Discussion 5:     (9/19/2018)

1- iid vs id
2- VC dimension
3- handle missing

_____

— iid   ( independent & identical distribution )

$X_1, X_2 \cdots, X_n$  they are generated from an iid
source, the it should satisfy the following

$$P(X_1, X_2, \ldots, X_n) = \prod_{i=1}^{n} P(X_i)$$

$P(X_i)$   is   the  same   $\forall \, i \in \{1, \ldots, n\}$

— id ( independent distribution )

$Y_i = X_i + \varepsilon$

$\quad \longrightarrow \varepsilon \perp X_i \, , \, \varepsilon \sim N(0,1)$

if $X_i$'s are iid, what is the distribution
of $Y_i$'s given $X_i$'s are observed.

$$P(Y_1, \ldots, Y_n \mid X_1 \cdots X_n) = \prod_{i=1}^{n} P(Y_i \mid X_i)$$

$P(Y_i \mid X_i) \sim N(X_i, 1) \rightsquigarrow = \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-1}{2} (y_i - x_i)^2\right)$

(1)

$P(Y_i|X_i)$ is a function $x_i \Rightarrow$ it is Not identical for all $i$, but $P(Y_i|X_i)$ are indep. from each other.

## VC dimension:

we have $\overset{seen}{\checkmark}$ how complexity a model can affect the performance of the model
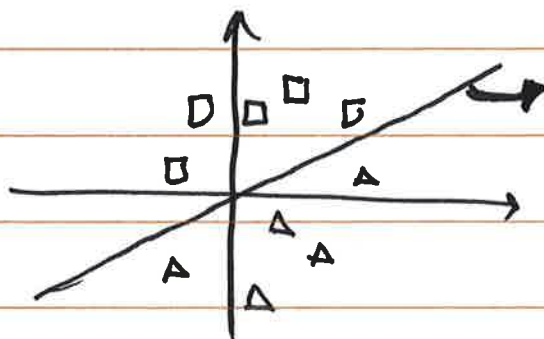
need to define complexity of model

VC dimension can be used to define complexity.

representation power is the ability of model to learn a wide variety of input-output relationship / or ability to memorize — to fit data

Example: perceptron for 2D

$$\hat{c}(\underline{x}) = Sgn\left( \theta_1 x_1 + \theta_2 x_2 + \theta_0 \right)$$

choose a different form of classifier

$\Rightarrow$ results in different model & performance
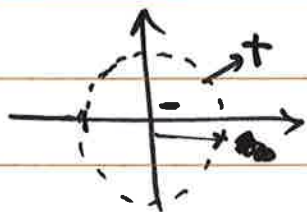
$$\hat{c}(\underline{x}) = Sgn(\theta_1 x_1 + \theta_2 x_2)$$

$\downarrow$

classifier (learner) should pass through

origin while classify data points.

Example: $\hat{c}(\underline{x}) = Sgn\left(\underbrace{\dfrac{\underline{x}^T \underline{x}}{\|x\|^2}}_{} - \theta_0\right)$

if value of $\|x\|$ is large $\longrightarrow$ it $\checkmark$ classified as $+1$

otherwise $\qquad\qquad\qquad \longrightarrow$ " " " $-1$



$\|x\|^2 > \theta_0 \longrightarrow +1$

$\|x\|^2 < \theta_0 \longrightarrow -1$

Tradeoff:

the more representation power we can have for classifiers, the more accurate model on training data we have, but this is at the cost of overfitting

the less rep. power $\downarrow$ $\Rightarrow$ poor performance on training

How can we quantify representation power?

(VC dim.)

Let's have some preliminary on VC:

— let's assume our training data are iid from dist. $p(x)$

— let's define "risk" & "emperical risk"
These are just "long term" test & "observed" training errors

$$R(\theta) = \text{Test Error} = E\left( \delta \left( c \neq c(x; \theta) \right) \right)$$

$$R^{emp}(\theta) = \text{Train Error} = \frac{1}{m} \sum_i \delta \left( c^{(i)} \neq c(x^{(i)}; \theta) \right)$$

(41

Q: How are these errors related?

- underfitting domain:
  Test error is similar to train error
- overfitting domain:
  Test error is be way worse!

VC dim. tells us about risk:

- Given some classifier, let $H$ be its VC dim.
- with "high prob" $1-y$, it can be shown that

$$\text{Test error} \leq \text{Train Error} + \sqrt{\frac{H \log\left(\frac{2m}{H}\right) + H - \log\frac{y}{4}}{m}}$$

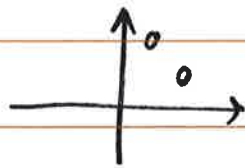($m$ is # of training data)

$H \uparrow \Rightarrow$ overfitting

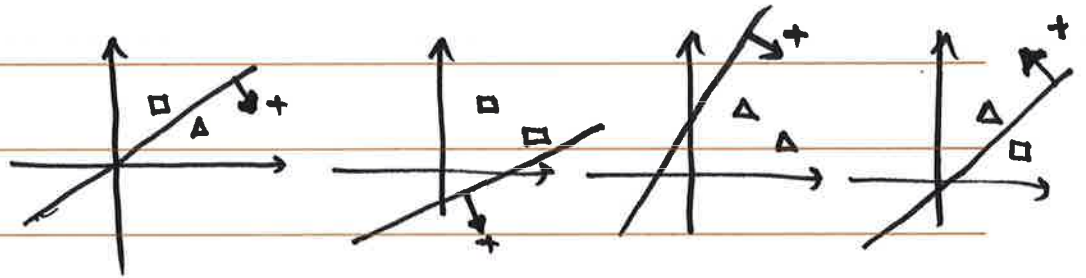$H \downarrow \Rightarrow$ Test & Trains are similar

Shattering:

We say a classifier $f(x)$ can shatter points $x^{(1)}...x^{(h)}$ iff for all $y^{(1)}...y^{(h)}$ $f(x)$ can achieve zero error on training data.
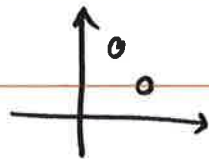
Example: (2D)   Can   $f(x;\theta) = Sgn(\theta_1 x_1 + \theta_2 x_2 + \theta_3)$
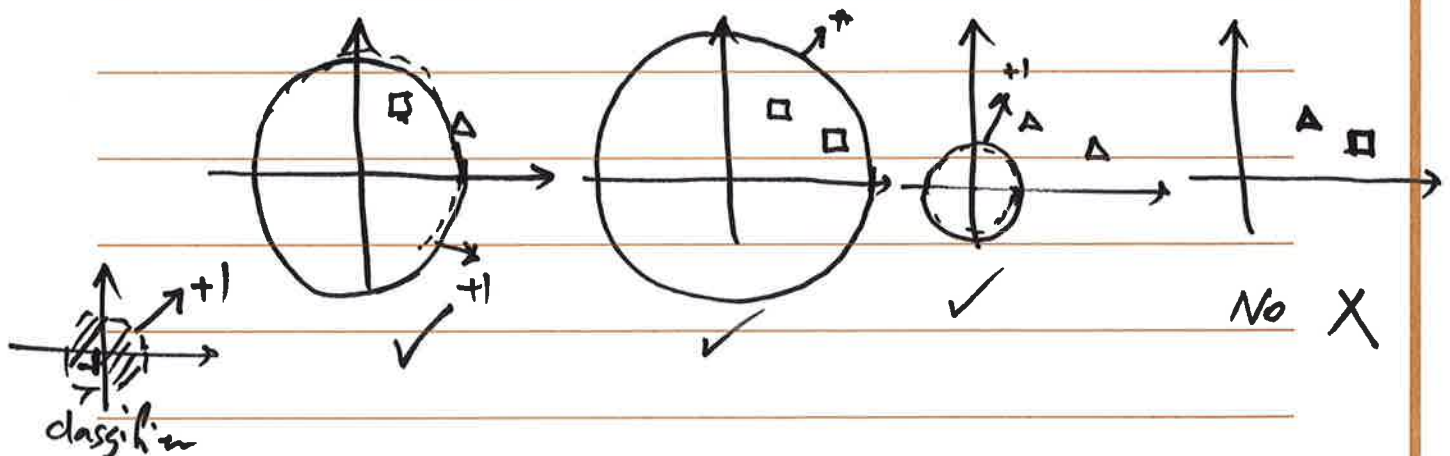shatter the following points?   Yes



$$\square \rightarrow -1$$
$$\triangle \rightarrow +1$$



Example: (2D)   Can   $f(x;\theta) = Sgn(\bullet\, x^T x - \theta)$
shatter the following points?   NO



$$\square \rightarrow -1$$
$$\triangle \rightarrow +1$$



✓           ✓           ✓           No  X

classifier

(6)

VC dim :

The VC dim H is defined as
the max. # of points h that
can be arranged so that $f(x)$ can
shatter all of them.

A game :
- fix the def $f(x; \theta)$
- Player 1 : choose location of $x^{(1)}...x^{(h)}$
- " 2 : " target labels of $x^{(1)}...x^{(h)}$
  $(y^{(1)}...y^{(h)})$

- Player 1 : choose value $\theta$

if $f(x; \theta)$ can reproduce the target labels,
Player 1 wins

$$\exists \{x^{(i)}\}_{i=1}^{h} \text{ s.t. } \forall \{y^{(i)}\}_{i=1}^{h} \quad \exists \theta \text{ s.t.}$$
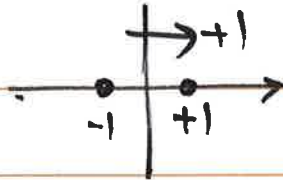
$$f(x^{(i)}; \theta) = y^{(i)} \quad \forall i \in \{1,...,h\}$$

what is the
Example: $\lor$ VC dimension of ~~1D~~ 1D

linear perception

$$\hat{c}(x) = sgn(\theta_1 x + \theta_0)$$

3 points

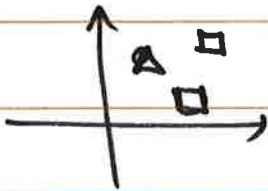$$\Rightarrow \quad \text{VC of 1D linear perceptron is 2.}$$

Example:    VC dim. of 2D
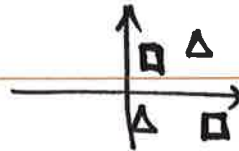
$$\hat{c}(x) = \text{Sgn}(\theta_1 x_1 + \theta_2 x_2 + \theta_0)$$

2 points ✓

3 points ✓

(8 possible labling)

4 points

VC dimension of this problem is 3.

\* VC dimension of $d$ dimensional linear perceptron class is equal to $d+1$.

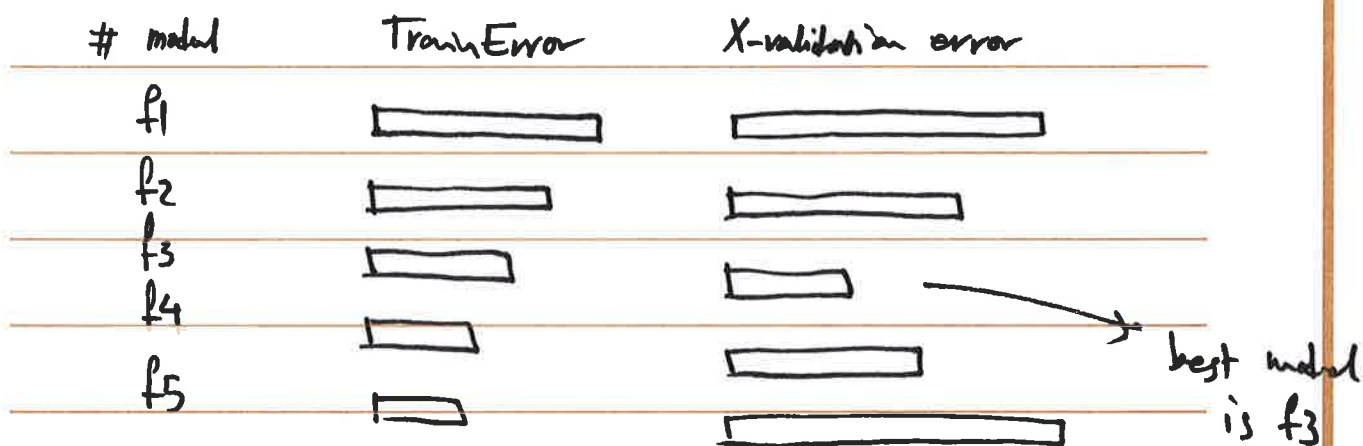VC dimension measures the power of the learner (NOT necessarily equal to # of parameters)

$\downarrow$

So # of parameters does NOT necessarily equals to the complexity of the model

$$\left( VC \left( 1NN \right) = \infty \right)$$

~~Trade off~~ Analysis with VC :

Cross-Validation

| # model | Train Error | X-validation error |
|---------|-------------|--------------------|
| $f_1$ | | |
| $f_2$ | | |
| $f_3$ | | |
| $f_4$ | | |
| $f_5$ | | |

best model is $f_3$

VC dimension

| # model | Train Error | VC Term | Sum (Train Error & VC term) |
|---------|-------------|---------|------------------------------|
| $f_1$ | | | |
| $f_2$ | | | |
| $f_3$ | | | |
| $f_4$ | | | |
| $f_5$ | | | |