

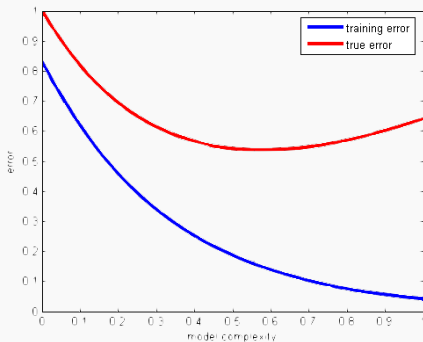
CS-GY 6923: Lecture 7

Taste of Learning Theory, PAC learning

NYU Tandon School of Engineering, Prof. Christopher Musco

THE FUNDAMENTAL CURVE OF ML

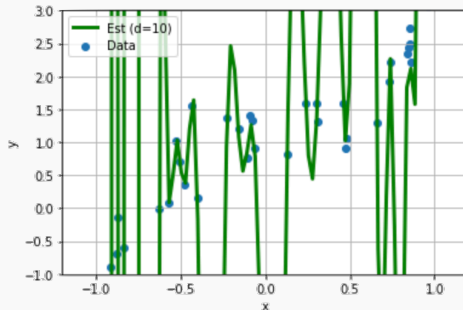
Key Observation: Due to overfitting, more complex models do not always lead to lower test error.



The more complex a model is, the more training data we need to ensure that we do not overfit.

EXAMPLE: POLYNOMIAL REGRESSION

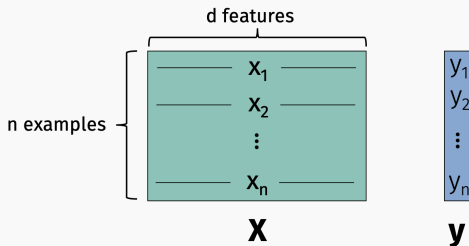
If we want to learn a degree q polynomial model, we will perfectly fit our training data if we have $n \leq q$ examples.



Need $n > q$ samples to ensure good generalization. How much more?

EXAMPLE: LINEAR REGRESSION

If we want to fit a multivariate linear model with d features, we will perfectly fit our training data if we have $n \leq d$ examples.



Need $> d$ samples to ensure good generalization.

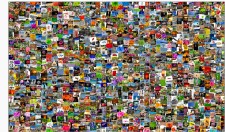
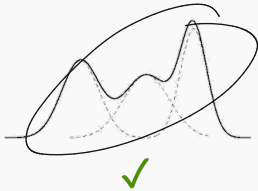
How much more?

Major goal in learning theory:

Formally characterize how much training data is required to ensure good generalization (i.e., good test set performance) when fitting models of varying complexity.

Statistical Learning Model:

- Assume each data example is randomly drawn from some distribution $(\mathbf{x}, y) \sim \mathcal{D}$



For today: We will only consider (classification problems) so assume that $y \in \{0, 1\}$.

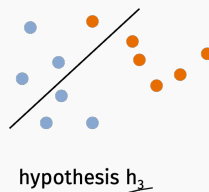
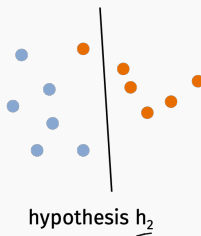
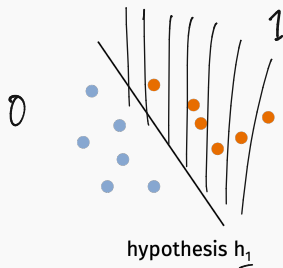
Statistical Learning Model:

- Assume each data example is randomly drawn from some distribution $(\mathbf{x}, y) \sim \mathcal{D}$.
- Assume we want to fit our data with a function h (a “hypothesis”) in some hypothesis class \mathcal{H} . For input \mathbf{x} , $\underline{h(\mathbf{x})} \rightarrow \underline{\{0, 1\}}$. $h \in \mathcal{H}$

You can think of h as a model, instantiated with a specific set of parameters. I.e. h is the same as $\underline{f_\theta}$.

EXAMPLE HYPOTHESIS CLASS

Linear threshold functions:

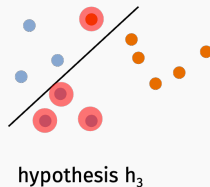
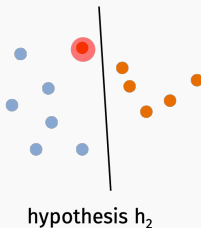
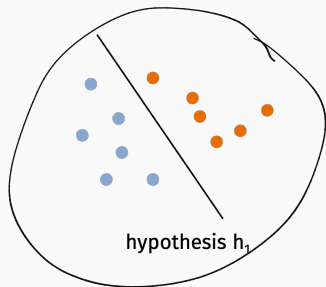


\mathcal{H} contains all functions of the form:

$$\underline{h}(\underline{x}) = \mathbb{1}[\underline{x}^T \underline{\beta} \geq \lambda]$$

EXAMPLE HYPOTHESIS CLASS

Linear threshold functions:



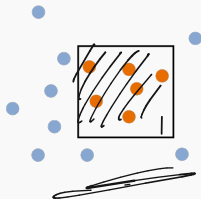
\mathcal{H} contains all functions of the form:

$$h(\mathbf{x}) = \mathbb{1}[\mathbf{x}^T \boldsymbol{\beta} \geq \lambda]$$

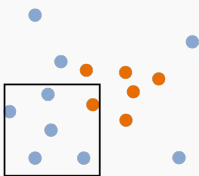
EXAMPLE HYPOTHESIS CLASS

Axis aligned rectangles:

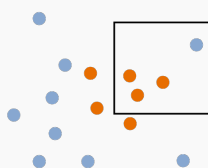
\emptyset



hypothesis h_1



hypothesis h_2



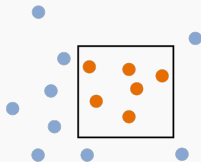
hypothesis h_3

\mathcal{H} contains all functions of the form:

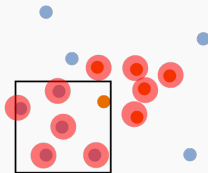
$$h(\mathbf{x}) = \mathbb{1}[l_1 \leq x_1 \leq u_1 \text{ and } l_2 \leq x_2 \leq u_2]$$

EXAMPLE HYPOTHESIS CLASS

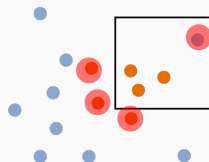
Axis aligned rectangles:



hypothesis h_1



hypothesis h_2



hypothesis h_3

\mathcal{H} contains all functions of the form:

$$h(\mathbf{x}) = \mathbb{1}[l_1 \leq x_1 \leq u_1 \text{ and } l_2 \leq x_2 \leq u_2]$$

EXAMPLE HYPOTHESIS CLASS

$$\mathbf{x} = (x_1, \dots, x_d)$$

Disjunctive Normal Form (DNF) formulas:

Assume $\mathbf{x} \in \{0,1\}^d$ is binary.

$$\bar{x}_5$$

$$x_5 = 1$$

$$\bar{x}_5 = 0$$

$$x_5 = 0$$

$$\bar{x}_5 = 1$$

\mathcal{H} contains functions of the form:

$$h(\mathbf{x}) = (\underbrace{x_1 \wedge \bar{x}_5 \wedge x_{10}}_{\text{false}}) \vee (\underbrace{\bar{x}_3 \wedge x_2}_{\text{true}}) \vee \dots \vee (\underbrace{\bar{x}_1 \wedge x_2 \wedge x_{10}}_{\text{false}})$$

\wedge = "and", \vee = "or"

$$d=10$$

k -DNF: Each conjunction has at most k variables.

POPULATION AND EMPIRICAL ERROR

Same as “population risk” for the zero one loss:

- Population (“True”) Error:

$$\mathbb{E} [R_{\text{emp}}(h)] = R_{\text{pop}}(h)$$

$$\underline{R_{\text{pop}}}(h) = \Pr_{(\underline{x}, \underline{y}) \sim \mathcal{D}} [\underline{h}(\underline{x}) \neq \underline{y}]$$

- Empirical Error: Given a set of samples

$$(\underline{x}_1, \underline{y}_1), \dots, (\underline{x}_m, \underline{y}_m) \sim \mathcal{D},$$

$$\underline{R_{\text{emp}}}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}[h(\underline{x}_i) \neq \underline{y}_i]$$

↗ # of mistakes
h makes
on
(x₁, y₁) ... (x_m, y_m)

Goal is to find h ∈ H that minimizes population error.

GENERALIZATION

Let $(\underline{x_1, y_1}), \dots, (\underline{x_n, y_n}) \sim \mathcal{D}$ be our training set and let h_{train} be the empirical error minimizer:

$$\textcircled{h_{train}} = \arg \min_{h \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \mathbb{1}[h(x_i) \neq y_i]$$

Let h^* be the population error minimizer:

$$\text{---} h^* = \arg \min_{h \in \mathcal{H}} \text{---} R_{pop}(h) = \arg \min_h \Pr_{(x,y) \sim \mathcal{D}} [h(x) \neq y]$$

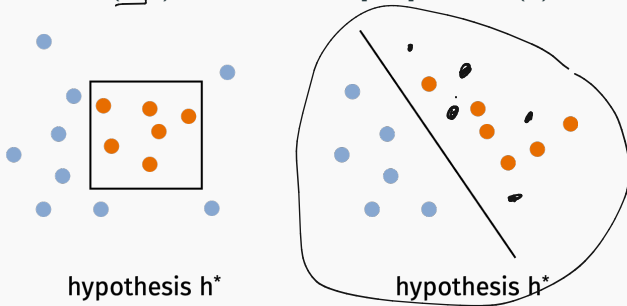
Goal: Ideally, for some small ϵ , $\text{---} \frac{R_{pop}(h_{train}) - R_{pop}(h^*)}{0} \leq \epsilon$.

$$R_{pop}(h_{train}) \leq \epsilon$$

SIMPLIFICATION

Simplification for today: Assume we are in the realizable setting, which means that $R_{pop}(h^*) = \underline{0}$. I.e. there is some hypothesis in our class \mathcal{H} that perfectly classifies the data.

Formally, for any (\underline{x}, y) such that $\Pr_{\mathcal{D}}[\underline{x}, y] > 0$, $h^*(\underline{x}) = y$.



Extending to the case when $R_{pop}(h^*) \neq 0$ is not hard, but the math gets a little trickier. And intuition is roughly the same.

Probably Approximately Correct (PAC) Learning (Valiant, 1984):

For a hypothesis class \mathcal{H} , data distribution \mathcal{D} and training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, let $h_{train} = \arg \min_h \frac{1}{n} \sum_{i=1}^n \mathbb{1}[h(\mathbf{x}_i) \neq y_i]$.

In the realizable setting, how many training samples n are required so that, with probability $1 - \delta$, \rightarrow *probably*

$$R_{pop}(h_{train}) \leq \underline{\epsilon?} \rightarrow \text{approximately.}$$

The number of samples n will depend on ϵ, δ and the complexity of the hypothesis class \mathcal{H} . Perhaps surprisingly, it will not depend at all on \mathcal{D} .

COMPLEXITY OF HYPOTHESIS CLASS

$$2^{O(d^3)}$$

of functions in H

Many ways to measure complexity of a hypothesis class. Today we will start with the simplest measure: the number of hypotheses in the class, $|\mathcal{H}|$.

Example: What is the number of hypothesis in the class of 3-DNF formulas on d dimensional inputs

$$\mathbf{x} = [x_1, \dots, x_d] \in \{0, 1\}^d$$

$$(_ \wedge _ \wedge _)$$

$$h(\mathbf{x}) = (x_1 \wedge \bar{x}_5 \wedge x_{10}) \vee (\bar{x}_3 \wedge x_2) \vee \dots \vee (\bar{x}_1 \wedge x_2 \wedge x_{10})$$

$$\binom{2d+1}{3}$$

$$O(d^3)$$

$$O(d^3) \left\{ \begin{array}{l} (x_1 \wedge x_2 \wedge x_3) \checkmark \\ (\bar{x}_1 \wedge x_2 \wedge x_3) \checkmark \\ \vdots \end{array} \right.$$

COMPLEXITY OF HYPOTHESIS CLASS

$$|\mathcal{H}| = C^{d+1} = O(2^d)$$

Caveat: Many hypothesis classes are infinitely sized. E.g. the set of linear thresholds

$$h(\mathbf{x}) = \mathbb{1}[\mathbf{x}^T \beta \geq \lambda]$$

But you could imagine approximating \mathcal{H} by a finite hypothesis class. E.g. take values in β, λ to lie on a finite grid of size C . Then how many hypothesis are there?

Formally moving from finite to infinite sized hypothesis classes is a huge area of learning theory (VC theory, Rademacher complexity, etc.)

MAIN RESULT

Consider the realizable setting with hypothesis class \mathcal{H} , data distribution \mathcal{D} , training data set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, and $h_{\text{train}} = \arg \min_h \frac{1}{n} \sum_{i=1}^n \mathbb{1}[h(\mathbf{x}_i) \neq y_i]$.

Theorem

If $n \geq \frac{1}{\epsilon} (\log |\mathcal{H}| + \log \frac{1}{\delta})$, then with probability $1 - \delta$,

$$R_{\text{pop}}(h_{\text{train}}) \leq \epsilon.$$

Roughly how many training samples are needed to learn 3-DNF formulas? To learn (discretized) linear threshold functions?

$$|\mathcal{H}| = O(2^{d^3})$$

$$\log(|\mathcal{H}|) = O(d^3)$$

$$|\mathcal{H}| = c^d$$

$$\log(|\mathcal{H}|) = \underline{d} \log(c)$$

$$O(d/\epsilon)$$

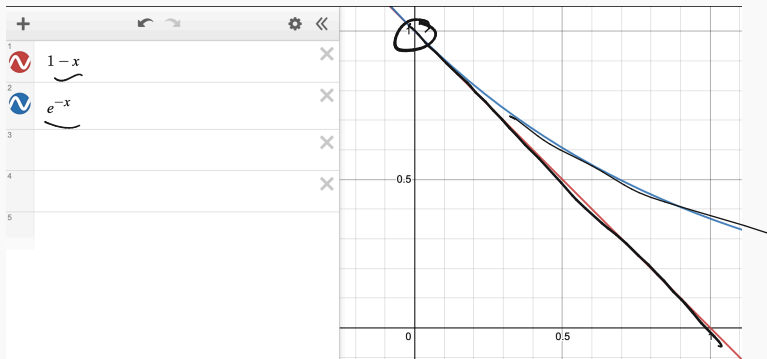
Two ingredients needed for proof:

1. For any $\epsilon \in [0, 1]$, $(1 - \epsilon) \leq e^{-\epsilon}$.
2. (Union bound) Basic but important inequality about probabilities.

ALGEBRAIC FACT

For any $\epsilon \in [0, 1]$ $(1 - \epsilon) \leq e^{-\epsilon}$

$$(e^{-\epsilon})^{1/\epsilon} = e^{-1} = 1/e$$



Raising both sides to $1/\epsilon$, we have the $(1 - \epsilon)^{1/\epsilon} \leq \frac{1}{e} \approx .37$.

~~The specific constant here won't be important.~~

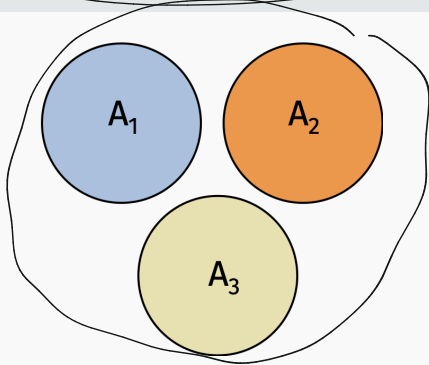
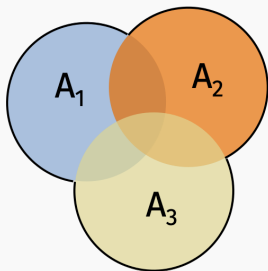
$$\lim_{\epsilon \rightarrow 0} (1 - \epsilon)^{1/\epsilon} = 1/e$$

UNION BOUND

Lemma (Union Bound)

For any random events $\underline{A_1}, \dots, \underline{A_k}$:

$$\Pr[\underline{A_1} \text{ or } \underline{A_2} \text{ or } \dots \text{ or } \underline{A_k}] \leq \underline{\Pr[A_1]} + \underline{\Pr[A_2]} + \dots + \underline{\Pr[A_k]}.$$



Proof by picture.

UNION BOUND

(1, 2, 3, 4, 5, 6)

What is the probability that a dice roll is odd, or that it is ≤ 3 ?

$$= \frac{4}{6}$$

$$\Pr(\text{roll is odd}) = 3/6$$

$$\Pr(\text{roll is } \leq 3) = 3/6 \quad + \quad = 6/6 = 1$$

$$\frac{4}{6} \leq 1$$

What is the probability that a dice roll is 1, or that it is ≥ 4 ?

$$= \frac{4}{6}$$

$$\Pr(\text{dice roll} = 1) = 1/6$$

$$\Pr(\text{dice roll} \geq 4) = 3/6 \quad + \quad = 4/6$$

MAIN RESULT

$$\delta = .001 \quad \log\left(\frac{1}{.001}\right) = \log(1000) = 10$$

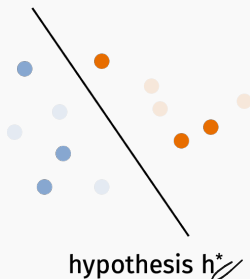
Consider the realizable setting with hypothesis class \mathcal{H} , data distribution \mathcal{D} , training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, and $h_{\text{train}} = \arg \min_h \frac{1}{n} \sum_{i=1}^n \mathbb{1}[h(\mathbf{x}_i) \neq y_i]$.

Theorem

If $n \geq \frac{1}{\epsilon} (\log |\mathcal{H}| + \log \frac{1}{\delta})$, then with probability $1 - \delta$,

$$\underline{R_{\text{pop}}(h_{\text{train}})} \leq \epsilon.$$

First observation: Note that because we are in the realizable setting, we always select and h_{train} with $R_{train}(h_{train}) = 0$. There is always at least one $h \in \mathcal{H}$ such that $h(\mathbf{x}_i) = y_i$ for all i .



Proof approach: Show that for any fixed hypothesis h^{bad} with $R_{pop}(h^{bad}) > \epsilon$, it is very unlikely that $R_{train}(h^{bad}) = 0$. So with high probability, we will not choose a bad hypothesis.

Let h^{bad} be a fixed hypothesis with $R_{pop}(h) > \epsilon$. For (x, y) drawn from \mathcal{D} , what is the probability that $h^{bad}(x) = y$?

$$< \underline{1 - \epsilon}$$

What is the probability that for a training set

$(x_1, y_1), \dots, (x_n, y_n)$ drawn from \mathcal{D} that $h^{bad}(x_i) = y_i$ for all i ? I.e. that $R_{train}(h^{bad}) = 0$.

$$(1 - \epsilon) (1 - \epsilon) \dots (1 - \epsilon) = (1 - \epsilon)^n$$

$$\leq \underline{e^{-\epsilon n}}$$

Claim

For any fixed hypothesis h with $R_{pop}(h^{bad}) > \epsilon$, the probability that $R_{train}(h) = 0$ can be bounded by:

$$\Pr[R_{train}(h^{bad}) = 0] < e^{-\epsilon n}.$$

$$e^{-\epsilon \cdot \frac{1}{\epsilon} \log(|\mathcal{H}|/\delta)} = \frac{\delta}{|\mathcal{H}|}$$

Set $n \geq \frac{1}{\epsilon} \log(|\mathcal{H}|/\delta)$. Then we have that for any fixed hypothesis h^{bad} with $R_{pop}(h^{bad}) > \epsilon$,

$$\Pr[R_{train}(h^{bad}) = 0] < \frac{\delta}{|\mathcal{H}|}$$

UNION BOUND APPLICATION

Let $h_1^{bad}, \dots, h_m^{bad}$ be all hypothesis in \mathcal{H} with $R_{pop}(h) > \epsilon$. How large can m be? Certainly no more than $|\mathcal{H}|$!

$$\begin{aligned} & \Pr[R_{train}(h_1^{bad}) = 0 \text{ or } \dots \text{ or } R_{train}(h_m^{bad}) = 0] \\ & \leq \Pr[R_{train}(h_1^{bad}) = 0] + \dots + \Pr[R_{train}(h_m^{bad}) = 0] \\ & \leq \frac{m \cdot \delta}{|\mathcal{H}|} < \delta \end{aligned}$$

So with probability $1 - \delta$ (high probability) no bad hypotheses have 0 training error. Accordingly, it must be that when we choose a hypothesis with 0 training error, we are choosing a good one. I.e. one with $R_{pop}(h) \leq \epsilon$.

THINGS WE DIDN'T COVER TODAY

- How to deal with the non-realizable setting? E.g. where $\min_h R_{pop} \neq 0$?
- How to deal with infinite hypothesis classes (most classes in ML are)?
- How to find $h_{train} = \arg \min_h \frac{1}{n} \sum_{i=1}^n \mathbb{1}[h(\mathbf{x}_i) \neq y_i]$ in a computationally efficient way?

HAVE A GOOD SPRING BREAK!