

LECTURE 1

• Big data:

- (a) # observations (data points) → Big
- (b) # dimensions (parameters describing each data pt.) → Big

• Examples:

1. Income of Kyiv population
= 3,000,000 observations, dimension = 1. low dim.

2. Avatars of people on FB

each avatar = 100 × 100 image

Each pixel = dimension ⇒ #dimensions = 10^4 .
high dim.

3. Other HD examples: text; sound; video;
genome; medical history;
chess games.

• Empirical Observation: it is exponentially harder

to deal with large # of dimensions
than with large # of observations.

↳ classical statistics, probability
(via limit thems)

↳ HDS, HDP (new).

WHY exponentially harder?

Example:

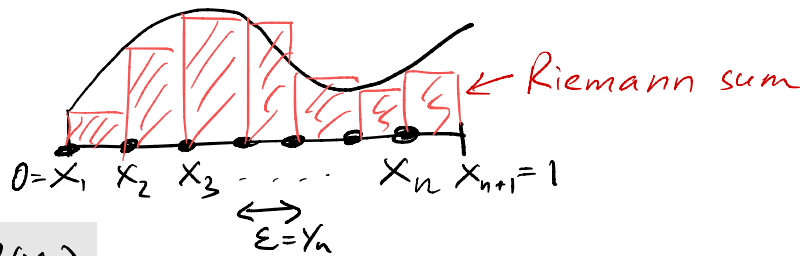
Problem (HD integration) Numerically compute the integral of a given function f

$$\int_0^1 \dots \int_0^1 f(x_1, \dots, x_d) dx_1 \dots dx_d = \int_{[0,1]^d} f(x) dx$$

\leftarrow parameters

(e.g. $f = \text{model of income}$,
 $\int f dx = \text{total income}$)

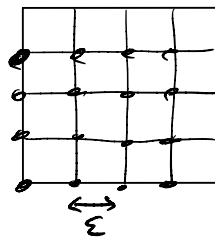
• If $d=1$: use the grid



$$\int_0^1 f(x) dx \approx \sum_{i=1}^n (x_{i+1} - x_i) f(x_i) = \frac{1}{n} \sum_{i=1}^n f(x_i)$$

"Resolution" $\epsilon = 1/n$

• If $d=2$, do similarly but use the grid



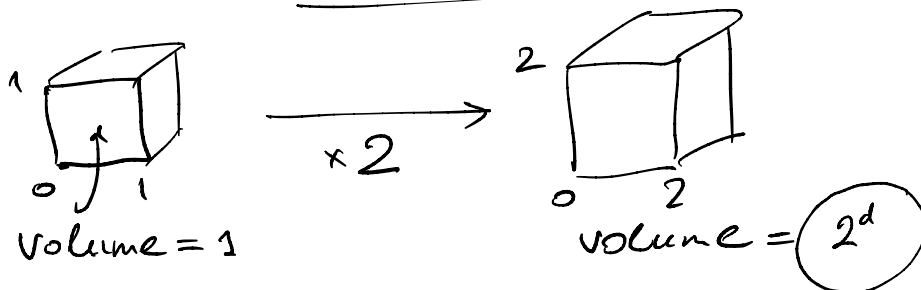
$$\Rightarrow n = \left(\frac{1}{\epsilon}\right)^2 \text{ pts.}$$

• For general dimension d , we need $n = \left(\frac{1}{\epsilon}\right)^d \text{ pts.}$

Exponential in d . Too large.

Complexity of many alg's is exponential in dimension

• Why? There is too much room in H.D.'s:



"THE CURSE OF DIMENSIONALITY"

Probability for rescue : Monte-Carlo method

- Instead of choosing x_i on the grid, choose them at random (independently, uniformly in $[0,1]^d$)
 $\Rightarrow f(x_i)$ are i.i.d. r.v.'s.

$$\frac{1}{d} \sum_{i=1}^d f(x_i) \stackrel{?}{\approx} \int_{[0,1]^d} f(x) dx$$

- Will use the following standard facts of probability theory:

① Def A r.v. X has density (= pdf = probability density function) $p(x)$ if
$$P\{X \in A\} = \int p(x) dx \quad \forall \text{ Borel } A \subset \mathbb{R}$$

In this case we say that X has a continuous distribution.

• $p(x) \geq 0 \quad \forall x$; $\int_{-\infty}^{\infty} p(x) dx = P\{X \in (-\infty, \infty)\} = 1$.

• Similarly for a random vector taking values in \mathbb{R}^n . Density $p: \mathbb{R}^n \rightarrow \mathbb{R}$

Example: uniform distribution. $X \sim \text{Unif}([0,1])$ if $p(x) = \begin{cases} 1, & x \in [0,1] \\ 0, & \text{otherwise} \end{cases}$
 $X \sim \text{Unif}([0,1]^d)$ if $p(x) = \begin{cases} 1, & x \in [0,1]^d \\ 0, & \text{otherwise} \end{cases}$

② Def The expected value (expectation) of a r.v. X with density $p(x)$ is

$$E[X] = \int_{-\infty}^{\infty} x \cdot p(x) dx$$

More generally, \forall function $f: \mathbb{R} \rightarrow \mathbb{R}$,

$$E[f(X)] = \int_{-\infty}^{\infty} f(x) \cdot p(x) dx$$

Similarly for a random vector X : $E[f(X)] = \int_{\mathbb{R}^n} f(x) dx$.

② Variance of a r.v. X :

$$\text{Var}(X) = \mathbb{E}(X - \mathbb{E}X)^2 = \mathbb{E}(X^2) - (\mathbb{E}X)^2$$

③ Linearity: (a) $\mathbb{E}(X_1 + \dots + X_n) = \mathbb{E}X_1 + \dots + \mathbb{E}X_n$

(b) If X_i are independent,

$$\text{Var}(X_1 + \dots + X_n) = \text{Var}(X_1) + \dots + \text{Var}(X_n)$$

④ (Strong) Law of Large Numbers (SLLN):

If X_1, X_2, X_3, \dots are independent and identically distributed (iid) random variables, then

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mathbb{E}X \text{ almost surely (a.s.)}$$

↑
i.e. with probability = 1.

Back to our situation: $X_i \sim \text{Unif}([0,1]^d)$ independent and identically distributed (iid)

$$\bullet \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n f(x_i) \right] = \frac{1}{n} \sum_{i=1}^n \underbrace{\mathbb{E} f(x_i)}_{\mathbb{E} f(x) \text{ by identical distribution}} = \mathbb{E} f(x)$$

$$= \int_{\mathbb{R}^d} f(x) p(x) dx, \text{ where density is } p(x) = \begin{cases} 1, & x \in [0,1]^d \\ 0, & \text{---} \end{cases}$$

$$= \int_{[0,1]^d} f(x) dx. \quad \text{😊}$$

⇒ we have an unbiased estimator of the integral

$$\bullet \text{SLLN} \Rightarrow \boxed{\frac{1}{n} \sum_{i=1}^n f(x_i) \rightarrow \int_{[0,1]^d} f(x) dx \text{ a.s.}}$$

• Rate of convergence? L^2 error ("MSE"):

$$\mathbb{E} \left(\underbrace{\frac{1}{n} \sum_{i=1}^n f(x_i)}_{\bar{Z}} - \underbrace{\int_{[0,1]^d} f(x) dx}_{\mathbb{E} Z} \right)^2 = \text{Var} \left(\frac{1}{n} \sum_{i=1}^n f(x_i) \right)$$

$$= \frac{1}{n^2} \sum_{i=1}^n \underbrace{\text{Var}(f(x_i))}_{\text{Var}(f(x)) \text{ by identical distribution}} = \frac{\text{Var}(f(x))}{n}$$

$$\leq \frac{1}{n} \text{ e.g. if } |f(x)| \leq 1 \quad \forall x.$$

• Taking square root ⇒ expected error = $\boxed{O(1/\sqrt{n})}$
regardless of dimension !! 😊