# Lecture 5: Learning from data.

Previously: we have a dataset $X$, can we learn about $X$?

e.g. nearest neighbor, heavy hitters.

Today: we have a dataset $X$ that is representative of some underlying process. Can we learn about this underlying process?

This is the basic question underlying ML!

e.g. image classification

$$X = \left\{ \left( \boxed{\text{🐱}}, \text{"cat"} \right), \right.$$

$$\left( \boxed{\text{👤}}, \text{"human"} \right),$$

$$\vdots$$

$$\left. \left( \boxed{\text{🐱}}, \text{"cat"} \right) \right\}$$

$$\Downarrow$$

$$f: \text{new image} \rightarrow \{ \text{"cat"}, \text{"human"} \}$$

Binary classification: $D$ is a distribution over $\mathbb{R}^d$
(Unknown) ground truth $f: \mathbb{R}^d \rightarrow \{0,1\}$.

$X_1, \cdots, X_n \sim D$ are independent draws from $D$.

Labels: $y_i = f(X_i)$, $i = 1, \cdots, n$.

Input: training set $\{ (X_1, y_1), \cdots, (X_n, y_n) \}$

Output: A predictor $g: \mathbb{R}^d \rightarrow \{0,1\}$

**Goal:** $g = f$, or at least, $g$ "looks like" $f$
for typical points in $D$. $\longrightarrow$ quantified via
generalization error.

**Def:**

$$\text{Gen Error}(g) = \Pr_{x \sim D}\left[ g(x) \neq f(x) \right]$$

How to evaluate gen. error? **use test set**

In addition to the training set, maintain separate
test set $\{(x_1', y_1'), \cdots, (x_m', y_m')\}$ that the
predictor cannot see (ie. is not trained on).
Then $\forall i = 1, \cdots, m$

$$\Pr\left[ g(x_i') \neq y_i' \right] = \Pr_{x \sim D}\left[ g(x) \neq f(x) \right]$$
$$= \text{Gen Error}(g)$$

$$\Rightarrow \text{Gen Error}(g) \underset{\uparrow}{\approx} \frac{1}{m} \#\{i : g(x_i') \neq y_i, \, i=1,\cdots,m\}.$$

$\qquad\qquad$ **Chernoff** $\qquad\qquad\underbrace{\phantom{xxxxxxxxxxx}}$

$\qquad\qquad\qquad\qquad\qquad\qquad$ "test error".

We can also define training error
$$\text{Train Error}(g) = \frac{1}{n} \#\{i : g(x_i) \neq y_i, \, i=1,\cdots,n\}$$

**Important: Train Error $\neq$ Test Error**

Overfitting: Train Error is small, but
$\qquad\qquad$ Test Error is large.
$\quad$ Q: When do models generalize or overfit?

# Finite, well-separated case:

$\mathcal{A}$ := finite set of predictors of size $k$

$$\mathcal{A} = \{f_1, \cdots, f_k\}$$

Assume: ① (Realizable) ground truth $f \in \mathcal{A}$

② (Well-separated) $\forall f_i \neq f$ in $\mathcal{A}$, we have

$$\text{GenError}(f_i) \geq \varepsilon.$$

$$\underset{X \sim D}{\Pr}[f_i(x) \neq f(x)] \geq \varepsilon.$$

Let $(x_1, y_1), \cdots, (x_n, y_n)$ be a training set

    ↳ $x_i \sim D$ independent,

       $y_i = f(x_i)$.

Q: How large does $n$ have to be to learn $f$?

Meta-Algo: Output any $g \in \mathcal{A}$ s.t.

$$\text{TrainError}(g) = 0.$$

Claim: For any $\delta > 0$, let $n \geq \frac{1}{\varepsilon}(\log k + \log^{1}/\delta)$
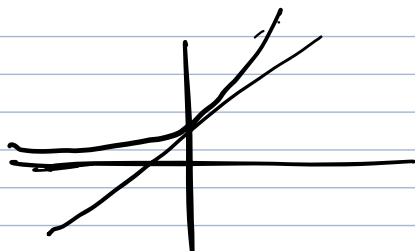
Then $\underset{x_1, \cdots, x_n}{\Pr}[g \neq f] \leq \delta.$

                                     tight.

pf: fix some $f_j \neq f$

$$\underset{x_1, \cdots, x_n}{\Pr}[\text{TrainError}(f_j) = 0] = \prod_{i=1}^{n} \underbrace{\Pr[f_j(x_i) = f(x_i)]}_{\leq (1-\varepsilon)}$$

$\forall x: (1-x) \leq e^x$



$$\leq (1-\varepsilon)^n$$

$$\leq e^{-\varepsilon n}$$

$$= e^{-\varepsilon \cdot \frac{1}{\varepsilon}(\log k + \log^{1}/\delta)}$$

$$= \delta/k.$$

Now, by union bound:

$$\Pr\left[\exists f_i \neq f : \text{TrainError}(f_i) = 0\right]$$

$$\leq \sum_{f_i \neq f} \Pr\left[\text{TrainError}(f_i) = 0\right] \leq \frac{\delta}{k} \cdot (k-1)$$

$$\leq \delta.$$

$\Rightarrow$ w.p. $\geq 1-\delta$, only $g \in A$ w/ TrainError$(g) = 0$
is $g = f$.

---

What about not well-separated?
  In general, can't hope to recover $f$!
But, maybe you can get something that acts
like $f$ for most points in D.

## Probably Approximately Correct (PAC) learning [Valiant '89]

Let $\varepsilon, \delta > 0$. Given training data, find $g \in A$

$$\text{GenError}(g) = \Pr_{x \sim D}\left[g(x) \neq f(x)\right] \leq \varepsilon$$

with probability $\geq 1-\delta$.

For this, the same proof works!

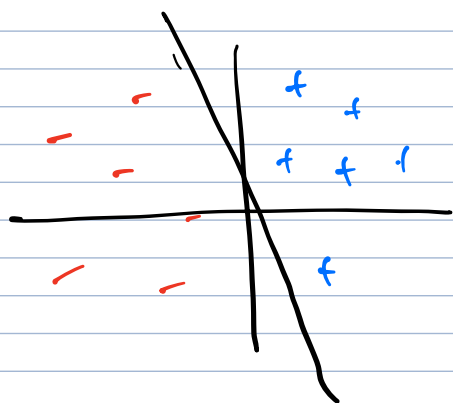$$n \geq \frac{1}{\varepsilon}\left(\log k + \log 1/\delta\right)$$

---

What about infinite families $A$?
  $A =$ any function?  Impossible!

need structure

## Linear classifiers:

parameter: $\theta \in \mathbb{R}^d$

$$f_\theta(x) = \text{sign}(\langle \theta, x \rangle)$$

$$= \begin{cases} 1 & \text{if } \langle \theta, x \rangle \geq 0 \\ 0 & \text{o.w.} \end{cases}$$



Obs 1: Since we only care about angles, wlog can assume $\|\theta\|_2 = 1$.

But there are still infinitely many unit vectors...

But "only" exponentially many distinct ones!

Recall curse of dimensionality: can fit exponentially many small balls into unit ball.

↰ tight

$$k = \# \text{ "distinct" directions } \approx C^d$$

$$\Rightarrow \quad n \geq \frac{1}{\varepsilon}(d \log C + \log 1/\delta) \quad \text{suffices!}$$

Can slightly improve this:

Vapnik-Chervonenkis theory: (VC theory).

$$n \geq \frac{1}{\varepsilon}(d + \log 1/\delta).$$

Rule of thumb: # samples ≥ # free parameters
→ good generalization

not necessary! e.g. deep networks.

What if $f \notin A$? "agnostic learning"
i.e. what if optimal error > 0 ?

Goal: Given $\varepsilon, \delta > 0$, output $g \in A$ s.t.
$$\Pr_{x \sim D}[g(x) \neq f(x)] \leq OPT + \varepsilon$$
$$OPT = \min_{g^* \in A} \Pr[g^*(x) \neq f(x)]$$

with probability $\geq 1 - \delta$.

Meta-Algo: empirical risk minimization (ERM)

output $g \in A$ with smallest training error

When $A$ = linear classifiers,        still linear in # params

$$n \geq \Omega\left(\frac{1}{\varepsilon^2}\left(d + \log 1/\delta\right)\right)$$

why $\varepsilon^2$?  How many coin flips do you need to
distinguish a coin w/ bias $p$ vs $p + \varepsilon$?

Realizable: $p = 0$ → just need to see 1 tails
occurs whp after $\geq 1/\varepsilon$ throws.

Agnostic: $p = \frac{1}{2}$ → need to see many tails
to distinguish. → $1/\varepsilon^2$ throws.

Algorithmic aspects: How do we compute the the best
fit classifier? say for linear classifiers.

Realizable: poly-time (e.g. support vector machine (SVM),
                                    linear programming).

Agnostic:    NP-hard    ..
             use heuristics!
                    ( e.g. linear programming, stochastic gradient descent)
                                    (SGD)

Generalizing beyond linear?
    there are many, many types of $A$ used in practice.

    However, many have a same basic structure:

    $k: \mathbb{R}^d \rightarrow \mathbb{R}^D$ ,    $D \gg d$,
    $k$ is "kernel" map , and apply linear classifier in $\mathbb{R}^D$
to lifted data.
                e.g.    $k =$ polynomials, neural networks.
Intuition:    $k$ maps data to a set of useful
            features, and just apply simple classifier on top.