

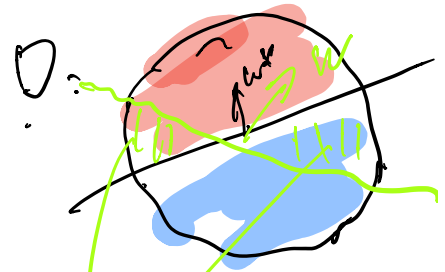
PAC learning: toy example



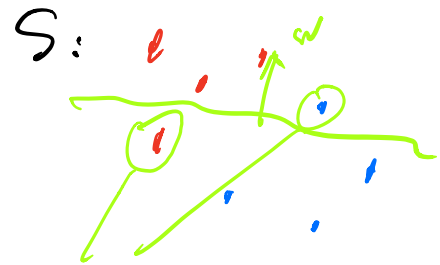
$$X = \mathbb{R}^2$$

$$y = \{\pm 1\}$$

$$h_w(x) = \text{sign}(w \cdot x)$$



$$\text{err}(h_w, D)$$



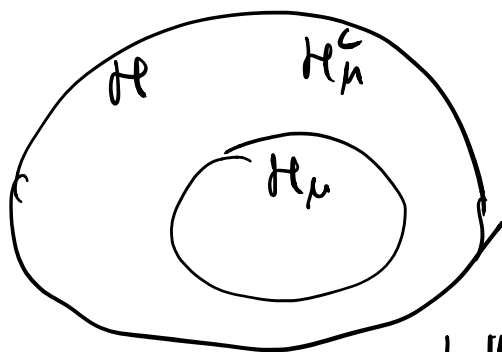
$$\begin{aligned} \text{err}(h_w, S) &= \frac{2}{8} \\ &= 0.25 \end{aligned}$$

Theorem: For finite \mathcal{H} . If the consistency alg. is given m iid training examples from \mathcal{D} realizable w.r.t. \mathcal{H} , then

with probability $1-\delta$. Its output \hat{h} is such that
 $\text{err}(\hat{h}, \mathcal{D}) \leq \underbrace{\sum_{\mathcal{R}}(m, \delta, |\mathcal{H}|)}_{\text{realizable}} = \frac{\ln|\mathcal{H}| + \ln\frac{1}{\delta}}{m} \leq \epsilon$

In other words. It ^{PAC} learns \mathcal{H} with sample complexity
 $f(m, \delta) = \frac{1}{\epsilon} (\ln|\mathcal{H}| + \ln\frac{1}{\delta})$.

1. error bound decreases polynomially with m .
 2. " " depends logarithmically with $|\mathcal{H}|$
 3. " " depends logarithmically with $\frac{1}{\delta}$
- $\delta = 10^{-5}$



$$H_\mu = \left\{ h \in H : \underline{\text{err}(h, D)} > \mu = \frac{\ln |H| + \ln \frac{1}{\delta}}{m} \right\}$$

Proof: Define event

$$E = \left\{ \text{for all } h \in H_\mu : \text{err}(h, S) > 0 \right\}$$

want to show $P(E) \geq 1 - \delta$. (*)

consistency alg. returns \hat{h} with $\text{err}(h, S) = 0$.

if E happens, $\hat{h} \in H_\mu^c \Rightarrow \text{err}(\hat{h}, D) \leq \mu$.

proof of (*):

$$E = \bigcap_{h \in H_\mu} \left\{ \text{err}(h, S) > 0 \right\}$$

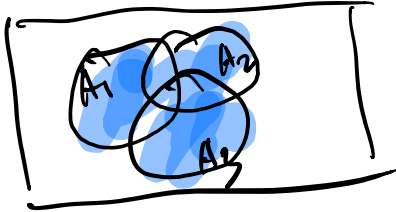
$$\bar{E} = \bigcup_{h \in H_\mu} \left\{ \text{err}(h, S) = 0 \right\}$$

↓
 A_h

Union bound:

for any events A_1, \dots, A_n ,

$$P\left(\bigcup_{i=1}^n A_i\right) \leq \sum_{i=1}^n P(A_i)$$



$$P(\bar{E}) \leq \sum_{h \in H} P(\text{err}(h, S) = 0)$$

given h , $\text{err}(h, D) \geq \mu$.

training example:
 (x_i, y_i)

$$P(\text{err}(h, S) = 0)$$

$$= P(\text{for all } i, h(x_i) = y_i) = \prod_{i=1}^m P(h(x_i) = y_i)$$

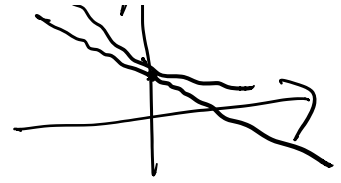
$$= \prod_{i=1}^m (1 - \underbrace{P(h(x_i) \neq y_i)}_{\text{err}(h, D)})$$

$$\leq \prod_{i=1}^m (1 - \mu) = (1 - \mu)^m$$

$$P(\bar{E}) \leq \sum_{h \in H} (1 - \mu)^m$$

$$\frac{1 - x \leq e^{-x}}{N}$$

$$\leq \sum_{h \in \mathcal{H}_\mu} e^{-m \cdot \mu}$$



$$= |\mathcal{H}_\mu| \cdot e^{-m \mu}$$

$$\leq |\mathcal{H}| \cdot e^{-m \mu} \quad \leftarrow \quad \mu = \frac{\ln |\mathcal{H}| + \ln \frac{1}{\delta}}{m}$$

$$= |\mathcal{H}| \cdot e^{-\left(\ln(|\mathcal{H}|/\delta)\right)}$$

$$= \cancel{|\mathcal{H}|} \cdot \frac{\delta}{\cancel{|\mathcal{H}|}} = \delta.$$

$$\Rightarrow P(E) = 1 - P(\bar{E}) \geq 1 - \delta. \quad \cancel{1}$$

The Agnostic PAC learning model

Def: We call an algorithm A agnostic PAC learns hypothesis class \mathcal{H} with sample complexity function $f: (0,1) \times (0,1) \rightarrow \mathbb{N}, \mathbb{R}$ for all distributions \mathcal{D} , $\epsilon > 0$, $\delta > 0$, when A

receives $m \geq f(\epsilon, \delta)$ i.i.d. training examples from D as input. It outputs a classifier \hat{h} , such that with prob. $\geq 1 - \delta$,

$$\text{err}(\hat{h}, D) \leq \min_{h \in \mathcal{H}} \text{err}(h, D) + \epsilon.$$

(A.10) ~~A PAC learns \mathcal{H}~~
~~↕ ↗~~
A agnostic PAC learns \mathcal{H} ?

Algorithm for agnostic PAC learning

$$|\mathcal{H}| < \infty.$$

Empirical risk minimization:

Return \hat{h} such that it has the smallest training error among \mathcal{H} :

$$\hat{h} = \underset{h \in \mathcal{H}}{\text{argmin}} \text{err}(h, S) \approx \text{err}(h, D)$$

Concentration of measure:

Given i.i.d. random variables X_1, \dots, X_n .

empirical mean concentrates around $\mathbb{E}[X]$

$$\frac{1}{n} \sum_{i=1}^n X_i$$

with large probability. $(-f(n))$

Given classifier $h \in \mathcal{H}$,

$$\underline{\text{err}(h, S)} \approx \text{err}(h, D)$$

$$Z_i = \mathbb{I}(h(x_i) \neq y_i) \sim \text{Bernoulli}(\text{err}(h, D))$$

$$\text{err}(h, S) = \frac{1}{m} \sum_{i=1}^m Z_i$$

Hoefding

Hoefding's Inequality:

Suppose that Z_1, \dots, Z_n ^{i.i.d.} such that for each i ,

$$\underline{Z_i \in [a, b]}, \quad \underline{\bar{Z}} = \frac{1}{n} \sum_{i=1}^n Z_i, \quad \mu = \mathbb{E}[Z_i].$$

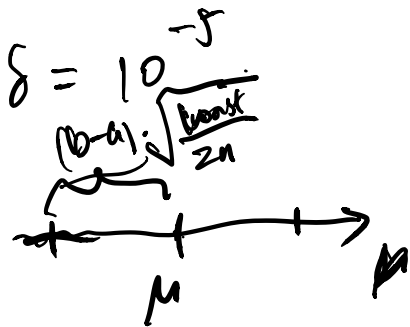
Then, for all $\epsilon > 0$,

$$P(|\bar{Z} - \mu| > \epsilon) \leq 2e^{-\frac{2n\epsilon^2}{(b-a)^2}} \quad (1)$$

~~Wasserstein~~ →

$\mu \in [a, b]$
Equivalently, $\forall \delta > 0$,

$$P(|\bar{Z} - \mu| > (b-a) \cdot \sqrt{\frac{\ln \frac{2}{\delta}}{2n}}) \leq \delta$$



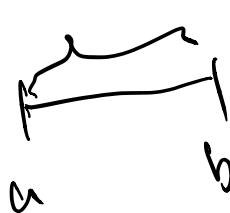
Comparison with Chebyshev's inequality:

$$\bar{Z}, \quad E[\bar{Z}] = \mu$$

$$E\left[\frac{1}{n}(X_1 + \dots + X_n)\right]$$

$$\text{Var}(\bar{Z}) = \frac{1}{n} \cdot \text{Var}(Z_1)$$

$$\leq \frac{1}{n} \cdot (b-a)^2$$



$$\mathbb{E} \left[\underline{\underline{Z_1 - \mathbb{E} Z_1}} \right]^2 \leq (b-a)^2$$

Chebyshev:

$$P(|\bar{Z} - \mu| \geq \epsilon) \leq \frac{\text{Var}(\bar{Z})}{\epsilon^2}$$

$$\leq \frac{(b-a)^2}{n\epsilon^2} = \delta.$$

Equivalently, $\forall \delta > 0,$

$$P(|\bar{Z} - \mu| > (b-a) \cdot \sqrt{\frac{1}{\delta}}) \leq \delta.$$

$\delta = 10^{-5}$. dependence on δ is much worse than Hoeffding.

Proof of Hoeffding's Ineq:

Method of moment generating functions

$$X \quad \phi_X(\lambda) = \mathbb{E}[e^{\lambda X}]$$

Def: random variable X is said to be σ^2 -subgaussian,

if $\forall \lambda \in \mathbb{R}$, $\mu = \mathbb{E}[X]$

$$\mathbb{E}[e^{\lambda(X-\mu)}] \leq e^{\frac{\sigma^2 \lambda^2}{2}}$$

$$(\text{If } X \sim N(\mu, \sigma^2), \mathbb{E}[e^{\lambda(X-\mu)}] = e^{\frac{\sigma^2 \lambda^2}{2}})$$

Next time: $[a, b]$

- bounded r.v's are subgaussian

- ^{independent} sum of subgaussian are subgaussian

- subgaussian \Rightarrow light probability tail