

14th January 2025 (Tuesday) →

Not AnkiFiel

④ Median →

$$\underline{\text{Ex:}} \quad 7, 9, 5, 10, 4, 3, 12 = x \\ \tilde{x} = \text{median.}$$

First, we need to order the data.

$$x = 4, 5, 7, 8, 9, 10, 12$$

$$\text{Median} = \frac{12 - 4}{2} + 4 = 4 + 4 = 8 \rightarrow \text{Work if } n \text{ is odd} \rightarrow \text{Does not work} \\ (\text{Why? What if data skewed})$$

If n is even, the median is any number b/w the middle two numbers.

$$\therefore \tilde{x} = \begin{cases} x_{\left(\frac{n+1}{2}\right)} & , n \text{ is odd} \\ \frac{1}{2} [x_{\left(\frac{n}{2}\right)} + x_{\left(\frac{n}{2}+1\right)}] & , n \text{ is even.} \end{cases}$$

$$\underline{\text{Ex}} \rightarrow \text{Mark} = 50, 41, 72, 89, 92, 72, 60, 100$$

$$\text{Ordered} = 41, 50, 60, 72, 72, 89, 92, 100$$

$n \rightarrow \text{even.}$

$$\therefore \tilde{x} = \frac{1}{2} [x_4 + x_5] \\ = \frac{1}{2} [(72) + (72)] = 72$$

Now, in a frequency table.

x_i	f_i	C.F.
x_1	f_1	f_1
x_2	f_2	$f_1 + f_2$
\vdots	\vdots	\vdots
x_k	f_k	n
$\sum f_k = n$		

If $n \rightarrow \text{odd}$

$$F_i \geq \frac{n}{2}$$

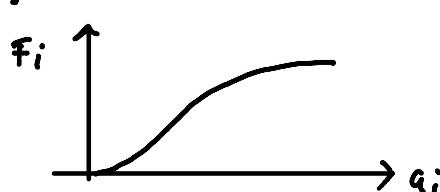
Now, if this is b/w two classes,

$$F_i \geq \frac{n}{2} \rightarrow x \Rightarrow \tilde{x} = \frac{x + x'}{2}$$

$$F_i \geq \frac{n}{2} + 1 \rightarrow x'$$

What if the data is continuous?

Or give, for example, look like —



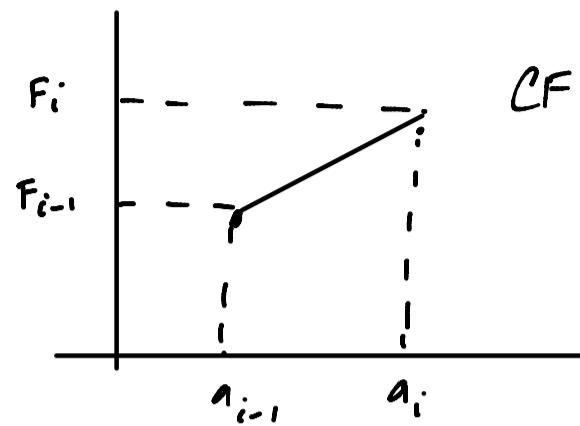
i.e. data is like,

Int	f_i	F
$[a_1, a_2]$	f_1	f_1
$(a_2, a_3]$	f_2	$f_1 + f_2$
\vdots	\vdots	\vdots

Say, we have,

$$F_i \geq \frac{n}{2} \rightarrow (a_{i-1}, a_i]$$

We assume that F_i grows linearly b/w i^{th} and i^{th} class.



CF changes linearly - assumption.

∴ Eqn of line,

$$\frac{y - F_{i-1}}{F_i - F_{i-1}} = \frac{x - a_{i-1}}{a_i - a_{i-1}}$$

$\underbrace{f_i}_{\text{f}_i}$ $\underbrace{h}_{\text{h}}$

We should also substitute, $y = \frac{n}{2}$

$$\therefore \frac{\frac{n}{2} - F_{i-1}}{f_i} = \frac{x - a_{i-1}}{h} \Rightarrow x = \frac{h}{f_i} \left[\frac{n}{2} - F_{i-1} \right] + a_{i-1}$$

How find out the median graphically? Simply, it the point of intersection of ogive with $F_i = \frac{n}{2}$.

Δ Formal defn of median: For any cont.. a.v. X , the median 'm' is defined as that value for which,

$$P(X \geq m) \geq \frac{1}{2} \text{ and } P(X \leq m) \geq \frac{1}{2}$$

Say, we do next on the a.v. values.

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$$

$$n = \text{Odd} = 2k+1 \rightarrow$$

$$\text{Let } m = k+1$$

$$\rightarrow \frac{\# \text{obs} \leq x_m}{\# \text{obs}} = \frac{\# \text{obs} \leq x_{k+1}}{2k+1} = \frac{k+1}{2k+1} \geq \frac{1}{2}$$

$$\text{Also, } \frac{\# \text{obs} \geq x_m}{\# \text{obs}} = \frac{k+1}{2k+1} \geq \frac{1}{2}$$

$$n = \text{Even} = 2k$$

$$\Rightarrow m = \frac{1}{2} [x_k + x_{k+1}]$$

$$\therefore \frac{\# \text{obs} \geq m}{\# \text{obs}} = \frac{k}{2k} = \frac{1}{2}$$

∴ Median defined for data satisfies the defn. for probability.

Merits	Demerits
--------	----------

- | | |
|---|--|
| <ol style="list-style-type: none"> ① Rigid Defn ② Easily understood ③ Not affected by outliers | <ol style="list-style-type: none"> ① no exact \bar{x} for $n=even$ ② algebraic treatment not possible ③ not based on all obs ④ Affected by sampling fluctuations |
|---|--|

What is the optimization provided by the median?

$$\tilde{x} = \operatorname{argmin}_c \left\{ \sum_i |x_{(i)} - c| \right\}$$

Proof: Say, data is,

$$x_{(1)} < x_{(2)} < \dots < x_{(n)}$$

Say, $n = 2k+1$ (odd)

Consider 2 terms,

$$|x_{(1)} - c| + |x_{(n)} - c|$$

When $c < x_{(1)}$,

$$\begin{aligned} |x_{(1)} - c| + |x_{(n)} - c| &= x_{(1)} - c + x_{(n)} - c + x_{(1)} - x_{(n)} \\ &= (x_{(n)} - x_{(1)}) + 2 \underbrace{(x_{(1)} - c)}_{+ve} \end{aligned}$$

If $c > x_{(n)}$,

$$\begin{aligned} |x_{(1)} - c| + |x_{(n)} - c| &= c - x_{(1)} + c - x_{(n)} + x_{(n)} - x_{(1)} \\ &= (x_{(n)} - x_{(1)}) + 2 \underbrace{(c - x_{(n)})}_{+ve} \end{aligned}$$

Now, $x_{(1)} \leq c \leq x_{(n)}$

$$\begin{aligned} \therefore |x_{(1)} - c| + |x_{(n)} - c| &= c - x_{(1)} + x_{(n)} - c \\ &= x_{(n)} - x_{(1)} \rightarrow \underline{\text{Minimal}} \end{aligned}$$

\Rightarrow for $c \in [x_{(1)}, x_{(n)}]$, we will minimize $|x_{(1)} - c| + |x_{(n)} - c|$

Similarly proceeding,

for $c \in [x_{(k)}, x_{(n-k+1)}]$, we will minimize $|x_{(k)} - c| + |x_{(n-k+1)} - c|$

16th January 2025 (Thursday) →

② Mode →

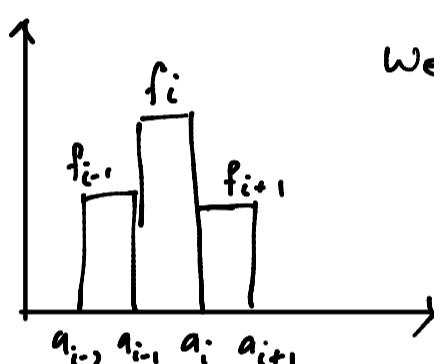
$$x = 10, 7, 8, 9, 10, 5, 12, 6, 10, 15, 9, 8, 10, 9, 12$$

$$\hat{x} = \text{Mode} = 10$$

What happens if there is more than one element having highest frequency?

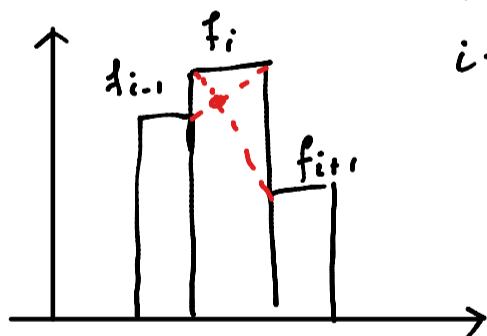
Mode is not defined here — Mode only defined for unimodal distributions.

What do we do with continuous data?



We just take the mode to be the middle point of $(a_{i-1}, a_i] \rightarrow i.e,$
 $\hat{x} = \frac{a_{i-1} + a_i}{2} = a_{i-1} + \frac{h}{2}$ → Assuming const class size.

What if we have →



Makes sense for the mode to be biased towards $i-1$ class.

∴ We find intersection of these interpolated linear curves. Line joining (a_{i-1}, f_{i-1}) to (a_i, f_i)

$$\frac{y - f_{i-1}}{f_i - f_{i-1}} = \frac{x - a_i}{a_i - a_{i-1}}$$
$$\Rightarrow y = f_{i-1} + \frac{(f_i - f_{i-1})(x - a_{i-1})}{h}$$

and the other line, joining (a_{i-1}, f_i) to (a_i, f_{i+1})

$$\frac{y - f_i}{f_{i+1} - f_i} = \frac{x - a_{i-1}}{a_i - a_{i-1}}$$
$$\Rightarrow y = f_i + \frac{(f_{i+1} - f_i)(x - a_{i-1})}{h}$$

Solving these,

$$\hat{x} = a_{i-1} + \frac{(f_i - f_{i-1})h}{2f_i - (f_{i+1} + f_{i-1})}$$

Empirical relation for nearly symmetric distribution \rightarrow

$$\bar{x} - \hat{x} = 3(\bar{x} - \tilde{x})$$

For perfectly symmetric dist,

$$\tilde{x} = \hat{x} = \bar{x}$$

Merits / Demerits of mode \rightarrow

Merits

- Fairly computed by inspection
- Not affected by outliers
- Unequal class-interval,
or open-ended or
class intervals

Demerits

- \tilde{x} not defined
- Not all observations
- Sampling fluctuation.
- Not mathematically treatable
(Chronic cancer)

$$\tilde{x} = \arg \min f(x_i)$$

Remarks \rightarrow

① Mean, median (except even n) rigidly defined.

Mode \rightarrow only for unimodal data.

② Determination of mode not possible if very few data points are given

③ Mean \rightarrow value passed by each unit if total volume ($\sum x_i$) were distributed equally amongst all units.

$\tilde{x}, \hat{x} \rightarrow$ missing data.

We may ignore top 25% and bottom 25% of data in \bar{x} calculation to deal

with the outlier effect. \rightarrow This is called a **trimmed mean**.
If do not delete, but replace top 25% with value just below it, bottom 25% with value just above it \rightarrow **Winsorized mean**.

* We thus generally take \bar{x} to be the best measure.

\rightarrow The minimization of sum of square deviations, accounting all of the data, etc.

Ex: Weights (kg) : 138, 143, 141, 139, 152, 148, 160, 267

$n=8, \bar{x}=161$ kg \rightarrow Not good central measure,
 \rightarrow lie below it

Why? 267 is extreme outlier.

Here, \tilde{x} is appropriate.

Eg: h values, $\{15, 30, 33, 36\}$

we sample (i.e. take out randomly with replacement)

		\bar{x}	\tilde{x}
1	15, 30, 33	26	30
2	15, 30, 36	27	30
3	15, 33, 36	28	33
4	30, 33, 36	33	33

} → Median is pretty stable,
mean fluctuates
(due to 15, outlier)

There are 2 more central tendencies, not generally used but useful in some distributions →

① Geometric Mean →

$$x_g = \left(\prod_{i=1}^n x_i \right)^{1/n}$$

$$\Rightarrow \log x_g = \frac{1}{n} \sum_{i=1}^n \log x_i \rightarrow \text{A.M. of log of values.}$$

Note: $\left\{ \prod_{i=1}^n \left(\frac{x_i}{y_i} \right) \right\}^{1/n} = \frac{\left\{ \prod_{i=1}^n x_i \right\}^{1/n}}{\left\{ \prod_{i=1}^n y_i \right\}^{1/n}}$ → Average price
Relative (Ratio)
useful here, the g.m. of ratio
is ratio of g.m.'s.

Also, say,

$$t_1 = 0, y = a$$

$$t_2 = t, y = a e^{rt}$$

To find value at mid point, we calculate G.M.

$$a e^{t/2} = \sqrt{a \cdot a e^{rt}}$$

② If we have any $x_i = 0$ or $x_i < 0$, we hit a problem with GM.

Done really? We may simply shift entire dataset by fixed value to ensure $x_i \neq 0$ and $x_i \neq 0 + i$, calculate GM, then shift back.

17th January 2025 (Friday) →

Harmonic Mean →

$$\bar{X}_n = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}} \rightarrow \text{AM of the reciprocals is the reciprocal of HM.}$$

→ x units per y units (?)?

Equal share of x units → HM

Equal share of y units → AM

① Weighed Data → (x_i, w_i) → Treating the weights as frequencies of the data.

$$\rightarrow \text{AM: } \bar{X}_A = \frac{\sum_i w_i x_i}{\sum_i w_i}$$

$$\rightarrow \text{GM: } \bar{X}_G = \left[\prod_i x_i^{w_i} \right]^{\frac{1}{\sum_i w_i}}$$

$$\rightarrow \text{HM: } \bar{X}_H = \frac{\sum_i w_i}{\sum_i \frac{w_i}{x_i}}$$

} Generalized
for weighed
data.