

Theory of Generalization

FACT

- with distribution in training data & testing data
- \Rightarrow low training error
- \nparallel
- low testing error

Def

Δ Training error:

$$E_{tr}(h) = \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n))$$

where x_1, \dots, x_N sampled from D

- h is determined by x_1, \dots, x_N

Δ Testing error:

$$E_{te}(h) = \frac{1}{N} \sum_{n=1}^N e(h(x_n), f(x_n))$$

where x_1, \dots, x_N sampled from D

- h is independent from x_1, \dots, x_N

Δ Generalization error

- G. error = Test error (on D)
- $E(h) = E_{x \sim D} [e(h(x), f(x))] = E_{te}(h)$

Δ Summary

if $E(h) = 0$

then $E(h) \approx E_{tr}(h) \rightarrow$ How?

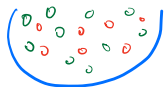
or

$E_{tr}(h) \approx 0 \rightarrow$ Training

Q: How do we make sure

$$E(h) \approx E_{tr}(h)$$

FACT Hoeffding's inequality



- $\Delta P[\text{pick red ball}] = \mu$
- $P[\text{pick green ball}] = 1 - \mu$
- \rightarrow we DO NOT know μ

- Δ by pick ball's independently we get fraction of V

- $\Delta V \rightarrow \mu$?
- perhaps

Δ Hoeffding's inequality

$$P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

anote: V & μ 的差距, 比 ϵ 大吗? P, 很小
多少? \rightarrow 比 $2e^{-2\epsilon^2 N}$ 还小

- Δ statement $\mu = V$ is probably approximately correct
- (PAC!)

FACT

$$\Delta P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

- valid for N
- $\epsilon > 0$
- independent from μ (real probability)



Δ in learning:

- given a function h
- we randomly draw x_1, \dots, x_N independent
- generalization error

$$E(h) = E_{x \sim D} [h(x) \neq f(x)] \Leftrightarrow \mu \quad \text{unknown}$$

sample data error

$$E_{tr}(h) = \frac{1}{N} \sum_{n=1}^N [h(x_n) \neq y_n] \Leftrightarrow V \quad \text{known}$$

$$P[|V - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

FACT

Δ for each h , h is a hypothesis

$$P[|E_{tr}(h) - E(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

Δ for all h , H is a hypothesis set

$$P[|E_{tr}(h_1) - E(h_1)| > \epsilon],$$

$$P[|E_{tr}(h_2) - E(h_2)| > \epsilon],$$

\vdots

$$P[|E_{tr}(h_{|H|}) - E(h_{|H|})| > \epsilon]$$

$$\leq P[\sup_{h \in H} |E_{tr}(h) - E(h)| > \epsilon]$$

$$\leq \sum_{m=1}^{|H|} P[|E_{tr}(h_m) - E(h_m)| > \epsilon] \leq 2|H|e^{-2\epsilon^2 N}$$

$$\text{from } P(\bigcup_{i=1}^{\infty} A_i) \leq \sum_{i=1}^{\infty} P(A_i)$$

Δ summary

$$P[|E_{tr}(h) - E(h)| > \epsilon] \leq P[\sup_{h \in H} |E_{tr}(h) - E(h)| > \epsilon] \leq 2|H|e^{-2\epsilon^2 N}$$

NOTE More on Hoeffding's inequality

D_1 Bad D_2 Bad \dots D_N Bad
 h_1 h_2 h_M Bad

$P_D[\text{BAD } D \text{ for } h_1] \leq \dots$
 $P_D[\text{BAD } D \text{ for } h_2] \leq \dots$
 $P_D[\text{BAD } D \text{ for } h_M] \leq \dots$

informed hypothesis (假设现在我的话的模型 i.e. 不是举出本的)
 对岸到到手上
 的首科 D_1, D_2, \dots, D_N
 informed 的 h on D
 可能導致 "Bad"
 存即 $E_{in}(h) \neq E_{out}(h)$

之所以要
 Hoeffding's ing.
 就是為量
 $P[\text{BAD}]$ 机率
 多高?

答覆是低值:
 bounded by $2e^{-2\epsilon^2 N}$

不過剛則是 1 個 h 啊...
 我 algo 都是在 1 個 h 啊...
 upper bound 是啥?

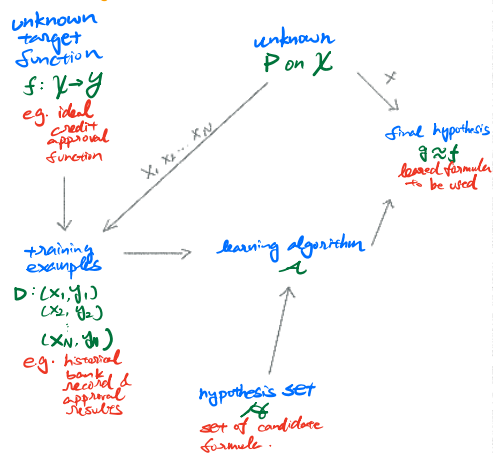
$\therefore P_D[\text{BAD } D]$

$= P_D[\text{BAD } D \text{ for } h_1 \text{ or BAD } D \text{ for } h_2 \dots \text{ or BAD } D \text{ for } h_M]$
 $\leq P_D[\text{BAD } D \text{ for } h_1] + P_D[\text{BAD } D \text{ for } h_2] + \dots + P_D[\text{BAD } D \text{ for } h_M]$
 (union bound)
 $\leq 2M e^{-2\epsilon^2 N} = 2|K| e^{-2\epsilon^2 N}$

- finite-bin version of Hoeffding
- & hope... $E_{in}(g) = E_{out}(g)$ is PAC.

$\rightarrow A$ will pick h_m w/ min. $E_{in}(h_m)$ as g

Summary: Statistical learning flow



& hope $E_{out}(g) \approx E_{in}(g) \approx 0$

- for batch & supervised learning, $g \approx f \Leftrightarrow E_{out}(g) \approx 0$ achieved through $E_{out}(g) \approx E_{in}(g)$ & $E_{in}(g) \approx 0$
- ① can we make sure $E_{out}(g) \approx E_{in}(g)$
- ② can we make $E_{in}(g)$ small enough

FACT $|K| = \infty$

Question: How do we deal with it?

- small $|K|$
 $P[\text{BAD}] \leq 2|K| e^{-2\epsilon^2 N}$
 small! great!
 but $|K|$ too little
 $E_{in}(g) \uparrow$
- large $|K|$
 $E_{in}(g) \rightarrow 0$
 small error! great!
 but $|K|$ too large
 $P[\text{BAD}] \uparrow$

FACT establish a finite quantity replace $|K|$

let $|K|$ replaced by m_K

s.t.

$$P[|E_{in}(g) - E_{out}(g)| > \epsilon] \leq 2m_K e^{-2\epsilon^2 N}$$

FACT $|K|$ is over-estimated for BAD events

- BAD events $B_m: |E_{in}(h_m) - E_{out}(h_m)| > \epsilon$
- over-lapping for similar hypothesis $h_1 \approx h_2$
- as ① $E_{out}(h_1) \approx E_{out}(h_2)$

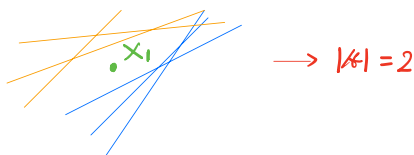
② for most D $E_{in}(h_1) \approx E_{out}(h_2)$

- should be instead of

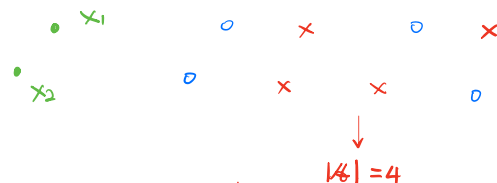


- so: can we group similar kinds?

eg. K in \mathbb{R}^2 $|K| \rightarrow \infty$



$\rightarrow |K| = 2$



$|K| = 4$

$N=3$ $|K|=8$ but if on same line, different
 $N=4$ $|K|=14$ but if on same line, different

FACT observation: effective $|K| \leq 2^N$
 perhaps can replace $|K|$ by effective $|K|$?
 need more rigorous proof

FACT Dichotomies: mini-hypotheses

- △ limited hypothesis: $K(x_1, x_2, \dots, x_N)$
- △ $|K(x_1, x_2, \dots, x_N)|$: depend on inputs (x_1, x_2, \dots, x_N)
- △ growth function:
 remove dependence by taking max of all possible (x_1, x_2, \dots, x_N)

$$m_K(N) = \max_{x_1, x_2, \dots, x_N \in X} |H(x_1, x_2, \dots, x_N)|$$

△ finite, upper-bounded by 2^N

Q: How to calculate growth function

FACT shattered

Δ if $m_H(N) = 2^N \Leftrightarrow$ exists N inputs that can be shattered

Δ eg. convex set

FACT summary of 4 growth function

- positive rays $N+1$
- positive intervals $C \binom{N+1}{2} + 1$
- convex sets 2^N
- 2D perceptrons $< 2^N$

polynomial good!

exponential bad!

FACT Break point \leadsto k 阈值, 无法被 shattered

Δ if no k inputs can be shattered by H
call k a break point for H

$\Delta m_H(k) < 2^k$

$\Delta k+1, k+2, k+3 \dots$ are all break points

Δ study minimum break point

eg. linear case break point $k=4$
note: 4 个 无法被 shattered

FACT conjecture:

Δ no break point: $m_H(N) = 2^N$

Δ break point k : $m_H(N) = O(N^{k-1})$

proof?

FACT $m_H(N) \leq$ maximum possible $m_H(N)$ given k
 $\leq \text{poly}(N)$

FACT Bounding function