

Lecture 1:

5 modules:

1) Estimating parameters/distributions of data
(5 - 6 lecs.)

2) Hypothesis testing (5 - 6 lecs)

Midterm

3) Supervised learning (4 - 5 lecs)

4) Unsupervised learning (4 lecs)

5) Introduction to statistical learning theory.
(4 lecs)

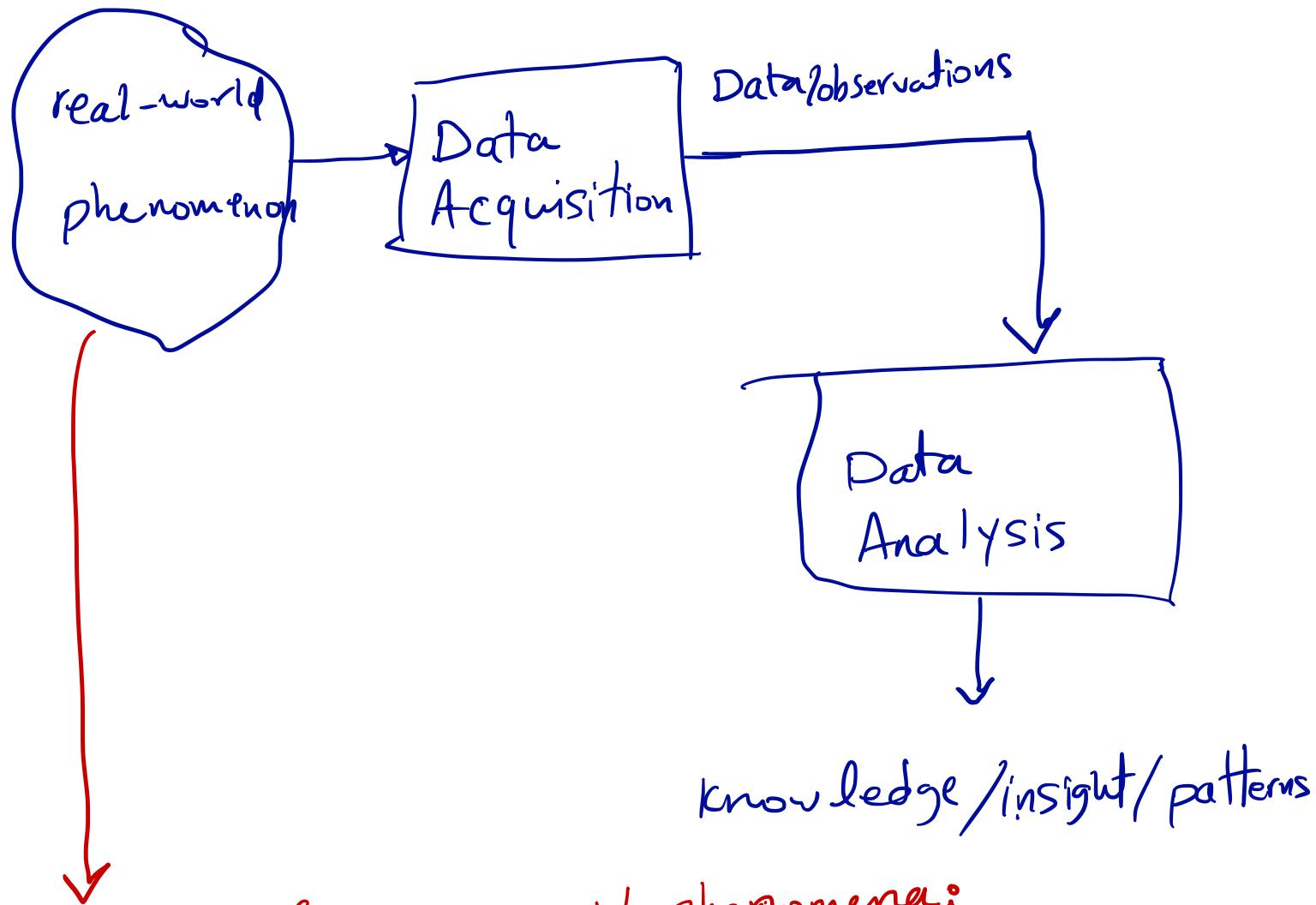
Grading :

Exercises : 40% \rightarrow 7-8 Hws

Exams (Midterm, Final) : 60%.

Text books \rightarrow See Canvas.

Lecture 2 :



Examples of real-world phenomena:

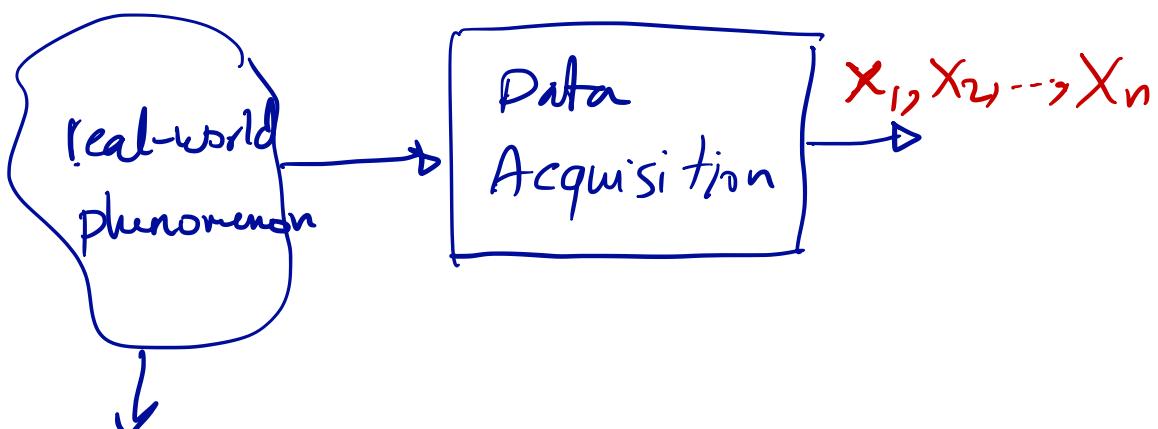
- finding the average height of women in the US
- Percentage of people who'd vote for a certain candidate in the state of Pennsylvania

- likelihood of a person to develop breast cancer within the next 5 years
- average speed of molecules in a specific gas
- likelihood of a user to click on a certain ad in her face book page

Data Science involves developing methods of acquisition/storage/analysis of data to effectively extract useful knowledge/information/insight /patterns.

Estimating Basic Statistics of Data:

Consider the task of estimation
the average height of women
across the US.



$\{y_1, \dots, y_N\}$
 $N \approx$ population size

$$\text{task} = \text{find } \mu = \frac{y_1 + y_2 + \dots + y_N}{N}$$

Typically, N is a very large number (e.g. $N = \text{tens of millions}$)

and the naive method that goes over the whole population is often inefficient or infeasible. So the basic idea of data acquisition is to sample from the real-world phenomenon.

Random Sampling: We will choose

uniformly at random a subset of the population.

Let X_i be the data corresponding to an individual chosen uniformly at random in the population.

$$X_i = \left\{ \begin{array}{l} y_1 \\ y_2 \\ \vdots \\ y_N \end{array} \right\}$$

w.p. $\frac{1}{N}$

w.p. $\frac{1}{N}$

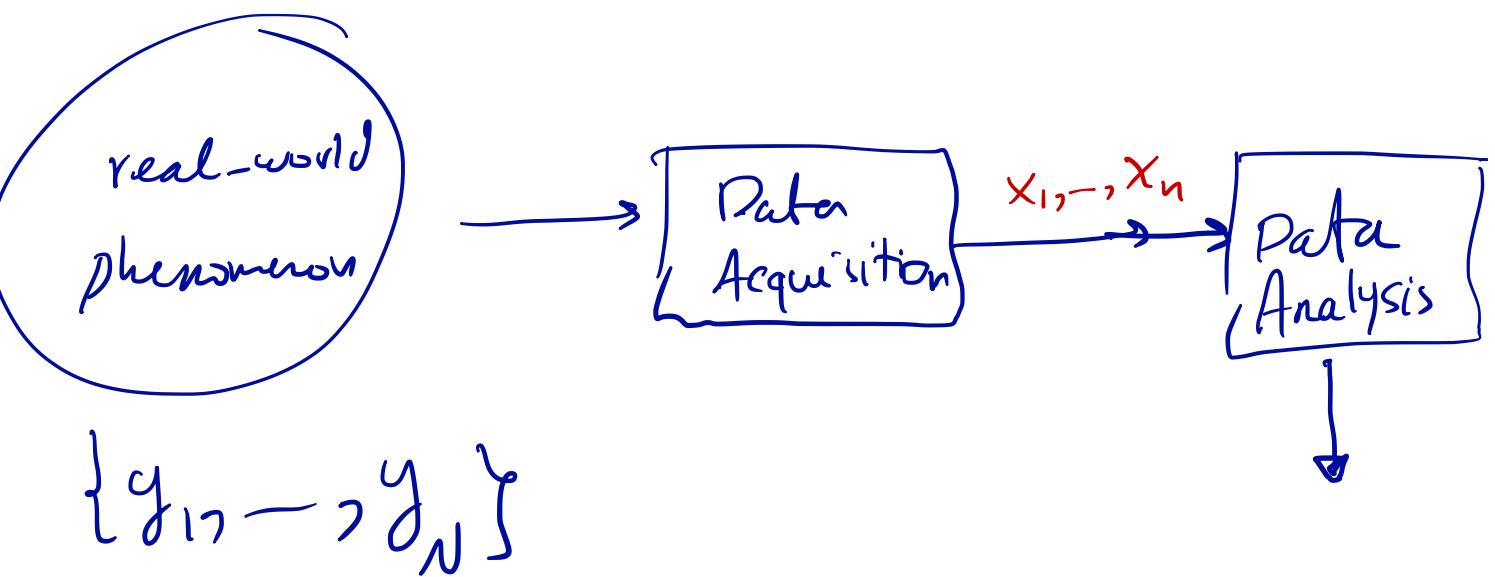
w.p. $\frac{1}{N}$

furthermore X_i 's are independently chosen.

Data: X_1, X_2, \dots, X_n

Typically n is much smaller

than N .

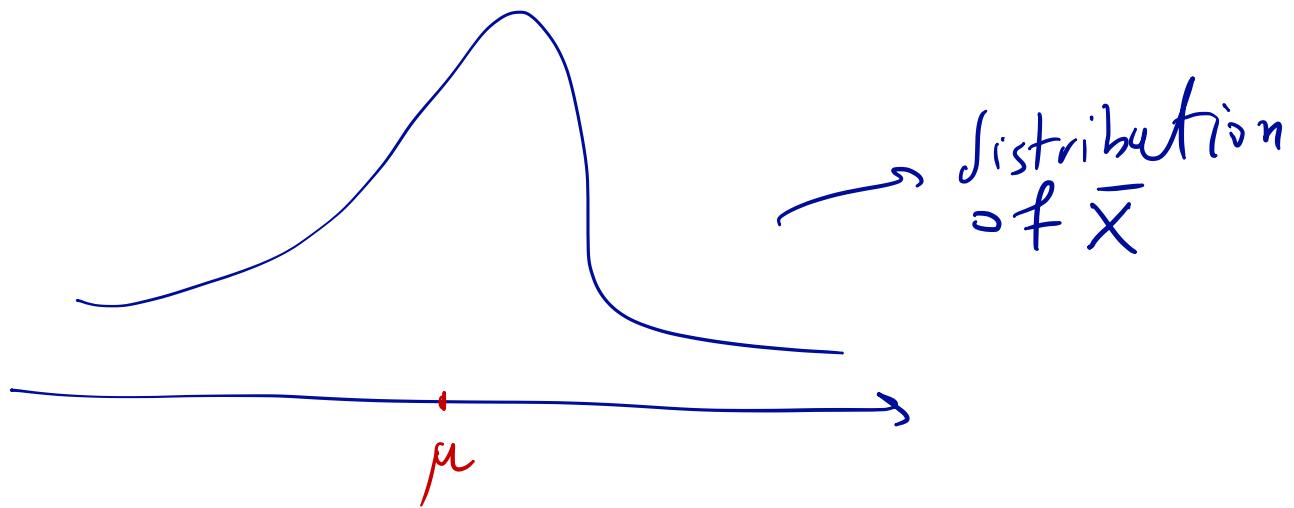


question : Estimate the average
of y_i 's (i.e. μ)
by only using the
data x_1, \dots, x_n .

Simplest answer:

$$\bar{x} = \frac{x_1 + \dots + x_n}{n} \quad (\text{sample mean})$$

Note that \bar{X} is a random object/variable.



Let's first compute the expected value of \bar{X} and how it is related to μ .

$$\begin{aligned} E[\bar{X}] &= E\left[\frac{x_1 + x_2 + \dots + x_n}{n}\right] \\ &= \frac{E[x_1] + \dots + E[x_n]}{n} \end{aligned}$$

Since X_i 's have exactly the same distribution, then

$$E[X_1] = \dots = E[X_n]$$

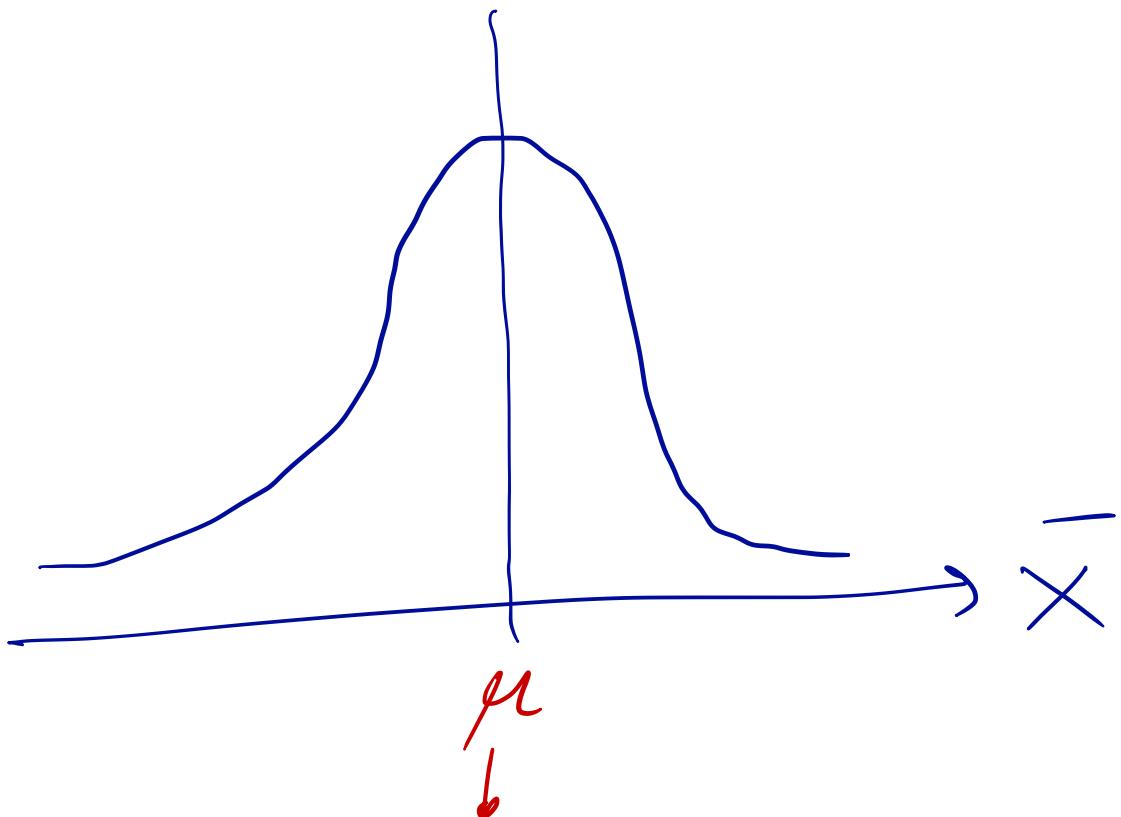
$$\bar{X}_i = \begin{cases} y_1 & \frac{1}{N} \\ y_2 & \frac{1}{N} \\ \vdots & \vdots \\ y_N & \frac{1}{N} \end{cases}$$

$$\Rightarrow E[\bar{X}] = \frac{E[X_1] + E[X_2] + \dots + E[X_n]}{n}$$

$$= E[X_1]$$

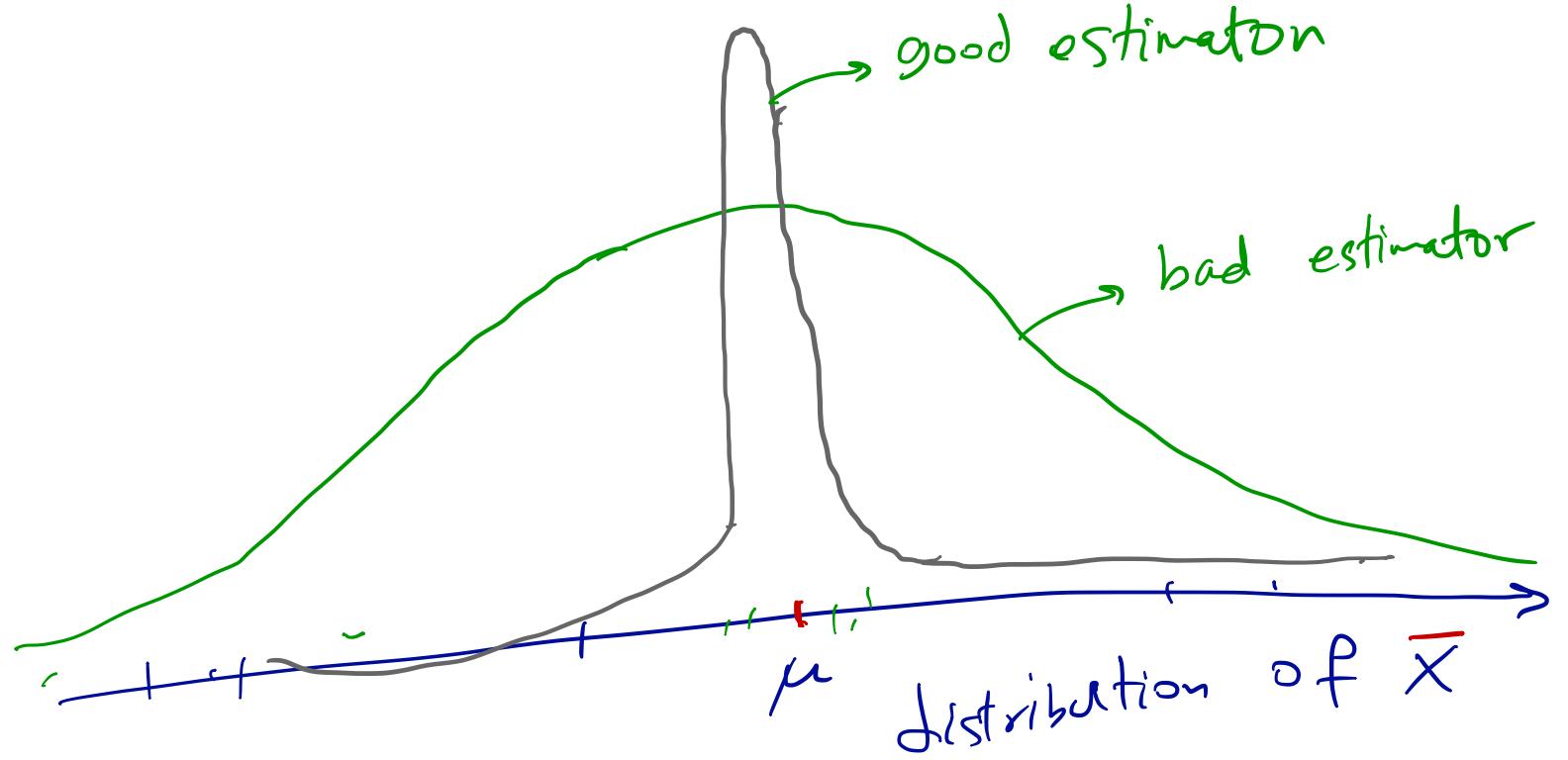
$$= \frac{1}{N} \cdot y_1 + \frac{1}{N} y_2 + \dots + \frac{1}{N} y_N$$

$$= \frac{y_1 + \dots + y_N}{N} = \mu$$



$$E[\bar{X}] = \mu$$

So, \bar{X} is on average a very good estimator because its expectation is μ . The next question is how concentrated is \bar{X} around its average μ .



In order to see how concentrated \bar{X} is around μ , we need to compute ...

$$E[(\bar{X} - \mu)^2].$$