

CSCI 699: Theory of ML

What are the goals for the class?

- How much data is needed for learning?
 - How much computation needed for learning?
 - How much memory need for learning?
 - Tradeoffs b/w these resources?
 - How to be robust to worst-case perturbations?
 - " to out-of-distribution data?
 - How to learn models which are fair?
 - How to understand these questions for deep nets?
- 

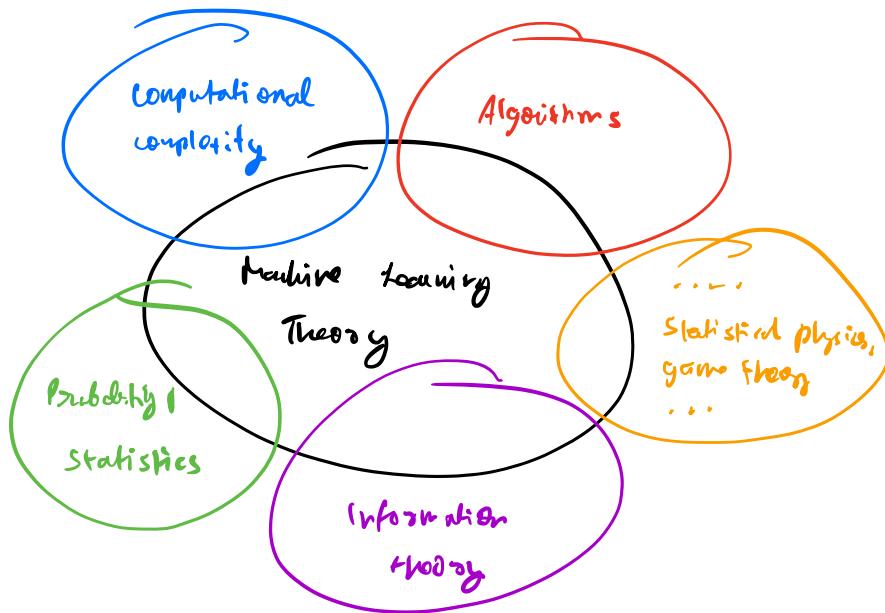
Goals for today's class:

- What is ML? What is ML theory?
- Define supervised learning
- PAC Learnability
- PAC Learnability results.

what is ML? What is ML theory?

ML: programs that can automate or improve based on data.

Goal of this class: understand the theory behind ML.



Why study ML theory?

- Inform practical algorithms, help in designing new algorithms / models for new settings.
- Boosting: Valiant + Kearns asked "Can you boost a weak learner algo. which gets 51% accuracy on binary classification to a strong learner which gets 95% accuracy?" led to Boosting, AdaBoost etc.
- Optimization algo. for ML.
- Modern applications (fairness, robustness, privacy)
- To develop a better understanding of learning, computation.

Supervised learning: A gentle introduction

Training set: input / output pairs

Learn some pattern to accurately predict outputs of unseen inputs
for e.g. image classifications
machine translation

SAT solving

X : input space

Y : output space

$X \subset \mathbb{R}^d$

$Y \in \{-1\}$

Come up with predictor $f(\cdot)$ which predicts output of x

Loss function: $l(f(\cdot), y)$

Classification problem: $l(f(\cdot), y) = \mathbb{I}(f(\cdot) \neq y)$ (0/1 loss)

Regression problem: $l(f(\cdot), y) = (f(\cdot) - y)^2$ (squared loss)

Distribution P over (x, y)

Risk of supervised learning problem:

$$R(f) = \mathbb{E}_{(x, y) \sim P} l(f(x), y)$$

We don't actually know P !

training / test paradigm:

i.i.d. data: get n samples drawn i.i.d from P .
Training set S .

$$f^* = \min_{\text{function } f: X \rightarrow Y} \mathbb{E}_{(x,y) \sim P} l(f(x), y)$$

Empirical risk minimizer :

Commit to a set of functions \mathcal{F} beforehand.

$$f_{S, \text{ERM}} = \underset{f \in \mathcal{F}}{\operatorname{arg\,min}} \frac{1}{n} \sum_i l(f(x_i), y_i)$$

Choose function f_S from training set S .

$$\begin{aligned} R(f_S) - R(f^*) &= R(f_S) - R(f_{S, \text{ERM}}) \quad \text{Computation error} \\ &\quad + \cancel{R(f_{S, \text{ERM}})} - \inf_{f \in \mathcal{F}} R(f) \quad \text{Estimation error} \\ &\quad + \cancel{\inf_{f \in \mathcal{F}} R(f) - R(f^*)} \quad \text{Approximation error} \end{aligned}$$

- Computation : finding the ERM might be expensive.
- Estimation : only have samples.
- Approximation : How much worse is it to work over \mathcal{F} .

An example : linear regression

$$X = \{x: x \in \mathbb{R}^d, \|x\|_2 \leq 1\} \quad (\text{unit ball})$$

$$Y = \mathbb{R} \quad \text{if } Y = \langle w^T, x \rangle, \text{ then approximation error} = 0,$$

$$\text{LC } f(A, Y) = \frac{1}{2} (f(x) - Y)^2$$

$$\underline{F = \{ f_w(x) = \langle w, x \rangle, w \in \mathbb{R}^d \}}$$

Computation error: Can compute ERM using closed-form solution.

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

$$\partial = \frac{\partial \hat{R}_s(w)}{\partial w} = \frac{1}{n} \sum_{i=1}^n x_i (\langle w, x_i \rangle - y_i)$$

$$\sum x_i y_i = \left(\sum_{i=1}^n x_i x_i^T w \right)$$

$$\underline{w = \left(\sum x_i x_i^T \right)^{-1} \sum x_i y_i}$$

$$\underline{w = (x^T)^{-1} x y}$$

$$X = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ 1 & 1 & \dots & 1 \end{pmatrix} \quad d+n$$

time complexity: $x x^T: O(nd^2)$ ($d \times n$ matrix
 $n \times d$ matrix)

$$(x x^T)^{-1}: O(d^3)$$

$$x y: O(nd)$$

$$(x x^T)^{-1} x y: O(d^2)$$

$$\therefore O(nd^2 + d^3)$$

* Gradient descent: $w_0 \leftarrow 0$

$$w_{t+1} \leftarrow w_t - \eta_t \nabla \hat{R}_s(w_t)$$

$$w_{t+1} = w_t - \eta_t \left(\frac{1}{n} \sum_{i=1}^n +_i (\langle w_t, +_i \rangle - y_i) \right)$$

per step complexity : $O(nd)$

in T steps, computation error is e^{-T} .

* stochastic gradient descent

Estimation error : How much data do you need to find a solution which has small error on the distribution?

Approximation error : If true model is not linear,

PAC learning (Probably Approximately Correct learning)

$$x \in \mathbb{R}^d$$

Label space: $Y = \{0, 1\}$

Loss function $l(f(x), y) = \mathbb{1}(f(x) \neq y)$

Crucial idea: no assumption over the dist. of x

but $y = h^*(x)$ for some $h \in \mathcal{H}$

[realizability assumption]

Learner knows \mathcal{H}

Def (PAC learnability):

with $n_{\mathcal{H}}(\epsilon, \delta)$ samples

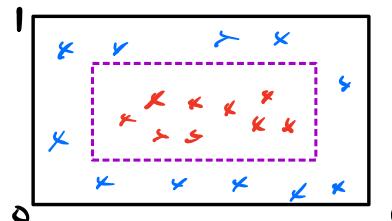
A hypothesis class \mathcal{H} is PAC-learnable if for a learning algorithm with the following property: for every $\epsilon, \delta \in (0, 1)$, every distribution D over X and every $h \in \mathcal{H}$, when the algorithm is given $n_{\mathcal{H}}(\epsilon, \delta)$ samples drawn from D & labelled by h , the alg. produces a hypothesis \hat{h} set. with probability $1-\delta$, $R(\hat{h}) \leq \epsilon$. (This prob. is over randomness in training sets & any interval alg. randomness.)

Example: Axis-aligned rectangles.

$$X = [0, 1] \times [0, 1]$$

\mathcal{H} = rectangles

D : uniform dist.



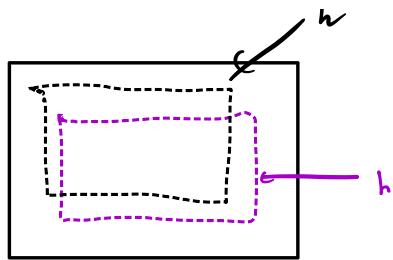
failure probability δ :

Can't get bad set of samples →



Approximation parameter ϵ :

Can only do approximately correct
(cannot recover h)



Any dist D:

Maybe the distribution
is such that we
cannot learn h , but
could still predict well

