# LECTURE 1

- Big data:
    (a)  # <u>observations</u> (data points)  $\longrightarrow$ Big

    (b)  # <u>dimensions</u> (parameters)  $\longrightarrow$ Big

- <u>Examples</u>:
    1. Income of Kyiv population
        = 3,000,000 observations,  dimension = 1.

    2. Images of people on FB
        each observation = $100 \times 100$ image
        Each pixel = dimension  $\Rightarrow$  # dimensions = $10^4$
        HD.

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

- <u>Empirical Observation</u>:  it is exponentially harder
    to deal with large # of dimensions
    than with lage # of observations.

        $\hookrightarrow$ classical statistics, probability
            (via limit thms)
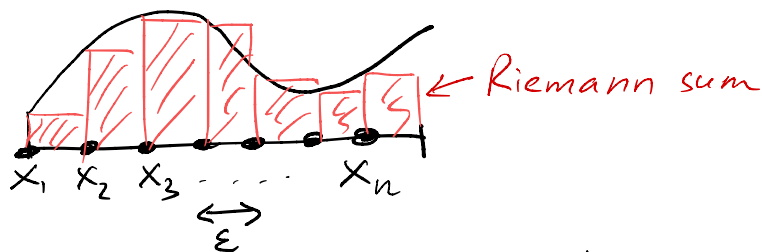
        $\hookrightarrow$ HDS, HDP (new).

WHY exponentially harder?          Example:

Problem (HD integration) | Numerically compute the integral

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

$\uparrow$ parameters

of a given function $f$          ( e.g. $f$ = model of income, $\int f \, dx$ = population income )

 ← Riemann sum

- If $d = 1$ : use the grid

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

$x_1, x_2, x_3 \cdots \cdots x_n$

resolution $= \varepsilon \implies n = 1/\varepsilon$ pts.

- If $d = 2$ : same $\nearrow$          $\implies n = \left(1/\varepsilon\right)^2$ pts.

$\overset{\varepsilon}{\longleftrightarrow}$

- For general $d$,          $N = \left(1/\varepsilon\right)^d$ pts.

  $\underline{\text{Exponential}}$ in $d \implies$ too large.

$\implies$ complexity of many alg's is $\underline{\text{exponential}}$ in $\underline{\text{dimension}}$

- There is too much room in H.D's :

volume $= 1$      $\xrightarrow{\times 2}$      volume $= \boxed{2^d}$ !

"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.   $\frac{1}{d} \sum\limits_{i=1}^{d} f(X_i) \overset{?}{\approx} \int\limits_{[0,1]^d} f(x)\, dx$

- Will use the following standard facts of probability theory:

① Expectation of a r.v. $X$ with density $p(x)$ :

$$\mathbb{E}X = \int\limits_{-\infty}^{\infty} x \cdot p(x)\, dx \qquad (def)$$

More generally,
$$\mathbb{E}f(x) = \int\limits_{-\infty}^{\infty} f(x)\, p(x)\, dx$$

More generally, if $X$ is a random vector in $\mathbb{R}^n$,
$$\mathbb{E}f(x) = \int\limits_{\mathbb{R}^n} f(x)\, p(x)\, dx$$

② Variance of a r.v. $X$ :
$$Var(X) = \mathbb{E}\left(X - \mathbb{E}X\right)^2$$

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

(b) If $X_i$ are independent,
$$Var(X_1 + \cdots + X_n) = Var(X_1) + \cdots + Var(X_n)$$

④ (Strong) Law of large numbers (SLLN):

If $X, X_1, X_2, \ldots$ are independent and identically distributed (iid) random variables, then
$$\frac{1}{N} \sum\limits_{i=1}^{N} X_i \longrightarrow \mathbb{E}X \quad \text{almost surely (a.s)}$$

$\uparrow$
i.e. with probability $= 1$.

In our situation,

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

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

(uniform distribution)

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

$\Rightarrow$ we have an <u>unbiased estimator</u> of the integral

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

- Rate of convergence? $L^2$ error:

$$\mathbb{E}\left(\underbrace{\frac{1}{n}\sum_{i=1}^{n} f(x_i)}_{\overset{\shortparallel}{Z} \text{ (red)}} - \underbrace{\int_0^1 f(x)\, dx}_{\mathbb{E}Z \text{ (red)}}\right)^2 = Var\left(\frac{1}{n}\sum_{i=1}^{n} f(x_i)\right)$$

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

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

- Taking square root $\Rightarrow$ expected error $= \boxed{O\left(\frac{1}{\sqrt{n}}\right)}$

<u>regardless of dimension</u> !! $\ddot\smile$