

Lecture 24

Data: Collection of numbers obtained as an observation of specific system

Considers • "free" data

y_1, y_2, \dots, y_n n numbers

Here, all you have is n observations

• "parameterized" data

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ n pair of numbers

or $(\underline{x}_1, \underline{y}_1), (\underline{x}_2, \underline{y}_2), \dots, (\underline{x}_n, \underline{y}_n)$ n pair of vectors

where $\underline{x}_i = (x_i^1, x_i^2, \dots, x_i^k)$ k element vector

$\underline{y}_i = (y_i^1, y_i^2, \dots, y_i^m)$ m element vector

Here, data y_i (or data vector \underline{y}_i) is obtained

under parameters x_i (or parameter vector \underline{x}_i)

I.e. each data have associated parameter

Example: (1.) flip coins and record head (0) or tail (1)

is (y_1, y_2, \dots, y_n)

where y_i is either 0 or 1

(2.) Consider m different coins manufactured (so each coin will be slightly different from another)

then
record

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

where

x_i is either coin 1, coin 2, ..., or coin m

y_i is either 0 or 1.

(3.) Measure drag coefficient by throwing an object of fixed shape & size n times and measuring C_d

$C_{d1}, C_{d2}, \dots, C_{dn}$ n observed drag coefficients

(4.) Measure drag coefficient by throwing an object of fixed shape & size n times

$$(m_1, C_{d1}), (m_2, C_{d2}), \dots, (m_n, C_{dn})$$

where

m_i = mass of an object in i^{th} experiment

Set containing data values

A set of numbers from which each observation value is drawn.

Discrete set: for coin flipping, the set is $\{0, 1\}$

Continuous/Continuum/Real set: for drag coefficient,

set is $\{x \geq 0\}$ of any positive number

Think of more examples

• Statistics of the data

Consider y_1, y_2, \dots, y_n n data

• Mean (arithmetic mean)

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$m_1 = m_2 = \dots = m_n = 1$$
$$\bar{C}_d = \frac{\sum_{i=1}^n C_{d_i}}{n}$$

for drug coefficients
 C_{d_1}, \dots, C_{d_n}

$$\bar{C}_d = \frac{1}{n} \sum_{i=1}^n C_{d_i}$$

but

$$\bar{C}_d = \frac{\sum_{i=1}^n m_i C_{d_i}}{\sum_{i=1}^n m_i}$$

↓
this is weighted mean

• Median (50th percentile of data)

arrange in increasing order

$$a_1 < a_2 < a_3 \dots < a_n$$

then if n is odd a_i where $i = \frac{n+1}{2}$
is median

if n is even $\frac{a_i + a_{i+1}}{2}$ where $i = \frac{n}{2}$

• **Mode** value in data that appears most frequently

Spread of data

while mean, mode, median etc. inform about the "center" or "key" value of data, we also want to know how large the values in data can vary from each other.

Example: two examples have same mean (zero mean)

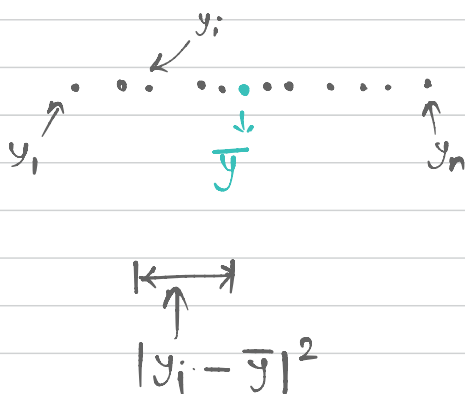
(i) 0, -0.1, 0.1, 0.2, -0.25, 0.35, 0.3

(ii) 0, -1, 1, 4, -6, 8, -6

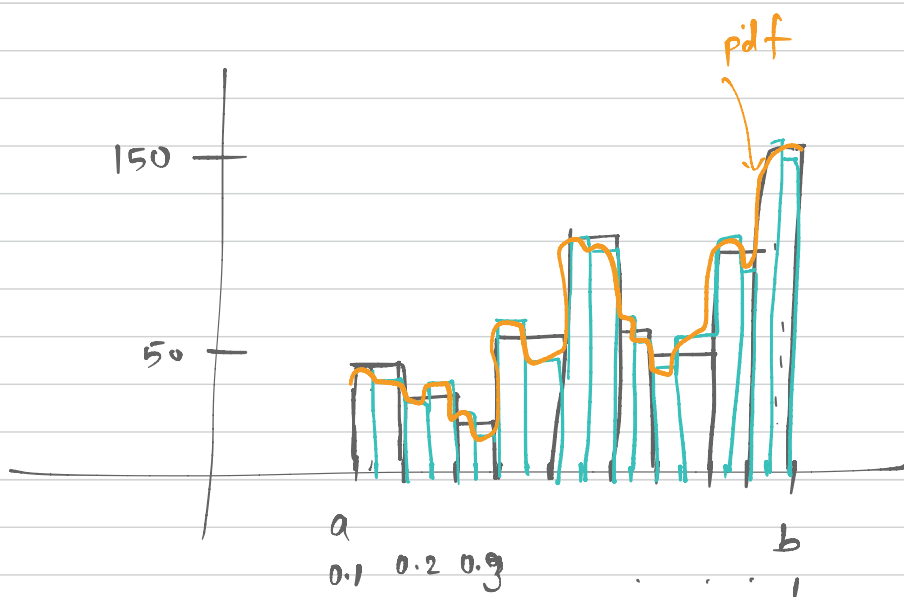
standard deviation (std)

$$s_y = \sigma = \sqrt{\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}}$$

Variance : $\sigma^2 = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}$



y_1, y_2, \dots, y_n



probability distribution function (pdf)

$$f: [a, b] \rightarrow [0, \infty)$$

$$x \in [a, b], f(x)$$

$$\int_a^b f(x) dx = 1$$

$$f: [a, b] \rightarrow [0, \infty)$$

$f(x) \approx$ how many times x will be observed in 1000 experiments

$$\bar{f}(x) = \frac{1}{\int_a^b f dx} f(x) \rightarrow$$

$c_{d_1}, c_{d_2}, \dots, c_{d_n}$

$[0.1, 1] \rightarrow [0.1, 0.2) \rightarrow \text{set 1}$

$[0.2, 0.3) \rightarrow \text{set 2}$

$[0.3, 0.4)$

\vdots

$[0.9, 1) \rightarrow \text{set 9}$

for each c_{d_i} , find $[a_i, b_i)$ s.t.

$$c_{d_i} \in [a_i, b_i)$$

1000 observations

set 1 appeared 50 times

set 9 appeared 150 times