

Lecture 1 Introduction

What this course is about:

How does one optimally extract information from data?

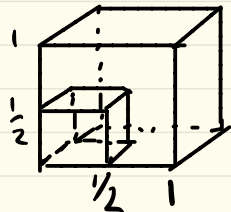
$$S_n = \{z_1, z_2, \dots, z_n\} \stackrel{\text{iid}}{\sim} \mathcal{P}$$

What makes this hard? $\xrightarrow{\text{probability measure}} \mathbb{R}^d$

Curse of dimensionality

- n is big \Rightarrow computation becomes harder
- d is big \Rightarrow statistical accuracy becomes worse

Problems become harder in high dimensions



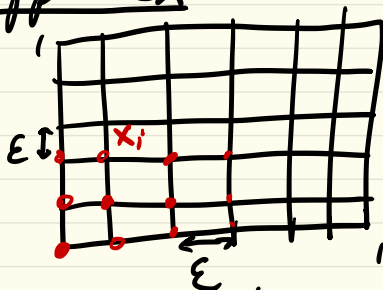
$\text{Vol}(2Q) = 2^d \text{Vol}(Q)$
 \Rightarrow need exponentially many small balls to cover a set in HD.

Ex: (Monte-Carlo Simulation)

How do we compute an integral

$$\int_{[0,1]^d} f(x) dx ?$$

Naive approach: Discretize



$(\frac{1}{\epsilon})^d$ points

If f is Lipschitz, then

$$\left| \epsilon^d \cdot \sum_i f(x_i) - \int_{[0,1]^d} f(x) dx \right| \leq O(\epsilon)$$

So you need $(\frac{1}{\epsilon})^d$ function eval
to estimate the integral to ϵ
accuracy.

Can do better with a randomized algo:

$$E \triangleq \int_{[0,1]} f(x) = E f(X)$$

where $X \sim \text{Unif}([0,1]^d)$

Draw $x_1, \dots, x_n \stackrel{iid}{\sim} \text{Unif}([0,1]^d)$
and set $I = \frac{1}{n} \sum_{i=1}^n f(x_i)$

Then

$$\begin{aligned} E(I - E)^2 &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n f(x_i)\right) \\ &\stackrel{\text{independence}}{=} \frac{1}{n} \underbrace{\text{Var}(f(X))}_{\text{bounded if } f \text{ is bounded.}} = \end{aligned}$$

So expected error is

$$E|I - E| \leq \frac{1}{\sqrt{n}}$$

Ex: (Machine Learning)

Suppose we see pairs

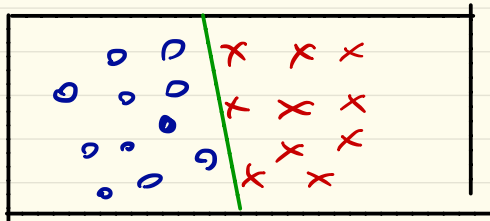
$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \in \mathbb{R}^d \times \{-1, 1\}$$

- x_i are called feature vectors
- y_i are called labels

Eg: x_i encodes properties of individual (incidence of smoking, distance of household to factory, genetics)

$$y_i = \begin{cases} 1 & \text{if cancer} \\ -1 & \text{if no cancer} \end{cases}$$

Goal: Find a predictive relationship between x and y .



Typical Approach:

- Fix a hypothesis class H

Eg: $H = \{h(x) = \langle v, x \rangle : v \in \mathbb{R}^d\}$

$H = \{\text{quadratics in } x\}$

$H = \{\text{neural networks}\}$ observed

- Fix a loss function $l(t, y)$.
prediction

Eg: $l(t, y) = \frac{1}{2} (t - y)^2$ MSE

$l(t, y) = \max\{0, 1 - ty\}$ Hinge

$l(t, y) = \log(1 + e^{-ty})$ Logistic

. Solve Empirical Risk Minimization

$$\star \min_{h \in \mathcal{H}} \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(h(x_i), y_i)}_{Z_h} \quad \uparrow \text{Optimization}$$

Is this any good?

If $(x_i, y_i) \stackrel{iid.}{\sim} \mathcal{P}$ we would like \star to approximate

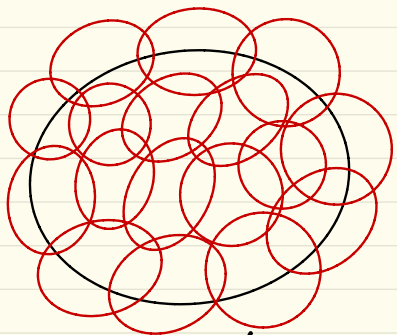
$$\min_{h \in \mathcal{H}} \underbrace{\mathbb{E}_{(x,y) \sim \mathcal{P}} \ell(h(x), y)}_{\mathbb{E} Z_h}$$

So we need to control

$$\sup_{h \in \mathcal{H}} |Z_h - \mathbb{E} Z_h|$$

How to bound the supremum of infinitely many random variables?

One idea: discretize \mathcal{H}



Cover \mathcal{H} with $\bigcup_{i=1}^k B_\epsilon(h_i)$

• Control $|Z_{h_i} - \mathbb{E}Z_{h_i}| \leq \delta$
with probability $1-p$

• Control $\sup_{i=1, \dots, k} |Z_{h_i} - \mathbb{E}Z_i| \leq \delta$
with probability $1-kp$ [union bound]

• Extend by regularity of the loss

$\sup_{h \in \mathcal{H}} |Z_h - \mathbb{E}Z_h| \leq \delta + L(\epsilon) \text{ w.p. } 1-kp$

But: $k \approx \left(\frac{1}{\epsilon}\right)^d$. So p needs to
be much smaller than ϵ^d
 \implies Concentration!