# 2. Measurement errors; statistical estimators

#### "Errors using inadequate data are much less than those using no data at all"

Charles Babbage

# Example

- Take one cuvette with bacterial culture
- Measure optical density (OD600)
- Result: 0.37
- Reading error
- Take five cuvettes and find mean OD600
- Results 0.42
- Sampling error
- These are examples of measurement errors





# Measurement errors

#### Systematic and random errors

#### Systematic errors

your mistakes

- Incorrect instrument calibration
- Change in experimental conditions
- Pipetting errors

#### **Random errors**

statistics sucks

- Reading errors
- Sampling errors
- Intrinsic variability



# YOU NEED REPLICATES

# Reading error

- The reading error is ±half of the smallest division
- Example: 23±0.5 mm from a ruler

- Beware of digital instruments that sometimes give readings much better than their real accuracy
- Read the instruction manual!

 Reading error does not take into account biological variability



### Random measurement error

- Determine the strength of oxalic acid in a sample
- Method: sodium hydroxide titration
- Uncertainties contributing to the final result
  - volume of the acid sample
  - judgement at which point acid is neutralized
  - volume of NaOH solution used at this point
  - accuracy of NaOH concentration
    - weight of solid NaOH dissolved
    - volume of water added



- Each of these uncertainties adds a random error to the final result
- Measurement errors are normally distributed

## Counting error

- Dilution plating of bacteria
- Found C = 10 colonies
- Counting statistics: Poisson distribution

 $\sigma=\sqrt{\mu}$ 

 Use standard deviation as error estimate to obtain the *standard error of the count*

$$S = \sqrt{C} = \sqrt{10} \approx 3$$

 $C = 10 \pm 3$ 



# Counting error

- Gedankenexperiment
- Measure counts on 10,000 plates

| $C_i$               | Count from plate <i>i</i> |
|---------------------|---------------------------|
| $S_i = \sqrt{C_i}$  | lts error                 |
| μ                   | Unknown population mean   |
| $\sigma=\sqrt{\mu}$ | Unknown population SD     |

- Counting errors, S<sub>i</sub>, are similar, but not identical, to σ
- $C_i$  is an estimator of  $\mu$
- $S_i$  is an estimator of  $\sigma$



#### Exercise: is Dundee a murder capital of Scotland?

- On 2 October 2013 *The Courier* published an article "Dundee is murder capital of Scotland"
- Data in the article (2012/2013):

| City      | Murders | Per 100,000 |
|-----------|---------|-------------|
| Dundee    | 6       | 4.1         |
| Glasgow   | 19      | 3.2         |
| Aberdeen  | 2       | 0.88        |
| Edinburgh | 2       | 0.41        |

- Compare Dundee and Glasgow
- Find errors on murder rates
- Hint: find errors on murder count first

### Exercise: is Dundee a murder capital of Scotland?



- and apply them to murder rate
  - $\Delta R_D = 4.1 \times 0.41 = 1.7$  $\Delta R_G = 3.2 \times 0.23 = 0.74$

#### Exercise: is Dundee a murder capital of Scotland?

| City      | Murders | Per 100,000 |
|-----------|---------|-------------|
| Dundee    | 6       | 4.1         |
| Glasgow   | 19      | 3.2         |
| Aberdeen  | 2       | 0.88        |
| Edinburgh | 2       | 0.41        |

95% confidence intervals (Lecture 4) p-values from chi-square test vs Dundee



#### Measurement errors: summary

- Random measurement errors are expected to be normally distributed
- Some errors can be estimated directly

   reading (scale, gauge, digital read-out)
   counting
- Other uncertainties require replicates (a sample)
   this introduces sampling error

# Population and sample

#### Population and sample





- Terms nicked from social sciences
- Most biological experiments involve sample selection
- Terms "population" and "sample" are not always literal

# What is a sample?

 The term "sample" has different meanings in biology and statistics

 Biology: sample is a specimen, e.g., a cell culture you want to analyse



• In these talks:  $x_1, x_2, \dots, x_n$ 



#### Population and sample



A **parameter** describes a population

#### A statistical estimator

(statistic) describes a sample

A statistical estimator approximates the corresponding parameter

# Sampling uncertainty



#### Sample size

Dilution plating experiment



What is the sample size?

n = 1

This sample consists of one measurement:  $x_1 = 10$ 

#### 10 colonies

# Statistical estimators

"The average human has one breast and one testicle"

Des MacHale

#### What is a statistical estimator?



"Right and lawful rood<sup>\*</sup>" from *Geometrei*, by Jacob Köbel (Frankfurt 1575) Stand at the door of a church on a Sunday and bid 16 men to stop, tall ones and small ones, as they happen to pass out when the service is finished; then make them put their left feet one behind the other, and the length thus obtained shall be a right and lawful rood to measure and survey the land with, and the 16th part of it shall be the right and lawful foot.

Over 400 years ago Köbel:

- introduced random sampling from a population
- required a representative sample
- defined standardized units of measure
- used 16 replicates to minimize random error
- calculated an estimator: the sample mean

# Example

- Weight of 7 mice
- This is a sample
- We can find
  - □ mean = 19.2 g
  - □ median = 18.7 g
  - $\square$  standard deviation = 4.4 g
  - $\square$  standard error = 1.7 g
  - $\Box$  interquartile range = 6.0 g
- These are examples of statistical estimators



| No. | Weight (g) |
|-----|------------|
| 1   | 13.6       |
| 2   | 16.1       |
| 3   | 25.1       |
| 4   | 24.8       |
| 5   | 16.6       |
| 6   | 19.8       |
| 7   | 18 7       |

#### Statistical estimators

 Statistical estimator is a sample attribute used to estimate a population parameter

From a sample  $x_1, x_2, \dots, x_n$  we can find

$$M = \frac{1}{n} \sum_{i=1}^{n} x_i \qquad \text{mean}$$

$$SD = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - M)^2}$$

standard deviation

median, proportion, correlation, ...



#### Mean vs median

#### Median

- More appropriate for skewed distributions
- Not sensitive to outliers

#### Mean

- Better estimate of the central value
- Statistical tests on the mean (e.g. t-test) are more power full than non-parametric tests

If your data are symmetric, use mean



# Standard deviation

- Standard deviation is a measure of spread of data points
- Idea:
  - calculate the mean
  - $\hfill\square$  find deviations from the mean
  - $\hfill\square$  get rid of negative signs
  - $\hfill\square$  combine them together



#### Standard deviation

 Standard deviation is a measure of spread of data points

Idea:

calculate the mean

□ find deviations from the mean□ get rid of negative signs

□ combine them together

Standard deviation of x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>n</sub>

$$SD_n = \sqrt{\frac{1}{n} \sum_{i} (x_i - M)^2}$$



$$SD_{n-1} = \sqrt{\frac{1}{n-1}\sum_{i}(x_i - M)^2}$$
  $\qquad SD_{n-1}^2$  estimates true variance better than  $SD_n^2$ 

# Sampling distribution

Population of mice with normal body weight:  $\mu = 20$  g,  $\sigma = 5$  g Draw lots of samples of size n = 5



#### Central limit theorem



#### Hypothetical experiment

- 100,000 samples of 5 mice
- Build a distribution of sample means
- Width of this distribution is the true uncertainty of the mean

$$\sigma_m = \frac{\sigma}{\sqrt{n}} = 2.2 \text{ g}$$

#### **Real experiment**

- 5 mice
- Measure body mass:

7.9, 14.4, 16.4, 21.7, 22.8 g

Find standard error

 $SE = \frac{SD}{\sqrt{n}} = 2.7 \text{ g}$ 

SE is an approximation of  $\sigma_m$ 



#### Hypothetical experiment

- 100,000 samples of 30 mice
- Build a distribution of sample means
- Width of this distribution is the true uncertainty of the mean

$$\sigma_m = \frac{\sigma}{\sqrt{n}} = 0.9 \text{ g}$$

#### **Real experiment**

- 30 mice
- Measure body mass:

9.9, 14.9, ..., 33.8 g

Find standard error

$$SE = \frac{SD}{\sqrt{n}} = 0.87 \text{ g}$$

#### SE is an approximation of $\sigma_m$





#### Standard deviation and standard error

| Standard deviation   | Standard error   |  |
|--|--|--|
| $SD = \sqrt{\frac{1}{n-1}\sum_{i}(x_i - M)^2}$                                 | $SE = \frac{SD}{\sqrt{n}}$   |  |
| Measure of dispersion in the sample  | Error of the mean  |  |
| Estimates the true standard deviation in the population, $\boldsymbol{\sigma}$ | Estimates the width (standard deviation) of the distribution of the sample means |  |
| Does not depend on sample size   | Gets smaller with increasing sample size   |  |

#### Correlation coefficient



• Two samples:  $x_1, x_2, \dots, x_n$  and  $y_1, y_2, \dots, y_n$ 

$$r = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{x_i - M_x}{SD_x} \right) \left( \frac{y_i - M_y}{SD_y} \right) = \frac{1}{n-1} \sum_{i=1}^{n} Z_{xi} Z_{yi}$$

where Z is a "Z-score"

#### Correlation doesn't mean causation!

r = 0.993



tylervigen.com

tylervigen.com

#### Statistical estimators

#### Central point

Mean

Geometric mean Harmonic mean **Median** 

Mode Trimemod mod

Trimmed mean

#### Dispersion

Variance Standard deviation Standard error Mean deviation Range Interquartile range Mean difference

#### Symmetry

Skewness

**Kurtosis** 

#### Dependence

**Pearson's correlation** Rank correlation

Distance

Hand-outs available at https://dag.compbio.dundee.ac.uk/training/Statistics\_lectures.html