

Understanding Machine Learning

Solution Manual

Written by Alon Gonen*
Edited by Dana Rubinfeld

November 17, 2014

2 Gentle Start

1. Given $S = ((\mathbf{x}_i, y_i))_{i=1}^m$, define the multivariate polynomial

$$p_S(\mathbf{x}) = - \prod_{i \in [m]: y_i = 1} \|\mathbf{x} - \mathbf{x}_i\|^2 .$$

Then, for every i s.t. $y_i = 1$ we have $p_S(\mathbf{x}_i) = 0$, while for every other \mathbf{x} we have $p_S(\mathbf{x}) < 0$.

2. By the linearity of expectation,

$$\begin{aligned} \mathbb{E}_{S|x \sim \mathcal{D}^m} [L_S(h)] &= \mathbb{E}_{S|x \sim \mathcal{D}^m} \left[\frac{1}{m} \sum_{i=1}^m \mathbb{1}_{[h(x_i) \neq f(x_i)]} \right] \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{x_i \sim \mathcal{D}} [\mathbb{1}_{[h(x_i) \neq f(x_i)]}] \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{P}_{x_i \sim \mathcal{D}} [h(x_i) \neq f(x_i)] \\ &= \frac{1}{m} \cdot m \cdot L_{(\mathcal{D}, f)}(h) \\ &= L_{(\mathcal{D}, f)}(h) . \end{aligned}$$

*The solutions to Chapters 13,14 were written by Shai Shalev-Shwartz

3. (a) First, observe that by definition, A labels positively all the positive instances in the training set. Second, as we assume realizability, and since the tightest rectangle enclosing all positive examples is returned, all the negative instances are labeled correctly by A as well. We conclude that A is an ERM.
- (b) Fix some distribution \mathcal{D} over \mathcal{X} , and define R^* as in the hint. Let f be the hypothesis associated with R^* a training set S , denote by $R(S)$ the rectangle returned by the proposed algorithm and by $A(S)$ the corresponding hypothesis. The definition of the algorithm A implies that $R(S) \subseteq R^*$ for every S . Thus,

$$L_{(\mathcal{D},f)}(R(S)) = \mathcal{D}(R^* \setminus R(S)) .$$

Fix some $\epsilon \in (0, 1)$. Define R_1, R_2, R_3 and R_4 as in the hint. For each $i \in [4]$, define the event

$$F_i = \{S|x : S|x \cap R_i = \emptyset\} .$$

Applying the union bound, we obtain

$$\mathcal{D}^m(\{S : L_{(\mathcal{D},f)}(A(S)) > \epsilon\}) \leq \mathcal{D}^m\left(\bigcup_{i=1}^4 F_i\right) \leq \sum_{i=1}^4 \mathcal{D}^m(F_i) .$$

Thus, it suffices to ensure that $\mathcal{D}^m(F_i) \leq \delta/4$ for every i . Fix some $i \in [4]$. Then, the probability that a sample is in F_i is the probability that all of the instances don't fall in R_i , which is exactly $(1 - \epsilon/4)^m$. Therefore,

$$\mathcal{D}^m(F_i) = (1 - \epsilon/4)^m \leq \exp(-m\epsilon/4) ,$$

and hence,

$$\mathcal{D}^m(\{S : L_{(\mathcal{D},f)}(A(S)) > \epsilon\}) \leq 4 \exp(-m\epsilon/4) .$$

Plugging in the assumption on m , we conclude our proof.

- (c) The hypothesis class of axis aligned rectangles in \mathbb{R}^d is defined as follows. Given real numbers $a_1 \leq b_1, a_2 \leq b_2, \dots, a_d \leq b_d$, define the classifier $h_{(a_1, b_1, \dots, a_d, b_d)}$ by

$$h_{(a_1, b_1, \dots, a_d, b_d)}(x_1, \dots, x_d) = \begin{cases} 1 & \text{if } \forall i \in [d], a_i \leq x_i \leq b_i \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

The class of all axis-aligned rectangles in \mathbb{R}^d is defined as

$$\mathcal{H}_{rec}^d = \{h_{(a_1, b_1, \dots, a_d, b_d)} : \forall i \in [d], a_i \leq b_i, \}.$$

It can be seen that the same algorithm proposed above is an ERM for this case as well. The sample complexity is analyzed similarly. The only difference is that instead of 4 strips, we have $2d$ strips (2 strips for each dimension). Thus, it suffices to draw a training set of size $\left\lceil \frac{2d \log(2d/\delta)}{\epsilon} \right\rceil$.

- (d) For each dimension, the algorithm has to find the minimal and the maximal values among the positive instances in the training sequence. Therefore, its runtime is $O(md)$. Since we have shown that the required value of m is at most $\left\lceil \frac{2d \log(2d/\delta)}{\epsilon} \right\rceil$, it follows that the runtime of the algorithm is indeed polynomial in $d, 1/\epsilon$, and $\log(1/\delta)$.

3 A Formal Learning Model

1. The proofs follow (almost) immediately from the definition. We will show that the sample complexity is monotonically decreasing in the accuracy parameter ϵ . The proof that the sample complexity is monotonically decreasing in the confidence parameter δ is analogous.

Denote by \mathcal{D} an unknown distribution over \mathcal{X} , and let $f \in \mathcal{H}$ be the target hypothesis. Denote by A an algorithm which learns \mathcal{H} with sample complexity $m_{\mathcal{H}}(\cdot, \cdot)$. Fix some $\delta \in (0, 1)$. Suppose that $0 < \epsilon_1 \leq \epsilon_2 \leq 1$. We need to show that $m_1 \stackrel{\text{def}}{=} m_{\mathcal{H}}(\epsilon_1, \delta) \geq m_{\mathcal{H}}(\epsilon_2, \delta) \stackrel{\text{def}}{=} m_2$. Given an i.i.d. training sequence of size $m \geq m_1$, we have that with probability at least $1 - \delta$, A returns a hypothesis h such that

$$L_{\mathcal{D}, f}(h) \leq \epsilon_1 \leq \epsilon_2 .$$

By the minimality of m_2 , we conclude that $m_2 \leq m_1$.

2. (a) We propose the following algorithm. If a positive instance x_+ appears in S , return the (true) hypothesis h_{x_+} . If S doesn't contain any positive instance, the algorithm returns the all-negative hypothesis. It is clear that this algorithm is an ERM.

(b) Let $\epsilon \in (0, 1)$, and fix the distribution \mathcal{D} over \mathcal{X} . If the true hypothesis is h^- , then our algorithm returns a perfect hypothesis.

Assume now that there exists a unique positive instance x_+ . It's clear that if x_+ appears in the training sequence S , our algorithm returns a perfect hypothesis. Furthermore, if $\mathcal{D}[\{x_+\}] \leq \epsilon$ then in any case, the returned hypothesis has a generalization error of at most ϵ (with probability 1). Thus, it is only left to bound the probability of the case in which $\mathcal{D}[\{x_+\}] > \epsilon$, but x_+ doesn't appear in S . Denote this event by F . Then

$$\mathbb{P}_{S|x \sim \mathcal{D}^m} [F] \leq (1 - \epsilon)^m \leq e^{-m\epsilon} .$$

Hence, $\mathcal{H}_{\text{Singleton}}$ is PAC learnable, and its sample complexity is bounded by

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{\log(1/\delta)}{\epsilon} \right\rceil .$$

3. Consider the ERM algorithm A which given a training sequence $S = ((\mathbf{x}_i, y_i))_{i=1}^m$, returns the hypothesis \hat{h} corresponding to the "tightest" circle which contains all the positive instances. Denote the radius of this hypothesis by \hat{r} . Assume realizability and let h^* be a circle with zero generalization error. Denote its radius by r^* .

Let $\epsilon, \delta \in (0, 1)$. Let $\bar{r} \leq r^*$ be a scalar s.t. $\mathcal{D}_{\mathcal{X}}(\{x : \bar{r} \leq \|x\| \leq r^*\}) = \epsilon$. Define $E = \{\mathbf{x} \in \mathbb{R}^2 : \bar{r} \leq \|\mathbf{x}\| \leq r^*\}$. The probability (over drawing S) that $L_{\mathcal{D}}(h_S) \geq \epsilon$ is bounded above by the probability that no point in S belongs to E . This probability of this event is bounded above by

$$(1 - \epsilon)^m \leq e^{-\epsilon m} .$$

The desired bound on the sample complexity follows by requiring that $e^{-\epsilon m} \leq \delta$.

4. We first observe that \mathcal{H} is finite. Let us calculate its size accurately. Each hypothesis, besides the all-negative hypothesis, is determined by deciding for each variable x_i , whether x_i , \bar{x}_i or none of which appear in the corresponding conjunction. Thus, $|\mathcal{H}| = 3^d + 1$. We conclude that \mathcal{H} is PAC learnable and its sample complexity can be bounded by

$$m_{\mathcal{H}}(\epsilon, \delta) \leq \left\lceil \frac{d \log 3 + \log(1/\delta)}{\epsilon} \right\rceil .$$

Let's describe our learning algorithm. We define $h_0 = x_1 \cap \bar{x}_1 \cap \dots \cap x_d \cap \bar{x}_d$. Observe that h_0 is the always-minus hypothesis. Let $((\mathbf{a}^1, y^1), \dots, (\mathbf{a}^m, y^m))$ be an i.i.d. training sequence of size m . Since

we cannot produce any information from negative examples, our algorithm neglects them. For each positive example a , we remove from h_i all the literals that are missing in a . That is, if $a_i = 1$, we remove \bar{x}_i from h and if $a_i = 0$, we remove x_i from h_i . Finally, our algorithm returns h_m .

By construction and realizability, h_i labels positively all the positive examples among $\mathbf{a}^1, \dots, \mathbf{a}^i$. From the same reasons, the set of literals in h_i contains the set of literals in the target hypothesis. Thus, h_i classifies correctly the negative elements among $\mathbf{a}^1, \dots, \mathbf{a}^i$. This implies that h_m is an ERM.

Since the algorithm takes linear time (in terms of the dimension d) to process each example, the running time is bounded by $O(m \cdot d)$.

5. Fix some $h \in \mathcal{H}$ with $L_{(\bar{\mathcal{D}}_m, f)}(h) > \epsilon$. By definition,

$$\frac{\mathbb{P}_{X \sim \mathcal{D}_1}[h(X) = f(X)] + \dots + \mathbb{P}_{X \sim \mathcal{D}_m}[h(X) = f(X)]}{m} < 1 - \epsilon .$$

We now bound the probability that h is consistent with S (i.e., that $L_S(h) = 0$) as follows:

$$\begin{aligned} \mathbb{P}_{S \sim \prod_{i=1}^m \mathcal{D}_i} [L_S(h) = 0] &= \prod_{i=1}^m \mathbb{P}_{X \sim \mathcal{D}_i} [h(X) = f(X)] \\ &= \left(\left(\prod_{i=1}^m \mathbb{P}_{X \sim \mathcal{D}_i} [h(X) = f(X)] \right)^{\frac{1}{m}} \right)^m \\ &\leq \left(\frac{\sum_{i=1}^m \mathbb{P}_{X \sim \mathcal{D}_i} [h(X) = f(X)]}{m} \right)^m \\ &< (1 - \epsilon)^m \\ &\leq e^{-\epsilon m} . \end{aligned}$$

The first inequality is the geometric-arithmetic mean inequality. Applying the union bound, we conclude that the probability that there exists some $h \in \mathcal{H}$ with $L_{(\bar{\mathcal{D}}_m, f)}(h) > \epsilon$, which is consistent with S is at most $|\mathcal{H}| \exp(-\epsilon m)$.

6. Suppose that \mathcal{H} is agnostic PAC learnable, and let A be a learning algorithm that learns \mathcal{H} with sample complexity $m_{\mathcal{H}}(\cdot, \cdot)$. We show that \mathcal{H} is PAC learnable using A .

Let \mathcal{D}, f be an (unknown) distribution over \mathcal{X} , and the target function respectively. We may assume w.l.o.g. that \mathcal{D} is a joint distribution over $\mathcal{X} \times \{0, 1\}$, where the conditional probability of y given x is determined deterministically by f . Since we assume realizability, we have $\inf_{h \in \mathcal{H}} L_{\mathcal{D}}(h) = 0$. Let $\epsilon, \delta \in (0, 1)$. Then, for every positive integer $m \geq m_{\mathcal{H}}(\epsilon, \delta)$, if we equip A with a training set S consisting of m i.i.d. instances which are labeled by f , then with probability at least $1 - \delta$ (over the choice of $S|_x$), it returns a hypothesis h with

$$\begin{aligned} L_{\mathcal{D}}(h) &\leq \inf_{h' \in \mathcal{H}} L_{\mathcal{D}}(h') + \epsilon \\ &= 0 + \epsilon \\ &= \epsilon . \end{aligned}$$

7. Let $x \in \mathcal{X}$. Let α_x be the conditional probability of a positive label given x . We have

$$\begin{aligned} \mathbb{P}[f_{\mathcal{D}}(X) \neq y | X = x] &= \mathbb{1}_{[\alpha_x \geq 1/2]} \cdot \mathbb{P}[Y = 0 | X = x] + \mathbb{1}_{[\alpha_x < 1/2]} \cdot \mathbb{P}[Y = 1 | X = x] \\ &= \mathbb{1}_{[\alpha_x \geq 1/2]} \cdot (1 - \alpha_x) + \mathbb{1}_{[\alpha_x < 1/2]} \cdot \alpha_x \\ &= \min\{\alpha_x, 1 - \alpha_x\}. \end{aligned}$$

Let g be a classifier¹ from \mathcal{X} to $\{0, 1\}$. We have

$$\begin{aligned} \mathbb{P}[g(X) \neq Y | X = x] &= \mathbb{P}[g(X) = 0 | X = x] \cdot \mathbb{P}[Y = 1 | X = x] \\ &\quad + \mathbb{P}[g(X) = 1 | X = x] \cdot \mathbb{P}[Y = 0 | X = x] \\ &= \mathbb{P}[g(X) = 0 | X = x] \cdot \alpha_x + \mathbb{P}[g(X) = 1 | X = x] \cdot (1 - \alpha_x) \\ &\geq \mathbb{P}[g(X) = 0 | X = x] \cdot \min\{\alpha_x, 1 - \alpha_x\} \\ &\quad + \mathbb{P}[g(X) = 1 | X = x] \cdot \min\{\alpha_x, 1 - \alpha_x\} \\ &= \min\{\alpha_x, 1 - \alpha_x\}, \end{aligned}$$

The statement follows now due to the fact that the above is true for every $x \in \mathcal{X}$. More formally, by the law of total expectation,

$$\begin{aligned} L_{\mathcal{D}}(f_{\mathcal{D}}) &= \mathbb{E}_{(x,y) \sim \mathcal{D}}[\mathbb{1}_{[f_{\mathcal{D}}(x) \neq y]}] \\ &= \mathbb{E}_{x \sim \mathcal{D}_X} \left[\mathbb{E}_{y \sim \mathcal{D}_{Y|x}}[\mathbb{1}_{[f_{\mathcal{D}}(x) \neq y]} | X = x] \right] \\ &= \mathbb{E}_{x \sim \mathcal{D}_X}[\alpha_x] \\ &\leq \mathbb{E}_{x \sim \mathcal{D}_X} \left[\mathbb{E}_{y \sim \mathcal{D}_{Y|x}}[\mathbb{1}_{[g(x) \neq y]} | X = x] \right] \\ &= L_{\mathcal{D}}(g) . \end{aligned}$$

¹As we shall see, g might be non-deterministic.

8. (a) This was proved in the previous exercise.
 (b) We proved in the previous exercise that for every distribution \mathcal{D} , the bayes optimal predictor $f_{\mathcal{D}}$ is optimal w.r.t. \mathcal{D} .
 (c) Choose any distribution \mathcal{D} . Then A is not better than $f_{\mathcal{D}}$ w.r.t. \mathcal{D} .
9. (a) Suppose that \mathcal{H} is PAC learnable in the one-oracle model. Let A be an algorithm which learns \mathcal{H} and denote by $m_{\mathcal{H}}$ the function that determines its sample complexity. We prove that \mathcal{H} is PAC learnable also in the two-oracle model.

Let \mathcal{D} be a distribution over $\mathcal{X} \times \{0, 1\}$. Note that drawing points from the negative and positive oracles with equal provability is equivalent to obtaining i.i.d. examples from a distribution \mathcal{D}' which gives equal probability to positive and negative examples. Formally, for every subset $E \subseteq \mathcal{X}$ we have

$$\mathcal{D}'[E] = \frac{1}{2}\mathcal{D}^+[E] + \frac{1}{2}\mathcal{D}^-[E].$$

Thus, $\mathcal{D}'[\{x : f(x) = 1\}] = \mathcal{D}'[\{x : f(x) = 0\}] = \frac{1}{2}$. If we let A an access to a training set which is drawn i.i.d. according to \mathcal{D}' with size $m_{\mathcal{H}}(\epsilon/2, \delta)$, then with probability at least $1 - \delta$, A returns h with

$$\begin{aligned} \epsilon/2 &\geq L_{(\mathcal{D}', f)}(h) = \mathbb{P}_{x \sim \mathcal{D}'}[h(x) \neq f(x)] \\ &= \mathbb{P}_{x \sim \mathcal{D}'}[f(x) = 1, h(x) = 0] + \mathbb{P}_{x \sim \mathcal{D}'}[f(x) = 0, h(x) = 1] \\ &= \mathbb{P}_{x \sim \mathcal{D}'}[f(x) = 1] \cdot \mathbb{P}_{x \sim \mathcal{D}'}[h(x) = 0 | f(x) = 1] \\ &+ \mathbb{P}_{x \sim \mathcal{D}'}[f(x) = 0] \cdot \mathbb{P}_{x \sim \mathcal{D}'}[h(x) = 1 | f(x) = 0] \\ &= \mathbb{P}_{x \sim \mathcal{D}'}[f(x) = 1] \cdot \mathbb{P}_{x \sim \mathcal{D}}[h(x) = 0 | f(x) = 1] \\ &+ \mathbb{P}_{x \sim \mathcal{D}'}[f(x) = 0] \cdot \mathbb{P}_{x \sim \mathcal{D}}[h(x) = 1 | f(x) = 0] \\ &= \frac{1}{2} \cdot L_{(\mathcal{D}^+, f)}(h) + \frac{1}{2} \cdot L_{(\mathcal{D}^-, f)}(h). \end{aligned}$$

This implies that with probability at least $1 - \delta$, both

$$L_{(\mathcal{D}^+, f)}(h) \leq \epsilon \text{ and } L_{(\mathcal{D}^-, f)}(h) \leq \epsilon.$$

Our definition for PAC learnability in the two-oracle model is satisfied. We can bound both $m_{\mathcal{H}}^+(\epsilon, \delta)$ and $m_{\mathcal{H}}^-(\epsilon, \delta)$ by $m_{\mathcal{H}}(\epsilon/2, \delta)$.

- (b) Suppose that \mathcal{H} is PAC learnable in the two-oracle model and let A be an algorithm which learns \mathcal{H} . We show that \mathcal{H} is PAC learnable also in the standard model.

Let \mathcal{D} be a distribution over \mathcal{X} , and denote the target hypothesis by f . Let $\alpha = \mathcal{D}[\{x : f(x) = 1\}]$. Let $\epsilon, \delta \in (0, 1)$. According to our assumptions, there exist $m^+ \stackrel{\text{def}}{=} m_{\mathcal{H}}^+(\epsilon, \delta/2), m^- \stackrel{\text{def}}{=} m_{\mathcal{H}}^-(\epsilon, \delta/2)$ s.t. if we equip A with m^+ examples drawn i.i.d. from \mathcal{D}^+ and m^- examples drawn i.i.d. from \mathcal{D}^- , then, with probability at least $1 - \delta/2$, A will return h with

$$L_{(\mathcal{D}^+, f)}(h) \leq \epsilon \wedge L_{(\mathcal{D}^-, f)}(h) \leq \epsilon .$$

Our algorithm B draws $m = \max\{2m^+/\epsilon, 2m^-/\epsilon, \frac{8\log(4/\delta)}{\epsilon}\}$ samples according to \mathcal{D} . If there are less than m^+ positive examples, B returns h^- . Otherwise, if there are less than m^- negative examples, B returns h^+ . Otherwise, B runs A on the sample and returns the hypothesis returned by A .

First we observe that if the sample contains m^+ positive instances and m^- negative instances, then the reduction to the two-oracle model works well. More precisely, with probability at least $1 - \delta/2$, A returns h with

$$L_{(\mathcal{D}^+, f)}(h) \leq \epsilon \wedge L_{(\mathcal{D}^-, f)}(h) \leq \epsilon .$$

Hence, with probability at least $1 - \delta/2$, the algorithm B returns (the same) h with

$$L_{(\mathcal{D}, f)}(h) = \alpha \cdot L_{(\mathcal{D}^+, f)}(h) + (1 - \alpha) \cdot L_{(\mathcal{D}^-, f)}(h) \leq \epsilon$$

We consider now the following cases:

- Assume that both $\alpha \geq \epsilon$. We show that with probability at least $1 - \delta/4$, the sample contain m^+ positive instances. For each $i \in [m]$, define the indicator random variable Z_i , which gets the value 1 iff the i -th element in the sample is positive. Define $Z = \sum_{i=1}^m Z_i$ to be the number of positive examples that were drawn. Clearly, $\mathbb{E}[Z] = \alpha m$. Using Chernoff bound, we obtain

$$\mathbb{P}[Z < (1 - \frac{1}{2})\alpha m] < e^{-\frac{m\alpha}{8}} .$$

By the way we chose m , we conclude that

$$\mathbb{P}[Z < m_+] < \delta/4 .$$

Similarly, if $1 - \alpha \geq \epsilon$, the probability that less than m^- negative examples were drawn is at most $\delta/4$. If both $\alpha \geq \epsilon$ and $1 - \alpha \geq \epsilon$, then, by the union bound, with probability at least $1 - \delta/2$, the training set contains at least m^+ and m^- positive and negative instances respectively. As we mentioned above, if this is the case, the reduction to the two-oracle model works with probability at least $1 - \delta/2$. The desired conclusion follows by applying the union bound.

- Assume that $\alpha < \epsilon$, and less than m^+ positive examples are drawn. In this case, B will return the hypothesis h^- . We obtain

$$L_{\mathcal{D}}(h) = \alpha < \epsilon .$$

Similarly, if $(1 - \alpha) < \epsilon$, and less than m^- negative examples are drawn, B will return h^+ . In this case,

$$L_{\mathcal{D}}(h) = 1 - \alpha < \epsilon .$$

All in all, we have shown that with probability at least $1 - \delta$, B returns a hypothesis h with $L_{(\mathcal{D},f)}(h) < \epsilon$. This satisfies our definition for PAC learnability in the one-oracle model.

4 Learning via Uniform Convergence

1. (a) Assume that for every $\epsilon, \delta \in (0, 1)$, and every distribution \mathcal{D} over $\mathcal{X} \times \{0, 1\}$, there exists $m(\epsilon, \delta) \in \mathbb{N}$ such that for every $m \geq m(\epsilon, \delta)$,

$$\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) > \epsilon] < \delta .$$

Let $\lambda > 0$. We need to show that there exists $m_0 \in \mathbb{N}$ such that for every $m \geq m_0$, $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] \leq \lambda$. Let $\epsilon = \min(1/2, \lambda/2)$. Set $m_0 = m_{\mathcal{H}}(\epsilon, \epsilon)$. For every $m \geq m_0$, since the loss is bounded

above by 1, we have

$$\begin{aligned}
\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] &\leq \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) > \lambda/2] \cdot 1 + \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) \leq \lambda/2] \cdot \lambda/2 \\
&\leq \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) > \epsilon] + \lambda/2 \\
&\leq \epsilon + \lambda/2 \\
&\leq \lambda/2 + \lambda/2 \\
&= \lambda .
\end{aligned}$$

(b) Assume now that

$$\lim_{m \rightarrow \infty} \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] = 0 .$$

Let $\epsilon, \delta \in (0, 1)$. There exists some $m_0 \in \mathbb{N}$ such that for every $m \geq m_0$, $\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] \leq \epsilon \cdot \delta$. By Markov's inequality,

$$\begin{aligned}
\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) > \epsilon] &\leq \frac{\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))]}{\epsilon} \\
&\leq \frac{\epsilon \delta}{\epsilon} \\
&= \delta .
\end{aligned}$$

2. The left inequality follows from Corollary 4.4. We prove the right inequality. Fix some $h \in \mathcal{H}$. Applying Hoeffding's inequality, we obtain

$$\mathbb{P}_{S \sim \mathcal{D}^m} [|L_{\mathcal{D}}(h) - L_S(h)| \geq \epsilon/2] \leq 2 \exp \left(-\frac{2m\epsilon^2}{(b-a)^2} \right) . \quad (2)$$

The desired inequality is obtained by requiring that the right-hand side of Equation (2) is at most $\delta/|\mathcal{H}|$, and then applying the union bound.

5 The Bias-Complexity Tradeoff

1. We simply follow the hint. By Lemma B.1,

$$\begin{aligned}
\mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) \geq 1/8] &= \mathbb{P}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) \geq 1 - 7/8] \\
&\geq \frac{\mathbb{E}[L_{\mathcal{D}}(A(S))] - (1 - 7/8)}{7/8} \\
&\geq \frac{1/8}{7/8} \\
&= 1/7 .
\end{aligned}$$

9 Linear Predictors

1. Define a vector of auxiliary variables $s = (s_1, \dots, s_m)$. Following the hint, minimizing the empirical risk is equivalent to minimizing the linear objective $\sum_{i=1}^m s_i$ under the following constraints:

$$(\forall i \in [m]) \quad \mathbf{w}^T \mathbf{x}_i - s_i \leq y_i \quad , \quad -\mathbf{w}^T \mathbf{x}_i - s_i \leq -y_i \quad (3)$$

It is left to translate the above into matrix form. Let $A \in \mathbb{R}^{2m \times (m+d)}$ be the matrix $A = [X \ -I_m; -X \ -I_m]$, where $X_{i \rightarrow} = x_i$ for every $i \in [m]$. Let $\mathbf{v} \in \mathbb{R}^{d+m}$ be the vector of variables $(w_1, \dots, w_d, s_1, \dots, s_m)$. Define $\mathbf{b} \in \mathbb{R}^{2m}$ to be the vector $\mathbf{b} = (y_1, \dots, y_m, -y_1, \dots, -y_m)^T$. Finally, let $\mathbf{c} \in \mathbb{R}^{d+m}$ be the vector $\mathbf{c} = (\mathbf{0}_d; \mathbf{1}_m)$. It follows that the optimization problem of minimizing the empirical risk can be expressed as the following LP:

$$\begin{aligned} \min \quad & \mathbf{c}^T \mathbf{v} \\ \text{s.t.} \quad & A\mathbf{v} \leq \mathbf{b} . \end{aligned}$$

2. Consider the matrix $X \in \mathbb{R}^{d \times m}$ whose columns are x_1, \dots, x_m . The rank of X is equal to the dimension of the subspace $\text{span}(\{x_1, \dots, x_m\})$. The SVD theorem (more precisely, Lemma C.4) implies that the rank of X is equal to the rank of $A = XX^T$. Hence, the set $\{x_1, \dots, x_m\}$ spans \mathbb{R}^d if and only if the rank of $A = XX^T$ is d , i.e., iff $A = XX^T$ is invertible.
3. Following the hint, let $d = m$, and for every $i \in [m]$, let $\mathbf{x}_i = \mathbf{e}_i$. Let us agree that $\text{sign}(0) = -1$. For $i = 1, \dots, d$, let $y_i = 1$ be the label of x_i . Denote by $\mathbf{w}^{(t)}$ the weight vector which is maintained by the Perceptron. A simple inductive argument shows that for every $i \in [d]$, $\mathbf{w}_i = \sum_{j < i} \mathbf{e}_j$. It follows that for every $i \in [d]$, $\langle \mathbf{w}^{(i)}, \mathbf{x}_i \rangle = 0$. Hence, all the instances $\mathbf{x}_1, \dots, \mathbf{x}_d$ are misclassified (and then we obtain the vector $w = (1, \dots, 1)$ which is consistent with x_1, \dots, x_m). We also note that the vector $\mathbf{w}^* = (1, \dots, 1)$ satisfies the requirements listed in the question.
4. Consider all positive examples of the form $(\alpha, \beta, 1)$, where $\alpha^2 + \beta^2 + 1 \leq R^2$. Observe that $\mathbf{w}^* = (0, 0, 1)$ satisfies $y \langle \mathbf{w}^*, \mathbf{x} \rangle \geq 1$ for all such (\mathbf{x}, y) . We show a sequence of R^2 examples on which the Perceptron makes R^2 mistakes.

The idea of the construction is to start with the examples $(\alpha_1, 0, 1)$ where $\alpha_1 = \sqrt{R^2 - 1}$. Now, on round t let the new example be such that the following conditions hold:

- (a) $\alpha^2 + \beta^2 + 1 = R^2$
- (b) $\langle \mathbf{w}_t, (\alpha, \beta, 1) \rangle = 0$

As long as we can satisfy both conditions, the Perceptron will continue to err. We'll show that as long as $t \leq R^2$ we can satisfy these conditions.

Observe that, by induction, $\mathbf{w}^{(t-1)} = (a, b, t-1)$ for some scalars a, b . Observe also that $\|\mathbf{w}_{t-1}\|^2 = (t-1)R^2$ (this follows from the proof of the Perceptron's mistake bound, where inequalities hold with equality). That is, $a^2 + b^2 + (t-1)^2 = (t-1)R^2$.

W.l.o.g., let us rotate $\mathbf{w}^{(t-1)}$ w.r.t. the z axis so that it is of the form $(a, 0, t-1)$ and we have $a = \sqrt{(t-1)R^2 - (t-1)^2}$. Choose

$$\alpha = -\frac{t-1}{a} .$$

Then, for every β ,

$$\langle (a, 0, t-1), (\alpha, \beta, 1) \rangle = 0 .$$

We just need to verify that $\alpha^2 + 1 \leq R^2$, because if this is true then we can choose $\beta = \sqrt{R^2 - \alpha^2 - 1}$. Indeed,

$$\begin{aligned} \alpha^2 + 1 &= \frac{(t-1)^2}{a^2} + 1 = \frac{(t-1)^2}{(t-1)R^2 - (t-1)^2} + 1 = \frac{(t-1)R^2}{(t-1)R^2 - (t-1)^2} \\ &= R^2 \frac{1}{R^2 - (t-1)} \\ &\leq R^2 \end{aligned}$$

where the last inequality assumes $R^2 \geq t$.

5. Since for every $w \in \mathbb{R}^d$, and every $x \in \mathbb{R}^d$, we have

$$\text{sign}(\langle w, x \rangle) = \text{sign}(\langle \eta w, x \rangle) ,$$

we obtain that the modified Perceptron and the Perceptron produce the same predictions. Consequently, both algorithms perform the same number of iterations.

6. In this question we will denote the class of halfspaces in \mathbb{R}^{d+1} by \mathcal{L}_{d+1} .

- (a) Assume that $A = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathbb{R}^d$ is shattered by \mathcal{B}_d . Then, $\forall \mathbf{y} = (y_1, \dots, y_d) \in \{-1, 1\}^d$ there exists $B_{\mu, r} \in \mathcal{B}$ s.t. for every i

$$B_{\mu, r}(\mathbf{x}_i) = y_i .$$

Hence, for the above μ and r , the following identity holds for every $i \in [m]$:

$$\text{sign} \left((2\mu; -1)^T(\mathbf{x}_i; \|\mathbf{x}_i\|^2) - \|\mu\|^2 + r^2 \right) = y_i , \quad (4)$$

where $;$ denotes vector concatenation. For each $i \in [m]$, let $\phi(x_i) = (x_i; \|x_i\|^2)$. Define the halfspace $h \in \mathcal{L}_{d+1}$ which corresponds to $w = (2\mu; -1)$, and $b = \|\mu\|^2 - r^2$. Equation (4) implies that for every $i \in [m]$,

$$h(x_i) = y_i$$

All in all, if $A = \{x_1, \dots, x_m\}$ is shattered by \mathcal{B} , then $\phi(A) := \{\phi(x_1), \dots, \phi(x_m)\}$, is shattered by \mathcal{L} . We conclude that $d + 2 = \text{VCdim}(\mathcal{L}_{d+1}) \geq VC(\mathcal{B}_d)$.

- (b) Consider the set C consisting of the unit vectors $\mathbf{e}_1, \dots, \mathbf{e}_d$, and the origin $\mathbf{0}$. Let $A \subseteq C$. We show that there exists a ball such that all the vectors in A are labeled positively, while the vectors in $C \setminus A$ are labeled negatively. We define the center $\mu = \sum_{e \in A} e$. Note that for every unit vector in A , its distance to the center is $\sqrt{|A|} - 1$. Also, for every unit vector outside A , its distance to the center is $\sqrt{|A|} + 1$. Finally, the distance of the origin to the center is $\sqrt{|A|}$. Hence, if $0 \in A$, we will set $r = \sqrt{|A|} - 1$, and if $0 \notin A$, we will set $r = \sqrt{|A|}$. We conclude that the set C is shattered by \mathcal{B}_d . All in all, we just showed that $\text{VCdim}(\mathcal{B}_d) \geq d + 1$.

10 Boosting

- Let $\epsilon, \delta \in (0, 1)$. Pick k ‘‘chunks’’ of size $m_{\mathcal{H}}(\epsilon/2)$. Apply A on each of these chunks, to obtain $\hat{h}_1, \dots, \hat{h}_k$. Note that the probability that $\min_{i \in [k]} L_{\mathcal{D}}(\hat{h}_i) \leq \min L_{\mathcal{D}}(h) + \epsilon/2$ is at least $1 - \delta_0^k \geq 1 - \delta/2$. Now, apply an ERM over the class $\hat{\mathcal{H}} := \{\hat{h}_1, \dots, \hat{h}_k\}$ with the training data being the last chunk of size $\left\lceil \frac{2 \log(4k/\delta)}{\epsilon^2} \right\rceil$. Denote the output hypothesis

11 Model Selection and Validation

- Let S be an i.i.d. sample. Let h be the output of the described learning algorithm. Note that (independently of the identity of S), $L_{\mathcal{D}}(h) = 1/2$ (since h is a constant function).

Let us calculate the estimate $L_V(h)$. Assume that the parity of S is

- Fix some fold $\{(\mathbf{x}, y)\} \subseteq S$. We distinguish between two cases:

- The parity of $S \setminus \{\mathbf{x}\}$ is 1. It follows that $y = 0$. When being trained using $S \setminus \{\mathbf{x}\}$, the algorithm outputs the constant predictor $h(\mathbf{x}) = 1$. Hence, the leave-one-out estimate using this fold is 1.
- The parity of $S \setminus \{\mathbf{x}\}$ is 0. It follows that $y = 1$. When being trained using $S \setminus \{\mathbf{x}\}$, the algorithm outputs the constant predictor $h(\mathbf{x}) = 0$. Hence, the leave-one-out estimate using this fold is 1.

Averaging over the folds, the estimate of the error of h is 1. Consequently, the difference between the estimate and the true error is $1/2$. The case in which the parity of S is 0 is analyzed analogously.

- Consider for example the case in which $\mathcal{H}_1 \subseteq \mathcal{H}_2 \subseteq \dots \subseteq \mathcal{H}_k$, and $|\mathcal{H}_i| = 2^i$ for every $i \in k$. Learning \mathcal{H}_k in the Agnostic-Pac model provides the following bound for an ERM hypothesis h :

$$L_D(h) \leq \min_{h \in \mathcal{H}_k} L_D(h) + \sqrt{\frac{2(k+1 + \log(1/\delta))}{m}}.$$

Alternatively, we can use model selection as we describe next. Assume that j is the minimal index which contains a hypothesis $h^* \in \operatorname{argmin}_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$. Fix some $r \in [k]$. By Hoeffding's inequality, with probability at least $1 - \delta/(2k)$, we have

$$|L_{\mathcal{D}}(\hat{h}_r) - L_V(\hat{h}_r)| \leq \sqrt{\frac{1}{2\alpha m} \log \frac{4}{\delta}}.$$

Applying the union bound, we obtain that with probability at least $1 - \delta/2$, the following inequality holds (simultaneously) for every $r \in$

[k]:

$$\begin{aligned}
L_{\mathcal{D}}(\hat{h}) &\leq L_V(\hat{h}) + \sqrt{\frac{1}{2\alpha m} \log \frac{4k}{\delta}} \\
&\leq L_V(\hat{h}_r) + \sqrt{\frac{1}{2\alpha m} \log \frac{4k}{\delta}} \\
&\leq L_{\mathcal{D}}(\hat{h}_r) + 2\sqrt{\frac{1}{2\alpha m} \log \frac{4k}{\delta}} \\
&= L_{\mathcal{D}}(\hat{h}_r) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}}.
\end{aligned}$$

In particular, with probability at least $1 - \delta/2$, we have

$$L_{\mathcal{D}}(\hat{h}) \leq L_{\mathcal{D}}(\hat{h}_j) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}}.$$

Using similar arguments⁸, we obtain that with probability at least $1 - \delta/2$,

$$\begin{aligned}
L_{\mathcal{D}}(\hat{h}_j) &\leq L_{\mathcal{D}}(h^*) + \sqrt{\frac{2}{(1-\alpha)m} \log \frac{4|\mathcal{H}_j|}{\delta}} \\
&= L_{\mathcal{D}}(h^*) + \sqrt{\frac{2}{(1-\alpha)m} \log \frac{4|\mathcal{H}_j|}{\delta}}
\end{aligned}$$

Combining the two last inequalities with the union bound, we obtain that with probability at least $1 - \delta$,

$$L_{\mathcal{D}}(\hat{h}) \leq L_{\mathcal{D}}(h^*) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}} + \sqrt{\frac{2}{(1-\alpha)m} \log \frac{4|\mathcal{H}_j|}{\delta}}.$$

We conclude that

$$L_{\mathcal{D}}(\hat{h}) \leq L_{\mathcal{D}}(h^*) + \sqrt{\frac{2}{\alpha m} \log \frac{4k}{\delta}} + \sqrt{\frac{2}{(1-\alpha)m} (j + \log \frac{4}{\delta})}.$$

Comparing the two bounds, we see that when the “optimal index” j is significantly smaller than k , the bound achieved using model selection is much better. Being even more concrete, if j is logarithmic in k , we achieve a logarithmic improvement.

⁸This time we consider each of the hypotheses in H_j , and apply the union bound accordingly.

3. Fix some $(\mathbf{x}, y) \in \{\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}'\|_2 \leq R\} \times \{-1, 1\}$. Let $\mathbf{w}_1, \mathbf{w}_2 \in \mathbb{R}^d$. For $i \in [2]$, let $\ell_i = \max\{0, 1 - y\langle \mathbf{w}_i, \mathbf{x} \rangle\}$. We wish to show that $|\ell_1 - \ell_2| \leq R\|\mathbf{w}_1 - \mathbf{w}_2\|_2$. If both $y\langle \mathbf{w}_1, \mathbf{x} \rangle \geq 1$ and $y\langle \mathbf{w}_2, \mathbf{x} \rangle \geq 1$, then $|\ell_1 - \ell_2| = 0 \leq R\|\mathbf{w}_1 - \mathbf{w}_2\|_2$. Assume now that $|\{i : y\langle \mathbf{w}_i, \mathbf{x} \rangle < 1\}| \geq 1$. Assume w.l.o.g. that $1 - y\langle \mathbf{w}_1, \mathbf{x} \rangle \geq 1 - y\langle \mathbf{w}_2, \mathbf{x} \rangle$. Hence,

$$\begin{aligned}
|\ell_1 - \ell_2| &= \ell_1 - \ell_2 \\
&= 1 - y\langle \mathbf{w}_1, \mathbf{x} \rangle - \max\{0, 1 - y\langle \mathbf{w}_2, \mathbf{x} \rangle\} \\
&\leq 1 - y\langle \mathbf{w}_1, \mathbf{x} \rangle - (1 - y\langle \mathbf{w}_2, \mathbf{x} \rangle) \\
&= y\langle \mathbf{w}_2 - \mathbf{w}_1, \mathbf{x} \rangle \\
&\leq \|\mathbf{w}_1 - \mathbf{w}_2\| \|\mathbf{x}\| \\
&\leq R\|\mathbf{w}_1 - \mathbf{w}_2\| .
\end{aligned}$$

4. (a) Fix a Turing machine T . If T halts on the input 0, then for every $h \in [0, 1]$,

$$\ell(h, T) = \langle (h, 1 - h), (1, 0) \rangle .$$

If T halts on the input 1, then for every $h \in [0, 1]$,

$$\ell(h, T) = \langle (h, 1 - h), (0, 1) \rangle .$$

In both cases, ℓ is linear, and hence convex over \mathcal{H} .

- (b) The idea is to reduce the halting problem to the learning problem⁹. More accurately, the following decision problem can be easily reduced to the learning problem described in the question: Given a Turing machine M , does M halt given the input M ? The proof that the halting problem is not decidable implies that this decision problem is not decidable as well. Hence, there is no computable algorithm that learns the problem described in the question.

13 Regularization and Stability

1. **From bounded expected risk to agnostic PAC learning:** We assume A is a (proper) algorithm that guarantees the following: If $m \geq m_{\mathcal{H}}(\epsilon)$ then for every distribution \mathcal{D} it holds that

$$\mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S))] \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon .$$

⁹Errata: “ T halts on the input 0” should be replaced (everywhere) by “ T halts on the input T ”

- Since $A(S) \in \mathcal{H}$, the random variable $\theta = L_{\mathcal{D}}(A(S)) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$ is non-negative.¹⁰ Therefore, Markov's inequality implies that

$$\mathbb{P}[\theta \geq \mathbb{E}[\theta]/\delta] \leq \frac{\mathbb{E}[\theta]}{\mathbb{E}[\theta]/\delta} = \delta .$$

In other words, with probability of at least $1 - \delta$ we have

$$\theta \leq \mathbb{E}[\theta]/\delta .$$

But, if $m \geq m_{\mathcal{H}}(\epsilon \delta)$ then we know that $\mathbb{E}[\theta] \leq \epsilon \delta$. This yields $\theta \leq \epsilon$, which concludes our proof.

- Let $k = \lceil \log_2(2/\delta) \rceil$.¹¹ Divide the data into $k + 1$ chunks, where each of the first k chunks is of size $m_{\mathcal{H}}(\epsilon/4)$ examples.¹² Train the first k chunks using A . Let h_i be the output of A for the i 'th chunk. Using the previous question, we know that with probability of at most $1/2$ over the examples in the i 'th chunk it holds that $L_{\mathcal{D}}(h_i) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) > \epsilon/2$. Since the examples in the different chunks are independent, the probability that for all the chunks we'll have $L_{\mathcal{D}}(h_i) - \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) > \epsilon/2$ is at most 2^{-k} , which by the definition of k is at most $\delta/2$. In other words, with probability of at least $1 - \delta/2$ we have that

$$\min_{i \in [k]} L_{\mathcal{D}}(h_i) \leq \min_{h \in \mathcal{H}} L_{\mathcal{D}}(h) + \epsilon/2 . \quad (7)$$

Next, apply ERM over the finite class $\{h_1, \dots, h_k\}$ on the last chunk. By Corollary 4.6, if the size of the last chunk is at least $\frac{8 \log(4k/\delta)}{\epsilon^2}$ then, with probability of at least $1 - \delta/2$, we have

$$L_{\mathcal{D}}(\hat{h}) \leq \min_{i \in [k]} L_{\mathcal{D}}(h_i) + \epsilon/2 .$$

Applying the union bound and combining with Equation (7) we conclude our proof. The overall sample complexity is

$$m_{\mathcal{H}}(\epsilon/4) \lceil \log_2(2/\delta) \rceil + 8 \left\lceil \frac{\log(4/\delta) + \log(\lceil \log_2(2/\delta) \rceil)}{\epsilon^2} \right\rceil$$

¹⁰One should assume here that A is a "proper" learner.

¹¹Note that in the original question the size was mistakenly $k = \lceil \log_2(1/\delta) \rceil$.

¹²Note that in the original question the size was mistakenly $m_{\mathcal{H}}(\epsilon/2)$.

2. **Learnability without uniform convergence:** Let \mathcal{B} be the unit ball of \mathbb{R}^d , let $\mathcal{H} = \mathcal{B}$, let $Z = \mathcal{B} \times \{0, 1\}^d$, and let $\ell : Z \times \mathcal{H} \rightarrow \mathbb{R}$ be defined as follows:

$$\ell(\mathbf{w}, (\mathbf{x}, \boldsymbol{\alpha})) = \sum_{i=1}^d \alpha_i (x_i - w_i)^2 .$$

- This problem is learnable using the RLM rule with a sample complexity that does not depend on d . Indeed, the hypothesis class is convex and bounded, the loss function is convex, non-negative, and smooth, since

$$\|\nabla \ell(\mathbf{v}, (\mathbf{x}, \boldsymbol{\alpha})) - \nabla \ell(\mathbf{w}, (\mathbf{x}, \boldsymbol{\alpha}))\|^2 = 4 \sum_{i=1}^d \alpha_i^2 (v_i - w_i)^2 \leq 4 \|\mathbf{v} - \mathbf{w}\|^2 ,$$

where in the last inequality we used $\alpha_i \in \{0, 1\}$.

- Fix some $j \in [d]$. The probability to sample a training set of size m such that $\alpha_j = 0$ ¹³ for all the examples is 2^{-m} . Since the coordinates of $\boldsymbol{\alpha}$ are chosen independently, the probability that for all $j \in [d]$ the above event will not happen is $(1 - 2^{-m})^d \geq \exp(-2^{-m+1}d)$. Therefore, if $d \gg 2^m$, the above quantity goes to zero, meaning that there is a high chance that for some j we'll have $\alpha_j = 0$ for all the examples in the training set.

For such a sample, we have that every vector of the form $\mathbf{x}_0 + \eta \mathbf{e}_j$ is an ERM. For simplicity, assume that \mathbf{x}_0 is the all zeros vector and set $\eta = 1$. Then, $L_S(\mathbf{e}_j) = 0$. However, $L_{\mathcal{D}}(\mathbf{e}_j) = \frac{1}{2}$. It follows that the sample complexity of uniform convergence must grow with $\log(d)$.

- Taking d to infinity we obtain a problem which is learnable but for which the uniform convergence property does not hold. While this seems to contradict the fundamental theorem of statistical learning, which states that a problem is learnable if and only if uniform convergence holds, there is no contradiction since the fundamental theorem only holds for binary classification problems, while here we consider a more general problem.

3. Stability and asymptotic ERM is sufficient for learnability:

¹³In the hint, it is written mistakenly that $\alpha_j = 1$

Proof. Let $h^* \in \operatorname{argmin}_{h \in \mathcal{H}} L_{\mathcal{D}}(h)$ (for simplicity we assume that h^* exists). We have

$$\begin{aligned}
& \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) - L_{\mathcal{D}}(h^*)] \\
&= \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) - L_S(A(S)) + L_S(A(S)) - L_{\mathcal{D}}(h^*)] \\
&= \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) - L_S(A(S))] + \mathbb{E}_{S \sim \mathcal{D}^m} [L_S(A(S)) - L_{\mathcal{D}}(h^*)] \\
&= \mathbb{E}_{S \sim \mathcal{D}^m} [L_{\mathcal{D}}(A(S)) - L_S(A(S))] + \mathbb{E}_{S \sim \mathcal{D}^m} [L_S(A(S)) - L_S(h^*)] \\
&\leq \epsilon_1(m) + \epsilon_2(m) .
\end{aligned}$$

□

4. Strong convexity with respect to general norms:

- (a) Item 2 of the lemma follows directly from the definition of convexity and strong convexity. For item 3, the proof is identical to the proof for the ℓ_2 norm.
- (b) The function $f(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_1^2$ is not strongly convex with respect to the ℓ_1 norm. To see this, let the dimension be $d = 2$ and take $\mathbf{w} = \mathbf{e}_1$, $\mathbf{u} = \mathbf{e}_2$, and $\alpha = 1/2$. We have $f(\mathbf{w}) = f(\mathbf{u}) = f(\alpha\mathbf{w} + (1 - \alpha)\mathbf{u}) = 1/2$.
- (c) The proof is almost identical to the proof for the ℓ_2 case and is therefore omitted.
- (d) *Proof.* ¹⁴
Using Holder's inequality,

$$\|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i| \leq \|\mathbf{w}\|_q \|(1, \dots, 1)\|_p ,$$

where $p = (1 - 1/q)^{-1} = \log(d)$. Since $\|(1, \dots, 1)\|_p = d^{1/p} = e < 3$ we obtain that for every \mathbf{w} ,

$$\|\mathbf{w}\|_q \geq \|\mathbf{w}\|_1 / 3 .$$

¹⁴Errata: In the original question, there's a typo: it should be proved that R is $(1/3)$ strongly convex instead of $\frac{1}{3 \log(d)}$ strongly convex.

Combining this with the strong convexity of R w.r.t. $\|\cdot\|_q$ we obtain

$$\begin{aligned} R(\alpha\mathbf{w} + (1-\alpha)\mathbf{u}) &\leq \alpha R(\mathbf{w}) + (1-\alpha)R(\mathbf{u}) - \frac{\alpha(1-\alpha)}{2} \|\mathbf{w} - \mathbf{u}\|_q^2 \\ &\leq \alpha R(\mathbf{w}) + (1-\alpha)R(\mathbf{u}) - \frac{\alpha(1-\alpha)}{2 \cdot 3} \|\mathbf{w} - \mathbf{u}\|_1^2 \end{aligned}$$

This tells us that R is $(1/3)$ -strongly-convex with respect to $\|\cdot\|_1$. \square

14 Stochastic Gradient Descent

1. Divide the definition of strong convexity by α and rearrange terms, to get that

$$\frac{f(\mathbf{w} + \alpha(\mathbf{u} - \mathbf{w})) - f(\mathbf{w})}{\alpha} \leq f(\mathbf{u}) - f(\mathbf{w}) - \frac{\lambda}{2}(1-\alpha)\|\mathbf{w} - \mathbf{u}\|^2 .$$

In addition, if \mathbf{v} be a subgradient of f at \mathbf{w} , then,

$$f(\mathbf{w} + \alpha(\mathbf{u} - \mathbf{w})) - f(\mathbf{w}) \geq \alpha \langle \mathbf{u} - \mathbf{w}, \mathbf{v} \rangle .$$

Combining together we obtain

$$\langle \mathbf{u} - \mathbf{w}, \mathbf{v} \rangle \leq f(\mathbf{u}) - f(\mathbf{w}) - \frac{\lambda}{2}(1-\alpha)\|\mathbf{w} - \mathbf{u}\|^2 .$$

Taking the limit $\alpha \rightarrow 0$ we obtain that the right-hand side converges to $f(\mathbf{w}) - f(\mathbf{u}) - \frac{\lambda}{2}\|\mathbf{w} - \mathbf{u}\|^2$, which concludes our proof.

2. Plugging the definitions of η and T into Theorem 14.13 we obtain

$$\begin{aligned} \mathbb{E}[L_{\mathcal{D}}(\bar{\mathbf{w}})] &\leq \frac{1}{1 - \frac{1}{1+3/\epsilon}} \left(L_{\mathcal{D}}(\mathbf{w}^*) + \frac{\|\mathbf{w}^*\|^2(1+3/\epsilon)\epsilon^2}{24B^2} \right) \\ &\leq (1+3/\epsilon)\epsilon/3 \left(L_{\mathcal{D}}(\mathbf{w}^*) + \frac{(1+3/\epsilon)\epsilon^2}{24} \right) \\ &= (1+\epsilon/3) \left(L_{\mathcal{D}}(\mathbf{w}^*) + \frac{\epsilon(\epsilon+3)}{24} \right) \\ &= L_{\mathcal{D}}(\mathbf{w}^*) + \frac{\epsilon}{3}L_{\mathcal{D}}(\mathbf{w}^*) + \frac{\epsilon}{3} \end{aligned}$$

Finally, since $L_{\mathcal{D}}(\mathbf{w}^*) \leq L_{\mathcal{D}}(\mathbf{0}) \leq 1$, we conclude the proof.

3. Perceptron as a sub-gradient descent algorithm:

- Clearly, $f(\mathbf{w}^*) \leq 0$. If there is strict inequality, then we can decrease the norm of \mathbf{w}^* while still having $f(\mathbf{w}^*) \leq 0$. But \mathbf{w}^* is chosen to be of minimal norm and therefore equality must hold. In addition, any \mathbf{w} for which $f(\mathbf{w}) < 1$ must satisfy $1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle < 1$ for every i , which implies that it separates the examples.
- A sub-gradient of f is given by $-y_i \mathbf{x}_i$, where $i \in \operatorname{argmax}\{1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle\}$.
- The resulting algorithm initializes \mathbf{w} to be the all zeros vector and at each iteration finds $i \in \operatorname{argmin}_i y_i \langle \mathbf{w}, \mathbf{x}_i \rangle$ and update $\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} + \eta y_i \mathbf{x}_i$. The algorithm must have $f(\mathbf{w}^{(t)}) < 0$ after $\|\mathbf{w}^*\|^2 R^2$ iterations. The algorithm is almost identical to the Batch Perceptron algorithm with two modifications. First, the Batch Perceptron updates with any example for which $y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle \leq 0$, while the current algorithm chooses the example for which $y_i \langle \mathbf{w}^{(t)}, \mathbf{x}_i \rangle$ is minimal. Second, the current algorithm employs the parameter η . However, it is easy to verify that the algorithm would not change if we fix $\eta = 1$ (the only modification is that $\mathbf{w}^{(t)}$ would be scaled by $1/\eta$).

4. **Variable step size:** The proof is very similar to the proof of Theorem 14.11. Details can be found, for example, in [1].

15 Support Vector Machines

1. Let \mathcal{H} be the class of halfspaces in \mathbb{R}^d , and let $S = ((\mathbf{x}_i, y_i))_{i=1}^m$ be a linearly separable set. Let $\mathcal{G} = \{(\mathbf{w}, b) : \|\mathbf{w}\| = 1, (\forall i \in [m]) y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 0\}$. Our assumptions imply that this set is non-empty. Note that for every $(\mathbf{w}, b) \in \mathcal{G}$,

$$\min_{i \in [m]} y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) > 0 .$$

On the contrary, for every $(\mathbf{w}, b) \notin \mathcal{G}$,

$$\min_{i \in [m]} y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \leq 0 .$$

It follows that

$$\operatorname{arg\,max}_{\substack{(\mathbf{w}, b): \\ \|\mathbf{w}\|=1}} \min_{i \in [m]} y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \subseteq \mathcal{G} .$$

Hence, solving the second optimization problem is equivalent to the following optimization problem:

$$\arg \max_{(\mathbf{w}, b) \in \mathcal{G}} \min_{i \in [m]} y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) .$$

Finally, since for every $(\mathbf{w}, b) \in \mathcal{G}$, and every $i \in [m]$, $y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) = |\langle \mathbf{w}, \mathbf{x}_i \rangle + b|$, we obtain that the second optimization problem is equivalent to the first optimization problem.

2. Let $S = ((x_i, y_i))_{i=1}^m \subseteq (\mathbb{R}^d \times \{-1, 1\})^m$ be a linearly separable set with a margin γ , such that $\max_{i \in [m]} \|x_i\| \leq \rho$ for some $\rho > 0$. The margin assumption implies that there exists $(\mathbf{w}, b) \in \mathbb{R}^d \times \mathbb{R}$ such that $\|w\| = 1$, and

$$(\forall i \in [m]) \quad y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \gamma .$$

Hence,

$$(\forall i \in [m]) \quad y_i(\langle \mathbf{w}/\gamma, \mathbf{x}_i \rangle + b/\gamma) \geq 1 .$$

Let $w^* = \mathbf{w}/\gamma$. We have $\|w^*\| = 1/\gamma$. Applying Theorem 9.1, we obtain that the number of iterations of the perceptron algorithm is bounded above by $(\rho/\gamma)^2$.

3. The claim is wrong. Fix some integer $m > 1$ and $\lambda > 0$. Let $\mathbf{x}_0 = (0, \alpha) \in \mathbb{R}^2$, where $\alpha \in (0, 1)$ will be tuned later. For $k = 1, \dots, m-1$, let $\mathbf{x}_k = (0, k)$. Let $y_0 = \dots = y_{m-1} = 1$. Let $S = \{(\mathbf{x}_i, y_i) : i \in \{0, 1, \dots, m-1\}\}$. The solution of hard-SVM is $\mathbf{w} = (0, 1/\alpha)$ (with value $1/\alpha^2$). However, if

$$\lambda \cdot 1 + \frac{1}{m}(1 - \alpha) \leq \frac{1}{\alpha^2} ,$$

the solution of soft-SVM is $\mathbf{w} = (0, 1)$. Since $\alpha \in (0, 1)$, it suffices to require that $\frac{1}{\alpha^2} > \lambda + 1/m$. Clearly, there exists $\alpha_0 > 0$ s.t. for every $\alpha < \alpha_0$, the desired inequality holds. Informally, if α is small enough, then soft-SVM prefers to “neglect” \mathbf{x}_0 .

4. Define the function $g : \mathcal{X} \rightarrow \mathbb{R}$ by $g(x) = \max_{y \in \mathcal{Y}} f(x, y)$. Clearly, for every $x \in \mathcal{X}$ and every $y \in \mathcal{Y}$,

$$g(x) \geq f(x, y) .$$

Hence, for every $y \in \mathcal{Y}$,

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} f(x, y) = \min_{x \in \mathcal{X}} g(x) \geq \min_{x \in \mathcal{X}} f(x, y) .$$

Hence,

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} f(x, y) \geq \max_{y \in \mathcal{Y}} \min_{x \in \mathcal{X}} f(x, y) .$$

16 Kernel Methods

1. Recall that \mathcal{H} is finite, and let $\mathcal{H} = \{h_1, \dots, h_{|\mathcal{H}|}\}$. For any two words u, v in Σ^* , we use the notation $u \preceq v$ if u is a substring of v . We will abuse the notation and write $h \preceq u$ if h is parameterized by a string v such that $v \preceq u$. Consider the mapping $\phi : \mathcal{X} \rightarrow \mathbb{R}^{|\mathcal{H}|+1}$ which is defined by

$$\phi(x)[i] = \begin{cases} 1 & \text{if } i = |\mathcal{H}| + 1 \\ 1 & \text{if } h_i \preceq x \\ 0 & \text{otherwise} \end{cases}$$

Next, each $h_j \in \mathcal{H}$ will be associated with $\mathbf{w}(h_j) := (2\mathbf{e}_j; -1) \in \mathbb{R}^{|\mathcal{H}|+1}$. That is, the first $|\mathcal{H}|$ coordinates of $\mathbf{w}(h_j)$ correspond to the vector \mathbf{e}_j , and the last coordinate is equal to -1 . Then, $\forall x \in \mathcal{X}$,

$$\langle \mathbf{w}, \phi(x) \rangle = 2[\phi(x)_j] - 1 = h_j(x) . \quad (8)$$

2. In the Kernelized Perceptron, the weight vector $\mathbf{w}^{(t)}$ will not be explicitly maintained. Instead, our algorithm will maintain a vector $\boldsymbol{\alpha}^{(t)} \in \mathbb{R}^m$. In each iteration we update $\boldsymbol{\alpha}^{(t)}$ such that

$$\mathbf{w}^{(t)} = \sum_{i=1}^m \alpha_i^{(t)} \psi(\mathbf{x}_i) . \quad (9)$$

Assuming that Equation (9) holds, we observe that the condition

$$\exists i \text{ s.t. } y_i \langle \mathbf{w}^{(t)}, \psi(\mathbf{x}_i) \rangle \leq 0$$

is equivalent to the condition

$$\exists i \text{ s.t. } y_i \sum_{j=1}^m \alpha_j^{(t)} K(\mathbf{x}_i, \mathbf{x}_j) \leq 0 ,$$

which can be verified while only accessing instances via the kernel function.

We will now detail the update $\alpha^{(t)}$. At each time t , if the required update is $\mathbf{w}_{t+1} = \mathbf{w}_t + y_i \mathbf{x}_i$, we make the update

$$\boldsymbol{\alpha}^{(t+1)} = \boldsymbol{\alpha}^{(t)} + y_i \mathbf{e}_i .$$

A simple inductive argument shows that Equation (9) is satisfied.

Finally, the algorithm returns $\alpha^{(T+1)}$. Given a new instance \mathbf{x} , the prediction is calculated using $\text{sign}(\sum_{i=1}^m \alpha_i^{(T+1)} K(\mathbf{x}_i, \mathbf{x}))$.

3. The representer theorem tells us that the minimizer of the training error lies in $\text{span}(\{\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_m)\})$. That is, the ERM objective is equivalent to the following objective:

$$\min_{\alpha \in \mathbb{R}^m} \lambda \left\| \sum_{i=1}^m \alpha_i \psi(\mathbf{x}_i) \right\|^2 + \frac{1}{2m} \sum_{i=1}^m \left(\left\langle \sum_{j=1}^m \alpha_j \psi(\mathbf{x}_j), \psi(\mathbf{x}_i) \right\rangle - y_i \right)^2$$

Denoting the gram matrix by G , the objective can be rewritten as

$$\min_{\alpha \in \mathbb{R}^m} \lambda \alpha^\top G \alpha + \frac{1}{2m} \sum_{i=1}^m (\langle \alpha, G_{\cdot, i} \rangle - y_i)^2. \quad (10)$$

Note that the objective (Equation (10)) is convex¹⁵. It follows that a minimizer can be obtained by differentiating Equation (10), and comparing to zero. Define $\lambda' = m \cdot \lambda$. We obtain

$$(\lambda' G + G G^\top) \alpha - G \mathbf{y} = 0$$

Since G is symmetric, this can be rewritten as

$$G(\lambda' I + G) \alpha = G \mathbf{y}.$$

A sufficient (and necessary in case that G is invertible) condition for the above to hold is that

$$(\lambda' I + G) \alpha = \mathbf{y}.$$

Since G is positive semi-definite and $\lambda' > 0$, the matrix $\lambda' I + G$ is positive definite, and thus invertible. We obtain that $\alpha^* = (\lambda' I + G)^{-1} \mathbf{y}$ is a minimizer of our objective.

4. Define $\psi : \{1, \dots, N\} \rightarrow \mathbb{R}^N$ by

$$\psi(j) = (\mathbf{1}^j; \mathbf{0}^{N-j}),$$

¹⁵ The term $\frac{\lambda}{2m} \sum_{i=1}^m (\langle \alpha, G_{\cdot, i} \rangle - y_i)^2$ is simply the least square objective, and thus it is convex, as we have already seen. The Hessian of $\alpha^\top G \alpha$ is G , which is positive semi-definite. Hence, $\alpha^\top G \alpha$ is also convex. Our objective is a weighted sum, with non-negative weights, of the two convex terms above. Thus, it is convex.

where $\mathbf{1}^j$ is the vector in \mathbb{R}^j with all elements equal to 1, and $\mathbf{0}^{N-j}$ is the zero vector in \mathbb{R}^{N-j} . Then, assuming the standard inner product, we obtain that $\forall (i, j) \in [N]^2$,

$$\langle \psi(i), \psi(j) \rangle = \langle (\mathbf{1}^i; \mathbf{0}^{N-i}), (\mathbf{1}^j; \mathbf{0}^{N-j}) \rangle = \min\{i, j\} = K(i, j) .$$

5. We will formalize our problem as an SVM (with kernels) problem. Consider the feature mapping $\phi : \mathcal{P}([d]) \rightarrow \mathbb{R}^d$ (where $\mathcal{P}([d])$ is the collection of all subsets of $[d]$), which is defined by

$$\phi(E) = \sum_{j=1}^d \mathbb{1}_{[j \in E]} \mathbf{e}_j .$$

In words, $\phi(E)$ is the indicator vector of E . A suitable kernel function $K : \mathcal{P}([d]) \times \mathcal{P}([d]) \rightarrow \mathbb{R}$ is defined by $K(E, E') = |E \cap E'|$. The prior knowledge of the manager implies that the optimal hypothesis can be written as a homogenous halfspace:

$$\mathbf{x} \mapsto \text{sign}(\langle 2\mathbf{w}, \phi(\mathbf{x}) \rangle - 1) ,$$

where $\mathbf{w} = \sum_{i \in I} \mathbf{e}_i$, where $I \subset [d]$, $|I| = k$, is the set of k relevant items. Furthermore, the halfspace defined by $(2\mathbf{w}, 1)$ has zero hinge-loss on the training set. Finally, we have that $\|(2\mathbf{w}, 1)\| = \sqrt{4k + 1}$, and $\|(\phi(\mathbf{x}), 1)\| \leq \sqrt{s + 1}$. We can therefore apply the general bounds on the sample complexity of soft-SVM, and obtain the following:

$$\begin{aligned} \mathbb{E}_S[L_D^{0-1}(A(S))] &\leq \min_{\mathbf{w}: \|\mathbf{w}\| \leq \sqrt{4k+1}} L_D^{\text{hinge}}(w) + \sqrt{\frac{8 \cdot (4k + 1) \cdot (s + 1)}{m}} \\ &= \sqrt{\frac{8 \cdot (4k + 1) \cdot (s + 1)}{m}} . \end{aligned}$$

Thus, the sample complexity is polynomial in $s, k, 1/\epsilon$. Note that according to the regulations, each evaluation of the kernel function K can be computed in $O(s \log s)$ (this is the cost of finding the common items in both carts). Consequently, the computational complexity of applying soft-SVM with kernels is also polynomial in $s, k, 1/\epsilon$.

6. We will work with the label set $\{\pm 1\}$.

Observe that

$$\begin{aligned}
h(\mathbf{x}) &= \text{sign}(\|\psi(\mathbf{x}) - c_-\|^2 - \|\psi(\mathbf{x}) - c_+\|^2) \\
&= \text{sign}(2\langle\psi(\mathbf{x}), c_+\rangle - 2\langle\psi(\mathbf{x}), c_-\rangle + \|c_-\|^2 - \|c_+\|^2) \\
&= \text{sign}(2(\langle\psi(\mathbf{x}), \mathbf{w}\rangle + b)) \\
&= \text{sign}(\langle\psi(\mathbf{x}), \mathbf{w}\rangle + b) .
\end{aligned}$$

(b) Simply note that

$$\begin{aligned}
\langle\psi(\mathbf{x}), \mathbf{w}\rangle &= \langle\psi(\mathbf{x}), c_+ - c_-\rangle \\
&= \frac{1}{m_+} \sum_{i:y_i=1} \langle\psi(\mathbf{x}), \psi(\mathbf{x}_i)\rangle + \frac{1}{m_-} \sum_{i:y_i=-1} \langle\psi(\mathbf{x}), \psi(\mathbf{x}_i)\rangle \\
&= \frac{1}{m_+} \sum_{i:y_i=1} K(\mathbf{x}, \mathbf{x}_i) + \frac{1}{m_-} \sum_{i:y_i=-1} K(\mathbf{x}, \mathbf{x}_i) .
\end{aligned}$$

17 Multiclass, Ranking, and Complex Prediction Problems

1. Fix some $(\mathbf{x}, y) \in S$. By our assumption about S (and by the triangle inequality),

$$\|\mathbf{x} - \boldsymbol{\mu}_y\| \leq r \quad , \quad (\forall y' \neq y) \|\mathbf{x} - \boldsymbol{\mu}_{y'}\| \geq 3r$$

Hence,

$$\|\mathbf{x} - \boldsymbol{\mu}_y\|^2 \leq r^2 \quad , \quad (\forall y' \neq y) \|\mathbf{x} - \boldsymbol{\mu}_{y'}\|^2 \geq 9r^2$$

It follows that for every $y' \neq y$,

$$\|\mathbf{x} - \boldsymbol{\mu}_{y'}\|^2 - \|\mathbf{x} - \boldsymbol{\mu}_y\|^2 = 2\langle\boldsymbol{\mu}_y, \mathbf{x}\rangle - \|\boldsymbol{\mu}_y\|^2 - (2\langle\boldsymbol{\mu}'_{y'}, \mathbf{x}\rangle - \|\boldsymbol{\mu}_{y'}\|^2) \geq 8r^2 > 0 .$$

Dividing by two, we obtain

$$\langle\boldsymbol{\mu}_y, \mathbf{x}\rangle - \frac{1}{2}\|\boldsymbol{\mu}_y\|^2 - (\langle\boldsymbol{\mu}'_{y'}, \mathbf{x}\rangle - \frac{1}{2}\|\boldsymbol{\mu}_{y'}\|^2) \geq 4r^2 > 0 .$$

Define \mathbf{w} as in the hint (that is, define $\mathbf{w} = [\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_k] \in \mathbb{R}^{(n+1)k}$, where each \mathbf{w}_i is defined by $\mathbf{w}_i = [\boldsymbol{\mu}_i, -\|\boldsymbol{\mu}_i\|^2/2]$). It follows that

$$\langle\mathbf{w}, \psi(\mathbf{x}, y)\rangle - \langle\mathbf{w}, \psi(\mathbf{x}, y')\rangle \geq 4r^2 > 0 .$$

Hence, $h_{\mathbf{w}}(\mathbf{x}) = y$, so $\ell(\mathbf{w}, (\mathbf{x}, y)) = 0$.

Then,

$$\begin{aligned}
\sum_{i=1}^r \tilde{v}_i y_i &= \sum_{i=1}^r \hat{v}_i y_i + (\tilde{v}_s - \hat{v}_s) y_s + (\tilde{v}_t - \hat{v}_t) y_t \\
&= \sum_{i=1}^r \hat{v}_i y_i + (s - \hat{v}_s) s + (\hat{v}_s - s) t \\
&= \sum_{i=1}^r \hat{v}_i y_i + (s - \hat{v}_s)(s - t) \\
&> \sum_{i=1}^r \hat{v}_i y_i .
\end{aligned}$$

Hence, we obtain a contradiction to the maximality of $\hat{\mathbf{v}}$.

5. (a) Averaging sensitivity and specificity, F1-score, F- β -score ($\theta = 0$):
Let $\mathbf{y}' \in \mathbb{R}^r$, and let $V = \{-1, 1\}^r$. Let $\hat{\mathbf{v}} = \operatorname{argmax}_{\mathbf{v} \in V} \sum_{i=1}^r v_i y'_i$.
We would like to show that $\hat{\mathbf{v}} = (\operatorname{sign}(y'_1), \dots, \operatorname{sign}(y'_r))$. Indeed,
for every $\mathbf{v} \in V$, we have

$$\sum_{i=1}^r v_i y'_i \leq \sum_{i=1}^r |v_i y'_i| = \sum_{i=1}^r \operatorname{sign}(y'_i) y'_i .$$

- (b) Recall at k , precision at k : The proof is analogous to the previous part.

18 Decision Trees

1. (a) Here is one simple (although not very efficient) solution: given h , construct a full binary tree, where the root node is $(x_1 = 0?)$, and all the nodes at depth i are of the form $(x_{i+1} = 0?)$. This tree has 2^d leaves, and the path from each root to the leaf is composed of the nodes $(x_1 = 0?)$, $(x_2 = 0?)$, \dots , $(x_d = 0?)$. It is not hard to see that we can allocate one leaf to any possible combination of values for x_1, x_2, \dots, x_d , with the leaf's value being $h(x) = h((x_1, x_2, \dots, x_d))$.
- (b) Our previous result implies that we can shatter the domain $\{0, 1\}^d$. Thus, the VC-dimension is exactly 2^d .
2. We denote by H the binary entropy.

- (a) The algorithm first picks the root node, by searching for the feature which maximizes the information gain. The information gain¹⁶ for feature 1 (namely, if we choose $x_1 = 0?$ as the root) is:

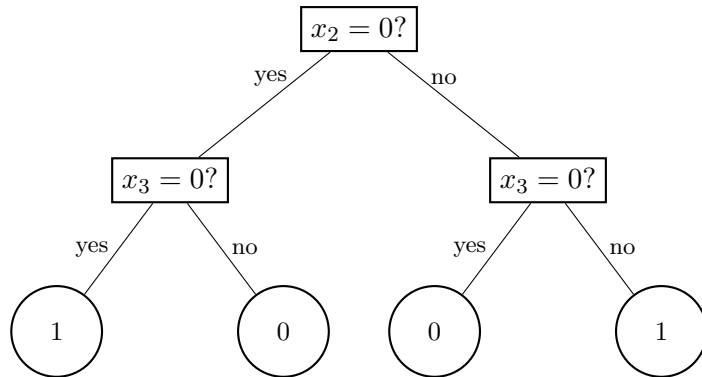
$$H\left(\frac{1}{2}\right) - \left(\frac{3}{4}H\left(\frac{2}{3}\right) + \frac{1}{4}H(0)\right) \approx 0.22$$

The information gain for feature 2, as well as feature 3, is:

$$H\left(\frac{1}{2}\right) - \left(\frac{1}{2}H\left(\frac{1}{2}\right) + \frac{1}{2}H\left(\frac{1}{2}\right)\right) = 0.$$

So the algorithm picks $x_1 = 0?$ as the root. But this means that the three examples $((1, 1, 0), 0)$, $((1, 1, 1), 1)$, and $((1, 0, 0), 1)$ go down one subtree, and no matter what question we'll ask now, we won't be able to classify all three examples perfectly. For instance, if the next question is $x_2 = 0?$ (after which we must give a prediction), either $((1, 1, 0), 0)$ or $((1, 1, 1), 1)$ will be mislabeled. So in any case, at least one example will be mislabeled. Since we have 4 examples in the training set, it follows that the training error is at least $1/4$.

- (b) Here is one such tree:



¹⁶Here we compute entropy where log is to the base of e . However, one can pick any other base, and the results will just change by a constant factor. Since we only care about which feature has the largest information gain, this won't affect which feature is picked.

19 Nearest Neighbor

1. We follow the hints for proving that for any $k \geq 2$, we have

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[\sum_{i: |C_i \cap S| < k} P[C_i] \right] \leq \frac{2rk}{m} .$$

The claim in the first hint follows easily from the linearity of the expectation and the fact that $\sum_{i: |C_i \cap S| \leq k} \mathbb{P}[C_i] = \sum_{i=1}^r \mathbb{1}_{[|C_i| \leq k]} \mathbb{P}[C_i]$. The claim in the second hint follows directly from Chernoff's bounds. The next hint leaves nothing to prove. Finally, combining the fourth hint with the previous hints,

$$\mathbb{E}_{S \sim \mathcal{D}^m} \left[\sum_{i: |C_i \cap S| < k} \mathbb{P}[C_i] \right] \leq \sum_{i=1}^r \max\{8/(me), 2k/m\} .$$

Since $k \geq 2$, our proof is completed.

2. The claim in the first hint follows from

$$\mathbb{E}_{Z_1, \dots, Z_k} \mathbb{P}_{y \sim p} [y \neq y'] = p \left(1 - \mathbb{P}_{Z_1, \dots, Z_k} [p' > 1/2] \right) + (1-p) \left(\mathbb{P}_{Z_1, \dots, Z_k} [p' > 1/2] \right) .$$

The next hints leave nothing to prove.

3. We need to show that

$$\mathbb{P}_{y \sim p} [y \neq y'] - \mathbb{P}_{y \sim p'} [y \neq y'] \leq |p - p'| \tag{15}$$

Indeed, if $y' = 0$, then the left-hand side of Equation (15) equals $p - p'$. Otherwise, it equals $p' - p$. In both cases, Equation (15) holds.

4. Recall that $\pi_j(\mathbf{x})$ is the j -th NN of \mathbf{x} . We prove the first claim. The probability of misclassification is bounded above by the probability that \mathbf{x} falls in a “bad cell” (i.e., a cell that does not contain k instances from the training set) plus the probability that x is misclassified given that \mathbf{x} falls in a “good cell”. The next hints leave nothing to prove.

20 Neural Networks

1. Let $\epsilon > 0$. Following the hint, we cover the domain $[-1, 1]^n$ by disjoint boxes such that for every x, x' which lie in the same box, we have $|f(\mathbf{x}) - f(\mathbf{x}')| \leq \epsilon/2$. Since we only aim at approximating f to an accuracy of ϵ , we can pick an arbitrary point from each box. By picking the set of representative points appropriately (e.g., pick the center of each box), we can assume w.l.o.g. that f is defined over the discrete set $[-1 + \beta, -1 + 2\beta, \dots, 1]^d$ for some $\beta \in [0, 2]$ and $d \in \mathbb{N}$ (which both depends on ρ and ϵ). From here, the proof is straightforward. Our network should have two hidden layers. The first layer has $(2/\beta)^d$ nodes which correspond to the intervals that make up our boxes. We can adjust the weights between the input and the hidden layer such that given an input \mathbf{x} , the output of each neuron is close enough to 1 if the corresponding coordinate of \mathbf{x} lies in the corresponding interval (note that given a finite domain, we can approximate the indicator function using the sigmoid function). In the next layer, we construct a neuron for each box, and add an additional neuron which outputs the constant $-1/2$. We can adjust the weights such that the output of each neuron is 1 if \mathbf{x} belongs to the corresponding box, and 0 otherwise. Finally, we can easily adjust the weights between the second layer and the output layer such that the desired output is obtained (say, to an accuracy of $\epsilon/2$).
2. Fix some $\epsilon \in (0, 1)$. Denote by \mathcal{F} the set of 1-Lipschitz functions from $[-1, 1]^n$ to $[-1, 1]$. Let $G = (V, E)$ with $|V| = s(n)$ be a graph such that the hypothesis class $\mathcal{H}_{V,E,\sigma}$, with σ being the sigmoid activation function, can approximate every function $f \in \mathcal{F}$, to an accuracy of ϵ . In particular, every function that belongs to the set $\{f \in \mathcal{F} : (\forall x \in \{-1, 1\}^n) f(\mathbf{x}) \in \{-1, 1\}\}$ is approximated to an accuracy ϵ . Since $\epsilon \in (0, 1)$, it follows that we can easily adapt the graph such that its size remains $\Theta(s(n))$, and $\mathcal{H}_{V,E,\sigma}$ contains all the functions from $\{-1, 1\}^n$ to $\{-1, 1\}$. We already noticed that in this case, $s(n)$ must be exponential in n (see Theorem 20.2).
3. Let $C = \{c_1, \dots, c_m\} \subseteq \mathcal{X}$. We have

$$\begin{aligned}
 |\mathcal{H}_C| &= |\{((f_1(\mathbf{c}_1), f_2(\mathbf{c}_2)), \dots, (f_1(\mathbf{c}_m), f_2(\mathbf{c}_m))) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}| \\
 &= |\{((f_1(\mathbf{c}_1), \dots, f_1(\mathbf{c}_m)), (f_2(\mathbf{c}_1), \dots, f_2(\mathbf{c}_m))) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}| \\
 &= |\mathcal{F}_{1C} \times \mathcal{F}_{2C}| \\
 &= |\mathcal{F}_{1C}| \cdot |\mathcal{F}_{2C}|.
 \end{aligned}$$

It follows that $\tau_{\mathcal{H}}(m) = \tau_{\mathcal{F}_2}(m)\tau_{\mathcal{F}_1}(m)$.

4. Let $C = \{\mathbf{c}_1, \dots, \mathbf{c}_m\} \subseteq \mathcal{X}$. We have

$$\begin{aligned} |\mathcal{H}_C| &= |\{f_2(f_1(\mathbf{c}_1)), \dots, f_2(f_1(\mathbf{c}_m)) : f_1 \in \mathcal{F}_1, f_2 \in \mathcal{F}_2\}| \\ &= \left| \bigcup_{f_1 \in \mathcal{F}_1} \{(f_2(f_1(\mathbf{c}_1)), \dots, f_2(f_1(\mathbf{c}_m))) : f_2 \in \mathcal{F}_2\} \right| \\ &\leq |\mathcal{F}_{1C}| \cdot \tau_{\mathcal{F}_2}(m) \\ &\leq \tau_{\mathcal{F}_1}(m)\tau_{\mathcal{F}_2}(m) . \end{aligned}$$

It follows that $\tau_{\mathcal{H}}(m) = \tau_{\mathcal{F}_2}(m)\tau_{\mathcal{F}_1}(m)$.

5. The hints provide most of the details. We skip to the conclusion part. By combining the graphs above, we obtain a graph $G = (V, E)$ with $V = O(n)$ such that the set $\{a_j^{(i)}\}_{(i,j) \in [n]^2}$ is shattered. Hence, the VC-dimension is at least n^2 .
6. We reduce from the k -coloring problem. The construction is very similar to the construction for intersection of halfspaces. The same set of points (namely $\{e_1, \dots, e_n\} \cup \{(e_i + e_j)/2 : \{i, j\} \in E, i < j\}$) is considered. Similarly to the case of intersection of halfspaces, it can be verified that the graph is k -colorable iff $\min_{h \in \mathcal{H}_{V,E,\text{sign}}} L_s(h) = 0$. Hence, the k -coloring problem is reduced to the problem of minimizing the training error. The theorem is concluded.

21 Online Learning

1. Let $\mathcal{X} = \mathbb{R}^d$, and let $\mathcal{H} = \{h_1, \dots, h_d\}$, where $h_j(\mathbf{x}) = \mathbb{1}_{[x_j=1]}$. Let $\mathbf{x}_t = \mathbf{e}_t$, $y_t = \mathbb{1}_{[t=d]}$, $t = 1, \dots, d$. The Consistent algorithm might predict $p_t = 1$ for every $t \in [d]$. The number of mistakes done by the algorithm in this case is $d - 1 = |\mathcal{H}| - 1$.
2. Let $d \in \mathbb{N}$, and let $\mathcal{X} = [d]$ and let $\mathcal{H} = \{h_A : A \subseteq [d]\}$, where

$$h_A(x) = \mathbb{1}_{[x \in A]} .$$

For $t = 1, 2, \dots$, let $x_t = t$, $y_t = 1$ (i.e., the true hypothesis corresponds to the set $[d]$). Note that at every time t ,

$$|\{h \in V_t : h(x_t) = 1\}| = |\{h \in V_t : h(x_t) = 0\}| ,$$

where the second equality follows from the fact that h_t only depends on x_1, \dots, x_{t-1} , and x_t is independent of x_1, \dots, x_{t-1} .

22 Clustering

1. Let $\mathbb{R}^2 \subseteq X = \{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4\}$, where $x_1 = (0, 0), x_2 = (0, 2), x_3 = (2\sqrt{t}, 0), x_4 = (2\sqrt{t}, 2)$. Let d be the metric induced by the ℓ^2 norm. Finally, let $k = 2$.

Suppose that the k -means chooses $\mu_1 = \mathbf{x}_1, \mu_2 = \mathbf{x}_2$. Then, in the first iteration it associates $\mathbf{x}_1, \mathbf{x}_3$ with the center $\mu_1 = (\sqrt{t}, 0)$. Similarly, it associates $\mathbf{x}_2, \mathbf{x}_4$ with the center $\mu_2 = (\sqrt{t}, 2)$. This is a convergence point of the algorithm. The value of this solution is $4t$. The optimal solution associates $\mathbf{x}_1, \mathbf{x}_2$ with the center $(0, 1)$, while $\mathbf{x}_3, \mathbf{x}_4$ are associated with the center $(2\sqrt{t}, 1)$. The value of this solution is 4.

2. The K-means solution is:

$$\mu_1 = 2, C_1 = \{1, 2, 3\}, \mu_2 = 4, C_2 = \{3, 4\} .$$

The value of this solution is $2 \cdot 1 = 2$. The optimal solution is

$$\mu_1 = 1.5, C_1 = \{1, 2\}, \mu_2 = 3.5, C_2 = \{3, 4\}$$

whose value is $1 = 4 \cdot (1/2)^2$.

3. For every $j \in [k]$, let $r_j = d(\mu_j, \mu_{j+1})$. Following the notation in the hint, it is clear that $r_1 \geq r_2 \geq \dots \geq r_k \geq r$. Furthermore, by definition of μ_{k+1} it holds that

$$\max_{j \in [k]} \max_{x \in \hat{C}_j} d(x, \mu_j) \leq r .$$

The triangle inequality implies that

$$\max_{j \in [k]} \text{diam}(\hat{C}_j) \leq 2r .$$

The pigeonhole principle implies now that at least 2 of the points μ_1, \dots, μ_{k+1} lie in the same cluster in the optimal solution. Thus,

$$\max_{j \in [k]} \text{diam}(C_j^*) \geq r .$$

4. Let $k = 2$. The idea is to pick elements in two close disjoint balls, say in the plane, and another distant point. The k -diam optimal solution would be to separate the distant point from the two balls. If the number of points in the balls is large enough, then the optimal center-based solution would be to separate the balls, and associate the distant point with its closest ball.

Here are the details. Let $\mathcal{X}' = \mathbb{R}^2$ with the metric d which is induced by the ℓ_2 -norm. A subset $\mathcal{X} \subseteq \mathcal{X}'$ is constructed as follows: Let $m > 0$, and let \mathcal{X}_1 be set of m points which are evenly distributed on the sphere of the ball $B_1((2, 0))$ (a ball of radius 1 around $(2, 0)$). Similarly, let \mathcal{X}_2 be a set of m points which are evenly distributed on the sphere of $B_1((-2, 0))$. Finally, let $\mathcal{X}_3 = \{(0, y)\}$ for some $y > 0$, and set $\mathcal{X} = \cup_{i=1}^3 \mathcal{X}_i$. Fix any monotone function $f : \mathbb{R}_+ \rightarrow \mathbb{R}_+$. We note that for large enough m , an optimal center-based solution must separate between \mathcal{X}_1 and \mathcal{X}_2 , and associate the point $(0, y)$ with the nearest cluster. However, for large enough y , an optimal k -diam solution would be to separate \mathcal{X}_3 from $\mathcal{X}_1 \cup \mathcal{X}_2$.

5. (a) Single Linkage with fixed number of clusters:
- i. Scale invariance: Satisfied. multiplying the weights by a positive constant does not affect the order in which the clusters are merged.
 - ii. Richness: Not satisfied. The number of clusters is fixed, and thus we can not obtain all possible partitions of \mathcal{X} .
 - iii. Consistency: Satisfied. Assume by contradiction that d, d' are two metrics which satisfy the required properties, but $F(\mathcal{X}, d') \neq F(\mathcal{X}, d)$. Then, there are two points x, y which belong to the same cluster in $F(x, d')$, but belong to different clusters in $F(\mathcal{X}, d)$ (we rely here on the fact that the number of clusters is fixed). Since d' assigns to this pair a larger value than d (and also assigns smaller values to pairs which belong to the same cluster), this is impossible.
- (b) Single Linkage with fixed distance upper bound:
- i. Scale invariance: Not satisfied. In particular, by multiplying by an appropriate scalar, we obtain the trivial clustering (each cluster consists of a single data point).
 - ii. Richness: Satisfied. Easily seen.
 - iii. Consistency: Satisfied. The argument is almost identical to the argument in the previous part: Assume by contradiction

that d, d' are two metrics which satisfy the required properties, but $F(\mathcal{X}, d') \neq F(\mathcal{X}, d)$. Then, either there are two points x, y which belong to the same cluster in $F(x, d')$, but belong to different clusters in $F(\mathcal{X}, d)$ or there are two points x, y which belong to different clusters in $F(x, d')$, but belong to the same cluster in $F(\mathcal{X}, d)$. Using the relation between d and d' and the fact that the threshold is fixed, we obtain that both of these events are impossible.

- (c) Consider the scaled distance upper bound criterion (we set the threshold r to be $\alpha \max\{d(x, y) : x, y \in \mathcal{X}\}$). Let us check which of the three properties are satisfied:
- i. Scale invariance: Satisfied. Since the threshold is scale-invariant, multiplying the metric by a positive constant doesn't change anything.
 - ii. Richness: Satisfied. Easily seen.
 - iii. Consistency: Not Satisfied. Let $\alpha = 0.75$. Let $\mathcal{X} = \{x_1, x_2, x_3\}$, and equip it with the metric d which is defined by: $d(x_1, x_2) = 3, d(x_2, x_3) = 7, d(x_1, x_3) = 8$. The resulting partition is $\mathcal{X} = \{x_1, x_2\}, \{x_3\}$. Now we define another metric d' by: $d(x_1, x_2) = 3, d(x_2, x_3) = 7, d(x_1, x_3) = 10$. In this case the resulting partition is $\mathcal{X} = \mathcal{X}$. Since d' satisfies the required properties, but the partition is not preserved, it follows that consistency is not satisfied.

Summarizing the above, we deduce that any pair of properties can be attained by one of the single linkage algorithms detailed above.

6. It is easily seen that Single Linkage with fixed number of clusters satisfies the k -richness property, thus it satisfies all the mentioned properties.

23 Dimensionality Reduction

1. (a) A fundamental theorem in linear algebra states that if V, W are finite dimensional vector spaces, and let T be a linear transformation from V to W , then the image of T is a finite-dimensional subspace of W and

$$\dim(V) = \dim(\text{null}(T)) + \dim(\text{image}(T)).$$