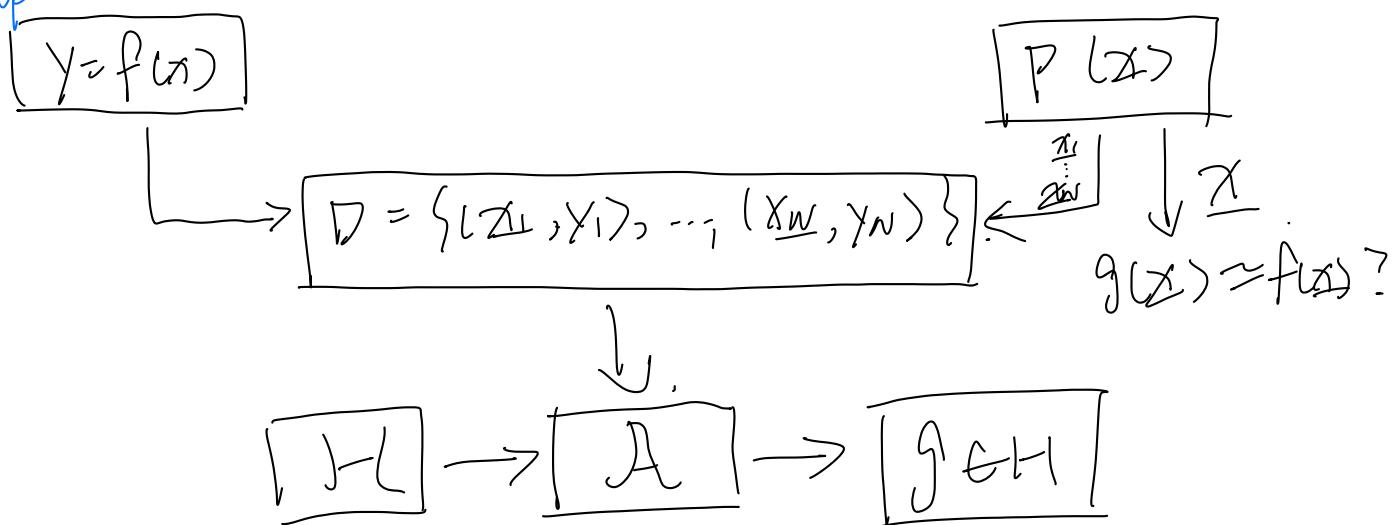


Lec 11 PAC Learning

① Recap:



For any given hypothesis $g \in \mathcal{H}$.

We define $e_n = \mathbb{E}[y_n, g(x_n)]$
 $= \mathbb{E}[\{y_n \neq g(x_n)\}]$

$$E_{in}(g) = \frac{1}{n} \left(\sum_{n=1}^N \mathbb{E}[\{y_n \neq g(x_n)\}] \right) - R.V.$$

Law of a seq of Bernoulli R.V.s

$$E_{out}(g) = \mathbb{E}[\{y \neq g(x)\}]$$

Expected val of Bernoulli R.V.s.

OR expectation of the sequence.

Some distribution for in-sample & out-sample Bernoulli R.V.s.

$$E_{in}(g) \approx E_{out}(g)$$

Min. in training. Care about in real-world.

Goal: $\Pr\{E_{in}(g) \approx E_{out}(g)\} = 1$ "PAC".

② Sequence of iid Bernoulli RVs.

e.g. Consider an urn containing infinitely many balls, each being either red or green.

Let μ be the fraction of red balls (also $\%$)

$$Z = \begin{cases} 1, & \text{if red ball is drawn} \\ 0, & \text{if green ball is drawn.} \end{cases}$$

$Z \sim \text{Bernoulli}(\mu)$.

$$\Pr\{Z=1\} = E(Z) = \mu.$$

Suppose we want to estimate μ .

Randomly draw N balls from the urn

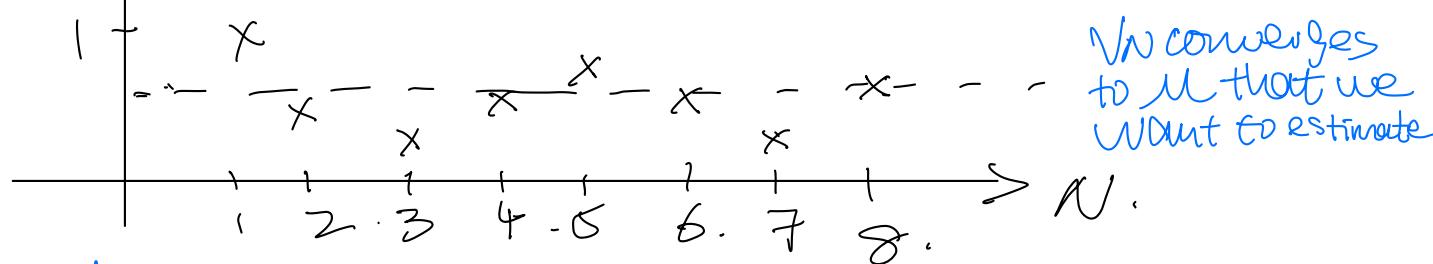
Z_1, Z_2, \dots, Z_N ~~iid~~ ^{"infinitely many balls."} $\sim \text{Bernoulli}(\mu)$.

e.g., Output sequence $(1, 0, 0, 1, 1, 1, 0, 0, 1)$

Define $V_N = \frac{1}{N} \sum_{n=1}^N Z_n$. (proportion of red ball(s)).

IS $V_N \approx \mu$? (same Q as $E_N \approx \text{true?}$)

e.g. V_N



Note.. V_N is a r.v.

$$E[V_N] = E\left[\frac{1}{N} \sum_{n=1}^N Z_n\right] = \mu. \quad \text{"Unbiased estimator!"}$$

$$\text{VAR}[V_N] = \text{VAR}\left[\frac{1}{N} \sum_{n=1}^N Z_n\right]$$

$$\stackrel{\text{indep}}{=} \frac{1}{N^2} \sum_{n=1}^N \text{VAR}(Z_n)$$

$$= \frac{1}{N^2} \cdot N(\mu(1-\mu))$$

$$= \frac{1}{N} \mu(1-\mu)$$

\therefore More balls \Rightarrow variance $\rightarrow 0$.
 Meaning: RV not longer random. \Rightarrow consistent estimator!
 $N \rightarrow \infty \Rightarrow \text{VAR} \rightarrow 0$ (converged).

$\Rightarrow V_N \rightarrow \mu$ as $N \rightarrow \infty$. (in the mean squared sense)

(3) Chebyshev Inequality. (Not enough!)

- It tells how far away RV can be from mean val.

$$\Pr(|V_N - \mu| > \varepsilon) \leq \frac{\text{Var}(V_N)}{\varepsilon^2}$$

"Upper bounded"

In Bernoulli case:

$$\Pr = \frac{\mu(1-\mu)}{N\varepsilon^2}$$

Notice: \Pr goes to 0 as $N \rightarrow \infty$ & $\varepsilon > 0$.

"We can get as close as we want if we have enough samples".

Can do better!

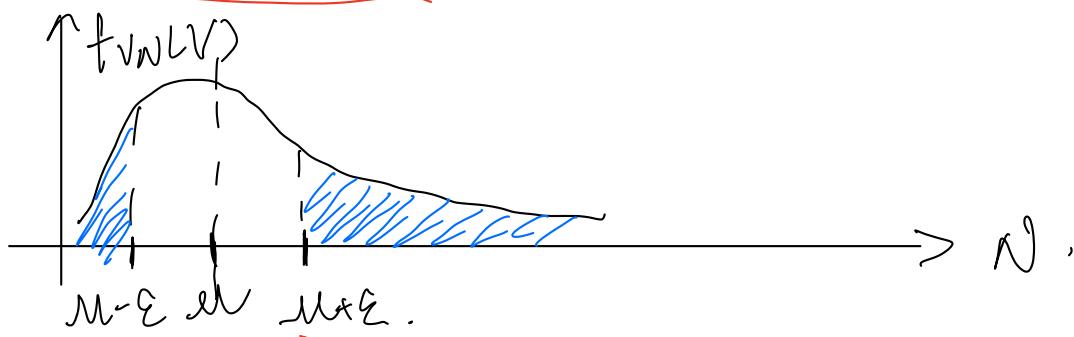
(4) Hoeffding's Inequality: (1963) \Rightarrow stronger bound.

$$\Pr(|V_N - \mu| > \varepsilon) \leq 2e^{-2N\varepsilon^2}$$

"tighter bound"

This is an exponential decrease in N .

Confidence Interval



Confidence Interval (ε -conf interval around mean)

- $\Pr \{ |\bar{X}_n - \mu| > \epsilon \} = \text{area of shaded (tail) region.}$
 $\triangleq S.$
- $1 - S = \text{area inside confidence interval}$
 L "confidence level".

Consider upper bound:

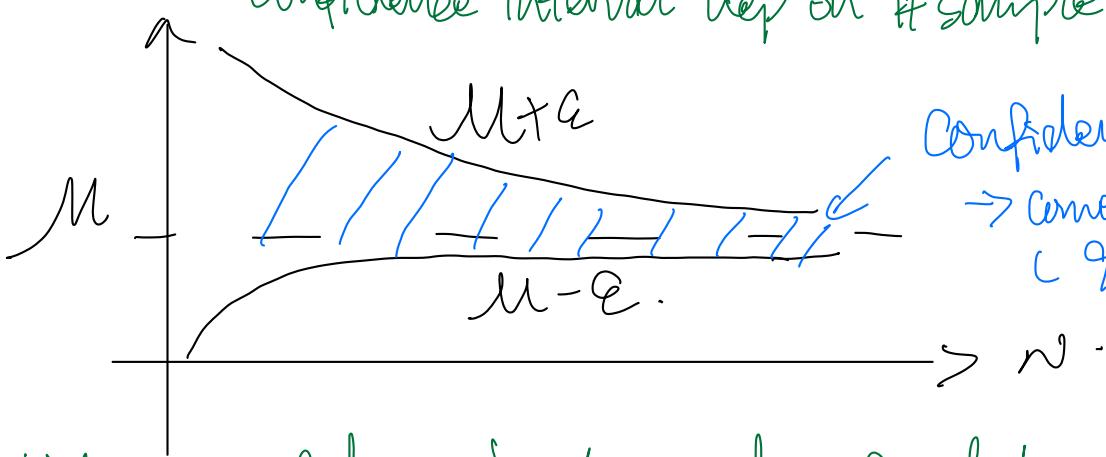
$$S = 2e^{-2N\epsilon^2}$$

$$\Rightarrow \epsilon = \sqrt{\frac{1}{2N} \ln \frac{2}{S}}$$

e.g., $S = 0.01$: what to be 99% sure / confident.

$$\epsilon = \sqrt{\frac{1}{2N} \ln 200} \approx \frac{1.63}{\sqrt{N}}$$

L Confidence interval dep on # sample points.



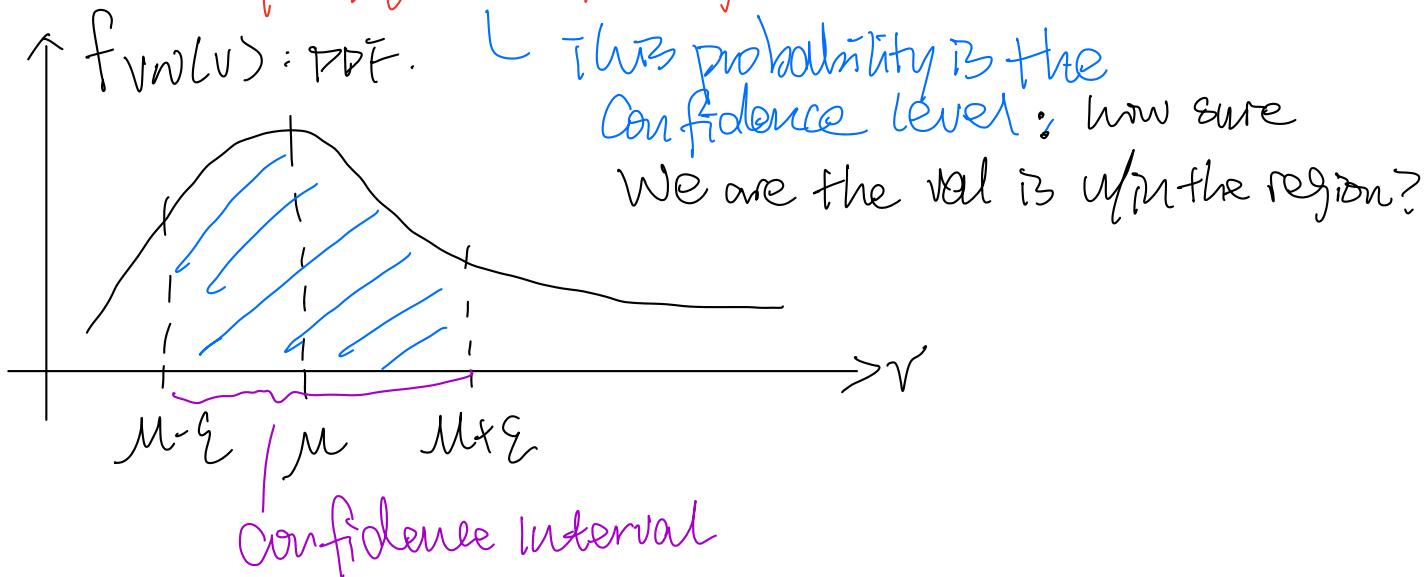
Confidence interval.
 \rightarrow Converges to M as $N \uparrow$.
 (99% Confid.).

Higher confidence level, wider confidence interval.
 \Rightarrow dep on how much error to tolerate.
 ∵ Big data = confident.

lec32 PAC learning & Confidence interval.

① Recap:

For iid Bernoulli RVs z_1, z_2, \dots, z_N w/
 $\Pr\{z_i\} = E[z] = \mu \& V_N = \frac{1}{N} \sum_{n=1}^N z_n$,
 we have: $\Pr\{|V_N - \mu| > \epsilon\} \leq 2e^{-2N\epsilon^2}$



② Back to PAC Learning:

$$E_{in}(h) = \frac{1}{N} \sum_{n=1}^N \mathbb{1}(y_n \neq h(x_n))$$

$$E_{out}(h) = \mathbb{E}[\mathbb{1}(y \neq h(x))]$$

Let $Z = \mathbb{1}(y \neq h(x))$.

$$Z_n = \mathbb{1}(y_n \neq h_n(x_n)), n = 1, 2, \dots, N$$

For any given $h \in \mathcal{H}$ (indep of training dataset D)

z, z_1, z_2, \dots, z_N are iid Bernoulli RVs.

$$E_{in}(h) = V_N$$

$$E_{out}(h) = \mu$$

Hoeffding Inequality:

$$\Rightarrow \Pr\{|E_{in}(h) - E_{out}(h)| > \epsilon\} \leq 2e^{-2N\epsilon^2} \forall \epsilon > 0$$

$$\text{Let } S = 2e^{-2n\alpha^2}$$

Then w/ probability $\geq 1 - \delta$.

$$|\bar{E}_{\text{in}(h)} - \bar{E}_{\text{out}(h)}| \leq \xi = \sqrt{\frac{1}{2n} \log \frac{2}{\delta}}.$$

$$\Rightarrow \bar{E}_{\text{out}(h)} \leq \bar{E}_{\text{in}(h)} + \sqrt{\frac{1}{2n} \log \frac{2}{\delta}}.$$

This creates an upper bound for test error

→ This is not a deterministic bound, but a probabilistic bound.

→ Not getting it for granted.

→ There's always a trade-off.

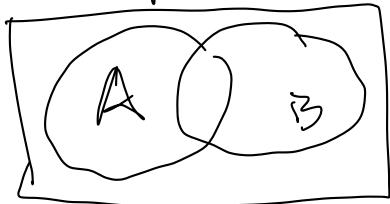
iid Bernoulli RV only if hypothesis h is indep of dataset
if we choose h based on dataset, these properties & straight forward analysis no longer apply!

Trouble! This only works for given h before you observe D (indep of D)

Our final hypothesis g is chosen based on D after observing D .

$\Rightarrow g(x_1), \dots, g(x_N)$ are not indep.

One useful thm: Union bound:



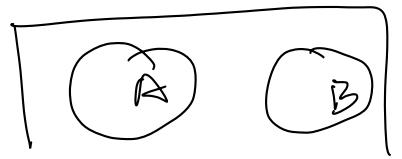
$$\begin{aligned} P\{A \cup B\} &= P\{A\} + P\{B\} - P\{A \cap B\} \\ &\leq P\{A\} + P\{B\}. \end{aligned}$$

- For arbitrary events A_1, A_2, \dots, A_M (M: # hypotheses)

$$P\{A_1 \cup \dots \cup A_M\} \leq \sum_{i=1}^M P\{A_i\}.$$

w/ equality iff A_1, \dots, A_M are mutually exclusive ("disjoint")

e.g.,



Since $\mathcal{G} \cdot \mathcal{H} = \{h_1, h_2, \dots, h_m\}$, the event that:

$$\{ |E_{in}(g) - E_{out}(g)| > \varepsilon \} \leq \bigcup_{i=1}^M \{ |E_{in}(h_i) - E_{out}(h_i)| > \varepsilon \}$$

g is one of the hypothesis
g is subset of all events.

Event: set of outcomes in a probability space.
g is a subset of event.

Probability of subset always \leq probability of set.

\Rightarrow the hypothesis we chose.

$$\Pr \{ |E_{in}(g) - E_{out}(g)| > \varepsilon \} \leq \Pr \{ \bigcup_{i=1}^M \{ |E_{in}(h_i) - E_{out}(h_i)| > \varepsilon \} \}$$

union bound

$$\leq \sum_{i=1}^M \Pr \{ |E_{in}(h_i) - E_{out}(h_i)| > \varepsilon \}.$$

$$\leq \sum_{i=1}^M 2e^{-2n\varepsilon^2} = 2Me^{-2n\varepsilon^2}.$$

h_i is indep of data
prior to training.

To get bound for final hypothesis we choose, we times M to Hoeffding's inequality

Note: If $g \in \mathcal{H}$ is selected after observing D (data silothing), then Hoeffding bound is weakened by a factor of M .

$$② S = 2Me^{-2n\varepsilon^2}$$

$$\varepsilon = \sqrt{\frac{1}{2n} \log \frac{2M}{S}}$$

④ Generalization Error:

Def: difference b/w E_{in} & E_{out} -

$$\Delta(g) = |E_{in}(g) - E_{out}(g)|.$$

w/ probability $1 - S$

$$\Delta(g) \leq \sqrt{\frac{1}{2n} \log \frac{2M}{S}}$$

related to model complexity.

M: #hypothesis,

N: # dataset.

C # samples.

Increase M

$\Delta(g)$

increases

Increase N

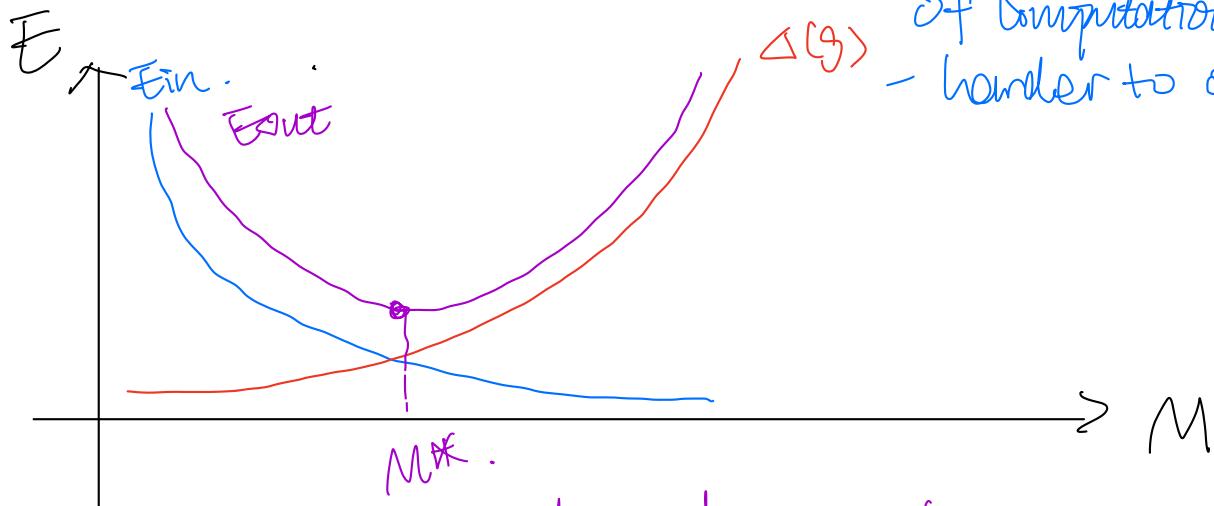
decreases

$E_{in}(g)$

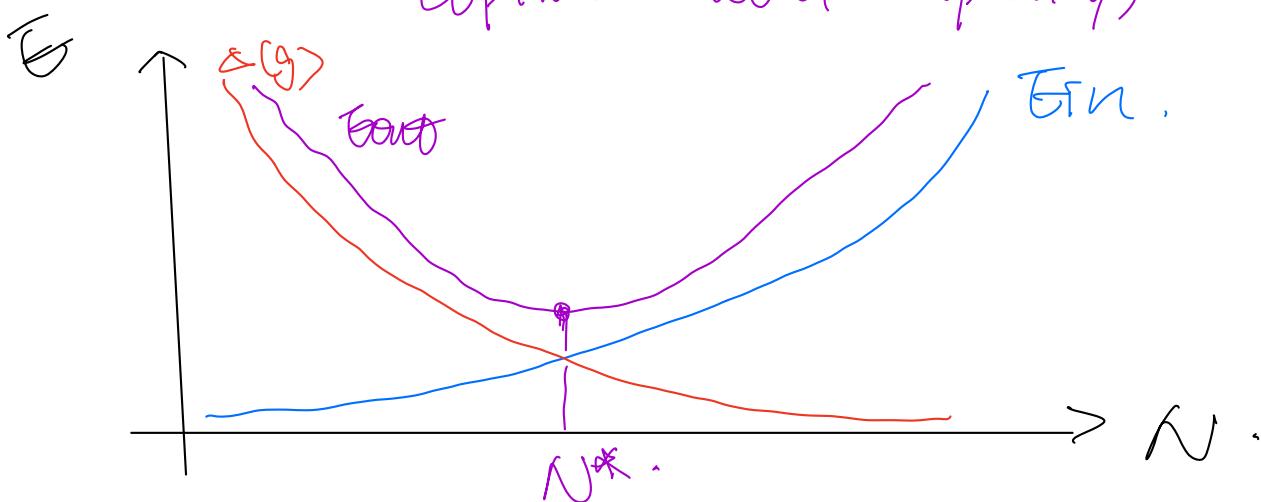
decreases

increases.

- from some amount of computation, E_{in} larger
- harder to optimize.



(Optimal model complexity)



Note: test error (E_{out}) is the only thing matters in the real world.

Note: Still trouble! Most useful/practical ML models in the world have ∞ M.

e.g., $y = \text{sign}(w^T x)$, $x \in \mathbb{R}^{d+1}$.

choices of w : infinite / uncountable

∴ Union bound Not Enough! (Loose bound).

We need effective number of hypotheses

Lec 33 PAC Learning for final hypothesis

① Recap: PAC learning for final hypothesis $g \in \mathcal{H}$
 $\mathcal{H} = \{h_1, h_2, \dots, h_m\}$.

Generalization error:

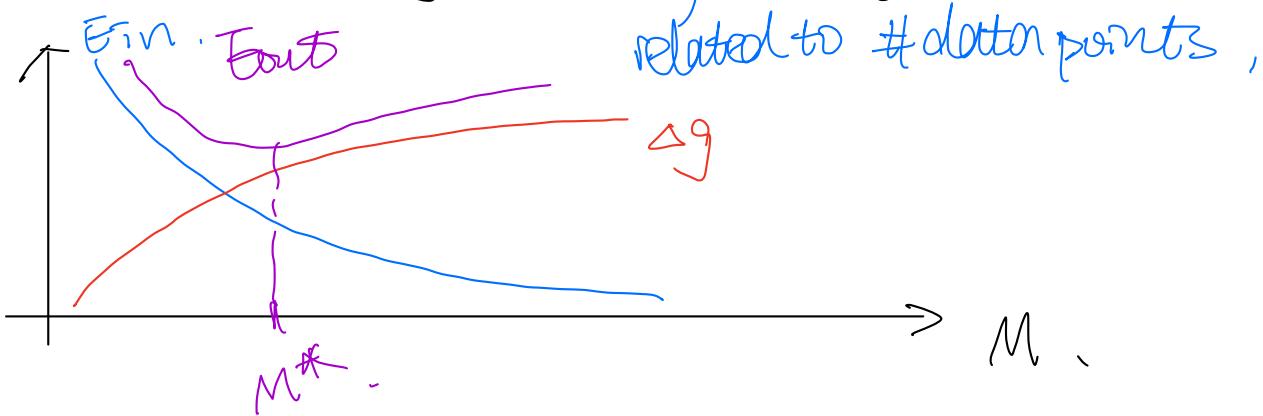
$$\Delta(g) = |E_{in}(g) - E_{out}(g)|$$

$$\Pr\{\Delta(g) \geq \epsilon\} \leq \Pr\left\{\bigcup_{i=1}^m \Delta(h_i) \geq \epsilon\right\}$$

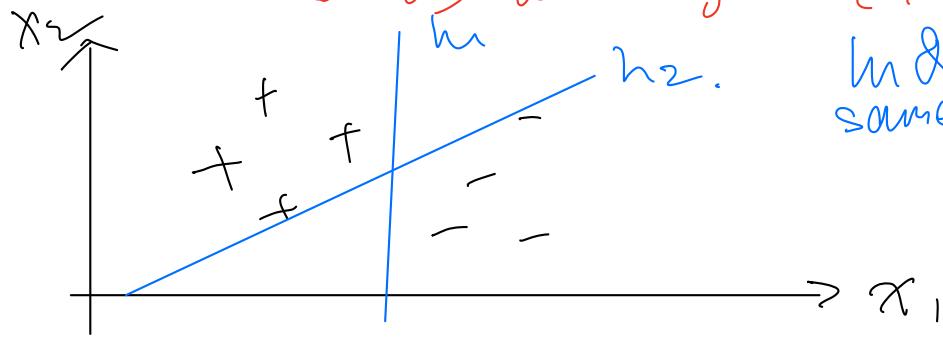
$$\text{Union bound} \leq 2M e^{-2N\epsilon^2}$$

\Rightarrow w/ probability $\geq 1 - \delta$

$$E_{out}(g) \leq E_{in}(g) + \sqrt{\frac{1}{2N} \log \frac{2M}{\delta}} \quad \begin{matrix} \text{related to} \\ \text{model complexity} \end{matrix}$$



Trouble! Union bound is loose in general.



② Idea: (chap 21.2.2)

Consider $E_{\text{in}}(h) = \frac{1}{N} \sum_{n=1}^N \mathbb{1}\{Y_n \neq h(\underline{x}_n)\}$ dep on h only through $\{h(\underline{x}_1), h(\underline{x}_2), \dots, h(\underline{x}_N)\}$.

\Rightarrow Replace M w/ an effective # hypothesis $m_H(N)$.

i if 2 h s are the same \Rightarrow essentially 1 line.

e.g. $h_1 \equiv h_2$.

Dichotomy:

let $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N \in \mathbb{R}^d$ be fixed points.
let $h \in H$ be some hypothesis, $h: \mathbb{R}^d \rightarrow \{+1, -1\}$.

Def: Dichotomy Vector

$(h(\underline{x}_1), h(\underline{x}_2), \dots, h(\underline{x}_N)) \in \{+1, -1\}^N$.

Def: Dichotomy set

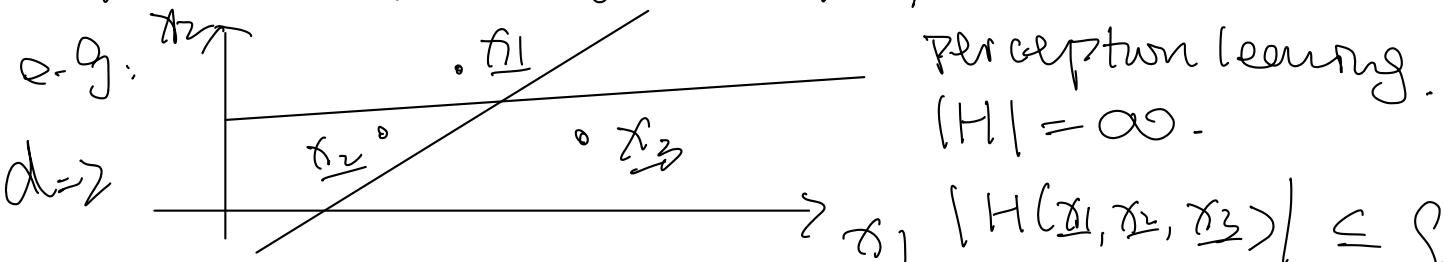
$H(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N) = \{(h(\underline{x}_1), h(\underline{x}_2), \dots, h(\underline{x}_N)) : h \in H\}$.

i.e., the collection of all binary vectors generated by H on $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N$.

Note: There are no repeated elements in a set.

$$|H(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)| \leq 2^N.$$

(No matter how large $M = |H|$) is !).



$$|H(\underline{x}_1, \underline{x}_2, \underline{x}_3)| \leq 2^3$$

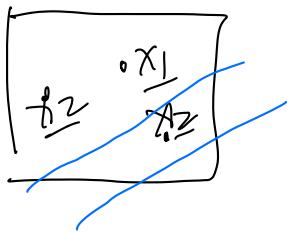
Def: the hypothesis set H shatters

$(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)$ if $|H(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)| = 2^N$.

e.g. linear classification in $d=2$.

$H = \{w \in \mathbb{R}^3 : w_0 + w_1 x_1 + w_2 x_2 = 0\}$.

$$Y = \text{sign}(w_0 + w_1 x_1 + w_2 x_2)$$

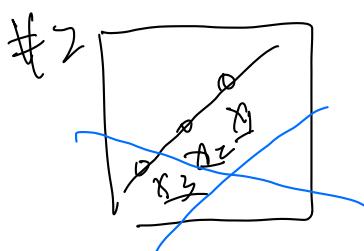


✓/8

| | $h(x_1)$ | $h(x_2)$ | $h(x_3)$ | |
|---|----------|----------|----------|---|
| + | + | + | + | ✓ |
| + | + | - | - | ✓ |
| + | - | - | + | ✓ |
| + | - | - | - | ✓ |
| - | - | - | - | ✓ |
| - | - | - | - | ✓ |
| - | - | - | - | ✓ |

by symmetry

Shattering



6/8

$$|H(\underline{x}_1, \underline{x}_2, \underline{x}_3)| = 6 \text{ (2^3)}$$

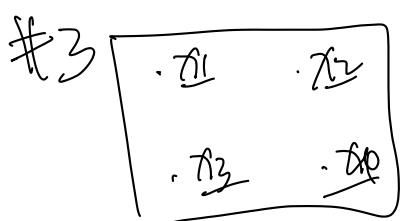
Not shattering

| | $h(x_1)$ | $h(x_2)$ | $h(x_3)$ | |
|---|----------|----------|----------|---|
| + | + | + | + | ✓ |
| + | + | - | - | ✓ |
| + | - | - | + | ✗ |
| + | - | - | - | ✓ |

not linearly separable

\Rightarrow has to go to quadratic classifier.

∴ Linear classifier can shatter points not in one linear line.



-- A linear classifier cannot shatter 4 points either.

$$|H(\underline{x}_1, \underline{x}_2, \underline{x}_3, \underline{x}_4)| = 14 \text{ (2^4 - 2)}$$

even though they are not co-linear

Def Growth Function

$$m_H(N) = \max_{\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N \in \mathbb{R}^d} |H(\underline{x}_1, \underline{x}_2, \dots, \underline{x}_N)|.$$

↳ [all possible location of points]

≡ Effective # of hypotheses

e.g., for linear classification in $d=2$

$$m_H(3) = 8 \quad m_H(4) = 14$$

Worst case (at most)
have 1 hypothesis
small complexity
to ∞

(Recall: Hypothesis represents model complexity wrt dataset).

Note: 1) $m_H(N) \leq 2^N$

2) A tighter bound can be found, using the VC dimension.

Def: Let k be an integer s.t. $m_H(k) < 2^k$, then k is a breakpoint of H .

e.g., for linear classifier, the ^(first) break point is 4.

Def Vapnik-Chervonenkis (VC) Dimension

Let N be an integer s.t. $m_H(N) = 2^N$ & $m_H(N+1) < 2^{N+1}$

i.e., $N+1$ is the first (smallest) breakpoint.

then VC dimension of H :

$$d_{VC}(H) = N.$$

Next week, VC dimension to bound growth func.