# VC dimension

We have seen PAC and Mistake bound algorithms for many concept classes.

In the case of halfspaces there was a $\frac{1}{\gamma^2}$ dependence on the margin $\gamma$.

In fact, one can make this $\log \frac{1}{\gamma}$.
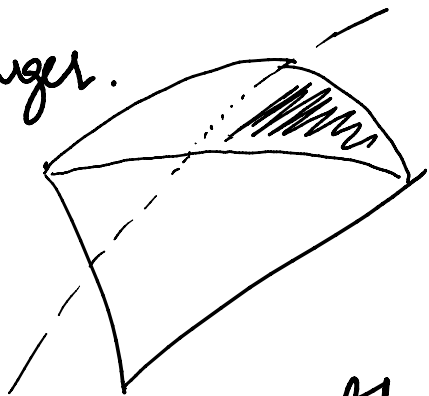
Suppose we predict majority of all surviving $w$.

i.e. suppose after examples $x^1, x^2, \ldots x^\ell$,

we have $W = \{ w : w^T x^i \geq 0, \|w\| \leq 1 \}$

as candidates and we consider which of

$$W \cap \{ w : w^T x^{\ell+1} \geq 0 \}, \quad W \cap \{ w : w^T x^{\ell+1} < 0 \}$$

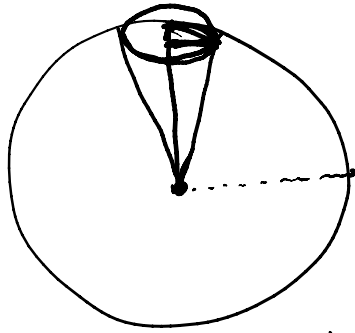is larger.



Predict according to that.

Then in each step we eliminate $\frac{1}{2}$ the volume.

$vol(W)$ starts at $vol(B)$.

vol($W$) starts at vol($0$).
At the end it is at least vol($\gamma$-cone)



$$\text{vol}(B) = \int_0^1 \left(\sqrt{1-t^2}\right)^{n-1} \text{vol}(B^{n-1})\, dt$$

$$\text{vol}(\gamma\text{-cap}) = \int_{\sqrt{1-\gamma^2}}^1 \left(\sqrt{1-t^2}\right)^{n-1} \text{vol}(B^{n-1})\, dt$$

$$\frac{\text{vol}(\gamma\text{-cap})}{\text{vol}(B)} = \frac{\int_{\sqrt{1-\gamma^2}}^1 (1-t^2)^{\frac{n-1}{2}}\, dt}{\int_0^1 (1-t^2)^{\frac{n-1}{2}}\, dt} \geq c.\gamma^n.$$

$$\therefore \#\text{mistakes} = O\left(n \log \frac{1}{\gamma}\right).$$

---

How to estimate volume fraction?
We can sample the current $W$ and take the majority vote of the sample. Even 1 sample suffices!

---

Alternatively we can use Linear Programming to find a feasible $w$ for all constraints so far.
To bound the number of examples needed we can use a more general theory.

a more general theory.

---

VC-dimension.    m points.

Concept class H.

How many distinct subsets of m points are defined by $h \in H$?    $H[m] \leq m^{VC-dim.}$

More precisely: $\underset{(H)}{VC\text{-dim}} = \max m$ s.t. $\exists$ m points that can be shattered, i.e. split in all possible ways by H.

e.g. intervals on a line     VC-dim = 2

rectangles in 2-d     VC-dim = 4

Halfspaces in $\mathbb{R}^d$     VC-dim = d+1

---

Thm1. For a concept class H of VC-dim d, # distinct ways to split m points using $h \in H$ is $\leq m^d$.

Thm2. # examples needed to $(\varepsilon, \delta)$-PAC learn H is $\leq \frac{2}{\varepsilon} \left( \log(2H[2m]) + \log \frac{1}{\delta} \right)$.

$$i_0 \leq \frac{2}{\varepsilon}\left(\log(2H[2m]) + \log \frac{1}{\delta}\right).$$

$$= O\left(\frac{1}{\varepsilon}\left(d\log \frac{1}{\varepsilon} + \log \frac{1}{\delta}\right)\right).$$

---

Pf (Thm 1). We will show $H[m] \leq \sum_{i=0}^{d} \binom{m}{i} = \binom{m}{\leq d}$.

Let $S$ be a set of $m$ points.

Induction on $m$. True for $m \leq d$.

Let $x \in S$. Consider $S \setminus \{x\}$.

By inductn $H(S \setminus \{x\}) \leq \binom{m-1}{\leq d}$.

Also note that

$$\binom{m}{\leq d} = \binom{m-1}{\leq d} + \binom{m-1}{\leq d-1}$$

So it suffices to show that

$$H(S) - H(S \setminus \{x\}) \leq \binom{m-1}{\leq d-1}$$

How can $H(S)$ be larger? there must be labelings $h$ and $h'$ s.t. they agree on all

labelings $h$ and $h'$ s.t. array agree
points except $x$.

Let $T = \{h \in H(S) : h(x) = 1, \; h' \in H(S)\}$.

Then we are interested in bounding $|T|$.

Let VC-dim $(T) = d'$. So $2^{d'}$ points can
be shattered by $T$. But then $d'+1$ points can
be shattered by $H$. So $d'+1 \leq d$.
i.e. $d' \leq d-1$.

Hence $H(T) \leq \binom{m-1}{\leq d-1}$.

---

Pf (Thm 2). We find a hypothesis $h_S$ that
correctly classifies $m$ points. We want to show
that with prob $\geq 1-\delta$

$$\Pr_{D}\left(h_S(x) \neq h(x)\right) \leq \varepsilon.$$

> Let $A$ be the
> complement of
> this event

Consider a different setting where we pick
2 subsets of size $m$, say

$S$ and $S'$. Let $B$ be the event that a hypothesis $h$ has zero error on $S$ and error $> \frac{\varepsilon}{2}$ on $S'$.

Claim. $\Pr(B) \geq \frac{1}{2} \Pr(A)$.

$$\Pr(B) = \Pr(A) \, \Pr(B/A)$$

$\Pr(B/A)$: $\Pr(h$ has error $\geq \frac{\varepsilon}{2}$ on $m$ points given that it has error $\geq \varepsilon$ on $D)$

This is a simple chernoff bound.

$$\Pr\left( \sum X_i - \mathbb{E}\left(\sum X_i\right) < \delta \, \mathbb{E}\left(\sum X_i\right)\right) \leq e^{-\frac{\delta^2 \mathbb{E}(\sum X_i)}{2}}.$$

$$\Pr\left( \sum X_i < \frac{\varepsilon}{2} m \right) \qquad \leq e^{-\frac{\varepsilon m}{8}}.$$

$$m \geq \frac{8}{\varepsilon} \implies \Pr(B/A) \geq \frac{1}{2}$$

___

So we want to show $\Pr(B) \leq \frac{8}{2}$.

For this we pick $2m$ points, $S''$ partition them randomly into two subsets $S, S'$ of $m$ points.

randomly into two subsets $S, S'$ of $m$ points.

Then we want to bound $Pr\left(err_h(S)=0, err_h(S') > \frac{\varepsilon}{2}\right)$.

Pair up the $2m$ points $(a_1, b_1), \ldots (a_m, b_m)$.

Fix hypothesis $h$.

if $h$ makes error on both $a_i$ and $b_i$,

then $Pr = 0$. (since no errors allowed on $S$).

Also at least $\frac{\varepsilon m}{2}$ indices $i$ must make an error.

So $Pr\left(\text{all } \frac{\varepsilon m}{2} \text{ errors fall in } S'\right) \leq \dfrac{1}{2^{\varepsilon m/2}}$.

# possible $h \leq \mathcal{H}(2m)$

$\therefore$ suffices to have $2^{-\varepsilon m/2} \mathcal{H}[2m] \leq \dfrac{\delta}{2}$

i.e. $m \geq \dfrac{2}{\varepsilon}\left(\log 2 \mathcal{H}[2m] + \log \frac{1}{\delta}\right)$.

---

## Chernoff

$X = \sum X_i$ independent $0/1$

$Pr\left(X \geq (1+\delta) \mathbb{E}(X)\right) < e^{-\frac{\delta^2}{2+\delta}\mathbb{E}(X)}$

$Pr\left(X \leq (1-\delta) \mathbb{E}(X)\right) < e^{-\frac{\delta^2 \mathbb{E}(X)}{2}}$

2

$\Pr(X \le (1-\delta)\mathbb{E}(x)) < c$

**Hoeffding** $a \le X_i \le b$

$$\Pr(X \ge \mathbb{E}(x) + t) < e^{-\frac{2t^2}{n(b-a)^2}}$$

$$\Pr(x \le \mathbb{E}(x) - t) < e^{-\frac{2t^2}{n(b-a)^2}}$$

---

Pf (Thm 3).   pairs $(a_i, b_i)$   randomly allocate to $S, S'$.

$$|err_h(S) - err_h(S')| \ge \frac{\varepsilon}{2}$$

$\Pr(S'$ gets $\frac{\varepsilon}{2}m$ more than $S)$

$\left( X_i = \begin{cases} 1 & \text{if } S' \\ -1 & \text{if } S \end{cases} \right.$   $\mathbb{E}(\Sigma X_i) = 0$

$\left. \Pr\left( \Sigma X_i \gg \frac{\varepsilon m}{2} \right) < e^{-2 \cdot \frac{\varepsilon^2 m}{4 \cdot 4}} \right.$

$$= e^{-\frac{\varepsilon^2 m}{8}}.$$

$$e^{-\frac{\varepsilon^2 m}{8}} \cdot \mathcal{H}(2m) \le \frac{\delta}{2}$$

$$\Rightarrow m \ge \frac{8}{\varepsilon^2}\left( \log 2\mathcal{H}(2m) + \log\frac{1}{\delta} \right)$$

$$\overline{\varepsilon^2}\,\big\rangle$$

suffices.

VC-dim $d$ : $\quad m = O\left(\dfrac{1}{\varepsilon^2}\left(d \log \dfrac{1}{\varepsilon} + \log \dfrac{1}{\delta}\right)\right)$

    suffices.