

REPORTS - Data generation philosophy

In some reports, you will **generate your own datasets**. This is a deliberate pedagogical choice. This approach helps you think more carefully about how different statistical metrics and methods behave and what they actually measure.

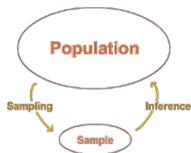
1

Statistics is about drawing conclusions from incomplete information. While errors can happen, statistical methods help us estimate the associated risks and uncertainties.



2

Statistics is based on samples!



The primary goal of statistics is to estimate (infer) an unknown characteristic of an entire population based on sample data.

These estimates, often derived from descriptive statistics, are then used to make informed decisions about the population.

3

Describing data

Samples and populations typically consist of numerous individual observational units along with their associated data (observations, variables).

To describe samples effectively, we use summary statistics (mean, median, variance, etc.), which serve as estimates of the corresponding statistics for the entire population.



4

Describing data

Samples and populations typically consist of numerous individual observational units along with their associated data (observations, variables).

To describe samples effectively, we use summary statistics (mean, median, variance, etc.), which serve as estimates of the corresponding statistics for the entire population.

Today: data summaries for each variable (separately).

Individual	Weight (kg)	Height (cm)
1	75.5	172
2	55.3	152
3	61.2	164
4	50.3	148
5	99.4	192
6	66.2	171
7	75.3	169
8	74.6	182
9	60.5	162
10	93.5	184
11	73.6	169



5

The primary goal of statistics is to **estimate (infer)** an unknown characteristic (e.g., height) of an entire population based on sample data.

We want to know about a large number of trees



Population

Population mean height (here the parameter of interest, i.e., unknown quantity)

Selected trees to measure height



Sample

Sample mean height

→
←

Inference

Inspired by <https://www.diffrnotes.com/study-guides/statistics/sampling/populations-samples-parameters-and-statistics>

6

Key Learning Objectives today

1. Differentiate between measures of location and measures of spread (or width).
2. Recognize that variability is not merely noise, but a fundamental parameter that can be estimated and interpreted.
3. Become familiar with the most commonly used descriptive statistics.
4. Understand when the mean or the median provides a more appropriate summary of location.
5. Summarize single variables using both location and spread measures (with extensions to multiple variables later in the course).

7

Scientific question: Did humans drive mammal extinctions in Australia?

↓

Statistical question: Are "victims" bigger than "survivors" and historical extinctions?

Australia

Number of Species vs Log Mass (g)

Legend: Survivors (grey), Victims (red), Historical Extinctions (blue)

Frequency distribution of mammal mass categorized into survivors, "victims" and older (historical) extinctions

Survivors (extant species, i.e., alive today).

Victims (late Pleistocene, i.e., past 50 000 years, 50 ka).

Historical extinctions (older than 50 ka) are based on samples (fossils).

We want to make inferences about all past and present mammals in Australia (i.e., statistical population are all mammal species, past or present, in Australia).

Study by Lyons et al. (2004; Evolutionary Ecology Research 6:339-358)

ka = kiloannus (1000); ~ 50 ka = "behavioural modernity" in humans.

8

Descriptive statistics or summary statistics are needed to make inferences

- **Location** tells us something about the average or typical individual units (i.e., where the observations are centered).
- **Spread** tells us how measurements vary among individual units (or observations), i.e., how widely scattered the values are around the center (location).

Australia

Number of Species vs Log Mass (g)

Legend: Survivors (grey), Victims (red), Historical Extinctions (blue)

location (vertical arrows)

spread (horizontal arrows)

Remember the jargon (lecture 2):

Individual units (of data) are called observation units (here each observational unit is a single species).

Study by Lyons et al. (2004; Evolutionary Ecology Research 6:339-358)

9

The most important location statistic: Arithmetic mean



"Flying" paradise tree snake (*Chrysopelea paradisi*). To better understand how lift is generated, Socha (2002) videotaped glides (from a 10-m tower) of 8 snakes. Rate of side-to-side undulation was measured in hertz (number of cycles per second). The values recorded were:

0.9, 1.4, 1.2, 1.2, 1.3, 2.0, 1.4, 1.6

The arithmetic mean is an algorithm = a process or set of rules to be followed in calculations - sum of all the observations in a sample divided by n , the number of observations.

$$\bar{Y} = \frac{0.9 + 1.2 + 1.2 + 2.0 + 1.6 + 1.3 + 1.4 + 1.4}{8} = 1.375 \text{ Hz.}$$

The sample mean is represented most often as Y or X said « Y bar » or « X bar »

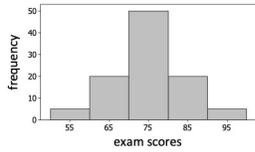
10

The concept of spread around the mean

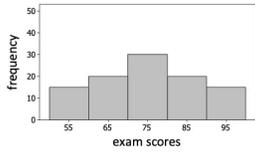
Deviation (from the mean) measures uncertainty. When deviation is small, values cluster near the mean, making the mean a good guess for a randomly chosen observation. When deviation is large, values are more spread out, increasing uncertainty in any guess.

Which class has the most variation in exam scores?

Class 1



Class 2



Note: scales (X and Y axis limits) are exactly the same

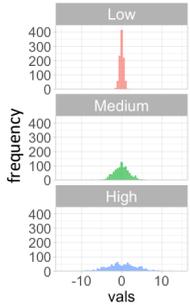
- Source: Cooper & Shore; Journal of Statistics Education (vol. 18, #2)

11

The concept of spread around the mean

Population variability should not be dismissed as mere noise around the mean; it holds biological significance in its own right.

As location values (e.g., mean, median), variation (spread) also has an inherent true value within a population (a parameter) that we aim to estimate through sampling.



frequency

400
300
200
100
0

vals

-10 0 10

Modified from: © 2020 W.H. Freeman and Company

12

The most important spread statistics: variance and standard deviation (the accompanying statistics of spread for the mean)

It indicates how far the different measurements typically are from the mean. The standard deviation is large if most observations are far from the mean, and it is small if most measurements lie close to the mean.

Quantities needed to calculate the standard deviation and variance of snake undulation rate ($\bar{Y} = 1.375$ Hz).

Observations (Y_i)	Deviations ($Y_i - \bar{Y}$)	Squared deviations ($(Y_i - \bar{Y})^2$)
0.9	-0.475	0.225625
1.2	-0.175	0.030625
1.2	-0.175	0.030625
1.3	-0.075	0.005625
1.4	0.025	0.000625
1.4	0.025	0.000625
1.6	0.225	0.050625
2.0	0.625	0.390625
Sum	0.000	0.735

Whitlock & Schluter, The Analysis of Biological Data, 3e © 2020 W. H. Freeman and Company

13

The most important spread statistics: variance and standard deviation (the accompanying statistics of spread for the mean)

Important measure: "Sum of Squared deviations from the mean"

Observations (Y_i)	Deviations ($Y_i - \bar{Y}$)	Squared deviations ($(Y_i - \bar{Y})^2$)
0.9	-0.475	0.225625
1.2	-0.175	0.030625
1.2	-0.175	0.030625
1.3	-0.075	0.005625
1.4	0.025	0.000625
1.4	0.025	0.000625
1.6	0.225	0.050625
2.0	0.625	0.390625
Sum	0.000	0.735

variance

$$s^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n - 1}$$

$s^2 = \frac{0.735}{7} = 0.11 \text{ Hz}^2$

standard deviation

$$s = \sqrt{\frac{\sum (Y_i - \bar{Y})^2}{n - 1}} = \sqrt{\frac{0.735}{7}} = 0.324037 \text{ Hz}$$

Variance is the average squared deviation of observations from the mean, measuring overall uncertainty or spread (units, here Hz, are squared, i.e., Hz²)

Square root of the variance (in the same unit as the original variable).

14

The most important spread statistics: variance and standard deviation

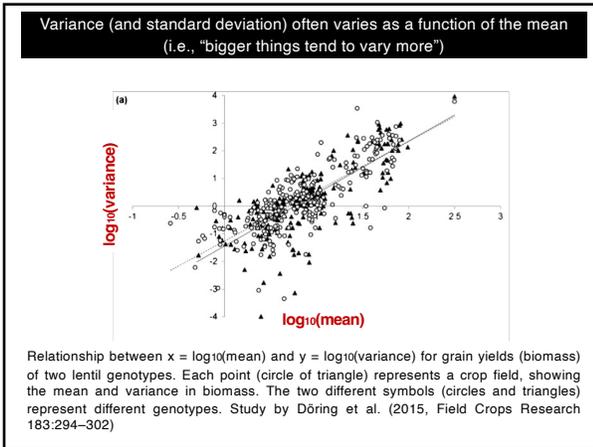
Why is the sum of the squared deviations from the mean divided by $n-1$ and not n (number of observations)?

variance

$$s^2 = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n - 1}$$

We will understand this in a couple of lectures!

15



16

A relative metric of spread: the coefficient of variation (CV) often important when comparing groups of individuals belonging to different classes or variables with different units.

CV tells you how *hard it is to guess* a typical value **relative to the size of what you're guessing**. A small CV means most values fall close to the target relative to its magnitude; a large CV means guesses are proportionally more uncertain.

coefficient of variation: Snake undulation data:

$$CV = \frac{s}{\bar{Y}} \times 100\% \quad CV = \frac{0.324}{1.375} 100\% = 24\%$$

CV answers: Are individuals similarly variable relative to what is typical (average)? even when values differ greatly in magnitude or units.

A 5-gram deviation matters a lot for a mouse but very little for an elephant. CV captures this by measuring variation relative to average size, not in absolute grams.

17

A relative metric of spread: the coefficient of variation (CV) often important when comparing groups of individuals belonging to different classes or variables with different units.

Making the coefficient of variation (CV) more obvious!

				X	s	CV
1	2	3	4	2.5	1.29	51.7%
31	32	33	34	32.5	1.29	3.97
204	205	206	207	205.5	1.29	0.63
1300	1301	1302	1303	1301.5	1.29	0.10

CV answers: Are individuals similarly variable relative to what is typical (average)? even when values differ greatly in magnitude or units.

18

When using the coefficient of variation, it is important to be clear about why variation (measured by the standard deviation) is being expressed relative to the mean.

Consider the case of species fluctuations in abundance through time:

Statistics	Species A	Species B
Mean (\bar{X})	1008.2	270.8
Standard deviation (s)	104.9	103.9
Coefficient of variation (CV)	10.4	38.4

Time Series of Abundance for Two Species

Species A: More abundance (mean) but the same variation (s) through time as species B.

$CV_B > CV_A =$ greater risk of extinction of Species B.

Species B: Less abundance (mean) but the same variation (s) through time as species A.

19

Before we go too far: a word on rounding numerical values

- When recording data, always retain as many significant digits (often involving decimals places) as your calculator or computer can provide.
- When presenting results, however, numbers should be rounded before being presented.
- There are no strict rules on the number of digits that should be retained when rounding.
- A common strategy, is to round descriptive statistics (e.g., means, standard deviations) to one decimal place MORE than the original measurements themselves.

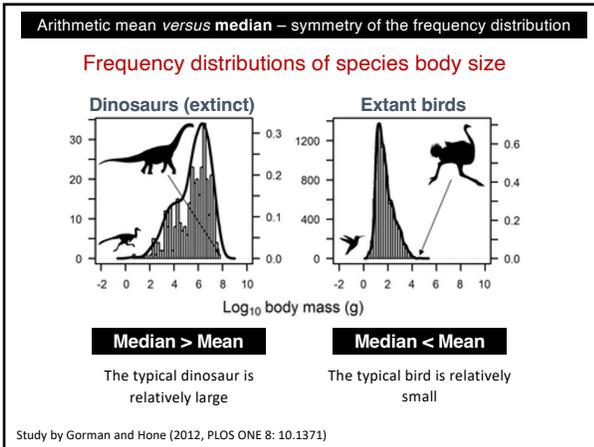
Example: the mean rate of undulation for the eight snakes (measured with a single decimal place; e.g., 0.9), calculated as 1.375 Hz, would be communicated as:

0.9, 1.4, 1.2, 1.2, 1.3, 2.0, 1.4, 1.6 $\bar{Y} = 1.38$ Hz.

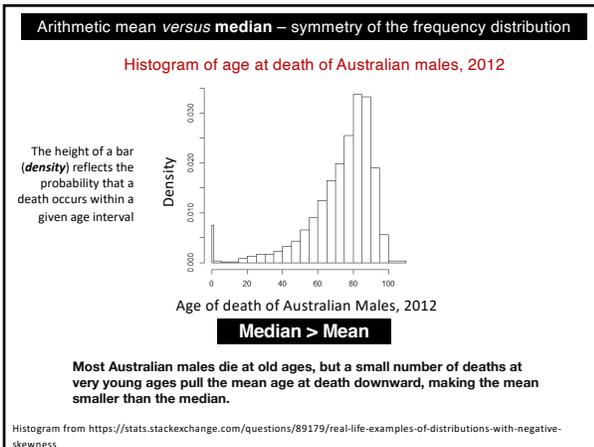
20

Let's take a break – 1 minute

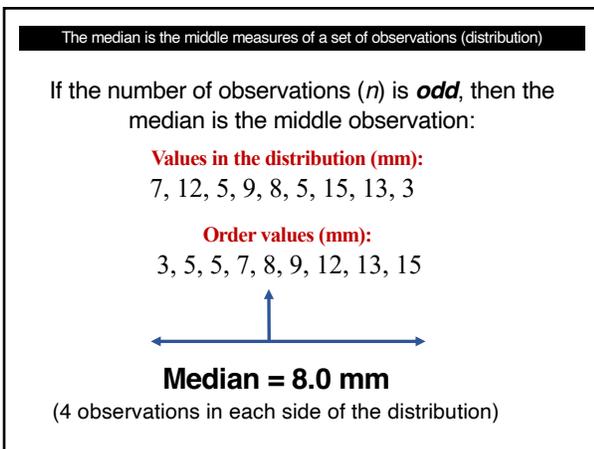
21



25



26



27

The median is the middle measures of a set of observations (distribution)

If the number of observations (n) is **odd**, then the median is the middle observation:

$$\text{Median} = Y_{((n+1)/2)} = Y_{((9+1)/2)} = Y_5$$

Ordered values (mm):
3, 5, 5, 7, 8, 9, 12, 13, 15

Median = $Y_5 = 8.0$ mm
(4 observations in each side of the distribution)

28

The median is the middle measures of a set of observations (distribution)

If the number of observations (n) is **even**, then the median is calculated differently:

It gives an "arm" (or a pedipalp that stores sperm) for a female spider.

Running speed (cm/s) of male *Tidarren* spiders before and after voluntary amputation of one pedipalp.

Spider	Speed before	Speed after	Spider	Speed before	Speed after
1	1.25	2.40	9	2.98	3.70
2	2.94	3.50	10	3.55	4.70
3	2.38	4.49	11	2.84	4.94
4	3.09	3.17	12	1.64	5.06
5	3.41	5.26	13	3.22	3.22
6	3.00	3.22	14	2.87	3.52
7	2.31	2.32	15	2.37	5.45
8	2.93	3.31	16	1.91	3.40

Oxyopes sallicus

29

The median is the middle measures of a set of observations (distribution)

Spider	Speed before	Speed after	Spider	Speed before	Speed after
1	1.25	2.40	9	2.98	3.70
2	2.94	3.50	10	3.55	4.70
3	2.38	4.49	11	2.84	4.94
4	3.09	3.17	12	1.64	5.06
5	3.41	5.26	13	3.22	3.22
6	3.00	3.22	14	2.87	3.52
7	2.31	2.32	15	2.37	5.45
8	2.93	3.31	16	1.91	3.40

For an **even** number of observations, the median is the average of the two central numbers. $n = 16$ in this study.

Median (speed before) = $M = 2.90$ cm/s

1.25 1.64 1.91 2.31 2.37 2.38 2.84 2.87 2.93 2.94 2.98 3.00 3.09 3.22 3.41 3.55

Median = $[Y_{(n/2)} + Y_{(n/2+1)}] / 2$.

Median = $(2.87 + 2.93) / 2 = 2.900$ cm/s

30

Arithmetic mean *versus* median – the second most common statistic to describe the location of a frequency distribution

					X	Median
1	2	3	4	5	3	3
1	2	3	4	489	99.8	3
1	2	3	4	6	3.2	3

Highlighting how extreme values have a greater impact on the mean than on the median!

31

Arithmetic mean *versus* median

The arithmetic mean is affected by the imbalance (asymmetry) in the distribution caused by the presence of extreme values.

Uniform: Mean equals the median
 Bell-shaped: Mean equals the median
 Asymmetric (skewed): Mean smaller than the median
 Bimodal: Mean greater than the median

Symmetric distributions: Mean Median
 Asymmetric distributions: Median Mean

32

Arithmetic mean *versus* median

The arithmetic mean is affected by the imbalance (asymmetry) in the distribution caused by the presence of extreme values.

Asymmetric distributions (skewed) can be either left or positive skewed

The rule comparing the mean to the median is particularly effective for large datasets (more than 30 observations).

Left (or Negative) skewed: Mean Median
 Right (or Positive) skewed: Median Mean

33

Mean or Median? Consider Skewness

- Few small values.
- > 1/2 of values exceed the mean.

- As many large as small values.
- ~ 1/2 of values exceed the mean.

- Few large values.
- > 1/2 of values are less than the mean.

Left Skewed (asymmetric) Not Skewed (symmetric) Right Skewed (asymmetric)

density

vals

For measures like income, **medians** are generally preferable to **means**. This is because income distributions are right-skewed (a few individuals are extremely wealthy), and we are usually more interested in what a typical ("average") person earns than in the arithmetic average of all incomes.

© 2020 W.H. Freeman and Company
