# CS 121.5 New Frontiers in ML Theory

## What is learning?

### Supervised Learning

- $X \rightarrow Y$
  input   labels
- given labeled ex.
  $(x_i, y_i)$

- Correct on examples close to those that were given
- Assume samples are iid from uniform $D$
  $(x, y) \sim D$

Given n iid samples $\{(x_i, y_i)\} \sim D$, $f \in F$   ← **what you're hoping to learn**

Goal   Output $\hat{f} : X \rightarrow Y$ w/ low

test error $L_D(\hat{f}) = \Pr_{(x_N) \sim D}[\hat{f}(x) \neq y] \leq \varepsilon$

   ↑ **what you actually learn**

$D_f = \{(x, f(x))\}$ ← label y guaranteed to be a function of x

### PAC Learning (Probably Approximately Correct)

Family $F$ is PAC-learnable if $\exists$ alg $A$ s.t. $\forall D_f \in D_F$
$\forall \varepsilon, \delta$ and $n = |S| = \text{poly}(\frac{1}{\varepsilon}, \frac{1}{\delta})$

$\hat{f} \leftarrow A(S)$, $S = \{(x_i, y_i)\} \sim D^n$

$\Pr[L_D(\hat{f}) \leq \varepsilon] \geq 1 - \delta \sim 0.999$

   ↑ accuracy   ↑ measurement of "good" dataset sample

no matter what dist of input, you'll get a good quality classifier (low value of loss function).

<u>Claim</u>. $|F|$ finite $\rightarrow$ F is PAL-learnable w/
$$n = \frac{\log(|F|/\delta)}{\varepsilon^2} \text{ samples}$$

<span style="color:blue">most important term</span>

w/o regards to efficiency

<span style="color:crimson">↑n means you can afford to shrink $\varepsilon$ and $\delta$</span>

<u>Empirical Risk Minimization (ERM)</u>
· Want $f$ s.t. $\mathcal{L}_D(f) = \underset{x,y \sim D}{\mathbb{E}}[\mathbb{1}\{f(x) \neq y\}]$ small

<span style="color:crimson">population loss</span>

$\underbrace{\qquad}_{l(f(x), y)}$

<u>Try</u>: $\underset{f \in F}{\text{argmin}} \; \mathcal{L}_S(f) = \frac{1}{n} \sum_{(x_i, y_i) \in S} l(f(x_i), y_i)$

<span style="color:crimson">empirical loss</span>

works w/ ==uniform convergence==$_{(\varepsilon)}$ (of F, D, n) =
will hold over $S \sim D^n$: $\{f \in F : |\mathcal{L}_D(f) - \hat{\mathcal{L}}_S(f)| \leq \varepsilon\}$

<u>Corollary</u>: if family has uniform convergence w/
$\varepsilon$, $\mathcal{L}_D(f_{erm}) \leq \underset{f \in F}{\min} \mathcal{L}_D(f) + \varepsilon$unit

<span style="color:blue">best classifier</span>

<u>Lemma</u>: $\underset{S \sim D^n}{\Pr}\left[\exists f \in F : |\mathcal{L}_D(f) - \hat{\mathcal{L}}_S(f)| > \varepsilon\right] \leq |F| e^{-\varepsilon^2 n}$

<span style="color:blue">$\delta$</span>

· 1D random variable variation from mean:
Chernoff Bounds (exp. small in $t^2$)

- Union bound over all function in family:
  probability that any function deviates ≤
  num functions · probability one deviates

EX: Binary Linear Classifier  $f_w(x) = I\{<w,x> \geq 0\}$

$$w \in \mathbb{R}^d, \quad |F| \sim 2^{O(d)}$$
$$n \sim O(d)$$

When does uniform convergence hold?
- not hold if function family is too big
- holds if small

**High Error:**
- too simple classifier (underfitting)
- minimum pop error of family is no good

- overfitting → classifier fits highly to train set

## Bias-complexity/variance Tradeoff



error/loss of (worst) ERM

aka generalization gap

$\epsilon_{unif}$, uniform convergn gap (best epsilon you can get)

— train error
— test error

complexity of F
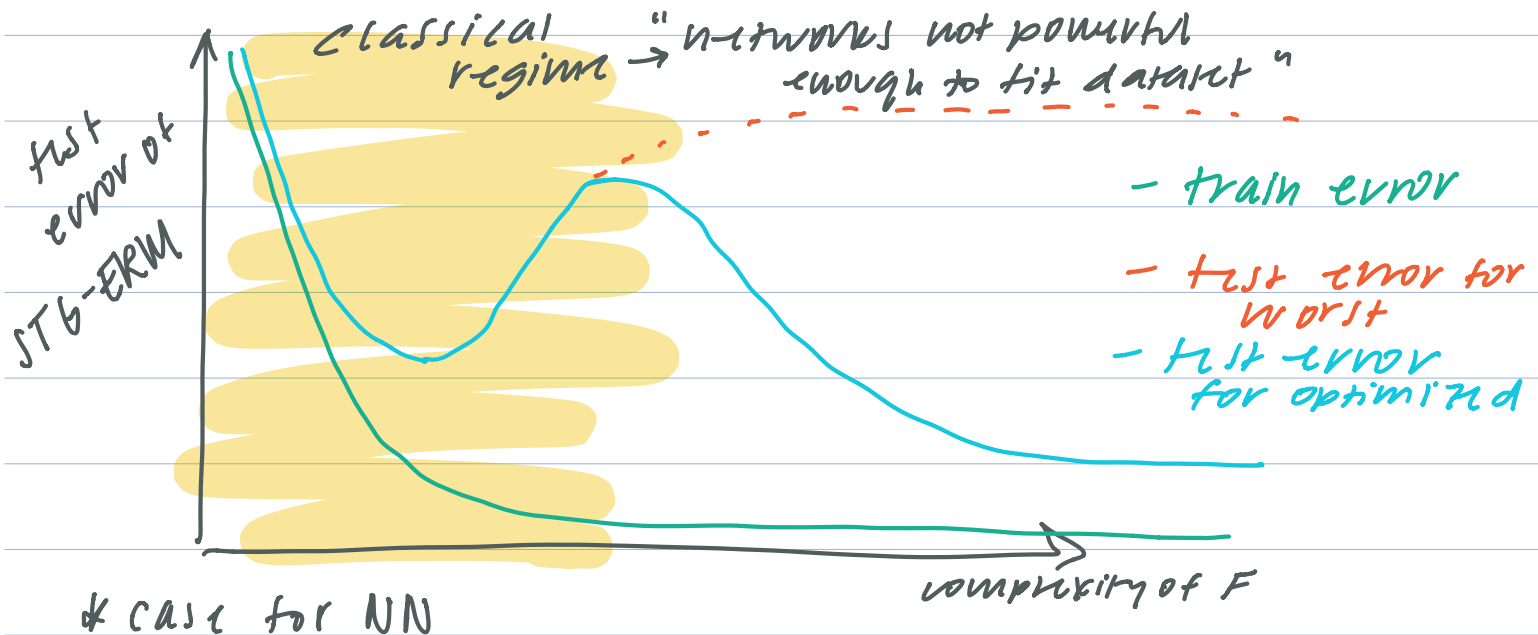
* n and D constant

<u>random matrix complexity</u>: how well $\widehat{F}$ can fit ← family of functions

n randomly samples from distribution

$R_{D,n}(F)$

· Can fit, family is too complex



test error of STG-ERM

classical regime → "networks not powerful enough to fit dataset"

— train error
— test error for worst
— test error for optimized

complexity of F

* case for NN

NN factors: dist D, model F, optimization algo, number samples n

<u>Training Algorithm A</u>: input $(X_i, y_i) \to$ output model M

Model Complexity$_D$ (A) := max n s.t. Train Error(A(S))

"$\approx 0$ for $S \sim D$"

For n train samples:

under-parameterized Model Complexity (A) < n

critically-parameterized Model Complexity (A) = n

over-parameterized Model Complexity (A) >> n