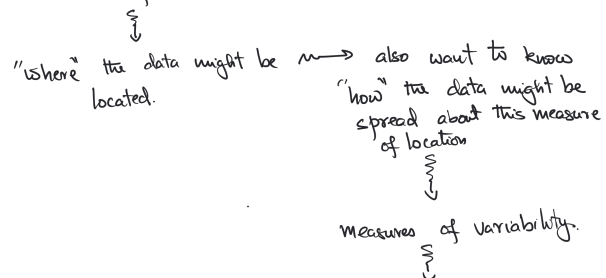


Measures of Variability

Recall: Measures of location



Note: choice of measure of variability depends on the measure of location under consideration.

I Standard deviation from mean

Suppose $\{x_1, x_2, \dots, x_n\}$ sample data. $n \rightarrow X$

$$\bar{x} = \text{sample mean} = \frac{\sum_{i=1}^n x_i}{n}$$

Want: Calculate the deviation of data points from \bar{x} .

Naive Approach:

Step 1: Calculate how each data point varies from \bar{x} .

$x_1 - \bar{x}$
 $x_2 - \bar{x}$
 $x_3 - \bar{x}$
 \vdots
 $x_n - \bar{x}$

Step 2: Add all these individual variations.

i.e. total variation = $(x_1 - \bar{x}) + (x_2 - \bar{x}) + (x_3 - \bar{x}) + \dots + (x_n - \bar{x})$.

Problems:

- Want the measure of spread to be positive.
- positive $(x_i - \bar{x})$ will cancel contribution of negative $(x_i - \bar{x})$.
 can give the impression of small variation if x_i 's are symmetrically distributed about \bar{x} .

Solution:

Step 2 corrected: Add the squares of individual variations.

i.e.

$$\text{Variation} = (x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + (x_3 - \bar{x})^2 + \dots + (x_n - \bar{x})^2$$

Still not a good measure!

If $|X|$ = size of data set $n \rightarrow$ can lead to larger value for Variation(X).

Solution: Divide by the size of the data set.

$$\begin{aligned} \text{Sample Variance } s^2 &= \frac{1}{n-1} ((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2) \\ &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \end{aligned}$$

Note: (n-1) instead of n

\downarrow
to make s^2 unbiased estimator for population variance (denoted by σ^2)

Units of s^2 are squared of units of x_i 's.
 \downarrow take square root

$$\text{Sample Std deviation } s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$$

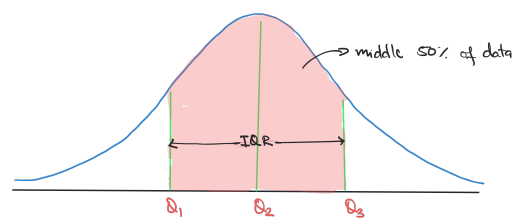
s^2 is a good measure of variability if data is

- symmetric
- unimodal
- does not have outliers

II Interquartile range for Median

Recall: $Q_1 \rightarrow 25^{\text{th}}$ percentile
 $Q_2 \rightarrow 50^{\text{th}}$ "
 $Q_3 \rightarrow 75^{\text{th}}$ percentile.

Interquartile := $Q_3 - Q_1$ = the range of the middle 50% of data.



Note: Like the median, IQR is robust to minor changes in data and to outliers.

III Summary Statistics

Location	Variation/spread
① Sample mean	① sample variance / std dev
② Sample median	② IQR
③ Quartiles: Q_1, Q_2, Q_3	③ range = $\max(X) - \min(X)$

IV Outliers

Define

upper fence = largest data value less than $Q_3 + 1.5 \times \text{IQR}$.

lower fence = smallest data value larger than $Q_1 - 1.5 \times \text{IQR}$.

An outlier \rightarrow any data value outside of the fences.

i.e. larger than upper fence or smaller than lower fence.

V Box plots

pictorial representation of summary statistics.

