

## CATEGORIES:

→ Active Learning

Ask / seek questions or labels during learning.

Goal: Reduce the number of queries.

## HOW TO MODEL SUPERVISED LEARNING:

INPUT: Some domain  $X$   
 $x$

eg: images, documents, browsing  
profiles, credit card history.

LABELS: cats / dogs / foxes

$L$   $\{0, 1\}$

$\mathbb{R}$

DATASET:  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$\swarrow \quad \searrow$   
 $x \quad L$

GOAL: Given a new  $x \in X$ ,

want to predict its label!

ASSUMPTION 1: There is an underlying ground truth distribution on training and test examples.

i.e. Both training, test similar cases.  
if we train on dog, cat and test on polar bear, makes no sense.

Distribution,  $D$  on  $x$ :  $x_1, \dots, x_n$  are i.i.d samples from  $D$ .

$x$  is also from  $D$ .

ASSUMPTION 2: Labels cannot be arbitrary.

There is a class of functions  $H: x \rightarrow L$ .

ASSUMPTION 3: Prediction cannot be accurate all the time.

eg: we say model is 60% accurate.

$$\rightarrow \Pr [y \neq f(x)] \leq \epsilon \rightarrow \text{"error"}$$

$x \sim D$   $\downarrow$   $\hookrightarrow$  my prediction

we need to relax " $\neq$ "  
depending on application.

**ASSUMPTION 4:** Cannot expect to succeed always!

Sometimes, the dataset training data when chosen  
comes out very bad  $\rightarrow$  like all dogs! Now  
there is no way it can predict cats.

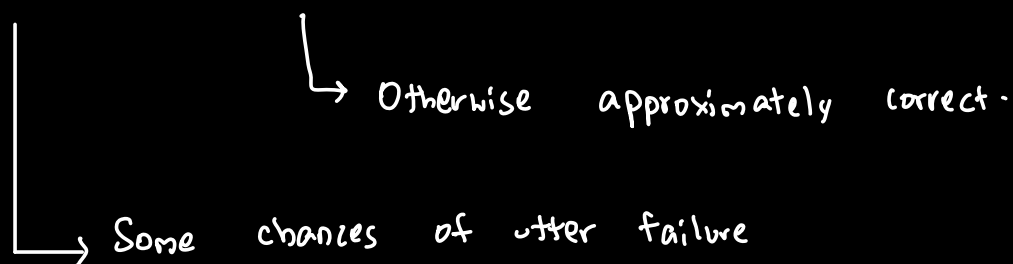
$\rightarrow$  "Probability of failure"  $\delta$ .

This is the chance of utter failure

$\rightarrow$  Probability of being able to build a  
model with 60% accuracy.

PAC MODEL:

(PROBABLY APPROXIMATELY CORRECT)



An algorithm  $(\epsilon, \delta)$ -PAC learns a hypothesis class  $H$  with sample complexity  $n(\epsilon, \delta, H)$  if

INPUT:  $(x_1, f^*(x_1)), \dots, (x_n, f^*(x_n))$  where  
 $x_i \leftarrow D$ ,  $f^* \in H$ .

OUTPUT: Some predictor  $h$

$\forall$  with probability  $1 - \delta$  over  $x_1, \dots, x_n$

$$\Pr [h(x) = f^*(x)] \geq 1 - \epsilon$$

$$x \leftarrow D$$

PAC is a very strong requirement! Very difficult to find  $h$  even with 51% accuracy. So use, may be, the structure in the distribution to simplify.

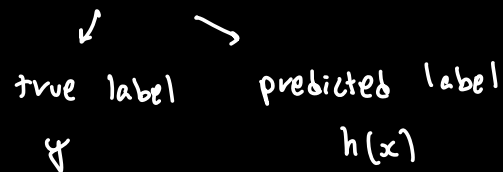
→ OFTEN BECOMES INTRACTABLE.

## LEARNING AS OPTIMIZATION:

→ We have an underlying  $D$  on  $X$ . Labels  $L$ .

→ Dataset  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

→ Loss function  $l: L \times L \rightarrow \mathbb{R}$



Goal: Find a hypothesis  $h$  such that

$$\sum_{i=1}^n l(h(x_i), y_i) \text{ is small}$$

Problem: Just memorize all  $(x_i, y_i)$  and this will satisfy.

PARAMETERIZED HYPOTHESIS CLASS:  $\mathcal{H}$

Restrict, say memory available, this will prevent memorization. This restricted space (subspace) is  $\mathcal{H}$ , parameterized by  $\theta$ .

$$\min_{\theta \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n \ell(h_{\theta}(x_i), y_i)$$

is the predictor as specified by the "parameters"  $\theta$ .

This is "Empirical Risk Minimization" (ERM).

We are using loss on training data.

IDEAL:  $\mathbb{E}_{x \sim \mathcal{D}} [\ell(h_{\theta}(x), y)]$   
(testing data)

EXAMPLE 1 (ERM):  $\mathcal{X} \equiv \mathbb{R}^d$  :  $d \equiv \mathbb{R}$  ↗  $d$ -dimensional vector.

Parametric family : Linear predictors (predictors are linear functions)  
 $(\mathcal{H}) \equiv \mathbb{R}^d$

$$\begin{aligned} h_{\theta}(x) &= \sum_{i=1}^d \theta_i x_i = \underbrace{\langle \theta, x \rangle}_{\text{inner product}} \\ &= \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d. \end{aligned}$$

(More generally,  $h_{\theta, b}(x) = \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_d x_d + b$ ).

LEAST SQUARES REGRESSION:

$$h_{\theta}(x) = \langle \theta, x \rangle$$

$$l : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$$

$$l(a, b) = (a - b)^2$$

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n l(h_{\theta}(x_i), y_i)$$

$$= \frac{1}{n} \sum_{i=1}^n (\langle \theta, x_i \rangle - y_i)^2$$

ERM: Find  $\theta$  to minimize  $L(\theta)$

EXAMPLE 2:  $l(a, b) = |a - b|$  "L1-REGRESSION"

$$\text{ERM} \quad \min L(\theta) = \frac{1}{n} \sum_{i=1}^n |\langle \theta, x_i \rangle - y_i|$$

EXAMPLE 3 : What if labels are discrete?

$$L = \{0, 1\}$$

PARAMETRIC FAMILY  $h_{\theta}(x) = \begin{cases} 1 & \text{if } \langle \theta, x \rangle > 0 \\ 0 & \text{if } \langle \theta, x \rangle < 0 \end{cases}$

$$l(a, b) = \begin{cases} 1 & \text{if } a \neq b \\ 0 & \text{else} \end{cases}$$



LINEAR THRESHOLD

FUNCTIONS

(HALFSPACES)

$$L(\theta) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(\text{sign}(\langle \theta, x_i \rangle) \neq y_i)$$



0 if sign matches  $y_i$

1 else.