac{1}{m}$, then clearly $\widehat{R}_S(h) = 0$. Using the bound obtained in the previous question, if $e^{-2\gamma^2 T} < \frac{1}{m}$, the empirical error is zero. This can be rewritten as

$$T > \frac{\log m}{2\gamma^2}. \quad (\text{E.131})$$

(e) Using the bound for the consistent case,

$$\mathbb{P}[R(h) > \epsilon] \leq 2\Pi e(2m)2^{-\frac{m\epsilon}{2}} \leq 2\left(\frac{2em}{d}\right)^d 2^{-\frac{m\epsilon}{2}}. \quad (\text{E.132})$$

Setting the right-hand side to δ , with probability at least $1 - \delta$, the following bound holds for that consistent hypothesis:

$$\text{error}_{\mathcal{D}}(\mathcal{H}) \leq \frac{2}{m} \left(d \log_2 \frac{2em}{d} + \log_2 \frac{2}{\delta} \right), \quad (\text{E.133})$$

with $d = 2(s+1)T \log_2(eT)$ and $T = \left\lfloor \frac{\log m}{2\gamma^2} \right\rfloor + 1$.

The bound is vacuous for $\gamma(m) = O(\sqrt{\frac{\log m}{m}})$. This could suggest overfitting.

7.9 AdaBoost example.

(a) At $t = 1$ we have:

$$\mathbf{d}_1^\top \mathbf{M} = \left(\frac{\sqrt{5}-1}{2}, 0, \frac{3-\sqrt{5}}{2}, \frac{3\sqrt{5}-1}{12}, \frac{3\sqrt{5}-1}{12}, \frac{3\sqrt{5}-1}{12}, \frac{1}{2}, \frac{11-3\sqrt{5}}{12} \right),$$

so we pick weak classifier 1. Now, the distribution at round two is:

$$\mathbf{d}_2 = \left(\frac{1}{4}, \frac{1}{4}, \frac{\sqrt{5}-1}{12}, \frac{\sqrt{5}-1}{12}, \frac{\sqrt{5}-1}{12}, \frac{3-\sqrt{5}}{8}, \frac{3-\sqrt{5}}{8}, 0 \right)^\top,$$

and the edges at round 2 are:

$$\mathbf{d}_2^\top \mathbf{M} = \left(0, \frac{3-\sqrt{5}}{2}, \frac{\sqrt{5}-1}{2}, \frac{4-\sqrt{5}}{6}, \frac{4-\sqrt{5}}{6}, \frac{4-\sqrt{5}}{6}, \frac{\sqrt{5}-1}{4}, \frac{5+\sqrt{5}}{12} \right),$$

so we pick weak classifier 3. Continuing this process, we then pick weak classifier 2 in round 3. However, now we observe that $\mathbf{d}_4 = \mathbf{d}_1$; hence, we have found a cycle, in which we repeatedly select classifiers 1, 3, 2, 1, 3, 2, ...

(b) $r_t = \frac{\sqrt{5}-1}{2}, t \in 1, 2, 3$. Thus, the coefficients used to combine classifiers in our example are: $[\frac{1}{3}, \frac{2}{3}, \frac{1}{3}, 0, 0, 0, 0, 0]$ and the margin equals the minimum value in the following vector: $\mathbf{M}[\frac{1}{3}, \frac{1}{3}, \frac{1}{3}, 0, 0, 0, 0, 0]^\top$, which is $\frac{1}{3}$.

- (c) $\frac{1}{16} \mathbf{M}[2, 3, 4, 1, 2, 2, 1, 1]^\top = \frac{3}{8}$ for all training points. This margin is greater than the one generated by AdaBoost. Therefore AdaBoost does NOT always maximize the L_1 norm margin.

7.10 Boosting in the presence of unknown labels

- (a) Say a ‘boosting-style algorithm’ is just AdaBoost with a possibly different step size α_t . Recall these definitions from the description of AdaBoost: The final hypothesis is $f(x) = \sum_t \alpha_t h_t(x)$ and the normalization constant in round t is $Z_t = \sum_i \mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i))$. We proved in the chapter that

$$\frac{1}{m} \sum_i 1_{y_i f(x_i) < 0} \leq \frac{1}{m} \sum_i \exp(-y_i f(x_i)) = \prod_t Z_t$$

and that AdaBoost’s step size can be derived by minimizing this objective in each round t . Taking that same approach, observe that

$$Z_t = \sum_i \mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i)) = \epsilon_t^0 + \epsilon_t^- \exp(\alpha_t) + \epsilon_t^+ \exp(-\alpha_t).$$

Differentiating the right-hand side with respect to α_t and setting equal to zero shows that Z_t is minimized by letting $\alpha_t = \frac{1}{2} \log \left(\frac{\epsilon_t^+}{\epsilon_t^-} \right)$.

- (b) One possible assumption is $\frac{\epsilon_t^+ - \epsilon_t^-}{\sqrt{1 - \epsilon_t^0}} \geq \gamma > 0$. Informally, this assumption says that the difference between the accuracy and error of each weak hypothesis is non-negligible relative to the fraction of examples on which the hypothesis makes any prediction at all. In part (d) we will prove that this assumption suffices to drive the training error to zero.
- (c) 1. Given: Training examples $((x_1, y_1), \dots, (x_m, y_m))$.
 2. Initialize \mathcal{D}_1 to the uniform distribution on training examples.
 3. for $t = 1, \dots, T$:
 a. $h_t \leftarrow$ base classifier in H .
 b. $\alpha_t \leftarrow \frac{1}{2} \log \left(\frac{\epsilon_t^+}{\epsilon_t^-} \right)$.
 c. For each $i = 1, \dots, m$: $\mathcal{D}_{t+1}(i) \leftarrow \frac{\mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$, where $Z_t \leftarrow \sum_i \mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i))$ is the normalization constant.
 4. $f \leftarrow \sum_{t=1}^T \alpha_t h_t$.
 5. Return: $\text{sgn}(f)$.
- (d) Plug in the value of α_t from part (a) into $Z_t = \epsilon_t^0 + \epsilon_t^- \exp(\alpha_t) + \epsilon_t^+ \exp(-\alpha_t)$ to obtain $Z_t = \epsilon_t^0 + 2\sqrt{\epsilon_t^- \epsilon_t^+}$. Therefore

$$\frac{1}{m} \sum_i 1_{y_i f(x_i) < 0} \leq \prod_t Z_t = \prod_t \left(\epsilon_t^0 + 2\sqrt{\epsilon_t^- \epsilon_t^+} \right).$$

Moreover, if the weak learning assumption from part (b) is satisfied then

$$\begin{aligned} \epsilon_t^0 + 2\sqrt{\epsilon_t^- \epsilon_t^+} &= \epsilon_t^0 + \sqrt{(1 - \epsilon_t^0)^2 - (\epsilon_t^+ - \epsilon_t^-)^2} \\ &= \epsilon_t^0 + (1 - \epsilon_t^0) \sqrt{1 - \frac{(\epsilon_t^+ - \epsilon_t^-)^2}{(1 - \epsilon_t^0)^2}} \\ &\leq \sqrt{1 - \frac{(\epsilon_t^+ - \epsilon_t^-)^2}{1 - \epsilon_t^0}} \\ &\leq \sqrt{1 - \gamma^2}. \end{aligned}$$

The first equality follows from $(\epsilon_t^+ + \epsilon_t^-)^2 - (\epsilon_t^+ - \epsilon_t^-)^2 = 4\epsilon_t^+ \epsilon_t^-$ (just multiply and gather terms) and $\epsilon_t^+ + \epsilon_t^- = 1 - \epsilon_t^0$. The first inequality follows from the fact that square root is

concave on $[0, \infty)$, and thus $\lambda\sqrt{x} + (1-\lambda)\sqrt{y} \leq \sqrt{\lambda x + (1-\lambda)y}$ for $\lambda \in [0, 1]$. The last inequality follows from the weak learning assumption.

Therefore we have $\frac{1}{m} \sum_i 1_{y_i f(x_i) < 0} \leq \left(\sqrt{1-\gamma^2}\right)^T \leq \exp\left(-\frac{\gamma^2 T}{2}\right)$, where we used $1+x \leq \exp(x)$.

7.11 HingeBoost

- (a) Since the hinge loss is convex, its composition with affine function of α is also convex and F is convex as a sum of convex functions.

For the existence of one-sided directional derivatives, one can use the fact that any convex function has one-sided directional derivatives or alternatively, that our specific function is the sum of piecewise affine functions, which are also known to have one-sided directional derivatives (think of one-dimensional hinge loss).

- (b) Distinguishing different cases depending on the value of $y_i f_{t-1}(x_i) = 1$, it is straightforward to derive the following expressions for all $j \in [N]$:

$$F'_+(\alpha_{t-1}, e_j) = \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f_{t-1}(x_i) < 1} + 1_{(y_i h_j(x_i) < 0) \wedge (y_i f_{t-1}(x_i) = 1)}]$$

$$F'_-(\alpha_{t-1}, e_j) = \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f_{t-1}(x_i) < 1} + 1_{(y_i h_j(x_i) > 0) \wedge (y_i f_{t-1}(x_i) = 1)}].$$

The key here is that when $y_i f_{t-1}(x_i) \neq 1$, each term in the sum will be either 0 or the affine function independent of $y_i h_j(x_i)$. On the other hand, when $y_i f_{t-1}(x_i) = 1$, the sign of $y_i h_j(x_i)$ determines whether the finite differences will extend into the 0 portion of the affine portion of the term.

- (c)

HINGEBOOST($S = ((x_1, y_1), \dots, (x_m, y_m))$)

```

1   $f \leftarrow 0$ 
2  for  $j \leftarrow 1$  to  $N$  do
3       $r \leftarrow \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f(x_i) < 1} + 1_{(y_i h_j(x_i) < 0) \wedge (y_i f(x_i) = 1)}]$ 
4       $l \leftarrow \sum_{i=1}^m -y_i h_j(x_i) [1_{y_i f(x_i) < 1} + 1_{(y_i h_j(x_i) > 0) \wedge (y_i f(x_i) = 1)}]$ 
5      if  $(l \leq 0) \wedge (r \geq 0)$  then
6           $d[j] \leftarrow 0$ 
7      elseif  $(l \leq r)$  then
8           $d[j] \leftarrow r$ 
9      else  $d[j] \leftarrow l$ 
10 for  $t \leftarrow 1$  to  $T$  do
11      $k \leftarrow \operatorname{argmin}_{j \in [N]} |d[j]|$ 
12      $\eta \leftarrow \operatorname{argmin}_{\eta \geq 0} G(f + \eta h_k) \triangleright \text{line search}$ 
13      $f \leftarrow f + \eta h_k$ 
14 return  $f$ 
```

7.12 Empirical margin loss boosting

(a) The following shows how $\hat{R}_{S,\rho}(f)$ can be upper-bounded:

$$\begin{aligned}\hat{R}_{S,\rho}(f) &= \frac{1}{m} \sum_{i=1}^m 1_{y_i f(x_i) \leq \rho} = \frac{1}{m} \sum_{i=1}^m 1_{y_i \sum_{t=1}^T \alpha_t h_t(x_i) - \rho \sum_{t=1}^T \alpha_t \leq 0} \\ &\leq \frac{1}{m} \sum_{i=1}^m \exp \left(-y_i \sum_{t=1}^T \alpha_t h_t(x_i) + \rho \sum_{t=1}^T \alpha_t \right).\end{aligned}$$

(b) Since \exp is convex and that composition with an affine function (of α) preserves convexity, each term of the objective is a convex function and G_ρ is convex as a sum of convex functions. The differentiability follows directly that of \exp .

(c) By definition of G_ρ , for any direction e_t we can write

$$G_\rho(\alpha_{t-1} + \eta e_t) = \frac{1}{m} \sum_{i=1}^m \exp \left(-y_i \sum_{s=1}^{t-1} \alpha_s h_s(x_i) + \rho \sum_{s=1}^{t-1} \alpha_s \right) e^{-y_i \eta h_t(x_i) + \rho \eta}.$$

Let \mathcal{D}_1 be the uniform distribution, that is $\mathcal{D}_1(i) = \frac{1}{m}$ for all $i \in [m]$ and for any $t \in \{2, \dots, T\}$, define \mathcal{D}_t by

$$\mathcal{D}_t(i) = \frac{\mathcal{D}_{t-1}(i) \exp(-y_i \alpha_{t-1} h_{t-1}(x_i))}{Z_{t-1}},$$

with $Z_{t-1} = \sum_{i=1}^m \mathcal{D}_{t-1}(i) \exp(-y_i \alpha_{t-1} h_{t-1}(x_i))$. Observe that $\mathcal{D}_t(i) = \frac{\exp \left(-y_i \sum_{s=1}^{t-1} \alpha_s h_s(x_i) \right)}{m \prod_{s=1}^{t-1} Z_s}$. Thus,

$$\begin{aligned}& \left. \frac{dG_\rho(\alpha_{t-1} + \eta e_t)}{d\eta} \right|_0 \\ &= \frac{1}{m} \sum_{i=1}^m [-y_i h_t(x_i) + \rho] \exp \left(-y_i \sum_{s=1}^{t-1} \alpha_s h_s(x_i) + \rho \sum_{s=1}^{t-1} \alpha_s \right) \\ &= \sum_{i=1}^m [-y_i h_t(x_i) + \rho] \mathcal{D}_t(i) \left(\prod_{s=1}^{t-1} Z_s \right) e^{\rho \sum_{s=1}^{t-1} \alpha_s} \\ &= \left[\sum_{i=1}^m -y_i h_t(x_i) \mathcal{D}_t(i) + \rho \right] \prod_{s=1}^{t-1} Z_s e^{\rho \sum_{s=1}^{t-1} \alpha_s} \\ &= (-(1 - \epsilon_t) + \epsilon_t + \rho) \left(\prod_{s=1}^{t-1} Z_s \right) e^{\rho \sum_{s=1}^{t-1} \alpha_s} \\ &= (2\epsilon_t - 1 + \rho) \left(\prod_{s=1}^{t-1} Z_s \right) e^{\rho \sum_{s=1}^{t-1} \alpha_s},\end{aligned}$$

where $\epsilon_t = \sum_{i=1}^m \mathcal{D}_t(i) 1_{y_i h_t(x_i) > 0}$. The direction selected is the one minimizing $(2\epsilon_t - 1 + \rho)$ that is ϵ_t . Thus, the algorithm selects at each round the base classifier with the smallest weighted error, as in the case of AdaBoost.

To determine the step η selected at round t , we solve the following equation

$$\begin{aligned}
 \frac{dG_\rho(\alpha_{t-1} + \eta e_t)}{d\eta} &= 0 \\
 \Leftrightarrow \sum_{i=1}^m \mathcal{D}_t(i) Z_t e^{-y_i \eta h_t(x_i) + \rho \eta} [-y_i h_t(x_i) + \rho] &= 0 \\
 \Leftrightarrow \sum_{y_i h_t(x_i) > 0} \mathcal{D}_t(i) Z_t e^{\eta(-1+\rho)} (-1 + \rho) + \sum_{y_i h_t(x_i) < 0} \mathcal{D}_t(i) Z_t e^{\eta(1+\rho)} (1 + \rho) &= 0 \\
 \Leftrightarrow (1 - \epsilon_t) e^{\eta(-1+\rho)} (-1 + \rho) + \epsilon_t e^{\eta(1+\rho)} (1 + \rho) &= 0 \\
 \Leftrightarrow (1 - \epsilon_t) e^{-\eta} (-1 + \rho) + \epsilon_t e^{\eta} (1 + \rho) &= 0 \\
 \Leftrightarrow e^{2\eta} = \frac{(1 - \rho)}{(1 + \rho)} \frac{(1 - \epsilon_t)}{\epsilon_t} \\
 \Leftrightarrow \eta = \frac{1}{2} \log \frac{(1 - \epsilon_t)}{\epsilon_t} - \frac{1}{2} \log \frac{(1 + \rho)}{(1 - \rho)}.
 \end{aligned}$$

Thus, the algorithm differs from AdaBoost only by the choice of the step, which it chooses more conservatively: the step size is smaller than that of AdaBoost by the additive constant $\frac{1}{2} \log \frac{(1+\rho)}{(1-\rho)}$.

- (d) For the coordinate descent algorithm to make progress at each round, the step size selected along the descent direction must be non-negative, that is

$$\begin{aligned}
 \frac{(1 - \rho)}{(1 + \rho)} \frac{(1 - \epsilon_t)}{\epsilon_t} > 1 &\Leftrightarrow (1 - \rho)(1 - \epsilon_t) > \rho \epsilon_t + \epsilon_t \\
 &\Leftrightarrow \epsilon_t < \frac{1 - \rho}{2}.
 \end{aligned}$$

Thus, the error of the base classifier chosen must be at least $\rho/2$ better than one half.

- (e) The normalization factor Z_t can be expressed in terms of ϵ_t and ρ using its definition:

$$\begin{aligned}
 Z_t &= \sum_{i=1}^m \mathcal{D}_t(i) \exp(-y_i \alpha_t h_t(x_i)) \\
 &= e^{-\alpha_t} (1 - \epsilon_t) + e^{\alpha_t} \epsilon_t \\
 &= \sqrt{\frac{1 + \rho}{1 - \rho}} (1 - \epsilon_t) \epsilon_t + \sqrt{\frac{1 - \rho}{1 + \rho}} (1 - \epsilon_t) \epsilon_t \\
 &= \sqrt{\epsilon_t (1 - \epsilon_t)} \left[\sqrt{\frac{1 + \rho}{1 - \rho}} + \sqrt{\frac{1 - \rho}{1 + \rho}} \right] \\
 &= \sqrt{\epsilon_t (1 - \epsilon_t)} \left[\frac{2}{\sqrt{1 - \rho^2}} \right] \\
 &= 2 \sqrt{\frac{\epsilon_t (1 - \epsilon_t)}{1 - \rho^2}}.
 \end{aligned}$$

```

 $\mathcal{A}_\rho(S = ((x_1, y_1), \dots, (x_m, y_m)))$ 
1  for  $i \leftarrow 1$  to  $m$  do
2       $\mathcal{D}_1(i) \leftarrow \frac{1}{m}$ 
3  for  $t \leftarrow 1$  to  $T$  do
4       $h_t \leftarrow$  base classifier in  $H$  with small error  $\epsilon_t = \mathbb{P}_{i \sim \mathcal{D}_t}[h_t(x_i) \neq y_i]$ 
5       $\alpha_t \leftarrow \frac{1}{2} \log \frac{(1-\epsilon_t)}{\epsilon_t} - \frac{1}{2} \log \frac{(1+\rho)}{(1-\rho)}$ 
6       $Z_t \leftarrow 2 \left[ \frac{\epsilon_t(1-\epsilon_t)}{1-\rho^2} \right]^{\frac{1}{2}} \triangleright$  normalization factor
7      for  $i \leftarrow 1$  to  $m$  do
8           $\mathcal{D}_{t+1}(i) \leftarrow \frac{\mathcal{D}_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$ 
9   $f \leftarrow \sum_{t=1}^T \alpha_t h_t$ 
10 return  $h = \text{sgn}(f)$ 

```

Note that \mathcal{A}_0 coincides exactly with AdaBoost.

- i. In view of the definition of \mathcal{D}_t and the bound derived in the first part of this exercise, we can write

$$\begin{aligned}
 \hat{R}_{S,\rho}(f) &\leq \frac{1}{m} \sum_{i=1}^m \exp \left(-y_i \sum_{t=1}^T \alpha_t h_t(x_i) + \rho \sum_{t=1}^T \alpha_t \right) \\
 &= \frac{1}{m} \sum_{i=1}^m \left(m \prod_{t=1}^T Z_t \right) \mathcal{D}_t(i) \exp \left(\rho \sum_{t=1}^T \alpha_t \right) \\
 &= \exp \left(\rho \sum_{t=1}^T \alpha_t \right) \left(\prod_{t=1}^T Z_t \right).
 \end{aligned}$$

- ii. The expression of Z_t was already given above. Plugging in that expression in the bound of the previous question and using the expression of α_t gives

$$\begin{aligned}
 \hat{R}_{S,\rho}(f) &\leq \left(\prod_t e^{\alpha_t} \right)^\rho \left(\prod_{t=1}^T \sqrt{\epsilon_t(1-\epsilon_t)} (u^{\frac{1}{2}} + u^{-\frac{1}{2}}) \right) \\
 &= \left(\sqrt{\frac{1-\rho}{1+\rho}} \right)^{\rho T} \left(\prod_t \sqrt{\frac{1-\epsilon_t}{\epsilon_t}} \right)^\rho \left(\prod_{t=1}^T \sqrt{\epsilon_t(1-\epsilon_t)} (u^{\frac{1}{2}} + u^{-\frac{1}{2}}) \right) \\
 &= \left(u^{\frac{1+\rho}{2}} + u^{-\frac{1-\rho}{2}} \right)^T \prod_{t=1}^T \sqrt{\epsilon_t^{1-\rho} (1-\epsilon_t)^{1+\rho}}.
 \end{aligned}$$

iii. Using the hint and the result of the previous question, we can write

$$\begin{aligned}
 \widehat{R}_{S,\rho}(f) &\leq \prod_t \left(1 - 2 \frac{\left(\frac{1-\rho}{2} - \epsilon_t\right)^2}{1 - \rho^2} \right) \\
 &\leq \prod_t \left(\exp \left(-2 \frac{\left(\frac{1-\rho}{2} - \epsilon_t\right)^2}{1 - \rho^2} \right) \right) \\
 &= \exp \left(-2 \frac{\left(\frac{1-\rho}{2} - \epsilon_t\right)^2}{1 - \rho^2} T \right) \\
 &\leq \exp \left(-\frac{2\gamma^2 T}{1 - \rho^2} \right).
 \end{aligned}$$

Thus, if the upper bound is less than $1/m$, then $\widehat{R}_\rho(f) = 0$ and every training point has margin at least ρ . The inequality $\exp \left(-\frac{2\gamma^2 T}{1 - \rho^2} \right) < 1/m$ is equivalent to $T > \frac{(\log m)(1 - \rho^2)}{2\gamma^2}$.

Chapter 8

8.1 Perceptron lower bound

Let w be the weight vector. Since each update is of the form $w \leftarrow w + y_i x_i$ and since the components of the sample points are integers, the components of w are also integers.

Let $n_1, \dots, n_N \in \mathbb{Z}$ denote the components of w . w correctly classifies all points iff $y_i(w \cdot x_i) > 0$ for $i = 1, \dots, m$, that is,

$$\left\{ \begin{array}{l} n_1 > 0 \\ n_1 - n_2 < 0 \\ -n_1 - n_2 + n_3 > 0 \\ \dots \\ (-1)^N(n_1 + n_2 + \dots + n_{N-1} - n_N) < 0 \end{array} \right. \Leftrightarrow \left\{ \begin{array}{l} n_1 > 0 \\ n_2 > n_1 \\ n_3 > n_1 + n_2 \\ \dots \\ n_N > n_1 + n_2 + \dots + n_{N-1} \end{array} \right.$$

These last inequalities show that the data is linearly separable with $w = (1, 2, \dots, 2^{N-1})$. They also imply that $n_1 \geq 1, n_2 \geq 2, n_3 \geq 4, \dots, n_N \geq 2^{N-1}$. Since each update can at most increment n_N by 1, the number of updates is at least $2^{N-1} = \Omega(2^N)$.

8.2 Generalized mistake bound

The bound is unaffected, as shown by the following, using the same definitions and steps as in this chapter:

$$\begin{aligned}
M\rho &\leq \frac{\mathbf{v} \cdot \sum_{t \in I} y_t \mathbf{x}_t}{\|\mathbf{v}\|} \\
&= \frac{\mathbf{v} \cdot \sum_{t \in I} (\mathbf{w}_{t+1} - \mathbf{w}_t)/\eta}{\|\mathbf{v}\|} \quad (\text{definition of updates}) \\
&= \frac{\mathbf{v} \cdot \mathbf{w}_{T+1}}{\eta \|\mathbf{v}\|} \\
&\leq \|\mathbf{w}_{T+1}\|/\eta \quad (\text{Cauchy-Schwarz ineq.}) \\
&= \|\mathbf{w}_{t_m} + \eta y_{t_m} \mathbf{x}_{t_m}\|/\eta \quad (t_m \text{ largest } t \text{ in } I) \\
&= \left[\|\mathbf{w}_{t_m}\|^2 + \eta^2 \|\mathbf{x}_{t_m}\|^2 + \underbrace{2\eta y_{t_m} \mathbf{w}_{t_m} \cdot \mathbf{x}_{t_m}}_{\leq 0} \right]^{1/2}/\eta \\
&\leq \left[\|\mathbf{w}_{t_m}\|^2 + \eta^2 R^2 \right]^{1/2}/\eta \\
&\leq \left[M\eta^2 R^2 \right]^{1/2}/\eta = \sqrt{M}R. \quad (\text{applying the same to previous } ts \text{ in } I).
\end{aligned}$$

8.3 Sparse instances

Clearly, it takes T updates and leads to $\mathbf{w} = \sum_{t=1}^T y_t \mathbf{x}_t$. Let $\mathbf{u} \in \mathbb{R}^T$ be a vector of norm 1 defining a separating hyperplane, thus $y_t \mathbf{u} \cdot \mathbf{x}_t = y_t u_t \geq 0$ for all $t \in [T]$. To obtain the maximum margin ρ , we seek a vector \mathbf{u} maximizing the minimum of $y_t u_t$ with $y_t u_t \geq 0$ for all t and $\|\mathbf{u}\| = 1$. By symmetry, all $y_t u_t$ s are equal, thus $u_t = y_t/\sqrt{T}$ for all $t \in [T]$ and $\rho = 1/\sqrt{T}$. Thus, Novikoff's bound gives $R^2/\rho^2 = 1/(1/T) = T$.

8.4 Tightness of lower bound

The lower bound is tight. This follows from the tightness of the Khintchine-Kahane inequality, which is the only inequality used in the proof.

8.5 On-line SVM algorithm

First we write the optimization in the following equivalent form:

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i)).$$

Then, using the general update rule in equation (8.22), we get the update rule,

$$\mathbf{w}_{t+1} \leftarrow \begin{cases} \mathbf{w}_t - \eta(\mathbf{w}_t - C y_t \mathbf{x}_t) & \text{if } y_t(\mathbf{w}_t \cdot \mathbf{x}_t) < 1, \\ \mathbf{w}_t - \eta \mathbf{w}_t & \text{if } y_t(\mathbf{w}_t \cdot \mathbf{x}_t) > 1, \\ \mathbf{w}_t & \text{otherwise,} \end{cases}$$

which corresponds exactly to the update in the pseudocode.

8.6 Margin Perceptron

- (a) By assumption, there exists $\mathbf{v} \in \mathbb{R}^N$ such that for all $t \in [T]$, $\rho \leq \frac{y_t(\mathbf{v} \cdot \mathbf{x}_t)}{\|\mathbf{v}\|}$, where ρ is the maximum margin achievable on S . Summing up these inequalities gives

$$\begin{aligned}
M\rho &\leq \frac{\mathbf{v} \cdot \sum_{t \in I} y_t \mathbf{x}_t}{\|\mathbf{v}\|} \leq \left\| \sum_{t \in I} y_t \mathbf{x}_t \right\| && (\text{Cauchy-Schwarz inequality}) \\
&= \left\| \sum_{t \in I} (\mathbf{w}_{t+1} - \mathbf{w}_t) \right\| && (\text{definition of updates}) \\
&= \|\mathbf{w}_{T+1}\| && (\text{telescoping sum, } \mathbf{w}_0 = 0).
\end{aligned}$$

(b) For any $t \in I$, by definition of the update, $\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \mathbf{x}_t$; thus,

$$\begin{aligned} \|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t\|^2 + \|\mathbf{x}_t\|^2 + 2y_t \mathbf{w}_t \cdot \mathbf{x}_t \\ &\leq \|\mathbf{w}_t\|^2 + \|\mathbf{x}_t\|^2 + \|\mathbf{w}_t\|\rho \quad (\text{def. of update condition}) \\ &\leq \|\mathbf{w}_t\|^2 + R^2 + \|\mathbf{w}_t\|\rho + \rho^2/4 \\ &= (\|\mathbf{w}_t\| + \rho/2)^2 + R^2. \end{aligned}$$

(c) In view of the previous result, $\|\mathbf{w}_{t+1}\|^2 - (\|\mathbf{w}_t\| + \rho/2)^2 \leq R^2$, that is

$$\begin{aligned} (\|\mathbf{w}_{t+1}\| - \|\mathbf{w}_t\| - \rho/2)(\|\mathbf{w}_{t+1}\| + \|\mathbf{w}_t\| + \rho/2) &\leq R^2 \\ \implies (\|\mathbf{w}_{t+1}\| - \|\mathbf{w}_t\| - \rho/2) &\leq \frac{R^2}{\|\mathbf{w}_{t+1}\| + \|\mathbf{w}_t\| + \rho/2} \\ \implies \|\mathbf{w}_{t+1}\| &\leq \|\mathbf{w}_t\| + \rho/2 + \frac{R^2}{\|\mathbf{w}_{t+1}\| + \|\mathbf{w}_t\| + \rho/2}. \end{aligned}$$

(d) If $\|\mathbf{w}_t\| \geq \frac{4R^2}{\rho}$ or $\|\mathbf{w}_{t+1}\| \geq \frac{4R^2}{\rho}$, then $\|\mathbf{w}_{t+1}\| + \|\mathbf{w}_t\| + \rho/2 \geq \frac{4R^2}{\rho}$, thus

$$\frac{R^2}{\|\mathbf{w}_{t+1}\| + \|\mathbf{w}_t\| + \rho/2} \leq \frac{R^2}{4R^2/\rho} = \frac{\rho}{4}.$$

In view of this, the inequality of the previous question implies

$$\begin{aligned} \|\mathbf{w}_{t+1}\| &\leq \|\mathbf{w}_t\| + \rho/2 + \frac{R^2}{\|\mathbf{w}_{t+1}\| + \|\mathbf{w}_t\| + \rho/2} \\ \implies \|\mathbf{w}_{t+1}\| &\leq \|\mathbf{w}_t\| + \rho/2 + \frac{\rho}{4} = \|\mathbf{w}_t\| + \frac{3}{4}\rho. \end{aligned}$$

(e) Since $\mathbf{w}_1 = y_1 \mathbf{x}_1$, $\|\mathbf{w}_1\| = \|\mathbf{x}_1\| \leq R$. The margin ρ is at most twice the radius R , thus, $\rho \leq 2R$ and $2R/\rho \geq 1$. This implies that $\|\mathbf{w}_1\| \leq R \leq 2R^2/\rho$. Since $\|\mathbf{w}_1\| \leq 2R^2/\rho$ and $\|\mathbf{w}_{T+1}\| \geq \frac{4R^2}{\rho}$, there must exist at least one update time $t \in I$ at which $\|\mathbf{w}_t\| \leq \frac{4R^2}{\rho}$ and $\|\mathbf{w}_{t+1}\| \geq \frac{4R^2}{\rho}$. The set of such times t is non empty and thus admits a largest element t_0 .

(f) By definition of t_0 , for any $t \geq t_0$, $\|\mathbf{w}_{t+1}\| \geq \frac{4R^2}{\rho}$. Thus, by the inequality of part (d), the following holds for any $t \geq t_0$,

$$\|\mathbf{w}_{t+1}\| \leq \|\mathbf{w}_t\| + \frac{3}{4}\rho.$$

This implies that

$$\begin{aligned} \|\mathbf{w}_{T+1}\| &\leq \|\mathbf{w}_{t_0}\| + |\{t_0, \dots, T\} \cap I| \frac{3}{4}\rho \\ &\leq \|\mathbf{w}_{t_0}\| + M \frac{3}{4}\rho \\ &\leq \frac{4R^2}{\rho} + M \frac{3}{4}\rho. \end{aligned}$$

By the first question $M\rho \leq \|\mathbf{w}_{T+1}\|$; therefore,

$$M\rho \leq \frac{4R^2}{\rho} + M \frac{3}{4}\rho \iff M\rho/4 \leq 4R^2/\rho \iff M \leq 16R^2/\rho^2.$$

8.7 Generalized RWM

(a) Observe that:

$$\begin{aligned} W_{t+1} &= \sum_{i=1}^N (1 - (1 - \beta)l_{t,i})w_{t,i} \\ &= W_t - W_t(1 - \beta)l_{t,i}w_{t,i}/W_t = W_t(1 - (1 - \beta)L_t). \end{aligned}$$

Thus, $W_{T+1} = N \prod_{t=1}^T (1 - (1 - \beta)L_t)$ and

$$\begin{aligned} \log W_{T+1} &= \log N + \sum_{t=1}^T \log(1 - (1 - \beta)L_t) \\ &\leq \log N + \sum_{t=1}^T -(1 - \beta)L_t \\ &= \log N - (1 - \beta)\mathcal{L}_T. \end{aligned}$$

(b) For all $i \in [N]$,

$$\begin{aligned} \log W_{T+1} &\geq \log w_{T+1,i} \\ &= \log \left(\prod_{t=1}^T (1 - (1 - \beta)l_{t,i}) \right) \\ &= \sum_{t=1}^T \log(1 - (1 - \beta)l_{t,i}) \\ &\geq \sum_{t=1}^T -(1 - \beta)l_{t,i} - (1 - \beta)^2 l_{t,i}^2 \\ &= -(1 - \beta)\mathcal{L}_{T,i} - (1 - \beta)^2 \sum_{t=1}^T l_{t,i}^2. \end{aligned}$$

(c) Comparing the lower and upper bounds gives:

$$\begin{aligned} -(1 - \beta)\mathcal{L}_{T,i} - (1 - \beta)^2 \sum_{t=1}^T l_{t,i}^2 &\leq \log N - (1 - \beta)\mathcal{L}_T \\ \implies \mathcal{L}_T &\leq \mathcal{L}_{T,i} + \frac{\log N}{(1 - \beta)} + (1 - \beta) \sum_{t=1}^T l_{t,i}^2. \end{aligned}$$

Clearly, for any $i \in [N]$, $\sum_{t=1}^T l_{t,i}^2 \leq T$. Thus, for all $i \in [N]$,

$$\mathcal{L}_T \leq \mathcal{L}_{T,i} + \frac{\log N}{(1 - \beta)} + (1 - \beta)T,$$

in particular, $\mathcal{L}_T \leq \mathcal{L}_T^{\min} + \frac{\log N}{(1 - \beta)} + (1 - \beta)T$. Differentiating with respect to β and setting the result to zero gives $\frac{\log N}{(1 - \beta)^2} - T = 0$, as in the case of the RWM algorithm. Thus, for $\beta = \max\{1/2, 1 - \sqrt{(\log N)/T}\}$, $\mathcal{L}_T \leq \mathcal{L}_T^{\min} + 2\sqrt{T \log N}$, that is $R_T \leq 2\sqrt{T \log N}$.

8.9 General inequality

(a) We analyze the function $f(x) = \log(1 - x) + x + \frac{x^2}{\alpha}$ and show that it is positive for $x \in [0, 1 - \frac{\alpha}{2}]$. First note that $f(0) = 0$, then note that

$$\begin{aligned} f'(x) \geq 0 &\iff -1 + 1 - x \frac{2}{\alpha} x(1 - x) \geq 0 \\ &\iff \frac{2}{\alpha} x(1 - x) \geq x \\ &\iff \frac{2}{\alpha} (1 - x) \geq 1 \quad (\text{for } x \geq 0) \\ &\iff x \leq 1 - \frac{\alpha}{2}. \end{aligned}$$

Thus, the derivative is only increasing for $x \in [1, 1 - \frac{\alpha}{2}]$, which implies that the function is positive for the same interval.

In order to apply the inequality, inequality the valid range of β is $[\frac{2}{\alpha}, 1)$.

- (b) The bound follows directly using the same steps as in the original proof, but with the general inequality. The optimal choice of β is $\max\{\frac{\alpha}{2}, 1 - \sqrt{\alpha(\log N)/T}\}$, which gives

$$R_T \leq \sqrt{\frac{\log(N)T}{\alpha}} + \sqrt{\alpha \log(N)T}.$$

- (c) Setting α close to 2 forces β close to 1, which results in an algorithm that downweights experts in a very conservative fashion. From the bound in part (b) we see that $\alpha = 1$, as is used in the chapter, is the optimal choice.

8.10 On-line to batch — non-convex loss

- (a) We use the following series of inequalities:

$$\begin{aligned} & \min_{i \in [T]} (R(h_i) + 2c_\delta(T - i + 1)) \\ & \leq \frac{1}{T} \sum_{i=1}^T (R(h_i) + 2c_\delta(T - i + 1)) \\ & = \frac{1}{T} \sum_{i=1}^T R(h_{i-1}) + \frac{2}{T} \sum_{i=0}^{T-1} \sqrt{\frac{1}{2(T-i)} \log \frac{T(T+1)}{\delta}} \\ & < \frac{1}{T} \sum_{i=1}^T R(h_{i-1}) + \frac{2}{T} \sum_{i=0}^{T-1} \sqrt{\frac{1}{2(T-i)} \log \left(\frac{(T+1)}{\delta} \right)^2} \\ & = \frac{1}{T} \sum_{i=1}^T R(h_{i-1}) + \frac{2}{T} \sum_{i=0}^{T-1} \sqrt{\frac{1}{(T-i)} \log \frac{(T+1)}{\delta}} \\ & \leq \frac{1}{T} \sum_{i=1}^T R(h_{i-1}) + 4\sqrt{\frac{1}{T} \log \frac{T+1}{\delta}}. \end{aligned}$$

The first inequality follows, since the minimum is always less than or equal to the average and the final inequality follows from $\sum_{i=0}^{T-1} \sqrt{1/(T-i)} = \sum_{i=1}^T \sqrt{1/i} \leq 2\sqrt{T}$.

- (b) Coupling the inequality of part (a) with the high probability statement of lemma 8.14 to bound $\frac{1}{T} \sum_{i=1}^T R(h_i)$ shows the desired bound.
- (c) The square-root terms in part (b) can be bounded further by $6\sqrt{\frac{1}{T} \log \frac{2(T+1)}{\delta}}$.

Now, note that for two events A and B that each occur with probability at least $1 - \delta$,

$$\begin{aligned} \mathbb{P}[\neg A \cup \neg B] & \leq \mathbb{P}[\neg A] + \mathbb{P}[\neg B] \leq 2\delta \\ \iff \mathbb{P}[A \wedge B] & \geq 1 - 2\delta. \end{aligned}$$

Thus, the probability that both bounds in (b) and (c) hold simultaneously is at least $1 - 2\delta$; substituting δ with $\delta/2$ everywhere completes the bound.

8.11 On-line to batch — kernel Perceptron margin bound

- (a) Let $\mathbf{u} \in \mathbb{H}$ with $\|\mathbf{u}\| = 1$. Observe that for any $t \in [T]$, we can write

$$1 - \left(1 - \frac{y_t(\mathbf{u} \cdot \Phi(x_t))}{\rho}\right)_+ \leq \frac{y_t(\mathbf{u} \cdot \Phi(x_t))}{\rho}.$$

Let \mathcal{J} be the set of $t \in [T]$ at which kernel Perceptron makes an update, and let M_T be the total number of updates made, then, summing up the previous inequalities over all such ts and using the Cauchy-Schwarz inequality yields

$$M_T - \sum_{t \in \mathcal{J}} \left(1 - \frac{y_t(\mathbf{u} \cdot \Phi(x_t))}{\rho}\right) \leq \sum_{t \in \mathcal{J}} \frac{y_t(\mathbf{u} \cdot \Phi(x_t))}{\rho} \leq \frac{\|\sum_{t \in \mathcal{J}} y_t \Phi(x_t)\|}{\rho}.$$

In view of the proof for separable case of the perceptron algorithm (theorem 8.8), the right-hand side can be bounded by $\frac{\sqrt{\sum_{t \in \mathcal{J}} K(x_t, x_t)}}{\rho}$. Thus, for any $\mathbf{u} \in \mathbb{H}$ with $\|\mathbf{u}\| = 1$,

$$M_T \leq \sum_{t \in \mathcal{J}} \left(1 - \frac{y_t(\mathbf{u} \cdot \Phi(x_t))}{\rho} \right)_+ + \frac{\sqrt{\sum_{t \in \mathcal{J}} K(x_t, x_t)}}{\rho}.$$

- (b) Plugging in the result from part (a) into (8.31) gives the desired bound. This bound is with respect to expected error of best in class, while corollary 6.13 is with respect to empirical error of selected hypothesis.

Chapter 9

9.5 Decision trees. A binary decision tree with n nodes has exactly $n+1$ leaves. Each node can be labeled with an integer from $\{1, \dots, N\}$ indicating which dimension is queried to make a binary split and each leaf can be labeled with ± 1 to indicate the classification made at that leaf. Fix an ordering of the nodes and leaves and consider all possible labelings of this sequence. There can be no more than $(N+2)^{2n+1}$ distinct binary trees and, thus, the VC-dimension of this finite set of hypotheses can be no larger than $(2n+1) \log(N+2) = O(n \log N)$.

Chapter 11

11.1 Pseudo-dimension and monotonic functions

If for some $m > 0$, there exists (t_1, \dots, t_m) and a set of points (x_1, \dots, x_m) that \mathcal{H} shatters, then $\phi \circ \mathcal{H}$ can also shatter it. To see that, note that if for some $h \in \mathcal{H}$,

$$h(x_i) \geq t_i,$$

then by the monotonic property of ϕ ,

$$\phi(h(x_i)) \geq \phi(t_i).$$

A similar argument holds for the case $h(x_i) < t_i$. Thus, $\phi \circ \mathcal{H}$ can shatter the set of points (x_1, \dots, x_m) with thresholds $(\phi(t_1), \dots, \phi(t_m))$, and this proves that $\text{Pdim}(\phi \circ \mathcal{H}) \geq \text{Pdim}(\mathcal{H})$.

Since ϕ is strictly monotonic, it is invertible, and a similar argument with ϕ^{-1} can be used to show $\text{Pdim}(\mathcal{H}) \geq \text{Pdim}(\phi \circ \mathcal{H})$.

11.2 Pseudo-dimension of linear functions

From equation (11.3) we have that

$$\text{Pdim}(\mathcal{H}) = \text{VCdim} \left(\left\{ (\mathbf{x}, t) \mapsto 1_{(\mathbf{w}^\top \mathbf{x} - t) > 0} : h \in \mathcal{H} \right\} \right).$$

Note that $1_{(\mathbf{w}^\top \mathbf{x} - t) > 0} = \frac{1 + \text{sgn}(\mathbf{w}^\top \mathbf{x} - t)}{2}$ is a linear separator with *fixed* offset. It is easy to show that the VC-dimension of such a hypothesis class is d (as oppose to $d+1$ in the case of linear separators with a free offset parameter).

11.3 Linear regression

- (a) In order for the matrix $\mathbf{X}\mathbf{X}^\top$ to be invertible, we need the number of (linearly independent) examples to outnumber the number of features used to represent each example.
- (b) For any $\mathbf{v} \in \mathbb{R}^m$ we can choose $\mathbf{w} = (\mathbf{X}\mathbf{X}^\top)^\dagger \mathbf{X}\mathbf{y} + (\mathbf{I} - (\mathbf{X}^\dagger)^\top \mathbf{X}^\top) \mathbf{v}$. To see this, observe that

$$\begin{aligned} \mathbf{X}^\top (\mathbf{I} - (\mathbf{X}^\dagger)^\top \mathbf{X}^\top) &= \mathbf{X}^\top - \mathbf{X}^\top (\mathbf{X}^\dagger)^\top \mathbf{X}^\top \\ &= \mathbf{X}^\top - (\mathbf{X}\mathbf{X}^\dagger \mathbf{X})^\top \\ &= \mathbf{X}^\top - \mathbf{X}^\top = \mathbf{0}, \end{aligned}$$

and hence, we have $\mathbf{X}^\top \mathbf{w} = \mathbf{X}^\dagger \mathbf{X} \mathbf{y}$.

11.4 Perturbed kernels

- (a) Using the closed form solutions $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$ and the fact $\mathbf{M}'^{-1} - \mathbf{M}^{-1} = -\mathbf{M}'^{-1}(\mathbf{M}' - \mathbf{M})\mathbf{M}^{-1}$ (this can be verified by simply expanding the right-hand side), we have

$$\begin{aligned} \boldsymbol{\alpha}' - \boldsymbol{\alpha} &= ((\mathbf{K}' + \lambda \mathbf{I})^{-1} - (\mathbf{K} + \lambda \mathbf{I})^{-1}) \mathbf{y} \\ &= ((\mathbf{K}' + \lambda \mathbf{I})^{-1} (\mathbf{K}' + \lambda \mathbf{I} - \mathbf{K} - \lambda \mathbf{I}) (\mathbf{K} + \lambda \mathbf{I})^{-1}) \mathbf{y} \\ &= ((\mathbf{K}' + \lambda \mathbf{I})^{-1} (\mathbf{K}' - \mathbf{K}) (\mathbf{K} + \lambda \mathbf{I})^{-1}) \mathbf{y}. \end{aligned}$$

- (b) Using the fact that for any vector \mathbf{v} and matrix \mathbf{A} , $\|\mathbf{A}\mathbf{v}\|^2 = \mathbf{v}^\top \mathbf{A}^\top \mathbf{A} \mathbf{v} \leq \|\mathbf{v}\|^2 \|\mathbf{A}^\top \mathbf{A}\| = \|\mathbf{v}\|^2 \|\mathbf{A}\|^2$ we have

$$\|\boldsymbol{\alpha}' - \boldsymbol{\alpha}\| \leq \|(\mathbf{K}' + \lambda \mathbf{I})^{-1}\| \|\mathbf{K}' - \mathbf{K}\| \|(\mathbf{K} + \lambda \mathbf{I})^{-1}\| \|\mathbf{y}\|.$$

Since $\|\mathbf{y}\| \leq M$, we have $\|\mathbf{y}\| \leq \sqrt{m}M$ and we can use the observation $\|(\mathbf{K} + \lambda \mathbf{I})^{-1}\| = \lambda_{\max}((\mathbf{K} + \lambda \mathbf{I})^{-1}) = \lambda_{\min}(\mathbf{K} + \lambda \mathbf{I})^{-1} \leq 1/\lambda$, where $\lambda_{\max}(A)$ and $\lambda_{\min}(A)$ are the maximum and minimum eigenvalues of A , respectively.

11.5 Huber loss

The primal function can be written as:

$$\begin{aligned} \min_{\mathbf{w}, b} \quad & \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (L_c(\xi_i) + L_c(\xi'_i)) \\ \text{s.t.} \quad & \mathbf{w} \cdot \Phi(\mathbf{x}_i) + b - y_i \leq \xi_i, \quad \forall i \in [m] \\ & y_i - \mathbf{w} \cdot \Phi(\mathbf{x}_i) - b \leq \xi'_i, \quad \forall i \in [m] \\ & \xi_i, \xi'_i \geq 0, \quad \forall i \in [m] \end{aligned}$$

The Lagrangian is written as follows,

$$\begin{aligned} L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\xi}', \boldsymbol{\alpha}, \boldsymbol{\alpha}', \beta, \beta') &= \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (L_c(\xi_i) + L_c(\xi'_i)) \\ &+ \sum_{i=1}^m \left(\alpha_i (\mathbf{w} \cdot \Phi(\mathbf{x}_i) + b - y_i - \xi_i) + \alpha'_i (y_i - \mathbf{w} \cdot \Phi(\mathbf{x}_i) - b - \xi'_i) \right) \\ &- \sum_{i=1}^m (\beta_i \xi_i + \beta'_i \xi'_i), \end{aligned}$$

and the associated KKT conditions are:

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} + \sum_{i=1}^m (\alpha_i - \alpha'_i) \Phi(\mathbf{x}_i) = 0,$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m (\alpha_i - \alpha'_i) = 0,$$

$$\frac{\partial L}{\partial \xi_i} = 0 \iff \begin{cases} C\xi_i = \alpha_i + \beta_i, & \text{if } \xi_i \leq c \\ Cc = \alpha_i + \beta_i, & \text{otherwise} \end{cases}, \quad (\text{E.134})$$

$$\frac{\partial L}{\partial \xi'_i} = 0 \iff \begin{cases} C\xi_i = \alpha'_i + \beta'_i, & \text{if } \xi_i \leq c \\ Cc = \alpha'_i + \beta'_i, & \text{otherwise} \end{cases}, \quad (\text{E.135})$$

$$\beta_i \xi_i = 0, \forall i \in [m], \quad (\text{E.136})$$

$$\beta'_i \xi'_i = 0, \forall i \in [m]. \quad (\text{E.137})$$

The first two conditions are the same as in SVR the standard ϵ -insensitive loss, and using them to simplify the Lagrangian gives several familiar terms. The novel conditions involve ξ and ξ' . Collecting all terms in the Lagrangian that depend on ξ_i , we have the following

equality (regardless of which condition holds in (E.134))

$$CL_c(\xi_i) - \xi_i(\alpha_i + \beta_i) = -\frac{(\alpha_i + \beta_i)^2}{2C}. \quad (\text{E.138})$$

If it is the case that $\xi_i > 0$, then $\beta_i = 0$ by condition (E.136). In the case $\xi_i = 0$, then $\alpha_i + \beta_i = 0$ by the first case in condition (E.134). However, this also implies $\alpha_i = \beta_i = 0$, since both dual variables are constrained to be positive. Thus, in either case, we can simplify (E.138) to $-\frac{\alpha_i^2}{2C}$. Similar arguments can be used to simplify the ξ'_i terms, resulting in the final dual problem:

$$\begin{aligned} \max_{\boldsymbol{\alpha}} \quad & \mathbf{y}(\boldsymbol{\alpha}' - \boldsymbol{\alpha}) - \frac{1}{2}((\boldsymbol{\alpha}' - \boldsymbol{\alpha})^\top \mathbf{K}(\boldsymbol{\alpha}' - \boldsymbol{\alpha}) + \frac{1}{C}\boldsymbol{\alpha}^\top \mathbf{1} + \frac{1}{C}\boldsymbol{\alpha}'^\top \mathbf{1}) \\ \text{s.t.} \quad & (\boldsymbol{\alpha}' - \boldsymbol{\alpha})^\top \mathbf{1} = 0 \\ & 0 \leq \alpha_i, \alpha'_i \leq cC, \forall i \in [m]. \end{aligned}$$

11.8 Optimal kernel matrix

- (a) Using the closed-form solution for the inner maximization problem $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}$, simplifies the joint optimization to a simpler minimization:

$$\min_{\mathbf{K} \succeq 0} \mathbf{y}^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}, \quad \text{s.t.} \quad \|\mathbf{K}\|_2 \leq 1.$$

Note that for any invertible matrix \mathbf{A} , $\mathbf{y}^\top \mathbf{A}^{-1} \mathbf{y} \geq \|\mathbf{y}\|^2 \lambda_{\min}(\mathbf{A}^{-1}) = \|\mathbf{y}\|^2 \lambda_{\max}(\mathbf{A})^{-1}$. Thus, it is easy to see that $\min_{\mathbf{K} \succeq 0} \mathbf{y}^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y} \geq \frac{\|\mathbf{y}\|^2}{1+\lambda}$ since $\|\mathbf{K}\|_2 = \lambda_{\max}(\mathbf{K}) \leq 1$.

We now show $\mathbf{K} = \frac{1}{\|\mathbf{y}\|^2} \mathbf{y} \mathbf{y}^\top$ achieves this lower bound. First, note that $(\frac{1}{\|\mathbf{y}\|^2} \mathbf{y} \mathbf{y}^\top + \lambda \mathbf{I}) \mathbf{y} = (1 + \lambda) \mathbf{y}$, so \mathbf{y} is an eigenvector of the matrix with eigenvalue $(1 + \lambda)$. Since the matrix is invertible, it can be shown that \mathbf{y} is also an eigenvector of $(\frac{1}{\|\mathbf{y}\|^2} \mathbf{y} \mathbf{y}^\top + \lambda \mathbf{I})^{-1}$ with eigenvalue $\frac{1}{1+\lambda}$ (for example, consider the eigen decomposition of the matrix).

- (b) The kernel matrix alone is not useful for classifying future unseen points x , which requires computing $\sum_{i=1}^m K(x_i, x)$ and needs access to an underlying kernel *function* that is consistent with the kernel matrix. Finding such a kernel function may be difficult in general, and furthermore the choice of function may not be unique.

11.9 Leave-one-out error

- (a) Note that the hypothesis $h_{S'_i}$ will make zero error on the i th point of S'_i and is defined as the minimizer with respect to the remainder of the points. Thus, $h_{S'_i}$ is also the minimizer with respect to the set S'_i .
- (b) Using part (a) and the definition of the KRR hypothesis with respect to the dual variables we have $h_{S'_i}(x_i) = h_{S'_i}(x_i) = \boldsymbol{\alpha}_{S'_i}^\top \mathbf{K} \mathbf{e}_i$, where $\boldsymbol{\alpha}_{S'_i}$ is the optimal set of dual variable for KRR trained with S'_i . Noting that the closed-form solution is $\boldsymbol{\alpha} = (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{y}_i$ proves the equality.

- (c) Using part (b) we can write

$$\begin{aligned} h_{S'_i}(x_i) - y_i &= \mathbf{y}_i^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K} \mathbf{e}_i - y_i \\ &= (\mathbf{y} - y_i \mathbf{e}_i + h_{S'_i}(x_i) \mathbf{e}_i)^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K} \mathbf{e}_i - y_i \\ &= h_S(x_i) - y_i + (h_{S'_i}(x_i) - y_i) \mathbf{e}_i^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K} \mathbf{e}_i, \end{aligned}$$

which implies $h_{S'_i}(x_i) - y_i = (h_S(x_i) - y_i) / (\mathbf{e}_i^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K} \mathbf{e}_i)$. Thus, we can write

$$\hat{R}_{\text{LOO}}(\mathcal{A}) = \frac{1}{m} \sum_{i=1}^m \left(\frac{h_S(x_i) - y_i}{\mathbf{e}_i^\top (\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K} \mathbf{e}_i} \right)^2.$$

- (d) In this case the two losses differ only by the factor $\frac{1}{\gamma^2}$. Thus, if $\gamma = \sqrt{m}$, the two performance measures coincide.

Chapter 14

14.1 Tighter stability bounds

- (a) No, even as $\beta \rightarrow 0$ the generalization bound of theorem 14.2 only guarantees $R(h_S) - \widehat{R}_S(h_S) \leq M \sqrt{\frac{\log \frac{1}{\delta}}{2m}} = O(1/\sqrt{m})$.
- (b) In this case, $M = C/\sqrt{m}$ and $M \sqrt{\frac{\log \frac{1}{\delta}}{2m}} = O(1/m)$; thus, it would suffice to have $\beta = O(1/m^{3/2})$ in order to guarantee an $O(1/m)$ generalization bound.

14.2 Quadratic hinge loss stability

We first show that the loss function is σ -admissible. Consider three cases:

- Both $h(x)$ and $h'(x)$ are correct with margin greater than 1, then

$$|L(h(x), y) - L(h'(x), y)| = 0.$$

- Only one hypothesis is correct with large enough margin. Without loss of generality assume $h'(x)$ is correct, then

$$\begin{aligned} |L(h(x), y) - L(h'(x), y)| &= (1 - h(x)y)^2 \\ &\leq ((1 - h(x)y) - (1 - h'(x)y))^2 = (h'(x) - h(x))^2 \\ &\leq 4\sqrt{M}|h(x) - h'(x)|. \end{aligned}$$

The first inequality follows from the assumption $1 - h'(x)y \leq 0$, and the second inequality follows from the bounded loss assumption, which implies $\forall h \in \mathcal{H}, |h(x)| \leq \sqrt{M} + 1 \leq 2M$.

- Finally, we consider the case where both $h(x)$ and $h'(x)$ incur a loss. Without loss of generality assume $(1 - h(x)y) \geq (1 - h'(x)y)$, then

$$\begin{aligned} |L(h(x), y) - L(h'(x), y)| &= (1 - h(x)y)^2 - (1 - h'(x)y)^2 \\ &= ((1 - h(x)y) + (1 - h'(x)y))((1 - h(x)y) - (1 - h'(x)y))) \\ &\leq |2 - y(h(x) + h'(x))||y(h(x) - h'(x))| \leq 6\sqrt{M}|h(x) - h'(x)|. \end{aligned}$$

Thus, the quadratic hinge loss is σ -admissible with $\sigma = 6\sqrt{M}$. By proposition 14.4, SVM with quadratic hinge loss is stable with $\beta = \frac{36r^2M}{m\lambda}$, and using theorem 14.2 gives the following bound:

$$R(h_S) \leq \widehat{R}_S(h_S) + \frac{36r^2M}{m\lambda} + \left(\frac{72r^2M}{\lambda} + M\right) \sqrt{\frac{\log \frac{1}{\delta}}{m}}.$$

14.3 Stability of linear regression

- (a) Observing corollary 14.6, if $\lambda \rightarrow 0$ then the bound becomes vacuous. No generalization is guaranteed.
- (b) Observe the following simple counter example in one dimension with $|x| \leq 1$ and $|y| \leq 1$. Let $S = ((0, 0), (1, 1))$ define one training sample of (x, y) pairs, so the optimal linear predictor is $h(x) = x$. Then we can define $S' = ((0, 0), (\alpha, 1))$ for $0 < \alpha \leq 1$ and the optimal linear predictor is then $h'(x) = \frac{1}{\alpha}x$. Thus, the two predictors can vary by an arbitrary amount as $\alpha \rightarrow 0$, which shows linear regression is not stable.

14.4 Kernel stability

First we note,

$$\begin{aligned} |h'(x) - h(x)| &= |(\alpha'^\top - \alpha^\top)\mathbf{k}_x| \\ &\leq \|\alpha' - \alpha\| \|\mathbf{k}_x\|. \end{aligned}$$

Using the bound on K , we have $\|\mathbf{k}_x\| \leq \sqrt{mr}$, and using the result from exercise 10.3 completes the bound.

14.5 Stability of relative-entropy regularization

- (a) This result follows from the property $|\int g(x) dx| \leq \int |g(x)| dx$ for integrable function g and the bound on L :

$$\begin{aligned} |H(g, z) - H(g', z)| &\leq \left| \int_{\Theta} L(h_{\theta}(x), y)(g(\theta) - g'(\theta)) d\theta \right| \\ &\leq \int_{\Theta} |L(h_{\theta}(x), y)(g(\theta) - g'(\theta))| d\theta \leq M \int_{\Theta} |g(\theta) - g'(\theta)| d\theta. \end{aligned}$$

- (b) i. We can write

$$\begin{aligned} \left(\int_{\Theta} |g(\theta) - g'(\theta)| d\theta \right)^2 &= \frac{1}{2} \left(\int_{\Theta} |g(\theta) - g'(\theta)| d\theta \right)^2 + \frac{1}{2} \left(\int_{\Theta} |g'(\theta) - g(\theta)| d\theta \right)^2 \\ &\leq K(g, g') + K(g', g). \end{aligned}$$

Thus, it suffices to show $K(g, g') = B_{K(\cdot, f_0)}(g \| g') + \int_{\Theta} g(\theta) - g'(\theta) d\theta$ in order to show the result. Note that $[\nabla_{g'} K(g', f_0)]_{\theta} = \log \frac{g'(\theta)}{f_0(\theta)} + 1$, thus

$$\begin{aligned} B_{K(\cdot, f_0)}(g \| g') &= K(g, f_0) - K(g', f_0) - \langle g - g', \nabla_{g'} K(g', f_0) \rangle \\ &= K(g, f_0) - K(g', f_0) - \int_{\Theta} (g(\theta) - g'(\theta)) \left(\log \frac{g'(\theta)}{f_0(\theta)} + 1 \right) d\theta \\ &= \int_{\Theta} g(\theta) \left(\log \frac{g(\theta)}{f_0(\theta)} - \frac{g'(\theta)}{f_0(\theta)} \right) - g(\theta) + g'(\theta) d\theta \\ &= \int_{\Theta} g(\theta) \left(\log \frac{g(\theta)}{g'(\theta)} \right) - g(\theta) + g'(\theta) d\theta \\ &= K(g, g') + \int_{\Theta} g'(\theta) - g(\theta) d\theta. \end{aligned}$$

- ii. To show the first inequality, note that $B_{F_S} = B_{\hat{R}_S} + \lambda B_{K(\cdot, f_0)}$, and, since Bregman divergences are positive, it implies $B_{K(\cdot, f_0)} \leq B_{F_S}$. Then, using the definition of g and g' as minimizers,

$$\begin{aligned} B_{K(\cdot, f_0)}(g \| g') + B_{K(\cdot, f_0)}(g' \| g) &\leq \frac{1}{\lambda} B_{F_S'}(g \| g') + B_{F_S}(g' \| g) \\ &= \frac{1}{\lambda} \hat{R}_{S'}(g) - \hat{R}_{S'}(g') + \hat{R}_S(g') - \hat{R}_S(g) \\ &= \frac{1}{m\lambda} |H(g', z_m) - H(g, z_m) + H(g, z'_m) - H(g', z'_m)|. \end{aligned}$$

Where the last equality follows from the fact that only the m th point differs in S and S' . Finally, using the M -Lipschitz property from part (a) completes the question.

- iii. Combining part (i) and (ii), we have

$$\begin{aligned} \left(\int_{\Theta} |g(\theta) - g'(\theta)| d\theta \right)^2 &\leq \frac{2M}{m\lambda} \int_{\Theta} |g(\theta) - g'(\theta)| d\theta \\ \iff \int_{\Theta} |g(\theta) - g'(\theta)| d\theta &\leq \frac{2M}{m\lambda} \end{aligned}$$

Thus, using the M -Lipschitz property, we have

$$|H(g, z) - H(g', z)| \leq \frac{2M^2}{m\lambda}.$$

Chapter 15

15.1 PCA and maximal variance

(a)

$$\begin{aligned}
\text{variance} &= \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i^\top \mathbf{u} - \bar{\mathbf{x}}^\top \mathbf{u})^2 = \frac{1}{m} \sum_{i=1}^m ((\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{u})^2 \\
&= \frac{1}{m} \sum_{i=1}^m ((\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{u})^\top ((\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{u}) \\
&= \frac{1}{m} \sum_{i=1}^m \mathbf{u}^\top (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{u} = \mathbf{u}^\top \mathbf{C} \mathbf{u}.
\end{aligned}$$

- (b) We want to maximize $\mathbf{u}^\top \mathbf{C} \mathbf{u}$ subject to $\mathbf{u}^\top \mathbf{u} = 1$. We can directly apply the Rayleigh quotient (section A.2.3) to show that \mathbf{u}^* is the largest eigenvector (or equivalently the largest singular vector) of \mathbf{C} , which proves that PCA with $k = 1$ projects the data onto the direction of maximal variance. Alternatively, we can prove this by introducing a Lagrange multiplier α to enforce the constraint $\mathbf{u}^\top \mathbf{u} = 1$. To maximize the Lagrangian $\mathbf{u}^\top \mathbf{C} \mathbf{u} + \alpha(1 - \mathbf{u}^\top \mathbf{u})$, we set the derivative w.r.t. \mathbf{u} to $\mathbf{0}$, which gives us $\mathbf{C} \mathbf{u} = \alpha \mathbf{u}$. This equality implies that \mathbf{u}^* is an eigenvector of \mathbf{C} . Finally, we observe that among all eigenvectors of \mathbf{C} , the eigenvector associated with the largest eigenvalue of \mathbf{C} maximizes $\mathbf{u}^\top \mathbf{C} \mathbf{u}$.

15.2 Double centering

- (a) Observe that $\|\mathbf{x}_i - \mathbf{x}_j\|^2 = (\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j) = \mathbf{x}_i^\top \mathbf{x}_i + \mathbf{x}_j^\top \mathbf{x}_j - 2\mathbf{x}_i^\top \mathbf{x}_j$ and rearrange terms.
- (b) Noting that $\mathbf{X}^* = \mathbf{X} - \frac{1}{m} \mathbf{X} \mathbf{1} \mathbf{1}^\top$ and plugging into the equation $\mathbf{K}^* = \mathbf{X}^{*\top} \mathbf{X}^*$ yields the result.
- (c) Note that the scalar form of the equation in (b) is

$$\mathbf{K}_{ij}^* = \mathbf{K}_{ij} - \frac{1}{m} \sum_{k=1}^m \mathbf{K}_{ik} - \frac{1}{m} \sum_{k=1}^m \mathbf{K}_{kj} + \frac{1}{m^2} \sum_k \sum_l \mathbf{K}_{k,l}.$$

Substituting with the equation $\mathbf{D}_{ij}^2 = \mathbf{K}_{ii} + \mathbf{K}_{jj} - 2\mathbf{K}_{ij}$ from (a) and simplifying yields the result.

- (d) We first observe that $-\frac{1}{2} \mathbf{H} \mathbf{D} \mathbf{H} = -\frac{1}{2} (\mathbf{D} - \frac{1}{m} \mathbf{D} \mathbf{1} \mathbf{1}^\top - \frac{1}{m} \mathbf{1} \mathbf{1}^\top \mathbf{D} + \frac{1}{m^2} \mathbf{1} \mathbf{1}^\top \mathbf{D} \mathbf{1} \mathbf{1}^\top)$. By inspection, the matrix expression on the RHS corresponds to the scalar expression with four terms on the RHS of the equation in (c).

15.3 Laplacian eigenmaps

$$\begin{aligned}
\operatorname{argmin}_{\mathbf{y}'} \sum_{i,j} \mathbf{W}_{ij} \|y'_i - y'_j\|_2^2 &= \operatorname{argmin}_{\mathbf{y}'} \sum_{i,j} \mathbf{W}_{ij} (y'^2_i + y'^2_j - 2y'_i y'_j) \\
&= \operatorname{argmin}_{\mathbf{y}'} 2 \sum_i \mathbf{D}_{ii} y'^2_i - 2 \sum_{i,j} \mathbf{W}_{ij} y'_i y'_j \\
&= \operatorname{argmin}_{\mathbf{y}'} \mathbf{y}'^\top \mathbf{D} \mathbf{y}' - \mathbf{y}'^\top \mathbf{W} \mathbf{y}' \\
&= \operatorname{argmin}_{\mathbf{y}'} \mathbf{y}'^\top \mathbf{L} \mathbf{y}'.
\end{aligned}$$

15.4 Nyström method

- (a) For the first part of question, note that \mathbf{W} is SPSD if $\mathbf{x}^\top \mathbf{W} \mathbf{x} \geq 0$ for all $\mathbf{x} \in \mathbb{R}^l$. This condition is equivalent to $\mathbf{y}^\top \mathbf{K} \mathbf{y} \geq 0$ for all $\mathbf{y} \in \mathbb{R}^m$ where $y_i = 0$ for $l+1 \leq i \leq m$. Since \mathbf{K} is SPSD by assumption, this latter condition holds. For the second part, we write $\tilde{\mathbf{K}}$

in block form as

$$\tilde{\mathbf{K}} = \begin{bmatrix} \mathbf{W} \\ \mathbf{K}_{21} \end{bmatrix} \mathbf{W}^\dagger \begin{bmatrix} \mathbf{W} & \mathbf{K}_{21}^\top \end{bmatrix} = \begin{bmatrix} \mathbf{W} & \mathbf{K}_{21}^\top \\ \mathbf{K}_{21} & \mathbf{K}_{21} \mathbf{W}^\dagger \mathbf{K}_{21}^\top \end{bmatrix}.$$

Comparison with the block form of \mathbf{K} then immediately yields the desired result.

- (b) Observe that $\mathbf{C} = \mathbf{X}^\top \mathbf{X}'$ and $\mathbf{W} = \mathbf{X}'^\top \mathbf{X}'$. Thus,

$$\tilde{\mathbf{K}} = \mathbf{C} \mathbf{W}^\dagger \mathbf{C}^\top = \mathbf{X}^\top \mathbf{X}' (\mathbf{X}'^\top \mathbf{X}')^\dagger \mathbf{X}'^\top \mathbf{X} = \mathbf{X}^\top \mathbf{U}_{\mathbf{X}'} \mathbf{U}_{\mathbf{X}'}^\top \mathbf{X} = \mathbf{X}^\top \mathbf{P}_{\mathbf{U}_{\mathbf{X}'}} \mathbf{X}.$$

- (c) Yes. Using the expression for $\tilde{\mathbf{K}}$ in (b) and the idempotency of orthogonal projection matrices, we can write $\tilde{\mathbf{K}} = \mathbf{X}^\top \mathbf{P}_{\mathbf{U}_{\mathbf{X}'}} \mathbf{X} = \mathbf{A}^\top \mathbf{A}$, where $\mathbf{A} = \mathbf{P}_{\mathbf{U}_{\mathbf{X}'}} \mathbf{X}$.

- (d) Since $\mathbf{K} = \mathbf{X}^\top \mathbf{X}$, $\text{rank}(\mathbf{K}) = \text{rank}(\mathbf{X}) = r$. Similarly, $\mathbf{W} = \mathbf{X}'^\top \mathbf{X}'$ implies $\text{rank}(\mathbf{W}) = \text{rank}(\mathbf{X}') = r$. The columns of \mathbf{X}' are columns of \mathbf{X} , and they thus span the columns of \mathbf{X} . Hence, $\mathbf{U}_{\mathbf{X}'}$ is an orthonormal basis for \mathbf{X} , i.e., $\mathbf{I}_N - \mathbf{P}_{\mathbf{U}_{\mathbf{X}'}} \in \text{Null}(\mathbf{X})$, and by part (b) of this exercise we have $\mathbf{K} - \tilde{\mathbf{K}} = \mathbf{X}^\top (\mathbf{I}_N - \mathbf{P}_{\mathbf{U}_{\mathbf{X}'}}) \mathbf{X} = \mathbf{0}$.

- (e) Storage of \mathbf{K} requires roughly 3200 TB, i.e.,

$$(20 \times 10^6)^2 \text{ entries} \times 8 \text{ bytes/entry} \times \frac{1 \text{ TB}}{10^{12} \text{ bytes}} = 3200 \text{ TB}.$$

Storage of \mathbf{C} requires roughly 160 GB, i.e.,

$$(20 \times 10^6 \times 10^3) \text{ entries} \times 8 \text{ bytes/entry} \times \frac{1 \text{ GB}}{10^9 \text{ bytes}} = 160 \text{ GB}.$$

Note that the computed numbers do not account for the symmetry of \mathbf{K} (doing so would change the storage requirements by less than a factor of two).

15.5 Expression for \mathbf{K}_{LLE}

We must find a SPSP matrix whose top k singular vectors are identical to the second through $(k+1)$ st smallest singular vectors of \mathbf{M} , since these singular vectors of \mathbf{M} form the solution to (15.10). If we define σ_{max} as the largest singular value of \mathbf{M} , then the largest $k+1$ singular vectors of $\tilde{\mathbf{M}} = \sigma_{max} \mathbf{I} - \mathbf{M}$ are identical to the bottom $k+1$ singular vectors of \mathbf{M} . Since $\tilde{\mathbf{M}}$ is SPSP, note that $\tilde{\mathbf{M}}$ is also SPSP. Moreover, its largest singular vector is simply the all ones vector properly normalized, i.e., $m^{-1/2} \mathbf{1}$. Hence, by projecting onto the subspace orthogonal to $\mathbf{1}$ LLE coincides with KPCA used with the kernel matrix $\mathbf{K}_{LLE} = (\mathbf{I} - \frac{1}{m} \mathbf{1} \mathbf{1}^\top) \tilde{\mathbf{M}} (\mathbf{I} - \frac{1}{m} \mathbf{1} \mathbf{1}^\top)$.

15.6 Random projection, PCA and nearest neighbors

(d,e) As shown in figure E.15, the random projection nearest neighbors become more accurate approximations of the true nearest neighbors as k increases. Moreover, in scenarios where we care about “near” neighbors, but not necessarily the “nearest” neighbors, the random projection approximation is quite good, as exhibited by the score_{50} performance.

(f) As shown in figure E.16, the nearest neighbor approximations generated via PCA are better than those generated via random projections, which is not surprising given that PCA aims to minimize reconstruction error, as shown in (15.1).

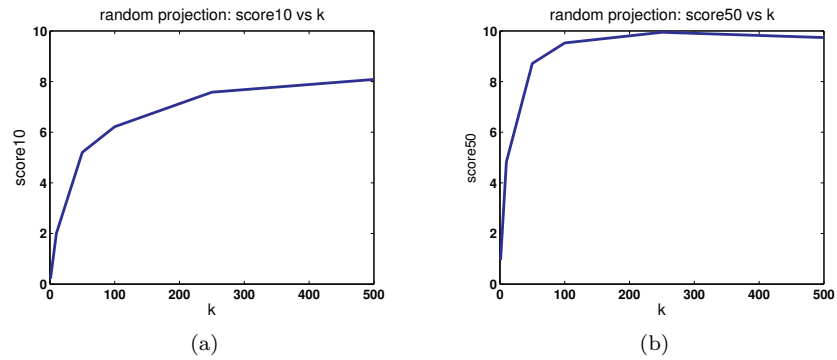
Chapter C

C.1 For any $\delta > 0$, let $t = f^{-1}(\delta)$. Plugging this in $\mathbb{P}[X > t] \leq f(t)$ yields $\mathbb{P}[X > f^{-1}(\delta)] \leq \delta$, that is $\mathbb{P}[X \leq f^{-1}(\delta)] \geq 1 - \delta$.

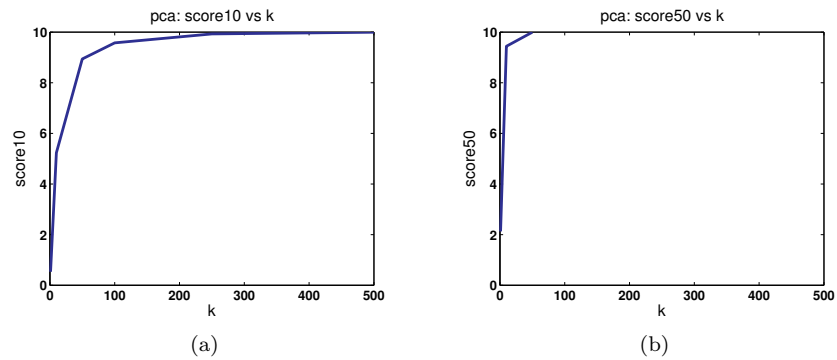
C.2 By definition of expectation and using the hint, we can write

$$\mathbb{E}[X] = \sum_{n \geq 0} n \mathbb{P}[X = n] = \sum_{n \geq 1} n (\mathbb{P}[X \geq n] - \mathbb{P}[X \geq n+1]).$$

Note that in this sum, for $n \geq 1$, $\mathbb{P}[X \geq n]$ is added n times and subtracted $n-1$ times, thus $\mathbb{E}[X] = \sum_{n \geq 1} \mathbb{P}[X \geq n]$.

**Figure E.15**

score₁₀ and score₅₀ as functions of k for random projection nearest neighbors.

**Figure E.16**

score₁₀ and score₅₀ as functions of k for PCA nearest neighbors.

More generally, by definition of the Lebesgue integral, for any non-negative random variable X , the following identity holds:

$$\mathbb{E}[X] = \int_0^{+\infty} \mathbb{P}[X \geq t] dt.$$

Chapter D

D.2 Estimating label bias. Let \hat{p}_+ be the fraction of positively labeled points in $\mathcal{S} = (x_1, \dots, x_m)$:

$$\hat{p}_+ = \frac{1}{m} \sum_{i=1}^m 1_{f(x_i)=+1}$$

Since the points are drawn i.i.d.,

$$\mathbb{E}[\hat{p}_+] = \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{s \sim \mathcal{D}^m} [1_{f(x_i)=+1}] = \mathbb{E}_{s \sim \mathcal{D}^m} [1_{f(x_1)=+1}] = \mathbb{E}_{x \sim \mathcal{D}} [1_{f(x)=+1}] = p_+.$$

Thus, by Hoeffding's inequality, for any $\epsilon > 0$,

$$\mathbb{P}[|p_+ - \hat{p}_+| > \epsilon] \leq 2e^{-2m\epsilon^2}.$$

Setting δ to match the right-hand side yields the result.

D.3 Biased coins

(a) By definition of the error of Oskar's prediction rule,

$$\begin{aligned} \text{error}(f_o) &= \mathbb{P}[f_o(S) \neq x] \\ &= \mathbb{P}[f_o(S) = x_A \wedge x = x_B] + \mathbb{P}[f_o(S) = x_B \wedge x = x_A] \\ &= \mathbb{P}\left[N(S) < \frac{m}{2} \middle| x = x_B\right] \mathbb{P}[x = x_B] + \\ &\quad \mathbb{P}\left[N(S) \geq \frac{m}{2} \middle| x = x_A\right] \mathbb{P}[x = x_A] \\ &= \frac{1}{2} \mathbb{P}\left[N(S) < \frac{m}{2} \middle| x = x_B\right] + \frac{1}{2} \mathbb{P}\left[N(S) \geq \frac{m}{2} \middle| x = x_A\right] \\ &\geq \frac{1}{2} \mathbb{P}\left[N(S) \geq \frac{m}{2} \middle| x = x_A\right]. \end{aligned}$$

(b) Note that $\mathbb{P}[N(S) \geq \frac{m}{2} | x = x_A] = \mathbb{P}[B(m, p) \geq k]$, with $p = 1/2 - \epsilon/2$, $k = \frac{m}{2}$, and $mp \leq k \leq m(1-p)$. Thus, by Slud's inequality (section D.5)

$$\text{error}(f_o) \geq \frac{1}{2} \mathbb{P}\left[N \geq \frac{m\epsilon/2}{\sqrt{1/4(1-\epsilon^2)m}}\right] = \frac{1}{2} \mathbb{P}\left[N \geq \frac{\sqrt{m}\epsilon}{\sqrt{1-\epsilon^2}}\right].$$

Using the second inequality of the appendix, we now obtain

$$\text{error}(f_o) \geq \frac{1}{4} \left(1 - \sqrt{1 - e^{-u^2}}\right),$$

with $u = \frac{\sqrt{m}\epsilon}{\sqrt{1-\epsilon^2}}$, which coincides with (D.29).

(c) If m is odd, since $\mathbb{P}\left[N(S) \geq \frac{m}{2} \middle| x = x_A\right] \geq \mathbb{P}\left[N(S) \geq \frac{m+1}{2} \middle| x = x_A\right]$, we can use the lower bound

$$\text{error}(f_o) \geq \frac{1}{2} \mathbb{P}\left[N(S) \geq \frac{m+1}{2} \middle| x = x_A\right].$$

Thus, in both cases we can use the lower bound expression with $\lceil m/2 \rceil$ instead of $m/2$.

(d) If $\text{error}(f_o)$ is at most δ , then $\frac{1}{4} \left[1 - \left[1 - e^{-\frac{2\lceil m/2 \rceil \epsilon^2}{1-\epsilon^2}}\right]^{\frac{1}{2}}\right] < \delta$, which gives

$$e^{-\frac{2\lceil m/2 \rceil \epsilon^2}{1-\epsilon^2}} < 1 - (1 - 4\delta)^2 = 4\delta(2 - 4\delta) = 8\delta(1 - 2\delta),$$

and

$$m > 2 \left\lceil \frac{1-\epsilon^2}{2\epsilon^2} \log \frac{1}{8\delta(1-2\delta)} \right\rceil.$$

The lower bound varies as $\frac{1}{\epsilon^2}$.

- (e) Let f be an arbitrary rule and denote by F_A the set of samples for which $f(S) = x_A$ and by F_B the complement. Then, by definition of the error,

$$\begin{aligned}
 error(f) &= \sum_{S \in F_A} \mathbb{P}[S \wedge x_B] + \sum_{S \in F_B} \mathbb{P}[S \wedge x_A] \\
 &= \frac{1}{2} \sum_{S \in F_A} \mathbb{P}[S|x_B] + \frac{1}{2} \sum_{S \in F_B} \mathbb{P}[S|x_A] \\
 &= \frac{1}{2} \sum_{\substack{S \in F_A \\ N(S) < m/2}} \mathbb{P}[S|x_B] + \frac{1}{2} \sum_{\substack{S \in F_A \\ N(S) \geq m/2}} \mathbb{P}[S|x_B] + \\
 &\quad \frac{1}{2} \sum_{\substack{S \in F_B \\ N(S) < m/2}} \mathbb{P}[S|x_A] + \frac{1}{2} \sum_{\substack{S \in F_B \\ N(S) \geq m/2}} \mathbb{P}[S|x_A].
 \end{aligned}$$

Now, if $N(S) \geq m/2$, clearly $\mathbb{P}[S|x_B] \geq \mathbb{P}[S|x_A]$. Similarly, if $N(S) < m/2$, clearly $\mathbb{P}[S|x_A] \geq \mathbb{P}[S|x_B]$. In view of these inequalities, $error(f)$ can be lower bounded as follows

$$\begin{aligned}
 error(f) &\geq \frac{1}{2} \sum_{\substack{S \in F_A \\ N(S) < m/2}} \mathbb{P}[S|x_B] + \frac{1}{2} \sum_{\substack{S \in F_A \\ N(S) \geq m/2}} \mathbb{P}[S|x_A] + \\
 &\quad \frac{1}{2} \sum_{\substack{S \in F_B \\ N(S) < m/2}} \mathbb{P}[S|x_B] + \frac{1}{2} \sum_{\substack{S \in F_B \\ N(S) \geq m/2}} \mathbb{P}[S|x_A] \\
 &= \frac{1}{2} \sum_{S: N(S) < m/2} \mathbb{P}[S|x_B] + \frac{1}{2} \sum_{S: N(S) \geq m/2} \mathbb{P}[S|x_A] \\
 &= error(f_o).
 \end{aligned}$$

Oskar's rule is known as the *maximum likelihood* solution.