

Week 2

Monday, 2 August 2021 23:36

BIOSTATISTICS

CONTINUOUS DATA MEASURES

NUMERICAL WAYS TO SUMMARIZE
CHARACTERISTICS OF SAMPLE

- MEASURE OF TENDENCY
- MEASURE OF VARIABILITY
- MEASURE OF LOCATION IN DISTRIBUTION
- COMPARING DISTRIBUTIONS
- DIFFERENCE IN SAMPLE SIZES

• USEFUL SUMMARY STATISTICS

CENTRAL TENDENCY MEDIAN, MEAN

VARIABILITY STANDARD DEVIATION

LOCATION PERCENTILES

$$\text{MEAN} = \frac{\sum_{i=1}^N x_i}{N} = \bar{x}$$

\bar{x} SAMPLE NOT μ POPULATION MEAN

μ CAN ONLY BE ESTIMATED BY \bar{x}

MEAN IN SMALL SAMPLE IS SENSITIVE TO EXTREME VALUE
JUST ONE VALUE CAN VARY IT BY A LOT.

MEDIAN MIDDLE VALUE IN ORDERED SET
OF CONTINUOUS DATA

50% PERCENTILE, II QUANTILE
LESS SENSITIVE TO EXTREMES

ONLY AFFECTED BY RELATIVE
POSITIONS OF THE VALUES

$$\text{MEDIAN} = \begin{cases} x_{(N/2)} & \frac{N}{2} \in \mathbb{N} \\ \frac{x_{\lfloor N/2 \rfloor} + x_{\lfloor N/2 \rfloor + 1}}{2} & \frac{N}{2} \notin \mathbb{N} \end{cases}$$

$$\text{VARIANCE} = \frac{\sum_{i=1}^N (x_i - \bar{x})^2}{N - 1} = S^2$$

MINIMAL IMPACT
IMPORTANT TO CORRECT SINCE
WE DON'T KNOW μ

CUMULATIVE SQUARE DISTANCE AVERAGED

STANDARD
DEVIATION $S = \sqrt{S^2}$

HOW FAR ON AVERAGE EVERY SINGLE
OBSERVATION IS FROM THE SAMPLE MEAN

$S = 0$ NO VARIABILITY $x_1 = x_2 = \dots = x_N$

S^2 BEST ESTIMATE OF σ^2

PERCENTILES $\gamma \in [0, 1]$

γ TH PERCENTILE THAT HAS
 $\gamma\%$ OF DATA LESS OR EQUAL

AND 1- $\delta\%$ MORE THAN VALUE

$$q_\delta = \begin{cases} \chi(m_\delta) & \text{IF } m_\delta \in \mathbb{N} \\ \frac{\chi(\lfloor m_\delta \rfloor) + \chi(\lfloor m_\delta \rfloor + 1)}{2} & \text{IF } m_\delta \notin \mathbb{N} \end{cases}$$

BEST ESTIMATES OF UNKNOWN UNDERLYING POPULATION VALUES

$$\begin{array}{l} \overline{X}_N \longrightarrow \mu \\ S_N^2 \longrightarrow \sigma^2 \\ N \longrightarrow \infty \end{array}$$

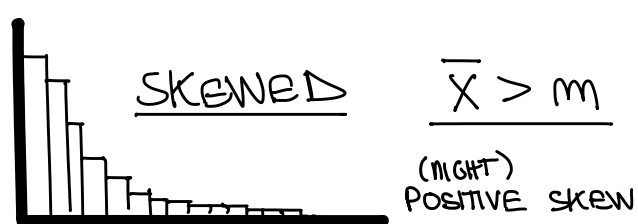
VISUAL DISPLAYS

HISTOGRAM } VISUALIZING DISTRIBUTION
BOX PLOT }

USEFUL TO COMPARE DATA FROM DIFFERENT SAMPLES

COUNTING # OF DATA IN A PREDEFINED RANGE

$\bar{X} \approx m$ NO SPECIAL TENDENCY (SYMMETRY) ALONG CENTER
MEAN \approx MEDIAN



$$\begin{aligned} \text{IQR} &= Q_3 - Q_1 & \text{MAX} &= Q_3 + \frac{3}{2} \text{IQR} \\ & & \text{MIN} &= Q_1 - \frac{3}{2} \text{IQR} \end{aligned}$$

INTER QUANTILE RANGE

SYMMETRIC, BELL-SHAPED



THE ROLE OF SAMPLE SIZE

RANDOM SAMPLE MIMICS CHARACTERISTICS OF THE UNDERLYING POPULATION

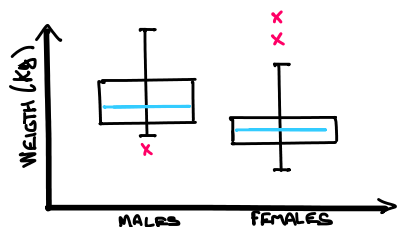
OF THE LARGER POPULATION

LARGER SAMPLE \rightarrow MINIMIZING VARIANCE

COMPARING CONTINUOUS DISTRIBUTIONS

VISUALLY

SIDE BY SIDE BOXPLOTS ARE EASIER TO INTERPRET THAN HISTOGRAMS



NUMERICALLY

$$\mu_{\text{MALES}} - \mu_{\text{FEMALES}} = \text{AVERAGE SHIFT}$$

STUDYING THE DIFFERENCE IN MEANS

$$\delta = \mu_1 - \mu_2 \approx \bar{X}_1 - \bar{X}_2$$

ON AVERAGE 1 DIFFERS FROM 2 BY δ (kg)

SKewed DISTRIBUTIONS ARE DIFFICULT TO ANALYZE EASILY

"STUDYING THE SHIFT OF MASS"

IF MORE THAN ONE GROUP WE CAN PICK ONE AS THE 'REFERENCE'

IMPORTANT

μ, δ PARAMETERS OF POPULATION
CAN'T BE DERIVED DIRECTLY.

\bar{X}_N, S_N ARE CALCULATED BASED ON THE CHOSEN SAMPLE (IMPERFECT REPRESENTATIONS)

\bar{X}_N, S_N, m CHANGE DEPENDING ON SAMPLE

THERE IS NO REAL WAY TO PREDICT HOW THEY WILL CHANGE GIVEN A DIFFERENT SAMPLE

S_N DECREASES IF N INCREASES \leftarrow
(ONLY WORKS IF $N_0 < 30$)