

Machine Learning Problemset 3

Aimee Barciauskas

26 February 2016

Problem 12

Consider the class \mathcal{A} of all sets of the form:

$$A_\alpha = \{x \in \mathbb{R} : \sin(\alpha x) > 0\}$$

where $\alpha > 0$. What is the VC dimension of \mathcal{A} ? (Note that \mathcal{A} has one free parameter.)

The VC dimension of \mathcal{A} is infinite.

For example, if $x_i = 2^{-i}, i = 1, \dots, m$ are assigned arbitrary labels $(y_1, \dots, y_m) \in \{-1, 1\}^m$, α may be chosen such that any label set is correctly classified:

$$\alpha = \pi(1 + \sum_{i=1}^m 2^i \frac{1 - y_i}{2})$$

Problem 13

Let $\mathcal{A}_1, \dots, \mathcal{A}_k$ be classes of sets, all of the with VC dimension at most V . Show that the VC-dimension of $\cup_{i=1}^k \mathcal{A}_i$ is at most $4V \log_2(2V) + 4k$. You may use the fact that for $a \geq 1$ and $b > 0$, if $x \geq 4a \log(2a) + 2b$ then $x \geq a \log x + b$.

Can you bound the VC dimension of the class of all sets of the form:

$$A_1 \cup \dots \cup A_k \text{ with } A_1 \in \mathcal{A}_1, \dots, A_k \in \mathcal{A}_k$$

Part 1 Solution:

$$s_{\mathcal{A}} = 2^u \text{ such that } u \text{ is the VC dimension of } \mathcal{A}$$

Sauer's Lemma:

$$s_{\mathcal{A}} \leq s_{\mathcal{A}_1} + s_{\mathcal{A}_2} + \dots + s_{\mathcal{A}_k} \leq k(u+1)^V$$

If there exists a u which satisfies:

$$s_{\mathcal{A}} = 2^u$$

$$2^u \leq k(u+1)^V$$

$$u \leq \log(k) + V \log(u+1)$$

$$= \log(k) + V\log(u+1)$$

$$\log(k) \leq 2k$$

$$= 2k + V\log(u+1)$$

Using $\log(u+1) \approx \log(u)$

$$u \leq 2k + V\log(u)$$

Using the hints from the problem where $a = V$, $x = u$ and $b = 2k$ this can be re-written as:

$$u \leq 4V\log(2V) + 4k$$

Part 2 Solution:

For the union of sets of $A_1 \cup \dots \cup A_k$:

$$s_{\mathcal{A}}(n) \leq s_{\mathcal{A}_1}(n) \times s_{\mathcal{A}_1}(n) \times \dots \times s_{\mathcal{A}_k}(n) \leq (n+1)^{V_k}$$

$$s_{\mathcal{A}}(u) = 2^u \leq (u+1)^{V_k}$$

$$u \leq V_k \log(u+1)$$

Let $a = V_k$, $b = \epsilon$, $x = u+1$, where ϵ is some small number:

$$u \leq 4V_k \log(2V_k)$$

Problem 14

$$\|w_t - w_*\|^2 \leq \|w_{t-1} - w_*\|^2 - 1$$

We have $w_t = w_{t-1} + \frac{Y_t X_t}{\|X_t\|}$, so the first term can be re-written as:

$$\begin{aligned} \|w_{t-1} - w_* + \frac{Y_t X_t}{\|X_t\|}\|^2 &= \|w_{t-1} - w_*\|^2 + \left(\frac{Y_t X_t}{\|X_t\|}\right)^2 + 2(w_{t-1} - w_*) \frac{Y_t X_t}{\|X_t\|} \\ &= \|w_{t-1} - w_*\|^2 + \left(\frac{Y_t X_t}{\|X_t\|}\right)^2 + 2w_{t-1} \frac{Y_t X_t}{\|X_t\|} - 2w_* \frac{Y_t X_t}{\|X_t\|} \end{aligned}$$

The first term above equals the first term in the RHS of our initial inequality to be proved.

The second term is 1, so we subtract it from the RHS and get an equality of the first two terms of the expanded LHS and RHS.

The last two terms formulate the inequality: We know the second to last term $2w_{t-1} \frac{Y_t X_t}{\|X_t\|} < 0$ when the perceptron makes no more updates. and the last term $2w_* \frac{Y_t X_t}{\|X_t\|} \geq 1$. Something negative minus something

positive is negative, so the whole term is negative. Adding this negative term to the other side we get the inequality:

$$\|w_t - w_*\|^2 \leq \|w_{t-1} - w_*\|^2 - 1$$

Using this inequality iteratively:

$$\|w_{t-1} - w_*\|^2 \leq \|w_{t-2} - w_*\|^2 - 2$$

$$\|w_{t-2} - w_*\|^2 \leq \|w_{t-3} - w_*\|^2 - 3$$

...

$$\|w_* - w_*\|^2 \leq \|w_0 - w_*\|^2 - k$$

where k is the number of steps, and the LHS is now 0:

$$k \leq \|w_0 - w_*\|^2$$

The number of steps is less than or equal to the $\|w_0 - w_*\|^2$.

Problem 15

Part 1

The expected risk of the data-dependent leave-one-out classifier:

$$\mathbb{E}\left\{R_n^D(g_n)\right\} = \mathbb{E}\left\{\frac{1}{n} \sum_{i=1}^n \mathbb{I}_{g_{n-1}(X_i, D_{n,i}) \neq Y_i}\right\}$$

The RHS is a random variable: it is the sum of the variations on the random data set being used to train the leave-one-out classifier. The law of iterated expectations shows this can be written as:

$$\mathbb{E}\left\{R_n^D(g_n)\right\} = \mathbb{E}\left\{\mathbb{E}\{R(g_{n-1})|D\}\right\} = \mathbb{E}\{R(g_{n-1})\}$$

Part 2

$$\mathbb{E}\left\{R_n^D(g_n)\right\} = \mathbb{E}\{R(g_{n-1})\}$$

can be used to bound the expected risk of the perceptron classifier. The expected risk of the leave-one-out classifier times n is the number of updates made during the perceptron classifier, (otherwise the x_i in the iteration is correctly classified by the algorithm and no update is made). The estimation of the risk based on the leave one out estimator is the same as M the number of iterations made by the perceptron algorithm. This M is upper-bounded (Novikoff) by $(\frac{R}{\gamma})^2$, where R is the distance to the furthest point and γ is the size of the margin. This is upperbounded by the number of steps bounded in Problem 14.

Problem 16

Consider the majority classifier:

$$g_n(x, D_n) = \begin{cases} 1, & \text{if } \sum_{i=1}^n Y_i \geq \frac{n}{2} \\ 0, & \text{otherwise} \end{cases}$$

(Thus, g_n ignores x and the X_i 's.) Assume that n is odd. What is the expected risk $\mathbb{E}R(g_n) = \mathbb{P}\{g_n(X) \neq Y\}$ of this classifier? Study the performance of the leave-one-out error estimate. Show that for some distributions $\text{Var}(R_n^D(g_n)) \geq c/\sqrt{n}$ for some constant c . *Hint:* Strang things happen when the number of 0's and 1's is about the same in the data.

Solution Part 1:

Let N_n be the number of $Y_i = 0$ in the sample. So N_n is binomial $(n, 1 - p)$ with p the $PY = 1$.

$$\mathbb{E}(R(g_n)) = p\mathbb{P}\left\{N_n \geq \frac{n}{2}\right\} + (1 - p)\mathbb{P}\left\{N_n < \frac{n}{2}\right\} \rightarrow \min(p, 1 - p)$$

Solution Part 2:

Minimal variance of the leave-one-out classifier can be estimated when the true probability of either class is $\mathbb{P}\{Y = 1\} = \mathbb{P}\{Y = 0\} = \frac{1}{2}$ and the probability the data set comes near the true probability $\frac{n+1}{2}$ (given n is odd). This probability can be represented as:

$$\mathbb{P}\left\{\text{Bin}\left(n, \frac{1}{2}\right) = \frac{n+1}{2}\right\}$$

This probability is (assuming $n \geq 1$):

$$= 2^{-n} \binom{n}{\frac{n+1}{2}}$$

Which is always greater than $\frac{c}{\sqrt{(n)}}$ where c is some positive constant, say $\frac{1}{n+1}$.