# Experimental design
# Statistical analysis and Interpretation of results

Konstadia (Dina) Lika

Biology Department

Univ of Crete

lika@uoc.gr

ΠΑΝΕΠΙΣΤΗΜΙΟ
ΚΡΗΤΗΣ

UNIVERSITY
OF CRETE

# Experiments using laboratory animals

- Well designed

- Efficiently executed

- Correctly analyzed

- Clearly presented

- Correctly interpreted

**Experiment**:
scientific procedure undertaken
- to test a hypothesis
  (confirmatory research) or
- to provide material for the
  generation of new hypotheses
  (exploratory research)

# Ethical consideration

Principles of the 3Rs (Russel & Burch, 1959)

- ## Replacement
    - Replace animals by less sentient ones or in vitro methods whenever possible

- ## Refinement
    - Experimental protocols should be refined to minimize adverse effects

- ## Reduction
    - Keep the number of animals to the minimum necessary to achieve the aims of the research
    - BUT NOT so few
        - miss biologically important effects or
        - require unnecessary repetition of experiments or
        - the results become statistically invalid

**The # of animals used can be reduced by**
- **good exp. design**
- **appropriate statistical procedures**

# Research description

Clearly state:

- the objectives of the research and the hypotheses to be tested

- the reason for choosing the particular animal model
  - the species, strain, source, and type of animal used

- the details of each experiment
  - experimental design
  - number of animals

- the statistical methods used for analysis

# The experimental unit (EU)

- the physical entity which can be assigned, at random, to a treatment

- It receives one level of treatment independently of other units

- Measurements from different units are assumed to be independent

- The number of experimental units determines the sample size

- commonly (not always) it is an individual animal

- it is also the unit of statistical analysis

# Example- mice in a cage

- 6 cages, 5 mice/cage
- 2 treatments, each applied to 3 cages



What is the EU? A mouse or a cage?

Answer: <u>A cage</u> of animals rather than the

individual animal

- because treatments are applied to whole cages
- mice in the cage can not have different treatments

# Principles of experimental design

- ## Replication
  - reduce variability and improve statistical reliability

- ## Randomization
  - avoid bias

- ## Controls
  - minimize confounding effects
  - measure treatment effects accurately

- ## Blinding
  - reduce bias and ensure objectivity
    - Animals, samples and treatments should be coded until the data are analyzed

# Replication

- Replication means applying each treatment to more than one independent experimental unit (EU).

- A replicate is an independent EU that receives the same treatment.

- Why replication matters
  - allows estimation of variability in the population
  - provides more *reliable* estimates of treatment effects

What is true replication in the mice in the cage example?

n=3

m=15 mice (pseudo-replication)

# Pseudo-replication

- Pseudo-replicates are not independent, while standard statistical tests assume independence

- This leads to incorrect estimates of variability and can invalidate your results

- Pseudo-replication can either be:
  - Temporal: repeated measures over time from the same mice or cage
  - Spatial: several measurements from the same vicinity or cage

- Solution: Use appropriate statistical methods (e.g., mixed models, repeated measures ANOVA) that account for dependence

# Randomization

- Randomization assigns each EU an equal chance of receiving any treatment

- It helps to:
  - Eliminate experimenter bias
  - Produce unbiased estimates of population parameters (e.g., means, treatment effects)
  - Ensure valid and reliable statistical inference

- Use random number generators (e.g., from statistical software) to assign EUs to treatment groups

# Experimental Size

How many EU (e.g., tanks, animals) do I need?

- Too few → may miss real effects (underpowered study)
- Too many → may waste resources and animals unnecessarily

**Goal:** Use **just enough** animals to detect a meaningful effect—if it exists—while minimizing waste.

How to Decide?

- **Power analysis** is the most common method for calculating appropriate sample size

# Formal experimental designs

- ## Completely randomized designs (CR)
  - EUs are assigned to treatments at random

- ## Completely randomized block designs (CRB)
  - EUs are grouped into blocks based on shared characteristics (e.g., age, weight, cage, batch)
  - Randomization occurs within each block.
  - Helps reduce unexplained variation without increasing the number of animals.

If blocks explain much of the variation in the response variable, the CRB may be more powerful than the CR

# Formal experimental designs

- ## Repeated measure designs
  - each EU measured repeatedly under different treatments and/or times

- ## Factorial designs
  - more than one type of factors (e.g. drug treatment and gender)
  - crossed designs: every level of one factor crossed with every level of a second factor
  - nested designs: different (randomly chosen) levels of a factor B nested in each level of a factor A

# Statistical Analysis

General aim is to extract all useful information present in the data

- **Descriptive statistics**
  - statistical methods to summarize, describe, and explore data
  - enable us to present the data in a more meaningful way, allowing a simpler interpretation of the data
  - provide insights without generalizing beyond the sample

- **Inferential statistics**
  - the process of drawing conclusions from data about a population
  - techniques allowing the use of samples to generalize about populations

# Descriptive Statistics

## Graphical exploration of data

- graphical methods summarize the data in a diagrammatic way

- qualitative and involve a degree of subjectivity

## Non-graphical method

Include summary statistics

– **Measures of central tendency** (e.g., mean, median)

– **Measures of variability** (e.g., variance, standard deviation)

- quantitative and objective, they do not give a full picture of the data

**Non-graphical and graphical methods complement each other**

# Inferential Statistics

**Estimating population parameters**

**Point estimate**

– single value estimate of parameter

**Confidence intervals**

– A **range of values** within which the true population parameter is expected to fall with a given probability

**Hypothesis testing**

– determine if observed differences are **statistically significant** or due to random chance

**Choosing appropriate statistical tests**

# Steps in Hypothesis testing

- Specify research question and null hypothesis

  e.g.

    Question: Do female and male mice have the same metabolic rate?

    $H_0$ : there is no difference between the sexes in the mean metabolic rate

- ## Choose test statistic

  – depends on the design

- ## Collect sample data

- ## Calculate test statistic relevant to hypothesis

  – Based on your sample data and chosen test

Common error: Collecting data before:
- Clearly defining the research question and hypothesis
- Choosing the appropriate statistical test

How likely is to get this result (this value of the statistic) under the null hypothesis?

- **Determine the p-value**

  Probability that any given experiment will produce *a value of the test statistic* that is equal to the one observed in our actual experiment or something more extreme, assuming that <u>the *null hypothesis* is true</u>

  – assumptions of the test are (reasonably well) met

- **Interpretation of p-value**

  – measures the "strength of evidence" against $H_0$

  - not the probability that $H_0$ is true!

  - not the probability of making a mistake by rejecting a true $H_0$

# Decision criterion

Compare p-value to the *a priori* significant level ($\alpha$)

  – if $p < \alpha$, conclude $H_0$ is "unlikely" to be true and reject it (statistically significant result)

   the observed effect is unlikely to have occurred by chance

  – else, conclude $H_0$ is "likely" to be true and do not reject it (statistically non- significant result)

Convention sets significant level $\alpha = 0.05$ (5%)

Arbitrary: other significant levels might be valid (e.g. 0.01, 0.001)

# Decision errors

Statistical hypothesis tests can produce decision errors

## Type I error
  – rejecting a true $H_0$
  – probability of making this error is set by significance level $\alpha$

## Type II error

– not rejecting a false $H_0$

– probability of making this error can only be determined if variability and desired detectable difference is known

- The risks of these two errors are inversely related

- Increasing sample size is the best way to minimise both errors

# Power of test

- Probability of detecting an effect if it exists

- Probability of rejecting incorrect $H_0$

- Complement to a Type II error
  - If $\beta$ is the probability of making a type II error,
    $1$-$\beta$ (the power) is the probability of not making a type II error

# Power analysis

Can be used

- Design an experiment
  - determine the <u>sample size required</u> to detect an effect of a given size with a given degree of confidence

- Make *a posteriori* assessment of the usefulness of an experiment
  - determine the <u>probability of detecting an effect</u> of a given size with a given level of confidence, under sample size constraints

# Statistical power

To calculate the power of a statistical test you need to specify:

- **Effect size (ES)**
  - differences between treatments
  - large effects easier to detect

- **Background variation**
  - variation between experimental units ($s^2$)
  - greater background variability; less likely to detect effects

# Statistical power

- Sample size (n) for each treatment group
  - increasing sample size makes effects easier to detect

- Significant level ($\alpha$)
  - Probability of Type I error
  - usually set at 5%, lower values sometimes specified
  - as $\alpha$ decreases, $\beta$ increases, power (1- $\beta$) decreases

- Alternative hypothesis

# Power analysis
## *a posteriori* power

For most types of analysis (t-test, ANOVA, regression), the **power** (1-β) is:

$$(1 - \beta) \propto \frac{ESa\sqrt{n}}{s}$$

**General shape**

# Sample size determination

To determine appropriate sample size, we need to:

- solve power equation for $n$

- know background variation

  (from pilot studies/previous literature)

- know the statistical power $(1-\beta)$ we want

- know what $ES$ we wish to be able to detect if it occurs

$$\sqrt{n} \propto \frac{s(1-\beta)}{aES}$$

# Statistical Power
## Conventions and Decisions

- Acceptable risk of a Type II error is often set at 1 in 5, i.e., a probability of 0.2 ($\beta$)

- "adequate" statistical power is therefore set at
  $1-\beta$ = 1 - 0.2 = 0.8

# Effect Size

- While **power (1-$\beta$)** and **significance level ($\alpha$)** are set irrespective of the data, the effect size is a property of the sample data

- *ES* formulae depend on statistical test

- Depending on the actual test, the *ES* may be expressed as
  - *d* (difference between two means),
  - *r* (correlation between two variables)
  - *f* (ANOVA test)
  - any other index related to specific test

Cohen, J., (1977). *Statistical power analysis for the behavioural sciences*. San Diego, CA: Academic Press.
Cohen, J., (1992). A Power Primer. *Psychological Bulletin* **112** 155-159.

# Calculating Cohen's *d*

Effect size $\qquad d = \dfrac{\bar{x}_1 - \bar{x}_2}{s_{Pooled}}$

$$s_{Pooled} = \sqrt{\dfrac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

$\bar{x}$: mean

$s$: standard deviation

$n$: sample size

Subscript refers to the two conditions compared

# How Do We Measure Effect Size?

- Use background information in the form of <u>preliminary/trial</u> data to get means and variation, then calculate effect size directly

- Use background information in the form of <u>similar studies</u> to get means and variation, then calculate effect size directly

- With no prior information, make an estimated guess on the expected effect size
  - Broad effect sizes categories are small, medium, and large
  - Different statistical tests will have different values of effect size for each category

# Cohen's Rules Of Thumb For Effect Size

| Effect size | Correlation coefficient | t-tests | ANOVA |
|---|---|---|---|
| "Small effect" | r = 0.1 | d = 0.2 | f=0.1 |
| "Medium effect" | r = 0.3 | d = 0.5 | f=0.25 |
| "Large effect" | r = 0.5 | d = 0.8 | f=0.4 |

Cohen's suggestions should be seen as rough guidelines.

# Required sample size for t-tests



Sample Size Estimation for independent t-test

# Required Sample size for one-way ANOVA



**Sample Size Estimation for ANOVA**

Sig=0.05 (k=3)

**Sample Size Estimation for ANOVA**

Sig=0.05 (k=5)

# Software for power analysis

- G*Power
- R  (e.g., package 'pwr')

# G*Power
# Main window

[Download GPower](#)

# Three basic steps:

◦ Select appropriate test

◦ Input parameters

◦ Determine effect size (can use background info or guess)

Determine effect size (can use background info or guess)

**Central and noncentral distributions**

Shows the distribution of the null hypothesis (red) and the alternative (blue)

**X-Y plot for a range of values**

Generates plots of one of the parameters α, effect size, power and sample size, depending on a range of values of the remaining parameters

# Exercise 1

The way productive animals are killed concerns the scientific community and society for both bioethical and productive reasons. The aim of the project is the evaluation of used methods of capturing and killing fish.
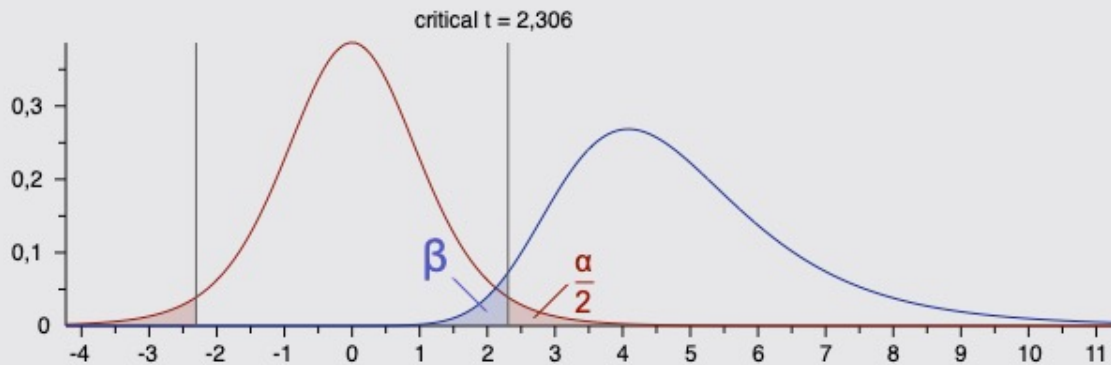
The evaluation of stress includes measuring cortisol concentrations in sea bass, *Dicentrarhus labrax* between:

**A) Two different ways of capturing the fish**

(i)   by using hook & line without any other handling (considered to be the fastest method of conception) and

(ii)  by the common way of capture (synchronization, netting, exposure to air)

You are interested in determining if the average cortisol concentration differs between treatments. (Determine the sample size you will use)

Cortisol preliminary data: (mean ± SD) Hook & line: 101.5 ± 30.6 ng/ml and Common way: 200.4 ± 40.0 ng/ml

**Answer:**
A total of 10 fish are needed
(5 per group)

# Exercise 2

**B) Four different ways of killing**

(i) by puncturing the spinal cord using Ikigun, (https://www.ikigun.com – considered to be the fastest way to kill) after anaesthesia with benzocaine;

(ii) by immersion in ice slurry (heat shock and anoxia) without prior use of anesthesia (the method used in fish farms);

(iii) with prolonged exposure to chemical anesthesia with benzocaine and

(iv) using electroanesthesia followed by immersion in ice water (recommended new humane method of killing in fish farming).

You are interested in determining if the average cortisol concentration differs between treatments. (Determine the sample size you will use)

(Cortisol preliminary data: (mean ± SD) ikigun: 617.1 ± 151.2; ice slurry: 478.7 ± 95.1; prolonged anesthesia: 531.2 ± 121.1; electroanesthesia: 389.6 ± 81.4)

Standard deviation within group = 140,25 ng/ml  group sizes of 5

**Answer:**
A total of 36 fish are needed (9 per group)
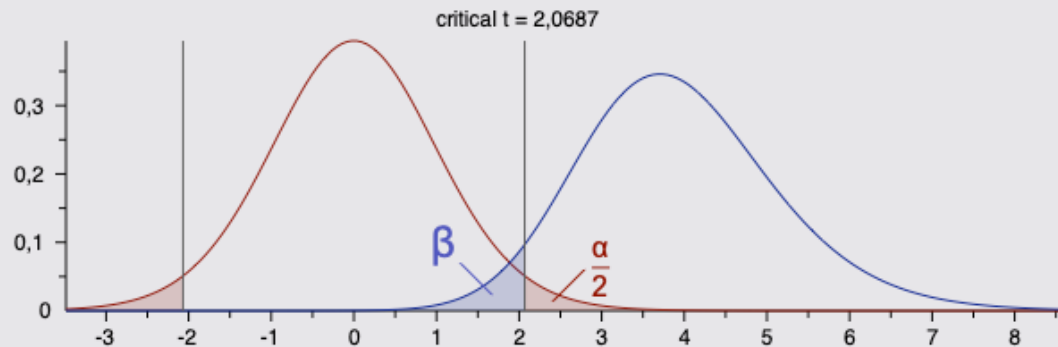
# Exercise 3

Behavioral experiment: The light-dark test to assess anxiety. Increased latency (in sec) means more anxious.

You are interested in determining

if there is difference between treatments.

(Determine the sample size you will use)

| mice | no stress | stress | difference |
|------|-----------|--------|------------|
| 1 | 28,00 | 10,20 | 17,80 |
| 2 | 64,00 | 28,00 | 36,00 |
| 3 | 31,60 | 23,00 | 8,60 |
| 4 | 46,00 | 30,00 | 16,00 |
| 5 | 19,20 | 26,51 | -7,31 |
| 6 | 68,93 | 3,83 | 65,10 |
| 7 | 13,00 | 17,85 | -4,85 |
| 8 | 24,70 | 8,34 | 16,36 |
| 9 | 28,81 | 16,22 | 12,59 |
| 10 | 29,47 | 17,52 | 11,95 |
| 11 | 13,98 | 12,90 | 1,08 |
| mean | 33,43 | 17,67 | 15,76 |
| SD | 18,70 | 8,51 | 20,24 |

**Answer**:
A sample size of 24 mice is needed to detect a difference with prob 0.95

# Power Plot

**t tests - Means: Difference between two dependent means (matched pairs)**
Tail(s) = Two, α err prob = 0,05, Effect size dz = 0,778656



Effect size dz

—○— = 0,778656

*Total sample size* (y axis)

*Power (1-β err prob)* (x axis): 0,8  0,81  0,82  0,83  0,84  0,85  0,86  0,87  0,88  0,89  0,9  0,91  0,92  0,93  0,94  0,95

## Parameters

Plot (on y axis)  [Total sample size]  ☑ with markers  ☐ displaying the values in the plot

as a function of  [Power (1-β err prob)]  from  0,8  in steps of  0,01  through to  0,95

Plot  [1]  graph(s)  [interpolating points]

with  [Effect size dz]  at  0,7786561

and  [α err prob]  at  0,05

**Draw plot**

# Assumptions

All statistical tests make assumptions about data

- parametric tests (like t-tests and ANOVA) make three main assumptions

  - probability distribution of errors (e.g. normal)

  - homogeneity of variance

  - independence of errors

- assumptions need to be checked before relying on the result of a test

# Assumptions not met

- Robust if equal sample sizes

- Transformations of data may be useful

- Nonparametric tests
  - rank transform tests
    - Mann-Whitney test or Wilcoxon rank sum test for comparing two groups (nonparametric equivalent of the t-test)
    - Kruskal-Wallis for comparing several groups (nonparametric equivalent of the one-way ANOVA)
    - Friedman test (nonparametric equivalent of the randomized block ANOVA)

# Thank you for your attention!

Dina Lika
lika@uoc.gr

# References

- G.P. Quinn and M.J. Keough, Experimental Design and Data Analysis for Biologists. Cambridge (2002)

- J. H. Zar,   Biostatistical Analysis, Prentice-Hall International, Inc. (1996)

- C. Dytham, Choosing and using statistics: A Biologist's guide. Blackwell science (1999)

- M.F.W. Festing and D.G. Altman. 2002. Guidelines for the design and statistical analysis of experiments using laboratory animals, *ILAR journal*, 43(4): 244-258