

# Mathematical Foundations of ML (CS 4785/5783)

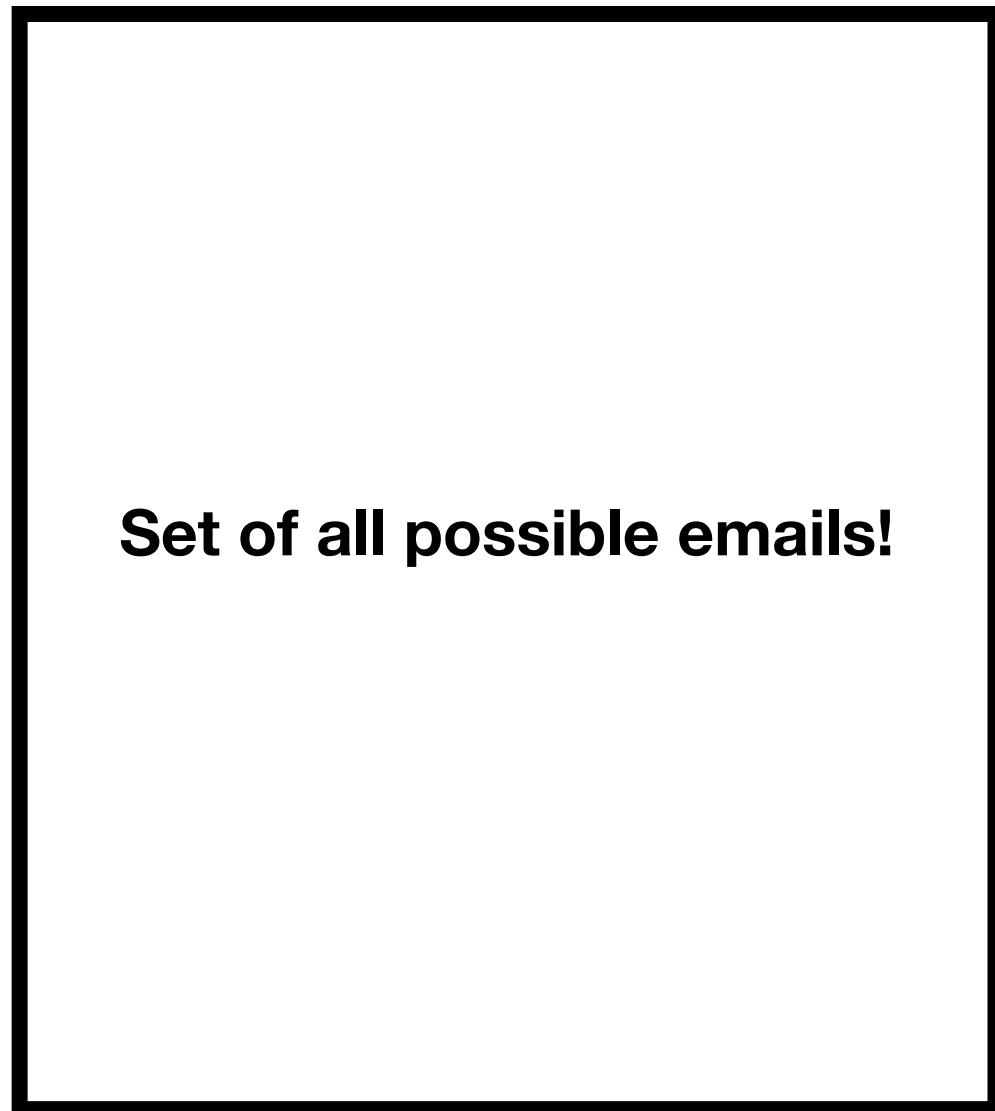
## Lecture 2

### Statistical Learning and Uniform Convergence

<http://www.cs.cornell.edu/Courses/cs4783/2022sp/notes02.pdf>

# SCENARIO II

**Universe of instances**



**Set of all possible emails!**

**U**

**On each round  $t$ :**

Email  $x_t$  is composed, possibly by spammer!

System classifies email as  $\hat{y}_t$

True label  $y_t = f_{i^*}(x_t)$  revealed

We get feedback every round. But spammer can pick next email.

Goal: Make as few mistakes as possible.

# SCENARIO II

**How about using the same algorithm from scenario 1 for each t (re-run)?**

**How many mistakes would it make?**

**Ans:  $N-1$**

# SCENARIO II

**Algorithm:**

Pick  $\mathcal{F}_t = \{f_i : i \in [N], \forall s < t, f_i(x_s) = y_s\}$

Set  $\hat{y}_t = \text{Majority}(\{f(x_t) : f \in \mathcal{F}_t\})$

**Mistake Bound:**

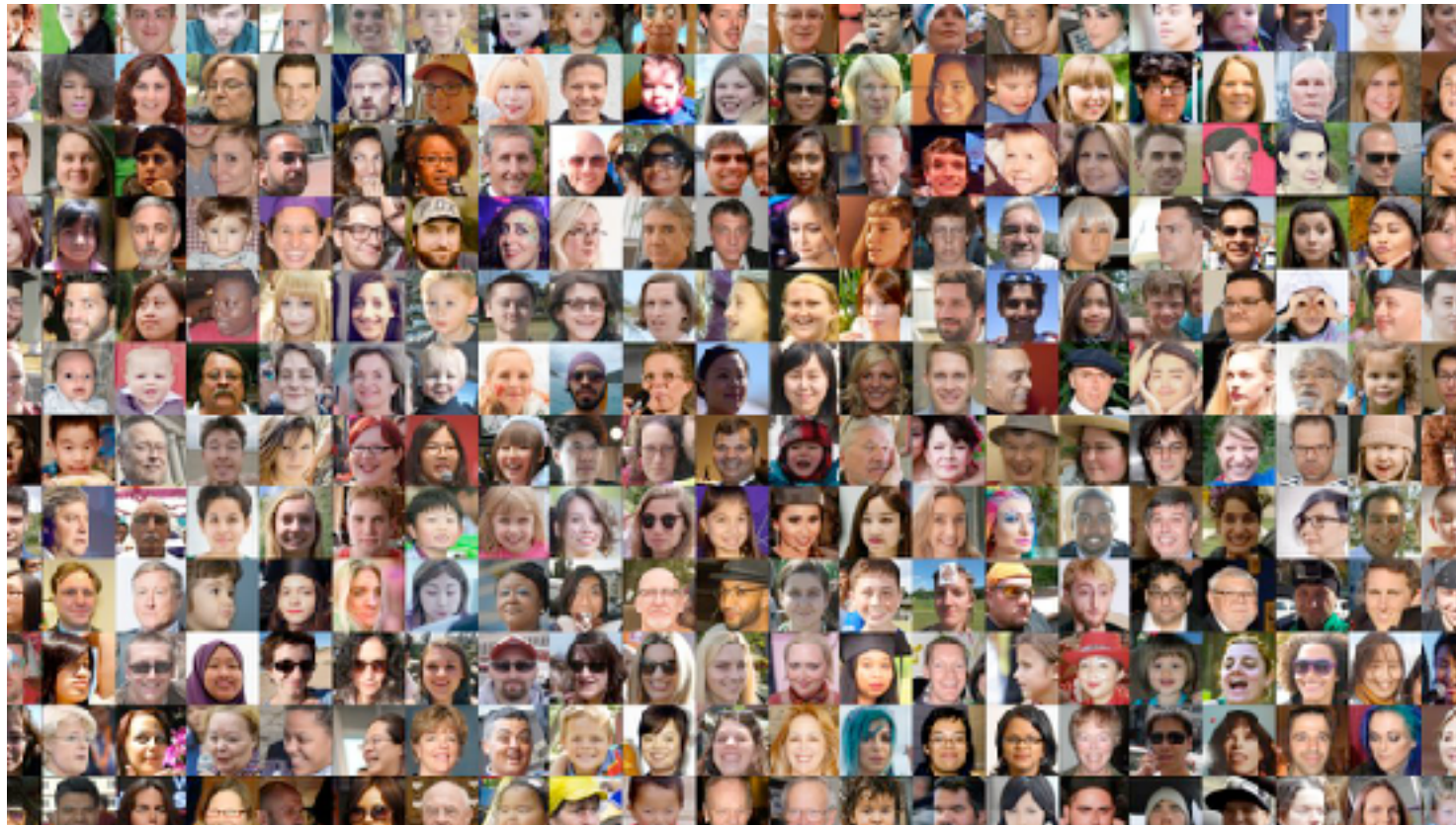
$$\sum_t \mathbf{1}\{\hat{y}_t \neq y_t\} \leq \log_2 N$$

**Why?**

# STATISTICAL LEARNING FRAMEWORK

**Eg: ML for Face recognition**

$\mathcal{X}$  : set of all images



**U**

**We don't have access to U, we just need the samples**

# STATISTICAL LEARNING FRAMEWORK

**When we deploy the system, do we really sample from  $U$  at random?**



**In summer**

**In winter**

**No assumption is right but some are useful!**

# STATISTICAL LEARNING FRAMEWORK

**D** is a distribution on  $\mathcal{X} \times \mathcal{Y}$

**D** captures the idea of this set **U**

Training sample  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$  Each  $(x_t, y_t) \sim \mathbf{D}$

Risk of a model  $g$  defined as  $L_{\mathbf{D}}(g) = \mathbb{E}_{(x,y) \sim \mathbf{D}} [\ell(g(x), y)]$

(Future instances drawn from **D**)

Excess risk of model  $g$  w.r.t. model class  $\mathcal{F}$  defined as

$$L_{\mathbf{D}}(g) - \min_{f \in \mathcal{F}} L_{\mathbf{D}}(f)$$

**Goal: provide an algorithm for which excess risk is small**



# TRAINING LOSS VS TEST LOSS

$$\text{Training loss: } \hat{L}_S(g) = \frac{1}{|S|} \sum_{(x,y) \in S} \ell(g(x), y)$$

Test loss: Draw fresh samples (not used by algorithm)  
and compute average error on that

Test loss is a good proxy for risk  
(provided we never use it in any sense for training/parameter tuning etc.)



# THE COMMON FALLACY

$\forall f \in \mathcal{F}, P \left( \left| L_{\mathbf{D}}(f) - \hat{L}_S(f) \right| \text{ is large} \right) \text{ is small}$

Algorithm picks  $\hat{f}_S \in \mathcal{F}$  and so

$P \left( \left| L_{\mathbf{D}}(\hat{f}_S) - \hat{L}_S(\hat{f}_S) \right| \text{ is large} \right) \text{ is small}$

**THIS IS FALSE IN GENERAL!**

**Breakout room 3 mins**

# THE COMMON FALLACY

- The issue with benchmark dataset like CIFAR and Imagenet
- Double edged sword

# EMPIRICAL RISK MINIMIZATION

**Pick a model in class that minimizes training error**

$$\hat{f}_{\text{ERM}} \in \arg \min_{f \in \mathcal{F}} \hat{L}_S(f)$$

- When does this succeed?
  - When model class is too complex, we already saw this can fail
  - When model class is say just one function, it succeeds due to law of large numbers (concentration)
  - In general how well does this algorithm do?

# ERM AND UNIFORM CONVERGENCE

$$\begin{aligned} P\left(L_{\mathbf{D}}(\hat{f}_{\text{ERM}}) - \min_{f \in \mathcal{F}} L_{\mathbf{D}}(f) > 2\epsilon\right) &= P\left(L_{\mathbf{D}}(\hat{f}_{\text{ERM}}) - \hat{L}_S(\hat{f}_{\text{ERM}}) + \hat{L}_S(\hat{f}_{\text{ERM}}) - \min_{f \in \mathcal{F}} L_{\mathbf{D}}(f) > 2\epsilon\right) \\ &= P\left(L_{\mathbf{D}}(\hat{f}_{\text{ERM}}) - \hat{L}_S(\hat{f}_{\text{ERM}}) + \max_{f \in \mathcal{F}} \left(\hat{L}_S(\hat{f}_{\text{ERM}}) - L_{\mathbf{D}}(f)\right) > 2\epsilon\right) \\ &\leq P\left(L_{\mathbf{D}}(\hat{f}_{\text{ERM}}) - \hat{L}_S(\hat{f}_{\text{ERM}}) + \max_{f \in \mathcal{F}} \left(\hat{L}_S(f) - L_{\mathbf{D}}(f)\right) > 2\epsilon\right) \\ &\leq P\left(\max_{f \in \mathcal{F}} \left|\hat{L}_S(f) - L_{\mathbf{D}}(f)\right| > \epsilon\right) \end{aligned} \tag{1}$$

# ERM OVER FINITE CLASS

If losses are bounded by 1 (in absolute) and  $|\mathcal{F}| < \infty$ , then, for any  $\delta > 0$  with probability at least  $1 - \delta$ ,

$$L_{\mathbf{D}}(\hat{f}_{\text{ERM}}) - \min_{f \in \mathcal{F}} L_{\mathbf{D}}(f) \leq \sqrt{\frac{8 \log(2|\mathcal{F}|/\delta)}{n}}$$

# ERM OVER FINITE CLASS

Hoeffding Inequality: Let  $Z_1, \dots, Z_n$  be a sequence of  $n$  random variables bounded by 1, drawn iid from a fixed distribution. Then:

$$P\left(\left|\frac{1}{n}\sum_{t=1}^n Z_t - \mathbb{E}Z\right| > \epsilon\right) \leq 2 \exp\left(-\frac{n\epsilon^2}{2}\right)$$



# ERM OVER FINITE CLASS

Hoeffding Inequality: Let  $Z_1, \dots, Z_n$  be a sequence of  $n$  random variables bounded by 1, drawn iid from a fixed distribution. Then:

$$P\left(\left|\frac{1}{n}\sum_{t=1}^n Z_t - \mathbb{E}Z\right| > \epsilon\right) \leq 2 \exp\left(-\frac{n\epsilon^2}{2}\right)$$

**Proof idea:**

For each  $f \in \mathcal{F}$  define  $Z_t^f = \ell(f(x_t), y_t)$

**Apply Hoeffding for each  $f$  individually**

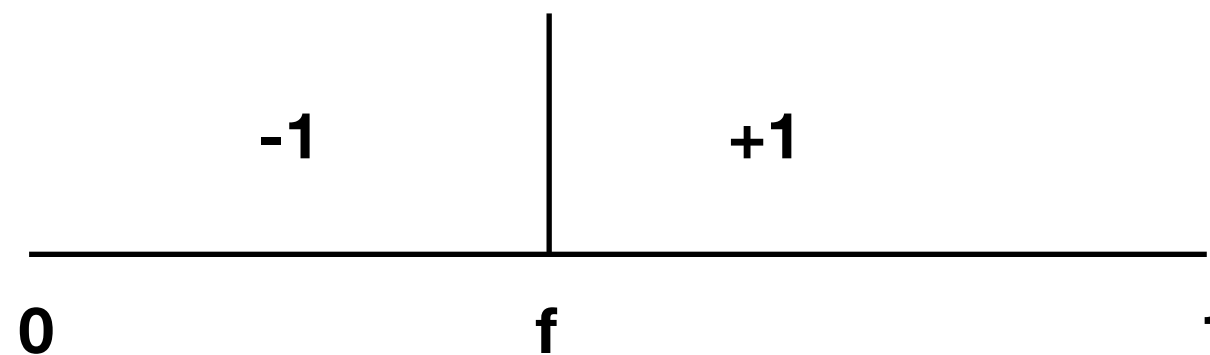
**Use union bound to move to uniform deviation**

# BEYOND FINITE MODEL CLASS

- Idea 1: Find a finite set  $\mathcal{F}'$  such that for any  $f \in \mathcal{F}$  there exists an  $f' \in \mathcal{F}'$  s.t.

$$\forall x, y, \quad |\ell(f'(x), y) - \ell(f(x), y)| < \Delta$$

- But this may not always work, consider the example of learning thresholds:



$$\mathcal{X} = [0, 1] \quad f(x) = \text{sign}(x - f)$$

$\mathcal{F}$  indexed by set  $[0, 1]$

For any  $\Delta < 1/2$ , this class cannot be approximated by a finite set.