

Generalization theory

①

① Big picture

② True error vs training error

H = set of classifiers, e.g. {linear separators in \mathbb{R}^d }

Training set of n points: $(x_1, y_1), \dots, (x_n, y_n)$

Training error: $\hat{\text{err}}(h) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}(h(x_i) \neq y_i)$

How does this relate to true error, $\text{err}(h)$?

Very roughly:

$$\text{err}(h) \leq \hat{\text{err}}(h) + \frac{c(H)}{n} + \sqrt{\hat{\text{err}}(h) \frac{c(H)}{n}}$$

where $c(H)$ is a measure of the size, or complexity, of H .

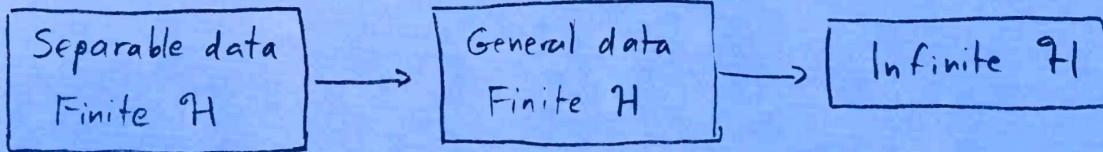
③ Occam's razor: simpler hypotheses are preferable.



ML version of this: "smaller" hypothesis classes generalize better.

We'll look at some sophisticated notions of size. [A far cry from parameter-counting.]

④ Cases of interest



② Separable data, finite \mathcal{H}

- (a) X : input space
 γ : label space
 \mathcal{H} : hypothesis class

Separable: $y = h^*(x)$ for some $h^* \in \mathcal{H}$

- (b) Distribution P over $X \times \gamma$.

True error of $h \in \mathcal{H}$:

$$\text{err}(h) = \Pr_{(x,y) \sim P} (h(x) \neq y)$$

- (c) Thm Pick any $0 < \delta < 1$. Suppose $(x_1, y_1), \dots, (x_n, y_n) \sim P$ i.i.d., and we pick $\hat{h} \in \mathcal{H}$ consistent with them. With probability at least $1 - \delta$ over the choice of training set, for any h consistent with the data,

$$\text{err}(h) \leq \frac{1}{n} \ln \frac{|\mathcal{H}|}{\delta}.$$

Pf. Define $\varepsilon = \frac{1}{n} \ln \frac{|\mathcal{H}|}{\delta}$.

- (i) Fix any $h \in \mathcal{H}$ with $\text{err}(h) > \varepsilon$.

Pick $(x_1, y_1), \dots, (x_n, y_n) \sim P$.

$$\begin{aligned} \Pr(h \text{ consistent with } (x_1, y_1), \dots, (x_n, y_n)) &< (1 - \varepsilon)^n \\ &< e^{-\varepsilon n} = \frac{\delta}{|\mathcal{H}|}. \end{aligned}$$

- (ii) $\therefore \Pr(\exists h \in \mathcal{H} : \text{err}(h) > \varepsilon \text{ and } h \text{ consistent with training set})$

$$\leq \sum_{\substack{h \in \mathcal{H} \\ \text{err}(h) > \varepsilon}} \Pr(h \text{ consistent with training set})$$

$$\leq |\mathcal{H}| \cdot \frac{\delta}{|\mathcal{H}|} = \delta. \quad \square$$

- (d) Example: \mathcal{H} = k-sparse disjunctions over d variables.

3 Non-separable case, finite H

a) Concentration of empirical averages

Fix any hypothesis $h \in H$.

How close is $\underbrace{(\text{training error on } n \text{ pts})}_{\text{empirical average}}$ to $\text{err}(h)$?

b) Central limit theorem

Informal version Suppose X_1, \dots, X_n are i.i.d. with mean μ and variance $\sigma^2 < \infty$.

Then the distribution of $\frac{X_1 + \dots + X_n - n\mu}{\sigma\sqrt{n}}$ approaches $N(0, 1)$

as $n \rightarrow \infty$.

Example A coin of bias p is flipped n times. Let \hat{p} be the fraction of heads observed. What is the distribution of \hat{p} for large n ?

. Define $X_i = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ coin heads} \\ 0 & \text{otherwise.} \end{cases}$

. Then $\mu = E X_i = p$ and $\sigma^2 = \text{var}(X_i) = p(1-p)$.

. By CLT:

$$\frac{n\hat{p} - np}{\sqrt{np(1-p)}} \sim N(0, 1)$$

$$\Rightarrow \hat{p} - p \sim \sqrt{\frac{p(1-p)}{n}} \cdot N(0, 1)$$

$$\Rightarrow \hat{p} \sim N(p, \frac{p(1-p)}{n})$$

$\therefore \hat{p}$ has roughly a Gaussian distribution with mean p and standard deviation $\sqrt{\frac{p(1-p)}{n}} \leq \frac{1}{2\sqrt{n}}$.

W.p. $\geq 95\%$, $\hat{p} = p \pm \frac{1}{\sqrt{n}}$. (for n large enough)

Hoeffding's inequality

Thm Let Z_1, \dots, Z_n be i.i.d. random variables with $a \leq Z_i \leq b$ and $\mu = \mathbb{E} Z_i$.

Then: $\Pr \left(\left| \frac{Z_1 + \dots + Z_n}{n} - \mu \right| \geq \varepsilon \right) \leq 2e^{-2\varepsilon^2 n / (b-a)^2}$

Another way to write it: With probability $\geq 1-\delta$,

$$\left| \frac{Z_1 + \dots + Z_n}{n} - \mu \right| \leq (b-a) \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}$$

(d) Application: relate true error to training error for finite \mathcal{H} .

Thm Suppose a training set $\{(x_i, y_i)\}$ is drawn i.i.d. from P . Define:

• $\hat{\text{err}}(h)$ = training error of h

$$= \frac{1}{n} \sum_{i=1}^n \mathbb{1}(h(x_i) \neq y_i)$$

• $\text{err}(h)$ = true error of h = $\Pr_{(x,y) \sim P}(h(x) \neq y)$

Then with prob $\geq 1-\delta$,

$$\max_h |\hat{\text{err}}(h) - \text{err}(h)| \leq \sqrt{\frac{1}{2n} \ln \frac{2|\mathcal{H}|}{\delta}}.$$

Pf. Fix any $h \in \mathcal{H}$. Apply Hoeffding's inequality with

$$Z_i = \mathbb{1}(h(x_i) \neq y_i).$$

Then take a union bound over \mathcal{H} . \otimes

In particular, let h^* be the best classifier:

$$h^* \text{ minimizes } \text{err}(h)$$

and let \hat{h} be the learned classifier:

$$\hat{h} \text{ minimizes } \hat{\text{err}}(h).$$

Then:

$$\text{err}(\hat{h}) \leq \text{err}(h^*) + \sqrt{\frac{2}{n} \ln \frac{2|\mathcal{H}|}{\delta}}.$$

(spell out the argument)

③ We have seen generalization bounds of the form

$$\max_h |\hat{\text{err}}(h) - \text{err}(h)| \leq \sqrt{\frac{c(H)}{n}}$$

where $c(H) = \log |H|$.

What if H is infinite?

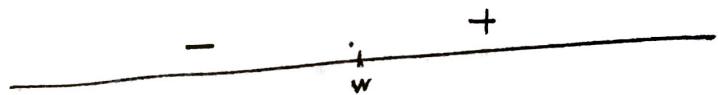
4] VC dimension

a) Example

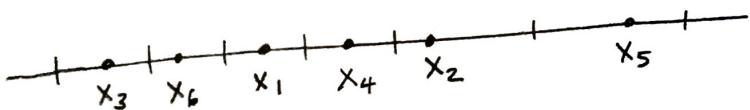
(i) $X = \mathbb{R}$

$H = \{\text{threshold classifiers } h_w : w \in \mathbb{R}\}$

$$h_w(x) = \begin{cases} + & \text{if } x \geq w \\ - & \text{if } x < w \end{cases}$$



(ii) Say we have n pts x_1, \dots, x_n .



Effectively there are just $n+1$ classifiers to consider, even though H is infinite.

\therefore Think of $|H|$ as $n+1$.

b) Effective hypothesis class

Pretend you get your training data in two phases:

(i) First, x_1, x_2, \dots, x_n

(ii) Then, the labels are revealed

After seeing (i), effective hypothesis class is

we have seen generalization bounds of the form

$$\max_n |\hat{\text{err}}(h) - \text{err}(h)| \leq \sqrt{\frac{c(H)}{n}}$$

where $c(H) = \log |H|$.

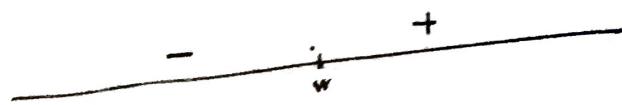
What if H is infinite?

VC dimension

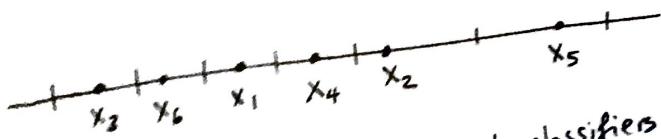
Example

$X = \mathbb{R}$
 $H = \{\text{threshold classifiers } h_w : w \in \mathbb{R}\}$

$$h_w(x) = \begin{cases} + & \text{if } x \geq w \\ - & \text{if } x < w \end{cases}$$



(ii) Say we have n pts x_1, \dots, x_n .



Effectively there are just $n+1$ classifiers to consider, even though H is infinite.
 \therefore Think of $|H|$ as $n+1$.

(b) Effective hypothesis class

Pretend you get your training data in two phases:

(i) First, x_1, x_2, \dots, x_n

(ii) Then, the labels are revealed

After seeing (i), effective hypothesis class is

$$\tilde{\mathcal{H}} = \{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\}$$

e.g. If \mathcal{H} = thresholds, $|\tilde{\mathcal{H}}| = n+1$.

We can use $c(\mathcal{H}) = \log |\tilde{\mathcal{H}}|$ in our generalization bounds. $\left[\sqrt{\frac{c(\mathcal{H})}{n}} \right]$

(c) How big is $\tilde{\mathcal{H}}$, in general?

• Clearly, $|\tilde{\mathcal{H}}| \leq 2^n$, but this gives a vacuous bound

• Any hypothesis class \mathcal{H} has an associated parameter called the Vapnik-Chervonenkis (VC) dimension, $VC(\mathcal{H})$.

$$\text{And: } |\tilde{\mathcal{H}}| \leq O(n^{VC(\mathcal{H})})$$

• Can use $c(\mathcal{H}) = VC(\mathcal{H})$, as a measure of the complexity of \mathcal{H} .

(d) VC dimension

(i) Def \mathcal{H} shatters a set of points $x_1, \dots, x_m \in X$ if it can realize all 2^m possible labelings of these points.

$$\text{Ex. } \mathcal{H} = \{\text{linear separators in } \mathbb{R}^2\}$$

\mathcal{H} shatters any 3 (non-collinear) points in \mathbb{R}^2 .

• •

•

(ii) Def The VC dimension of \mathcal{H} is the size of the largest set shattered by \mathcal{H} .

$$\text{Ex. } \mathcal{H} = \{\text{linear separators in } \mathbb{R}^2\}$$

\mathcal{H} cannot shatter any 4 pts

• •

$$\therefore VC(\mathcal{H}) = 3$$

• •

(7) Ex. $\mathcal{H} = \{\text{thresholds in } \mathbb{R}\}$. (iv) Can show:
 $\text{VC}(\text{linear separators in } \mathbb{R}^d) = d+1$.

$$\text{VC}(\mathcal{H}) = 1.$$

e) Sauer's lemma If $\text{VC}(\mathcal{H}) = V$ then for any $x_1, \dots, x_n \in X$,
 $\tilde{\mathcal{H}} = \{(h(x_1), \dots, h(x_n)) : h \in \mathcal{H}\}$
 has size at most $O(n^V)$.

f) Thm Suppose \mathcal{H} has VC dimension $V < \infty$. Suppose a training set of n pts is drawn i.i.d. from the underlying distribution p on $X \times Y$ and that $\hat{\text{err}}$ denotes training error and err true error. Then, with probability $\geq 1 - \delta$,

$$\text{For all } h \in \mathcal{H}: |\hat{\text{err}}(h) - \text{err}(h)| \leq C_0 \sqrt{\frac{V + \log \frac{1}{\delta}}{n}}$$

where C_0 is some constant.

5 More refined notions of hypothesis class size

a) $\mathcal{H} = \{\text{linear separators in } \mathbb{R}^d\} = \{h_w : w \in \mathbb{R}^{d+1}\} \quad h_w(x) = \text{sign}(w \cdot x)$

i) We have seen $\text{VC}(\mathcal{H}) = d+1$, so we get generalization bounds of the form $\sqrt{d/n}$. But what if we look at separators with large margin?

ii) Fix any $x_1, \dots, x_n \in \mathbb{R}^d$ and define
 $\tilde{\mathcal{H}}_\gamma = \{(h_w(x_1), \dots, h_w(x_n)) : w \in \mathbb{R}^{d+1}, \|w\| = 1, |w \cdot x_i| \geq \gamma \text{ for all } i\}$

What is $|\tilde{\mathcal{H}}_\gamma|$?

For any labeling in $\tilde{\mathcal{H}}_\gamma$, the Perceptron alg will converge within $1/\gamma^2$ updates. \therefore There are $\leq n^{1/\gamma^2}$ choices for the Perceptron's final classifier $\Rightarrow |\tilde{\mathcal{H}}_\gamma| \leq n^{1/\gamma^2}$