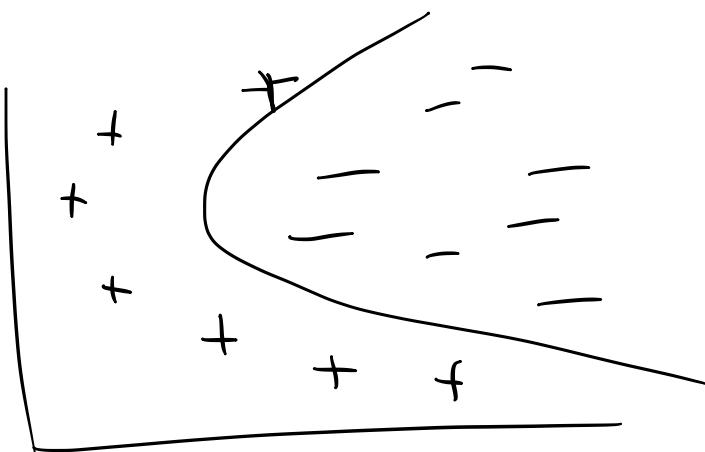


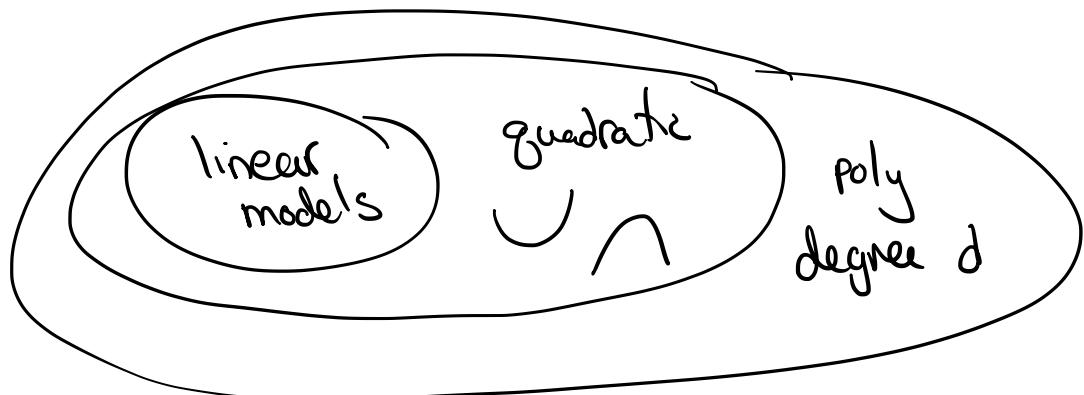
Statistical Learning Theory

Generalization = Data + Knowledge

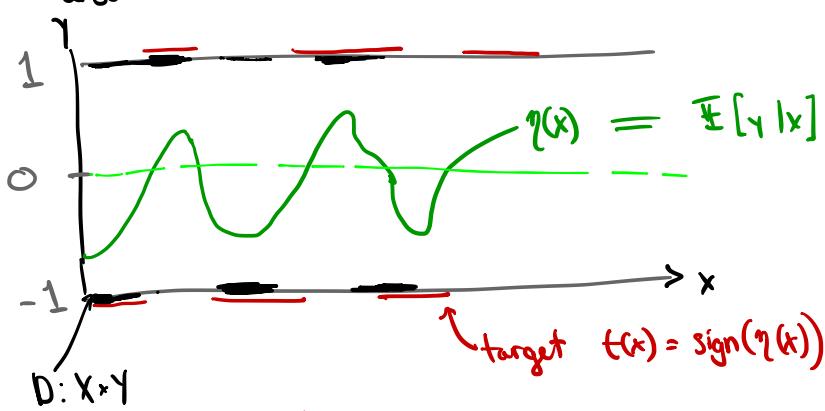
- cannot generalize from data alone. Need wise choice for function class.



- need "simplest" function class that describes the data well.
Occhan's Razor - choose the simplest possible explanation to explain the phenomenon



- can you measure the "complexity" of a function class by the number of parameters?
answer: sometimes yes, sometimes no!
- how do we measure the "complexity" of a function class?
answer: VC dim, Rademacher complexity, covering numbers



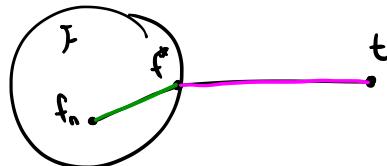
$$t \in \underset{\substack{f \\ \text{all functions}}}{\operatorname{argmin}} R^{\text{true}}(f)$$

$$\mathbb{E}_{(x,y) \sim D} \mathbb{1}_{[f(x) \neq y]}$$

$$R^* = R^{\text{true}}(t)$$

$$f_n \in \underset{f \in \mathcal{F}}{\operatorname{argmin}} R^{\text{emp}}(f) + C \|f\|^2$$

best-in-class is f^*
 $R^{\text{true}}(f^*) := \inf_{f \in \mathcal{F}} R^{\text{true}}(f)$



How bad is f_n ?

$$R^{\text{true}}(f_n) - R^* = \underbrace{[R^{\text{true}}(f^*) - R^*]}_{\text{approximation error}} + \underbrace{[R^{\text{true}}(f_n) - R^{\text{true}}(f^*)]}_{\text{estimation error}}$$

need to know something about D

want this to be as small as possible

Can't measure $R^{\text{true}}(f_n)$ what to do?

$$R^{\text{true}}(f_n) = \underbrace{R^{\text{emp}}(f_n)}_{\text{ok}} + \underbrace{[R^{\text{true}}(f_n) - R^{\text{emp}}(f_n)]}_{\text{stuff } (n, F)}$$

Hoeffding's Inequality (one-sided)

Z_1, \dots, Z_n iid rv

$h(z) \in [a, b]$ bounded function

$\forall \varepsilon > 0$,

$$P_{Z \sim D^n} \left[E_{Z \sim D^n} [h(Z)] - \frac{1}{n} \sum_{i=1}^n h(Z_i) \geq \varepsilon \right] \leq \exp \left(-\frac{2n\varepsilon^2}{(b-a)^2} \right)$$

If I set $h(z) = \mathbb{1}_{\{f(x) \neq y\}}$ then I get

$$P_{Z \sim D^n} \left[R^{\text{true}}(f) - R^{\text{emp}}(f) \geq \varepsilon \right] \leq \exp \left(-\frac{2n\varepsilon^2}{(b-a)^2} \right)$$

b-a is 1-0

$$P_{Z \sim D^n} \left[R^{\text{true}}(f) \geq R^{\text{emp}}(f) + \varepsilon \right] \leq \exp \left(-\frac{2n\varepsilon^2}{(b-a)^2} \right)$$

δ

reparameterize

$$\delta = \exp \left(-\frac{2n\varepsilon^2}{(b-a)^2} \right)$$

$$\ln \delta = -\frac{2n\varepsilon^2}{(b-a)^2}$$

$$\varepsilon^2 = -\frac{\ln \delta}{2n} (b-a)^2 = \frac{\ln(1/\delta)}{2n} (b-a)^2 \Rightarrow \varepsilon = \sqrt{\frac{\ln(1/\delta)}{2n}} (b-a)$$

$$P_{Z \sim D^n} \left[R^{\text{true}}(f) > R^{\text{emp}}(f) + \sqrt{\frac{\ln(1/\delta)}{2n}} (b-a) \right] \leq \delta$$

"inversion"

with prob $\geq 1 - \delta$,

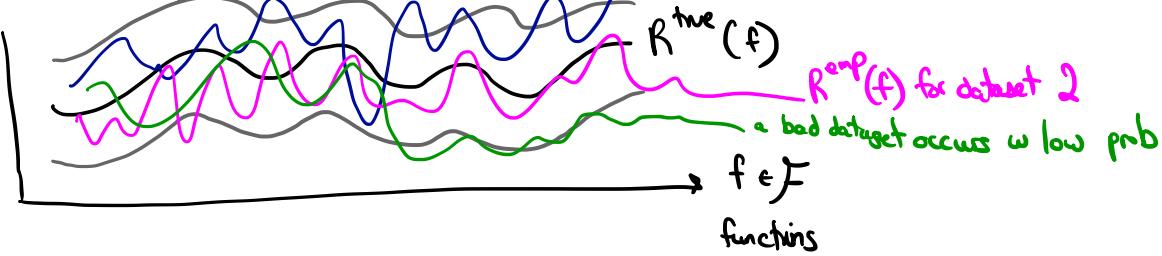
$$R^{\text{true}}(f) \leq R^{\text{emp}}(f) + \sqrt{\frac{\ln(1/\delta)}{2n}} (b-a)$$

← Bound for a single $f \in \mathcal{F}$

This bound says that for each f , there is a set of "good" datasets (S) where $R^{\text{true}}(f)$ is not too large and

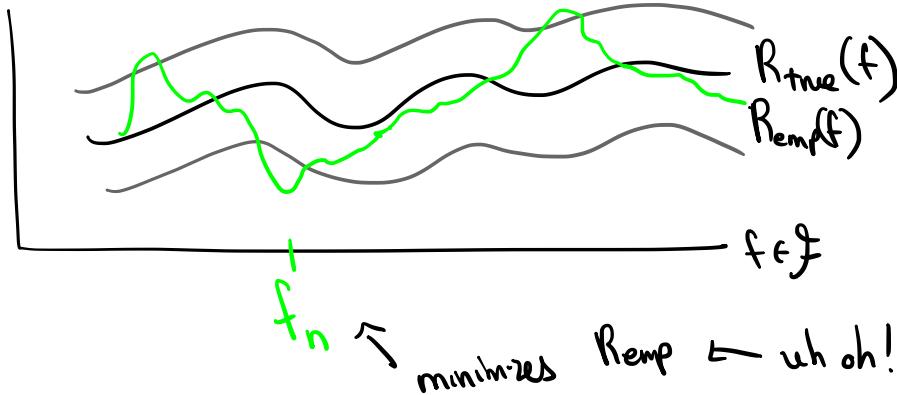
$$P_{Z \sim D^n} [Z \in S] \geq 1 - \delta$$

But f might depend on S !



Hoeffding's says that $R^{\text{true}}(f)$ & $R^{\text{emp}}(f)$ are close w.h.p. for any fixed f .

However, f is not fixed. The algorithm chooses f knowing the data. It tries to minimize R^{emp} .



- If \mathcal{F} is really large, this problem gets worse. Oy!

- Want a bound that holds uniformly over all $f \in \mathcal{F}$

$$R^{\text{true}}(f_n) - R^{\text{emp}}(f_n) \leq \sup_{f \in \mathcal{F}} (R^{\text{true}}(f) - R^{\text{emp}}(f))$$

Occhan's Razor Bound (Hoeffding's + Union Bound)

$$P[C_1^{\text{bad}} \cup C_2^{\text{bad}}] \leq P(C_1^{\text{bad}}) + P(C_2^{\text{bad}})$$

↑ ↑
prob we hit a bad dataset for either f_1 or f_2

From Hoeffding's inequality $P(C_1^{\text{bad}}) \leq \delta$, $P(C_2^{\text{bad}}) \leq \delta$ so $P[C_1^{\text{bad}} \cup C_2^{\text{bad}}] \leq 2\delta$

$$P[C_1^{\text{bad}} \cup \dots \cup C_M^{\text{bad}}] \leq M\delta$$

say \mathcal{F} is finite f_1, \dots, f_m .

$$P_{Z \sim D^n} [\exists f \in \mathcal{F} : R^{\text{true}}(f) - R^{\text{emp}}(f) \geq \varepsilon] \leq$$

$$\sum_{j=1}^m P_{Z \sim D^n} [R^{\text{true}}(f_j) - R^{\text{emp}}(f_j) \geq \varepsilon] \leq \sum_{j=1}^m e^{-2n\varepsilon^2} = M e^{-2n\varepsilon^2}$$

reparameterize $\delta = M e^{-2n\varepsilon^2} \Rightarrow \varepsilon^2 = \frac{-1}{2n} \log \frac{\delta}{M}$

$$\varepsilon = \sqrt{\frac{\log(M) + \log(1/\delta)}{2n}}$$

Plug: $P_{Z \sim D^n} [\exists f \in \mathcal{F} : R^{\text{true}}(f) - R^{\text{emp}}(f) \geq \sqrt{\frac{\log(M) + \log(1/\delta)}{2n}}] \leq \delta$

inverting: For all $\delta > 0$ w.p. at least $1-\delta$,

$$\forall f \in \mathcal{F} = \{f_1, f_2, \dots, f_m\}$$

$$R^{\text{true}}(f) \leq R^{\text{emp}}(f) + \sqrt{\frac{\log(M) + \log(1/\delta)}{2n}}$$

Occhan's
Raz or
Bound

Summary: Our bound can be applied to f_n even though it depends on the data.

Extra $\log M$ term \rightarrow says we want M bounds to hold simultaneously.

What if \mathcal{F} is infinite? Bound is vacuous.



Recap:

- Generalization = data + knowledge
restricting f to lie in restricted \mathcal{F}

- for a fixed f , for most datasets

$$R^{\text{the}}(f) - R^{\text{emp}}(f) \approx \frac{1}{\sqrt{n}}$$

- for most datasets if \mathcal{F} is finite, $|\mathcal{F}| = M$

$$\sup_{f \in \mathcal{F}} [R^{\text{the}}(f) - R^{\text{emp}}(f)] \approx \sqrt{\frac{\log M}{n}}$$

- Occam's Razor