

17/02/2021



MA2040: Probability, Stochastic Process & Statistics.

* INSTRUCTORS: Prof. S.H. Kulkarni
shk@iitpkd.ac.in
Prof. Ashok Kumar
ashokm@iitpkd.ac.in

* TAs : Abhinav Thakur
142002002@smail.iitpkd.ac.in
Shree Ganesha Sharma M S
142002015@smail.iitpkd.ac.in

* SLOTS : G1 Wed 1130-1230 hrs
T4 Thu 1400-1530 hrs

• Evaluation: 1 assignment (10 Marks)
4 Quizzes (60 Marks)
1 endsem exam (30 Marks)

⚠ Quizzes ~ 8 Mar, 26 Mar, 19 Apr, 7 May (0900-0945 hrs)

* Textbook: Introduction to Probability &
→ Statistics for Engineers & Scientists

- Sheldon. M. Ross
5th ed

• Probability, Random Variables & Stochastic Processes

- Athanasios Papoulis & S. Unnikrishna Pillai

⚠ Take own notes.

⚠ Discuss actively in the forum.

Freely ask doubts, keep your shyness
(& ego) aside:)

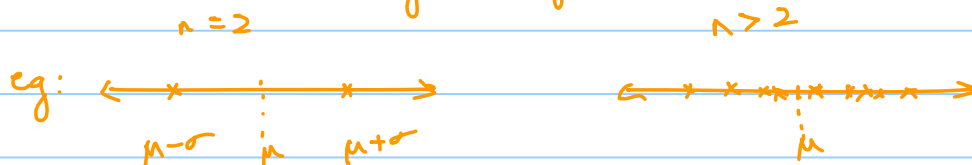
• Basic steps in Statistics

- a) Data Collection (& cleaning)
Outlier removal, standardization, ...
- b) Representing the data
Visualization - plot, histogram, ...
- c) Summarize the data
Central tendency, dispersion,
* Skewness, * Kurtosis
- d) Conclusion

△ Given mean, $\mu = 5.4$, } for a dataset X,
std deviation, $\sigma = 1.2$ }

→ what can you say about the distribution of the data (distrib)

→ Significance of sample size n of X in making inference about X



17/02/2021



18/02/2021

LECTURE 2

a) DATA COLLECTION

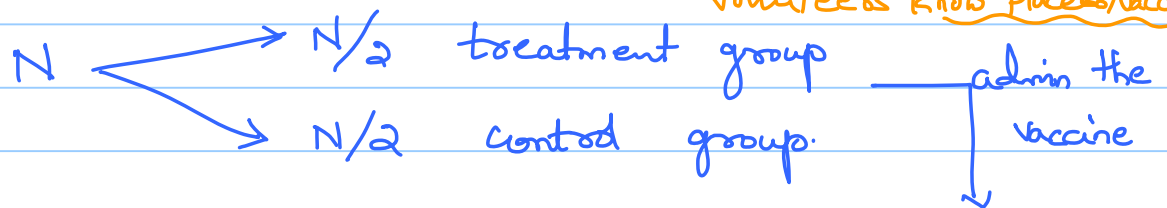
- eg: Clinical testing of efficacy of a vaccine.

Design of }
experiments }

N individuals/volunteers.
Treatment group (expt. treatment) Control group (Placebo)

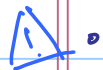
* Some objectives.

Neither doctors, nor volunteers know placebo/vaccine.



Expectation:

- More people in control group are infected.



		CONFUSION MATRIX	
Disease	Test	present	absent
		detected TP	not detected FP
		not detected FN	TN

https://en.wikipedia.org/wiki/False_positives_and_false_negatives



What's a good N? 100 / 1000 / 10000 ...
to make good decision.
Will find out in due course.

• Minimize N

F2 score etc.

• Still have good testing accuracy, efficiency, & confidence in inference.

b) REPRESENTING / DESCRIBING THE DATA

• MATLAB (optional, useful for assignment).

++ Can use MATLAB online using institute mail-id
@
matlab.mathworks.com.

- Setup: 30 students choose an integer b/w 1 & 10. (frequently used). Denoted by X .

$$X = [1, 2, 2, \dots, 10]$$

a) Frequency table

# written, x_i	1	2	3	4	5	6	7	8	9	10	Total
# of students, f_i	1	3	3	2	1	4	9	3	3	1	30
frequency.											(N)

discrete data, $x_i \in \{1, 2, \dots, 10\}$

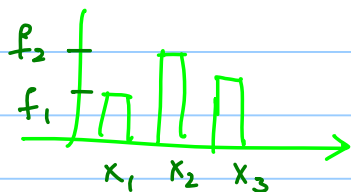
Plot

<https://en.wikipedia.org/wiki/Histogram>



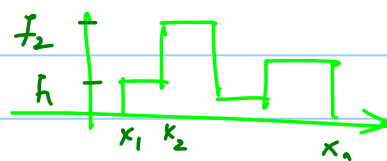
b) Bar chart

- for discrete observations



c) Histogram.

- Usually for grouped/continuous intervals.



- Normalized / relative frequency = proportion of students

$$f_i(\text{rel}) = f_i / \text{length}(X)$$

$$f_i(\text{rel}) \in [0, 1] \quad \forall i$$

- !! • area under this curve is 1 \rightsquigarrow Similar to a probability mass / density function.

- relative frequency is analogous to probability.

d) • pie-chart is a helpful tool / appropriate for this.

$\left. \begin{array}{l} \text{area of sector} \\ \text{angle of sector} \end{array} \right\} \propto \text{frequency of observation.}$

c) * Stem and leaf plot.

↳ Helps keep record of the data in a concise manner.

Say all observations are in lakhs or in fractions.
eg1. eg2.

phone #s, pincodes, ...
9886784591, 9901425926, ...
9886907843, 9482569380, ...

0.0011, 0.0025, 0.0019,
0.0042, 0.0069, ...

airtel
vodafone
BSNL

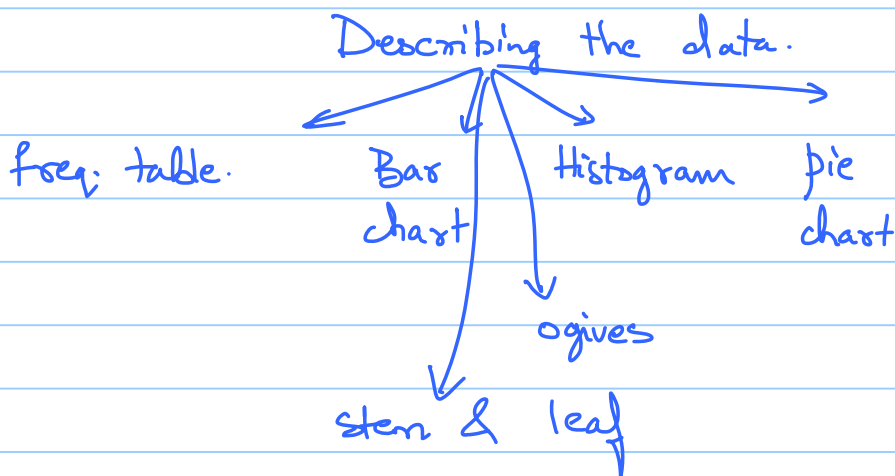
Stem	Leaf
← 990	1425926
← 988	6784591, 6907843
← 948	2569380

stem	leaf
0.001	1, 9
0.002	5
0.004	2
0.006	9

← first 3 digits
of phone #

→ rest of the
observation.

++ Stem & leaf plot helps us group/form clusters based on the stem.



& many more.

18/02/2021
↑

Textbook : Chapter 1
Chapter 2 till section 2.2 page 14.

24/02/2021

LECTURE 3

Summarizing the data.

- ++ Inference @ a glance.
- ++ Helps compare datasets.

Some measures - mean, median, mode, sd, variance, ...

a) Measures of Central tendency.

- ↳ The value (need not be an observation in our dataset)
↳ around which the dataset is spread.
- ↳ One value that "represents" the data.

given $\{X_i\}_{i \in I_n}$ // $I_n := \{1, 2, 3, \dots, n\}$

Sample Mean.

$$\bar{X} = \frac{\sum_i X_i}{n}$$

• arithmetic mean

• ++ mathematically tractable

• ++ simple to use, intuitive

• -- Susceptible to changes with extreme values.

• for frequency table, $\bar{X} = \frac{\sum X_i}{n} = \frac{X_1 f_1 + X_2 f_2 + \dots}{f_1 + f_2 + \dots}$

$\{(X_i : f_i)\}_{i=1}^l$
(X_i appears f_i times)

$$\bar{X} = \frac{\sum_{i=1}^l f_i X_i}{\sum_{i=1}^l f_i}$$

There are l distinct observations

* $\bar{X} = \sum_{i=1}^l X_i \left(\frac{f_i}{\sum f_i} \right)$ → relative frequency as weight.
→ \equiv EMPIRICAL probability/proportion of

times X_i appears in the data set.

* as $n \rightarrow \infty$ $f_i/n \rightarrow$ (theoretical) probability

Other means:
Averages

a) Geometric mean,

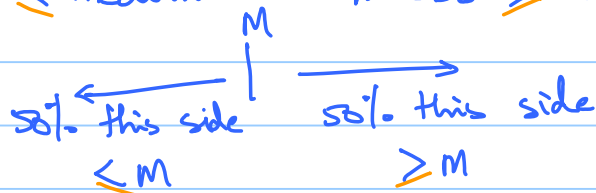
b) Harmonic mean,

c) \vdots

$$\bar{X} := \left(\prod_{i=1}^l X_i^{f_i} \right)^{1/n}$$
$$\bar{X} := \left(\frac{\sum X_i^2}{n} \right)^{-1}$$

Sample
median

- A locational / positional measure.
- divides the dataset into 2 halves of equal # of observations.
- $\# \text{ obs} \leq \text{median} = \# \text{ obs} \geq \text{median}$.



→ Not unique ++ less sensitive to outliers. ROBUST

Procedures to
find Median

① • sort & find the "middle" element.

↳ given X_1, X_2, \dots, X_n arrange as

$X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ () denotes indices may

get permuted & we're using order statistics.

↳ if n is odd:

$$\text{median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ observation}$$

else if n is even:

$$\text{median} = \left[\frac{\left(\frac{n}{2} \right)^{\text{th}} \text{ obs} + \left(\frac{n}{2} + 1 \right)^{\text{th}} \text{ obs}}{2} \right]$$

→ need not be in $\{X_i\}$.

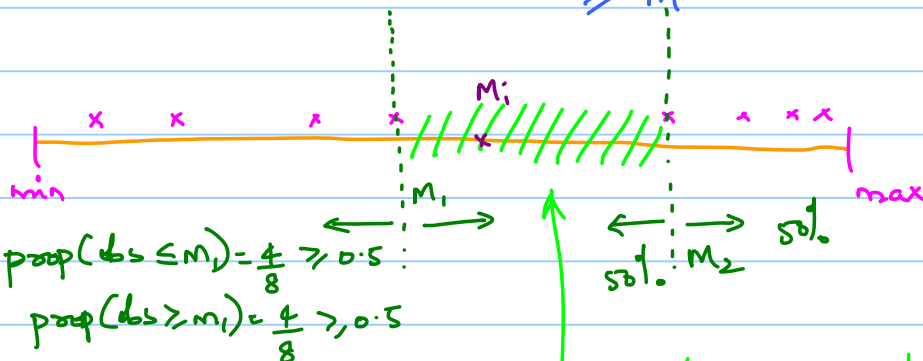
→ Need not be unique.

Definition

M is 'a' median of $\{X_i\}$ if

a) proportion of observations $\leq M$ is atleast 50%

b) $\geq M$



any M here can be a median.

△ Image proc. salt & pepper noise.

△ Median income, salary than mean.

LECTURE 4

Percentiles

→ Another positional measure.

let $0 \leq p \leq 1$

then $100p^{\text{th}}$ percentile is a value m_p , such that
 $\# \frac{\text{obs} \leq m_p}{n} \geq p$ &

$\# \frac{\text{obs} \geq m_p}{n} \geq 1-p$

- $m_{1/2}$ (m @ $p=0.5$) is the median.
- $m_{1/4}, m_{3/4}, m_{5/4}$ are called Quartiles. I, II & III quartile or Q_1, Q_2, Q_3 respectively.
- they are equivalent to 25th, 50th & 75th percentiles resp.

Mode

- observation ($\in \{X_i\}$) with the largest frequency.

Empirical relation

- $3\text{Median} \approx 2\text{Mean} + 1\text{Mode}$ (practically, for many datasets)

• Property of mean,
 given $\{X_i\}_{i \in I_n}, \bar{X}$
 $X_i \mapsto aX_i + b = Y_i$

\Downarrow

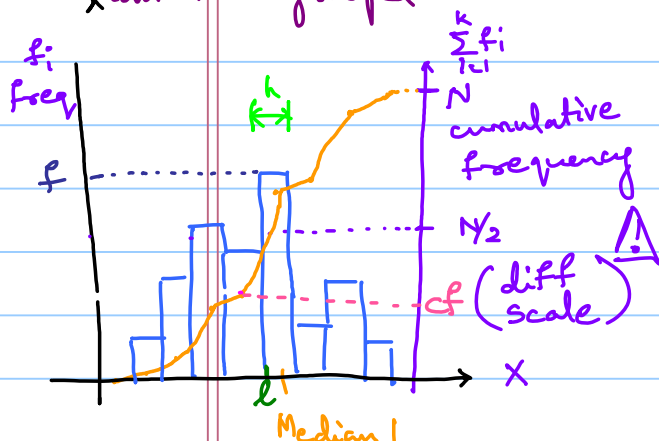
* $\bar{X} \mapsto a\bar{X} + b = \bar{Y}$

• change in origin & scale affects the mean.

• linear if $b=0$

Median for grouped data-

Interpolated median



center of median class

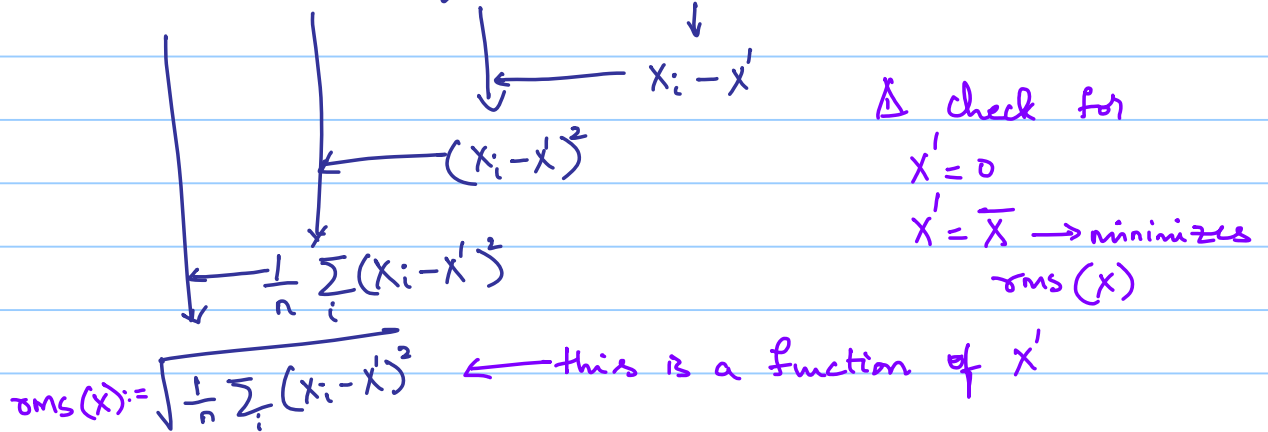
$$M := l + \frac{h}{f} \left(\frac{N}{2} - cf \right)$$

Annotations for the formula:

- l : lower limit of median class
- h : width of median class
- f : freq of median class
- $\frac{N}{2}$: cum freq, previous to median class
- cf : cum freq, previous to median class

b) Measures of deviation / spread / variability (of the data from a point X' . If $\frac{X'}{\bar{X}} = \bar{X}$, we call it central moments.)

4) Root Mean Square deviation.

$$\rightarrow \text{root}(\text{mean}(\text{square}(\text{deviation})))$$


Doesn't have the units of X

ii) Variance.

$$v(x) := \frac{1}{n-1} \sum_i (x_i - \bar{x})^2 =: \text{Var}(X), \sigma_x^2 \text{ etc}$$

(iii) Standard deviation := positive square root of $V(X)$

$$\sigma_x = \frac{|\sqrt{\text{Var}(X)}|}{\sqrt{\frac{1}{n-1} \sum (x_i - \bar{x})^2}}$$

(iv) Sample Standard deviation.

$$S_x := \sqrt{\frac{1}{n-1} \sum_i (x_i - \bar{x})^2}$$

$$= \sqrt{\frac{n}{n-1}} \cdot \sigma_x$$

$$\lim_{n \rightarrow \infty} S_x = \sigma_x$$

But why $n-1$?

! Bias of an estimator

Unbiased estimator

(v) Mean Absolute deviation. from 0 / Mean / arb. point (x')

$$\text{MAD}(x) = \frac{1}{n} \sum_i |x_i - x'| \quad // \text{L}_1 \text{ norm}$$

↑ Median minimizes MAD(x') !!

++ Not differentiable, mathematically less treatable.

- Requires no multiplication, less complex.

LECTURE 5

* Range := $X_{\max} - X_{\min}$

→ sample range \leq population range.

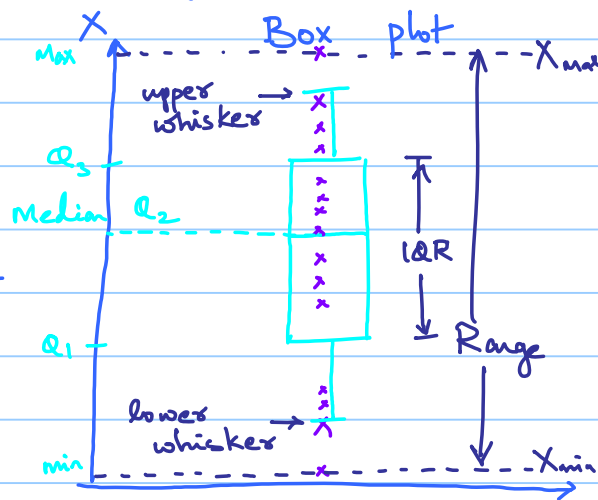
* Interquartile range (IQR)

$$IQR := Q_3 - Q_1$$

Outliers :- elements X_i such that

$$X_i > Q_3 + 1.5 IQR$$

$$X_i < Q_1 - 1.5 IQR$$



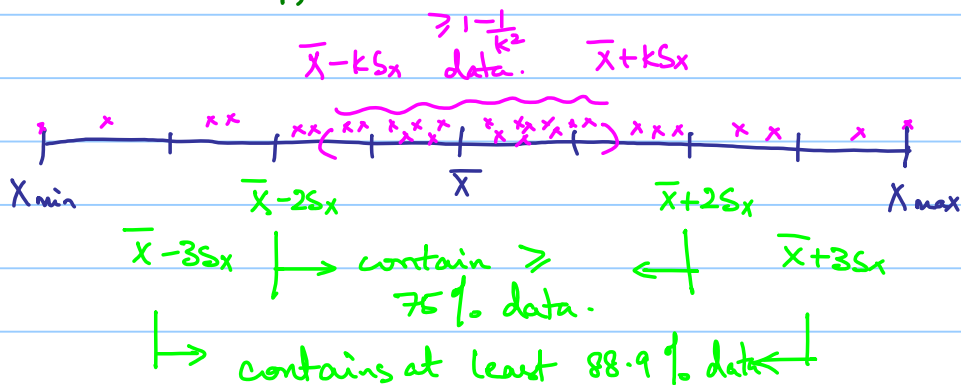
Chebyshev's inequality (in statistics)

Statement:

At least $100(1 - \frac{1}{k^2})\%$ of data points lie within k std. deviations from the mean. ($k > 1$)

$$\frac{\# \{X_i \mid |X_i - \bar{X}| \leq k S_x\}}{n} \geq 1 - \frac{1}{k^2}$$

$$\text{or } \frac{\# \{X_i \mid |X_i - \bar{X}| \geq k S_x\}}{n} \leq \frac{1}{k^2}$$



++ Empirical relationship, true for all kinds of data.

$\{S_x \neq 0 \text{ is not a constant sequence}\}$

⚠ → Caveats of Chebyshev's inequality in terms of probability, continuous Random Variables, integrability,

-- Chebyshev's bound is loose, better bounds are there.
(Chernoff, Markov, ...)

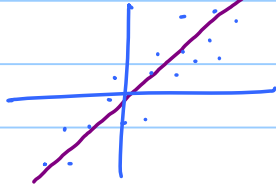
• Multivariate dataset.

eg in 2D $\{(X_i, Y_i)\}_{i \in I_n}$ ordered pairs. n points (X_i, Y_i) in 2D

CORRELATION
measure of LINEAR relationship

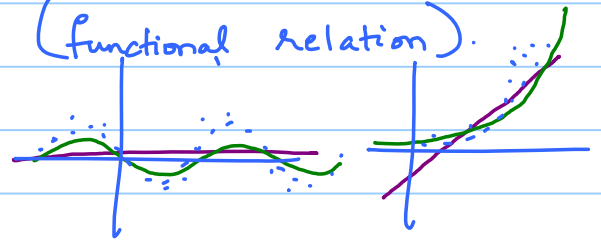
$$Y = ax + b$$

for some a & b



subset of

all kinds of relations b/w 2 variables.
eg: $Y = e^X$, $Y = \log X$,
 $X \equiv \text{time}$, $Y \equiv \text{acceleration}$,
(functional relation).

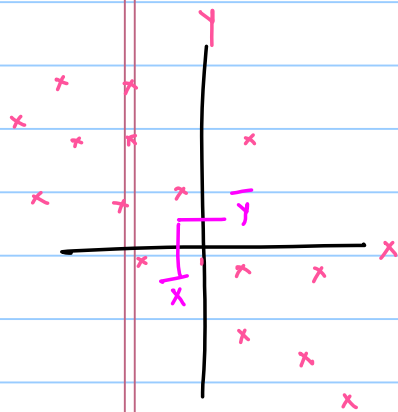
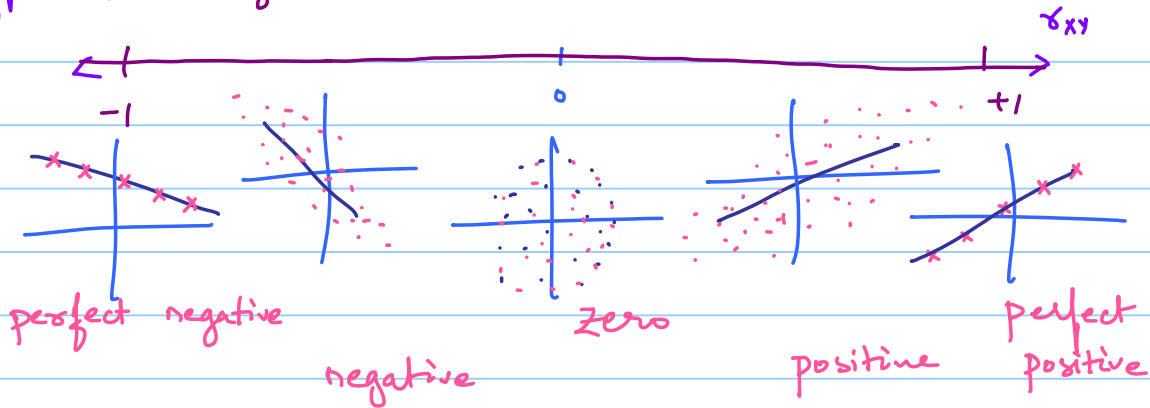


Regression: finding this function \rightarrow actual relation - $Y = f(X)$
linear relation -

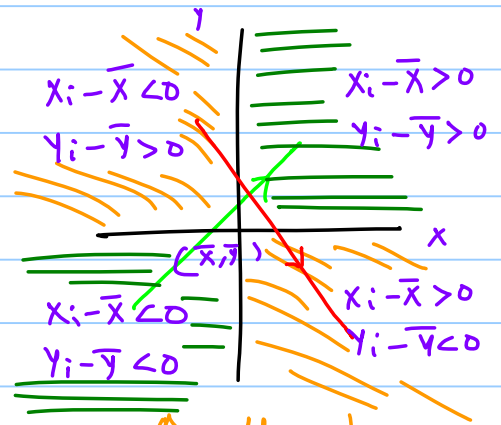
r_{xy} } Correlation \rightarrow sign
Coefficient \rightarrow magnitude

we want a measure like this
 $r_{xy} \in [-1, 1]$

Correlation



without loss of generality
 \rightarrow translate origin to the mean of data



$$\text{Cov}(X, Y) := \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n-1}$$

\rightarrow as one \uparrow , other \downarrow
 $\boxtimes (X_i - \bar{X})(Y_i - \bar{Y}) < 0$

$$\boxplus (X_i - \bar{X})(Y_i - \bar{Y}) > 0$$

-- X & Y can be heterogeneous. } if $X_i = \text{mass of } i^{\text{th}} \text{ body}$,
 Δ units of X, Y , unintuitive. } $Y_i = \text{force on } i^{\text{th}} \text{ body}$

$\text{Cov}(X, Y)$ has the dimension of force (N) or $[MLT^{-2}]$

Karl Pearson's Correlation coefficient, $r_{xy} = \frac{\text{Cov}(X, Y)}{S_X \cdot S_Y}$

++ r_{xy} is normalized, unitless. ++ Compare datasets \checkmark

gray areas

ROUGH !

!?

With an inner product processor L_2 norm is easy

+

$$X = \begin{bmatrix} x_1^T & x_2^T & \dots & x_n^T \end{bmatrix} \quad \text{vector, one instance.}$$

$$X = \begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix}_{d \times n} \quad \text{matrix}$$

$$XX^T = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} & & & \\ \vdots & & & \\ x_{d1} & & & x_{nd} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1d} \\ x_{21} \\ \vdots \\ x_{n1} & & & x_{nd} \end{bmatrix}$$

$$X = (x_1, x_2, \dots, x_n)^T$$

$$Y = (y_1, y_2, \dots, y_n)^T$$

$$\begin{aligned} \text{Var}(x) &= \frac{1}{n-1} \sum_i (x_i - \bar{x})^2 = \frac{1}{n-1} \sum_i (x_i^2 + \bar{x}^2 - 2x_i\bar{x}) = \frac{1}{n-1} \left\{ \begin{array}{l} x^T x + n\bar{x}^2 \\ - 2\bar{x} \cdot n\bar{x} \end{array} \right\} \\ &\rightarrow \frac{(x - \bar{x})(x - \bar{x})^T}{n-1} = \frac{1}{n-1} (x^T x - n\bar{x}^2) \end{aligned}$$

$$\text{Var}(y) := \frac{1}{n-1} (y - \bar{y})^T (y - \bar{y})$$

$$\text{Cov}(x, y) := \frac{1}{n-1} (x - \bar{x})^T (y - \bar{y})$$

$$r_{xy} = \frac{(x - \bar{x})^T (y - \bar{y})}{\sqrt{(x - \bar{x})(x - \bar{x})} \sqrt{(y - \bar{y})^T (y - \bar{y})}} \quad \left\{ \div \frac{n-1}{n-1} \right\}$$

Prob Set

bijections
functions
injections

SRSWOR
SRSWR