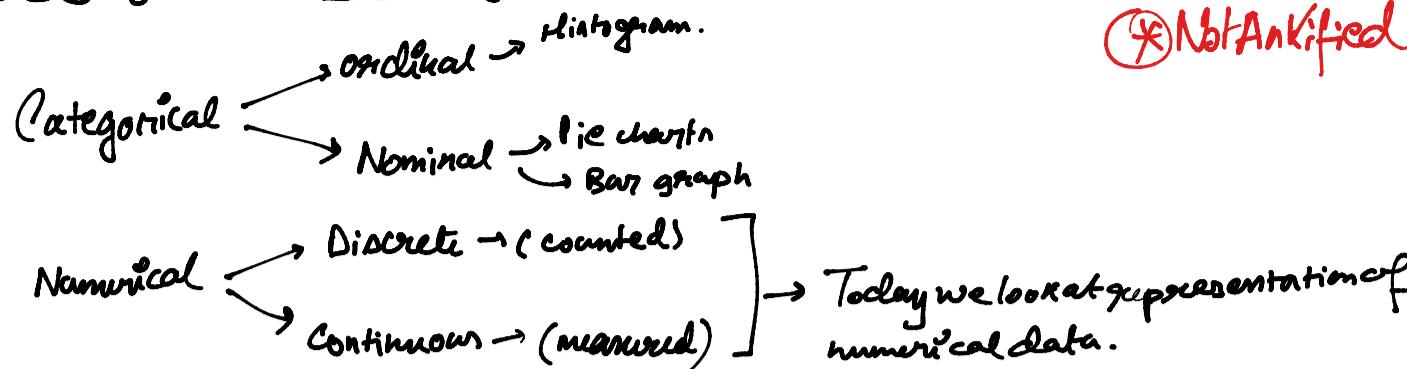


9th January 2025 (Thursday) →



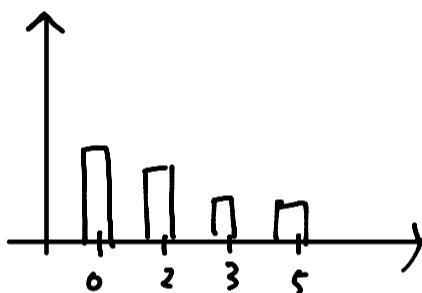
Discrete data →

No. of phone calls received in a week →

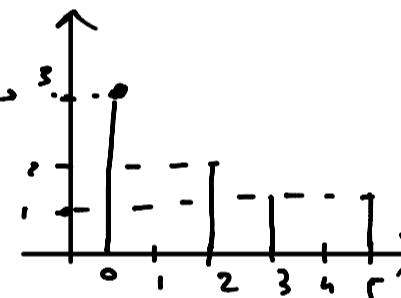
2, 0, 0, 3, 2, 5, 0

# Phone calls	# days this no. of calls were received
0	3
2	2
3	1
5	1

To represent this,



or, a line diagram →



But these are not good representation of data — especially if there is repetition.

② Stem-leaf Plot →

Marks → 78, 74, 82, 66, 94, 71, 64, 88, 55, 80, 91, 74, 82, 75, 96, 78, 84, 79, 71, 88

(Why did I copy all of this? No clue)

Stem → last digit

leaf → All other digits.

Stem	Leaves
7	8, 4, 1, 4, 5, 8, 9, 1
8	2, 8, 0, 2, 4, 3
6	6, 4
9	4, 1, 6
5	5

Stem	Leaves
5	5
6	6, 4
7	8, 4, 1, 4, 5, 8, 9, 1
8	2, 8, 0, 2, 4, 3
9	4, 1, 6

Now,
we
order
by
stem

Note: Requires # of digits does not vary a lot over the data.

Exercise → Do stem-leaves plot.

23, 26, 11, 18, 9, 21, 23, 30, 22, 11, 21, 20, 11, 13, 23, 11, 29, 25,
28, 26

	L
2	3, 6, 1, 3, 2, 1, 0, 3, 9, 5, 6, 6
1	1, 8, 1, 1, 3, 1
0	9
3	0

ordering,

S	L
0	9
1	1, 8, 1, 1, 3, 1
2	3, 6, 1, 3, 2, 1, 0, 3, 9, 5, 6, 6
3	0

* Tutorial: Tuesday, 7-8 pm (Probably from
(8AM class not going to be there))

L-S plot is called exploratory analysis.

① Representation of continuous data via histogram →

Suppose x_1, x_2, \dots, x_n are observations of heights of n individuals.

Let the data value lie in the interval $[a, b]$, where $a < b$ are two real numbers.

Divide the interval $[a, b]$ into K classes of equal length, say $h > 0$.

∴ The classes are,

$$[a, a_1], (a_1, a_2], (a_2, a_3], \dots, (a_{K-1}, b]$$

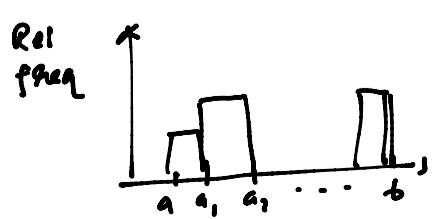
Then, we assign a frequency value to each class i , $i = 1, \dots, K-1$,

$$f_i = \#\{x_i \in (a_i, a_{i+1}]\}, \text{ with } a_0 = a, a_K = b$$

Class	Frequency	Relative freq	freq. density
$[a, a_1]$	f_1	f_1/n	$f_1/h \rightarrow \text{step size}$
$(a_1, a_2]$	f_2	f_2/n	f_2/h
$(a_2, a_3]$	f_3	f_3/n	f_3/h
\vdots	\vdots	\vdots	\vdots
$(a_{K-1}, b]$	f_K	f_K/n	f_K/h

$$n = \sum_i f_i$$

We may plot histogram →



Some what is use of freq. density?

It is an estimation of underlying frequency density function.

Remarks →

- ① Suppose X be a r.v. with CDF $F(x)$ and the observations are,
 x_1, x_2, \dots, x_n are iid copies of X

Define another r.v.,

$$Y_{i,j} = \begin{cases} 1, & x_j \in (a_i, a_{i+1}] \\ 0, & \text{otherwise.} \end{cases}$$

$$\Rightarrow Y_{i,j} \sim \text{Bernoulli}(1, p_i), \quad p_i = P(a_i < x_j \leq a_{i+1})$$

$$Y_i = \sum_{j=1}^n Y_{i,j} \rightarrow \text{No. of success in prob } p_i \text{ in } n \text{ trials.}$$

Obviously,

$$Y_i \sim \text{Binomial}(n, p_i)$$

$$\Rightarrow E(Y_i) = np_i$$

$$\text{Also, note, } Y_i = f_i$$

$$\Rightarrow E(f_i) = np_i \Rightarrow E\left(\frac{f_i}{nh}\right) = \frac{p_i}{h} = \frac{P(a_i < x \leq a_{i+1})}{h}$$

We see that,

$$\begin{aligned} f_x(x) &= \frac{d}{dx} F_x(x) \\ &= \lim_{h \rightarrow 0} \frac{F_x(x+h) - F_x(x)}{h} \\ &= \lim_{h \rightarrow 0} \frac{P(x \leq x+h) - P(x \leq x)}{h} &= \lim_{h \rightarrow 0} \frac{P(x < x < x+h)}{h} \end{aligned}$$

10th January 2025 (Friday) →

(*) Not Audited

Correction to L-S plot: The leaf is always the last digit, the rest is the stem.

(*) Freq distribution of the no. of pear in a pod →

Obviously, there are discrete values.

Data: A lot of numbers. I cannot fathom why writing this is useful.

I will still attempt, I guess.

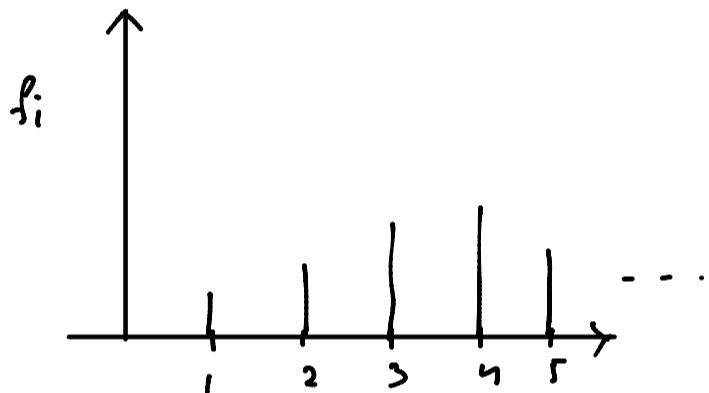
Data → 7, 1, 2, 4, 5, 2, 4, 5, 2, 4, 4, 4, 5, 3, 5, 3, 6, 3, 2, 2, 3, 4, 7,
1, 1, 1, 5, 3, 3, 4, 3, 2, 1, 6.

# Peas	Tally	f_i
1		5
2		6
3		:

Something. — Cannot finish.

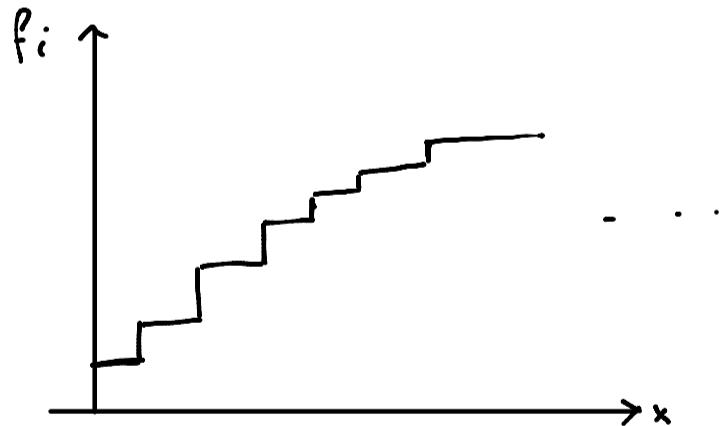
We calculate the cumulative freq $>$ and $<$ → The sum of all frequencies after and before some f_i

Freq Bar Diagram →



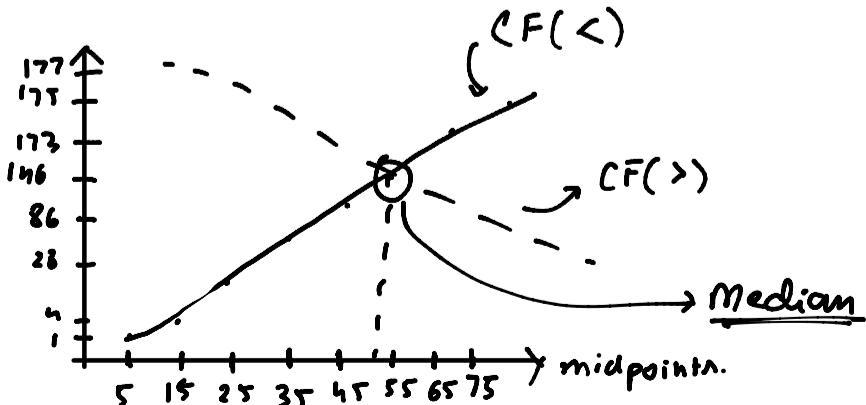
We may also draw the CF diagrams.

CF Bar diagram →



If this is not discrete, the continuous curve is called an **Ogive**
— credit to Niravra

We draw a few ogives →



Problems (Tutorial/Exercise) → (Next to next week)

Goon Gupta, Dangupta → Chapter - 3
3.8, 3.9, 3.10, 3.11, 3.12

① Measures of Central Tendency →

Tendency of data points to cluster around some mean point.

Δ Arithmetic Mean →

Given a dataset x_1, x_2, \dots, x_n ,

$$\text{A.M.}, \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

If the data is given in freq table,

$$\bar{x} = \frac{\sum_{i=1}^n f_i x_i}{\sum_{i=1}^n f_i}$$

How to do this for continuous data in freq table,

We simply associate class mean with the freq of the class, and calculate like discrete data in freq table.

This is an approximate mean.

Some properties →

$$① \bar{x} = \frac{\sum_{i=1}^n x_i}{n} \Rightarrow \sum_{i=1}^n x_i = n\bar{x}$$

$$\Rightarrow \sum_{i=1}^n (x_i - \bar{x}) = n\bar{x} - n\bar{x} = 0$$

⇒ Sum of the deviations of given values from its mean is necessarily 0.

② Mean of a variable who all have the same value must also be the same value.

Proof: $x_i = a$

$$\therefore \bar{x} = \frac{\sum_{i=1}^n a}{n} = \frac{na}{n} = a$$

③ If y is a linear function of x , then \bar{y} and \bar{x} are related in the same way as y and x

Proof $\rightarrow y = f(x)$, f is linear.

$$\text{i.e., } f(x_1 + x_2) = f(x_1) + f(x_2) \quad , \quad f(ax) = af(x)$$

$$\therefore \bar{y} = \frac{\sum_{i=1}^n y_i}{n} = \frac{\sum_{i=1}^n f(x_i)}{n} = f\left(\sum_{i=1}^n \frac{x_i}{n}\right) = f(\bar{x})$$

(1) If we have attached freq to each value, $(x_1, n_1), (x_2, n_2), \dots, (x_k, n_k)$

$$\therefore \bar{x} = \frac{\sum_{i=1}^n x_i n_i}{\sum_{i=1}^n n_i}$$

(2) $Z = ax + by$

$$\Rightarrow \bar{Z} = a\bar{x} + b\bar{y}$$

(3) \bar{x} solves the optimization problem given data x_1, x_2, \dots, x_n

$$\bar{x} \text{ minimizes } \sum_{i=1}^n (x_i - c)^2 \text{ w.r.t } c \in \mathbb{R}.$$

$$\bar{x} = \operatorname{argmin} \left\{ \sum_{i=1}^n (x_i - c)^2 \right\}$$

Proof Now,

$$\sum_{i=1}^n \{(x_i - \bar{x}) + (\bar{x} - c)\}^2 = \sum_{i=1}^n \left[\underbrace{(\bar{x} - \bar{x})^2}_{\geq 0} + \underbrace{(\bar{x} - c)^2}_{\geq 0} + 2(\bar{x} - \bar{x})(\bar{x} - c) \right]$$

for min.

$$\therefore \bar{x} = c$$

Merits of A.M. \rightarrow

- ① Rigid Definition (?)
- ② Easy to understand and compute.
- ③ Based on all observations.
- ④ Amenable to algebraic manipulation.
- ⑤ Least affected by sampling fluctuation.

Demerits of A.M. \rightarrow

- ① Location by inspection not possible
- ② Affected by outliers
- ③ Categorical Data not possible

④ Opened interval

⑤ Skewed dist.