

* General amalgamation from Sabano's notes for this week. Enjoy!

Measures of Dispersion →

△ Dispersion: It quantifies how far the datapoints are from their measure of central tendency.

There are 3 main measures of dispersion →

- ① Range
- ② Mean Absolute deviation (MAD)
- ③ Standard deviation.

Let $X = \{x_1, \dots, x_n\}$ be an ordered dataset.

Here,

$$\textcircled{1} \text{ Range} = x_n - x_1$$

$$\textcircled{2} \text{ } MAD_{\mu} = \frac{1}{n} \sum_{i=1}^n |x_i - \mu|$$

Measure of Central Tendency

$$\textcircled{3} \text{ } \sigma_{\mu} = \left[\sum_{i=1}^n \frac{(x_i - \mu)^2}{n} \right]^{1/2}$$

Measure of Central Tendency

If we have the central tendency used for MAD as the mean, it is minimized.

$$\begin{aligned} MAD_{\bar{x}} &= \frac{1}{n} \sum_i |x_i - \bar{x}| \\ &= \frac{1}{n} \sum_{x_i < \bar{x}} |x_i - \bar{x}| + \frac{1}{n} \sum_{x_i \geq \bar{x}} |x_i - \bar{x}| \\ &= \frac{1}{n} \sum_{x_i < \bar{x}} (\bar{x} - x_i) + \frac{1}{n} \sum_{x_i \geq \bar{x}} (x_i - \bar{x}) \end{aligned}$$

We use the fact that, $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) = 0$

$$\Rightarrow \sum_{x_i < \bar{x}} (x_i - \bar{x}) = - \sum_{x_i \geq \bar{x}} (x_i - \bar{x})$$

Putting this in,

$$MAD_{\bar{x}} = \frac{2}{n} \sum_{x_i < \bar{x}} (\bar{x} - x_i) \quad \text{OR} \quad MAD_{\bar{x}} = \frac{2}{n} \sum_{x_i \geq \bar{x}} (x_i - \bar{x})$$

Thus, we only have to compute for half the dataset if we know the mean.

Now, for standard deviation →

$$RMS_A = \left[\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 \right]^{1/2} \rightarrow \text{Root mean square deviation about } \mu.$$

When $\mu = \bar{x}$ (AM), RMS about \bar{x} is called Standard deviation

$$S.D = \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2} = \left[\sum_{i=1}^n \left(\frac{x_i^2}{n} \right) - \bar{x}^2 \right]^{1/2}$$

$$\text{Similar to, } \langle x - \bar{x} \rangle^2 = \langle x^2 \rangle - \langle x \rangle^2$$

For weighed data,

$$\begin{aligned} S.D &= \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 f_i \right]^{1/2}, \quad n = \sum_{i=1}^n f_i \\ &= \left[\frac{1}{n} \sum_{i=1}^n (x_i^2 f_i) - \bar{x}^2 \right]^{1/2} \end{aligned}$$

① Properties of standard deviation →

$$(1) x_i = a \quad \forall i \Rightarrow \bar{x} = \sum_{i=1}^n \frac{x_i}{n} = a \Rightarrow \sigma_{\bar{x}} = 0$$

\hookrightarrow Constant has 0 S.D.

$$(2) \text{ Let } y = a + bx \Rightarrow \sigma_{\bar{y}} = |b| \sigma_{\bar{x}} \rightarrow \text{Scaling of S.D.}$$

(3) Let there be 2 sets of x ,

$$x_1, \dots, x_{n_1} \rightarrow n_1 \text{ elements, } \bar{x}_1, \sigma_{\bar{x}_1}$$

$$x_{21}, \dots, x_{2n_2} \rightarrow n_2 \text{ elements, } \bar{x}_2, \sigma_{\bar{x}_2}$$

$$\Rightarrow \sigma_{\bar{x}}^2 = \frac{n_1 \sigma_{\bar{x}_1}^2 + n_2 \sigma_{\bar{x}_2}^2}{n_1 + n_2} + \frac{n_1 (\bar{x}_1 - \bar{x})^2 + n_2 (\bar{x}_2 - \bar{x})^2}{n_1 + n_2}$$

\downarrow
Mean of both
dist considered
to be one.

This, in general → (Pooled Variance)

$$\sigma_{\bar{x}}^2 = \frac{\sum_{i=1}^t n_i \sigma_{\bar{x}_i}^2}{\sum_{i=1}^t n_i} + \frac{\sum_{i=1}^t n_i (\bar{x}_i - \bar{x})^2}{\sum_{i=1}^t n_i}$$

\downarrow
Mean of t sets

* CT soon!!

Proof (1) trivial,

Proof of (2) \rightarrow

$$\sigma_{\bar{x}} = \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2}$$

Changing variables,

$$\begin{aligned}\sigma_{\bar{y}} &= \left[\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \right]^{1/2} = \left[\frac{1}{n} \sum_{i=1}^n (6x_i - 6\bar{x})^2 \right]^{1/2} \\ &= |6| \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2} \\ &= |6| \sigma_{\bar{x}}\end{aligned}$$

∴ True Proof of (3).

② Moments \rightarrow

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad , \quad \bar{x} = \frac{1}{n} \sum_i x_i f_i \text{ , where } n = \sum_i f_i$$

Let us consider the moments, defined as,

$$m_n = \frac{1}{n} \sum_{i=1}^n (x_i - A)^n \quad , \quad m'_n = \frac{1}{n} \sum_i (x_i - A)^n f_i$$

Some measure
of Central
Tendency

③ When there is grouped data, $x_i \rightarrow$ midpoint of a class.

When $A = \bar{x}$, it is called central moment.

$$\therefore m_n = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^n$$

$$\therefore m'_n = \frac{1}{n} \sum_i (x_i - \bar{x})^n f_i$$

Some values,

$$\textcircled{1} \quad m_0 = m'_0 = 1$$

$$\textcircled{2} \quad m_1 = 1, \quad m'_1 = \bar{x} - A$$

$$\text{Note, } (x_i - \bar{x})^n = [(x_i - A) - (\bar{x} - A)]^n$$

Using Binomial theorem,

$$\Rightarrow (x_i - \bar{x})^n = \sum_j \binom{n}{j} (x_i - A)^{n-j} (\bar{x} - A)^j (-1)^j$$

$$\Rightarrow m_n = \frac{1}{n} \sum_i \sum_j \binom{n}{j} (x_i - A)^{n-j} (\bar{x} - A)^j (-1)^j$$

$$= \frac{1}{n} \sum_j \binom{n}{j} (m'_{n-j}) (m')^j (-1)^j$$

) DCheck!

Now,

$$m_1 = m'_1 - m_1' = 0$$

$$m_2 = m'_2 - m_2'^2$$

$$m_3 = m'_3 - 3m'_2 m'_1 + 3m'_1^3 - m'_1^3$$

$$= m'_3 - 3m'_2 m'_1 + 2m'_1^3$$

$$m'_4 = m'_4 - 4m'_3 m'_1 + 6m'_2 m'_1^2 - 3m'_1^4$$

Note. $m'_1 = d$, $m'_2 = m_2 + d$, $m'_3 = m_3 + 3m_2 d + d^3$

$$m'_4 = m_4 + 4m_3 d + 6m_2 d^2 + d^4$$

① Sheppard's Correction → (for general moments)

c : Length of class interval.

$$m'_1(c) = m'_1$$

$$m'_2(c) = m'_2 - \frac{c^2}{12}$$

$$m'_3(c) = m'_3 - \frac{c^2 m_1}{4}$$

$$m'_4(c) = m'_4 - \frac{c^2}{2} m'_2 + \frac{7}{240} c^4$$

$$\Rightarrow m_1(c) = m_1$$

$$m_2(c) = m_2 - \frac{c^2}{12}$$

$$m_3(c) = m_3$$

$$m_4(c) = m_4 - \frac{c^2}{2} m_2 + \frac{7}{240} c^4$$

② Note how there are no corrections for $m_n(c)$, $\forall n = \text{odd}$.

When are these corrections valid?

① Observations from continuous data

② Freq curve of dist is tapered towards end (??)

③ Total freq. should be large

④ Class length should not be too small.

⑤ Relations b/w moments →

⊗ Cauchy-Schwarz Inequality for Expectations →

$$|E(f(x)g(x))| \leq \sqrt{E(f(x)^2)} \cdot \sqrt{E(g(x)^2)}$$

Assumption: $E(f(x)^2) < \infty$, $E(g(x)^2) < \infty$

$$E(f(x) \cdot g(x)) < \infty$$

⊗ $E(x) = \sum_i x_i P(x = x_i) = \bar{x}$

\downarrow
uniform dist

$$\text{Here, } \frac{1}{n} | \sum_i f(x_i) g(x_i) | \leq \sqrt{\frac{1}{n} \sum_i f^2(x_i)} \cdot \sqrt{\frac{1}{n} \sum_i g^2(x_i)}$$

$$\text{Set } f(x) = (x - \bar{x})^m, g = (x - \bar{x})$$

$$\Rightarrow \frac{|m_{m+1}|}{\sigma} \leq \sqrt{m}, \sigma : \text{standard deviation}$$

$$\text{Set, } f(x) = (x - \bar{x}_0)^{m+1}, g(x) = (x - \bar{x})^m$$

$$\Rightarrow |m_{m+1}| \leq \sqrt{m_{m+2}} \cdot \sqrt{m_m}$$

⑥ Jensen's Inequality →

If ϕ is a convex function, and $E(|\phi(x)|) < \infty$, then,

$$E(|\phi(x)|) \geq \phi(E(x))$$

In the discrete case,

$$\frac{1}{n} \sum_i \phi(x_i) \geq \phi\left(\sum_i \frac{x_i}{n}\right)$$

$$\text{Say, } \phi(x) = \frac{1}{x}$$

Now,

$$\frac{1}{n} \sum_i \frac{1}{x_i} \gg \frac{n}{\sum_i x_i}$$

$$\Rightarrow \bar{x} \gg \frac{1}{\frac{1}{n} \sum_i \frac{1}{x_i}} \rightarrow \text{Note: This is AM-HM inequality.}$$

Now, let $\phi(x) = -\log(x)$

$$\therefore \frac{1}{n} \sum_i \log(x_i) \leq \log\left(\frac{1}{n} \sum_i x_i\right)$$

$$\Rightarrow \log\left\{\left(\prod_i x_i\right)^{1/n}\right\} \leq \log\left(\frac{1}{n} \sum_i x_i\right)$$

④ \log is a monotonic function, so we may remove \log to get AM-GM inequality!

∴ In general,

$$\text{QM} \geq \text{AM} \geq \text{GM} \geq \text{HM}$$

To get this, take AM > GM
and then set
 $x_i \rightarrow \frac{1}{x_i}$

(?)

Surprise addition!

① Grouped data →

$$\Delta_1 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n f_i f_j |x_i - x_j|$$

So, here the data is all possible absolute differences b/w the data points

Why $f_i f_j$? f_i ways to select x_i , f_j ways to select x_j for same value of abs. diff.

$$\Delta_1 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2 \quad (?)$$

$$\Delta_2 = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n f_i f_j (x_i - x_j)^2$$

Now,

$$\begin{aligned}
 \sum_{i=1}^n \sum_{j=1}^n (x_i - x_j)^2 &= \sum_{i=1}^n \sum_{j=1}^n (x_i - \bar{x} + \bar{x} - x_j)^2 \\
 &= \sum_{i=1}^n \sum_{j=1}^n \left[(x_i - \bar{x})^2 + (x_j - \bar{x})^2 - 2(x_i - \bar{x})(x_j - \bar{x}) \right] \\
 &= \sum_{i=1}^n \sum_{j=1}^n \left[(x_i - \bar{x})^2 + \underbrace{(x_j - \bar{x})^2}_{\text{This term}} \right] \rightarrow 0 \\
 &= 2n^2 \sigma_{\bar{x}}^2
 \end{aligned}$$

① Quartile deviation →

Inter Quartile Range → IQR

Reminder: for any R.V. X , the median m is defined as value for which,

$$P(X \geq m) \geq \frac{1}{2}, \quad P(X \leq m) \geq \frac{1}{2}$$

Generalizing this,

we define p^{th} quartile of a R.V. X , denoted as z_p , as,

$$P(X \leq z_p) \geq p, \quad P(X \geq z_p) \geq 1-p$$

i.e. if $p = \frac{1}{4} \Rightarrow z_{1/4} \rightarrow Q_1$ (1st quartile)

$p = \frac{1}{2} \Rightarrow z_{1/2} \rightarrow Q_2$ (median)

$p = \frac{3}{4} \Rightarrow z_{3/4} \rightarrow Q_3$ (3rd Quartile)

$$\Delta IQR = Q_3 - Q_1$$

Classes	f_i	$C.F(c)$
$[a_0, a_1]$	f_1	$F_1 = f_1$
$(a_1, a_2]$	f_2	$F_2 = f_1 + f_2$
\vdots	\vdots	\vdots
$(a_{k-1}, a_k]$	f_k	$F_k = n$
	$\sum f_i = n$	

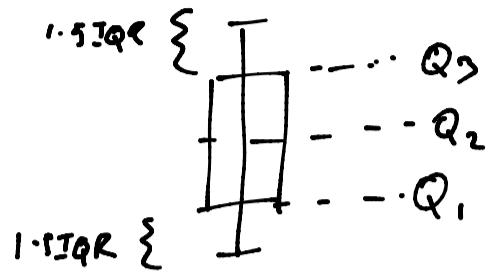
\Rightarrow for $F_{i+1} \geq \frac{n}{2}$ and $F_i < \frac{n}{4}$
 $\Rightarrow (a_i, a_{i+1}] \rightarrow \underline{1\text{st quartile}}$
and $F_{i+1} \geq \frac{3n}{4}$ and $F_i < \frac{3n}{4}$
 $\Rightarrow (a_i, a_{i+1}] \rightarrow \underline{3\text{rd quartile}}$

We use standard linearizing techniques to determine values of the quartiles.

$$\tilde{Q}_i = q_i + \frac{h}{f_{i+1}} \left(\frac{u}{h} - F_i \right)$$

$$\tilde{Q}_j = q_j + \frac{h}{f_{j+1}} \left(\frac{3u}{h} - F_j \right)$$

⊗ Box and whisker plot (?) →



I assume this is how one data point, with spread, is represented in a plot.

Something about 'coefficient of variation'?

Δ $V = 100 \cdot \frac{\sigma \bar{x}}{\bar{x}}$, given $\bar{x} \neq 0$

Some sort of confidence measure? Who knows.
