# Elements of experimental design

### … is it really that important that the experiment is well planned?

*Gdańsk, May 13, 2023*

**Peter Grešner, PhD.**

*Centre of Biostatistics and Bioinformatics Analyses*
*Division of Translational Oncology*
*Medical University of Gdansk*

*peter.gresner@gumed.edu.pl*

www.gumed.edu.pl

# Author's background…

1997 **programming class** high school (Slovak Republic)

**2002 biomedical physics (MSc)**
Faculty of Mathematics, Physics and Informatics,
Comenius University, Bratislava, Slovak Republic

**2006 biophysics (PhD)**
Faculty of Mathematics, Physics and Informatics,
Comenius University, Bratislava, Slovak Republic

...

**2006 medical statistician**
Military Teaching Hospital No.2, Medical University of Lodz,
Poland

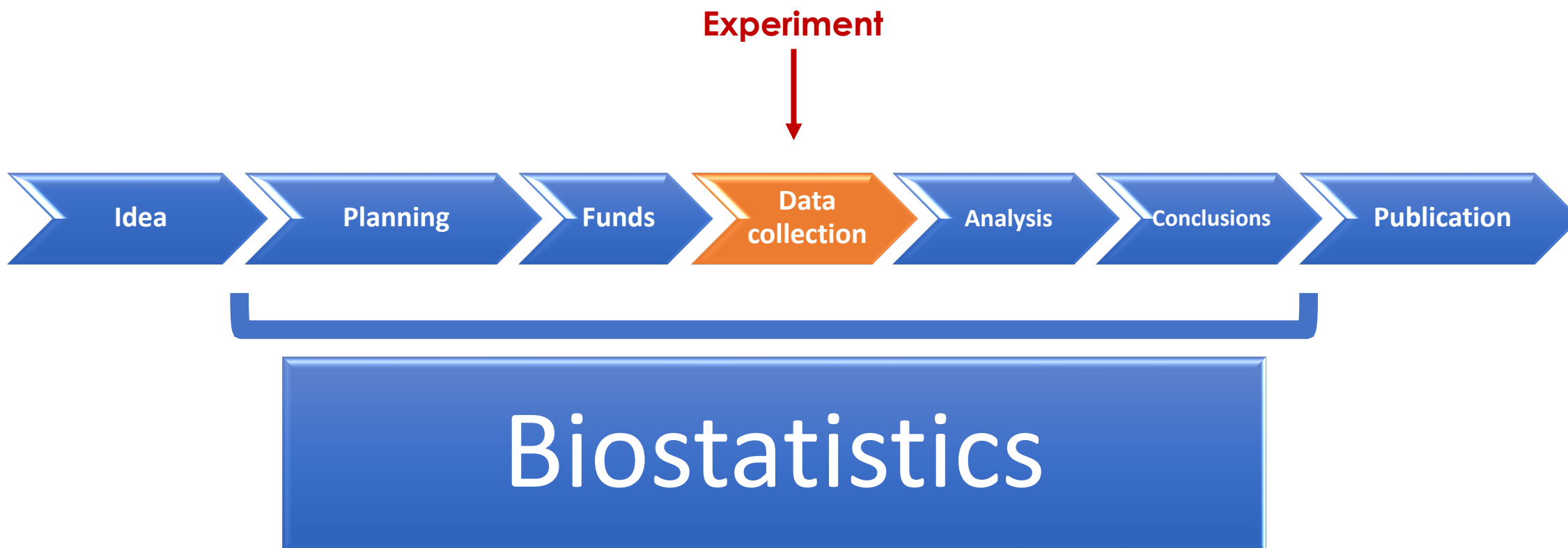**2006 – 2020  postdoctoral fellow […] adjunct associate professor**
Nofer Institute of Occupational Medicine, Lodz, Poland

**2020 – … senior biostatistician**
Medical University of Gdansk, Poland

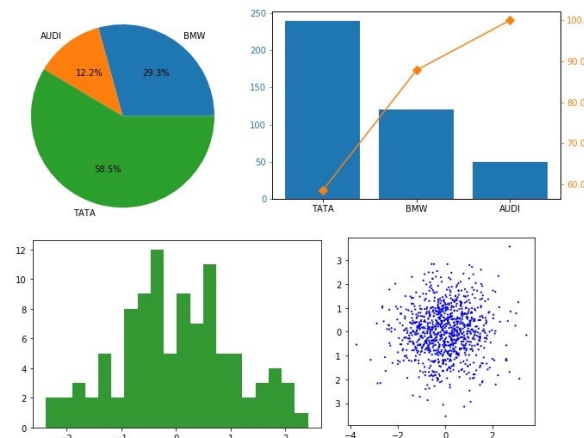# What exactly is (bio)statistics?

# Mostly, it's not done properly...

– approximately 60% of original papers in the field of biomed/pharm contain statistical errors(De Muth 1999)

- improper experimental design and **planning**
- poorly formulated research **hypothesis**
- incorrectly estimated **sample size** (or no estimation at all)
- misused **mean & SD**
- wrong selection of **parametric/non-parametric** tests
- incorrect use of **paired/unpaired** tests
- using standard error (**SE**) instead of standard deviation (**SD**)

- using multiple t-tests as an extension of the analysis of variance method
- using the chi² test and Fisher's exact test

*Watała: Biostatystyka - wykorzystanie metod statystycznych w pracy badawczej w naukach biomedycznych. α-medica press, Bielsko-Biała, 2002.*
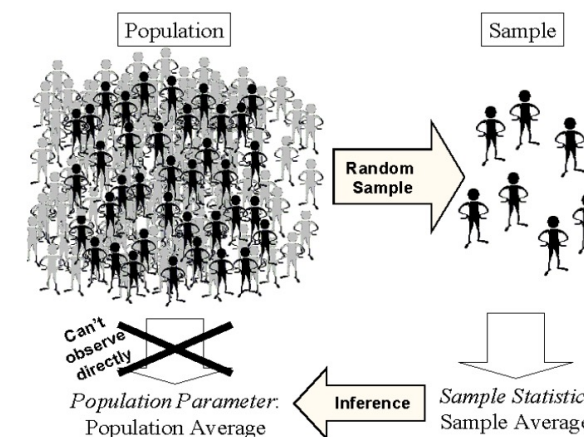
# Biostatistics – what are our options?



**Descriptive statistics**

- characteristics of collected dataset
- answers questions:
  - how to describe/present the collected data?
  - what is the most representative value?
  - what are the extremes?
  - what is the spread of data?
  - how to compare data to those from similar sets?
- not going "beyond the scope" of the collected data

**Statistical inference**

- a way of "generalizing"
- drawing conclusions about a population from which ONLY a small part (sample) has been analyzed

# Descriptive statistics

**Histogram**

**Measures of central tendency**
mode, mean ($\mu$), median

**Measures of spread**
range, quartiles (**Q**)
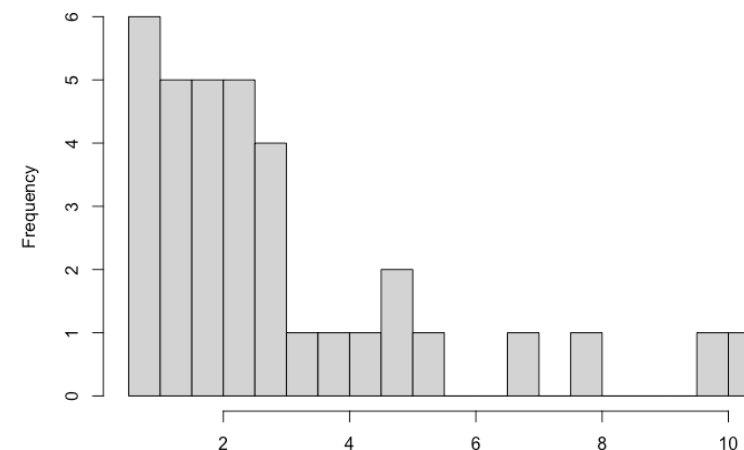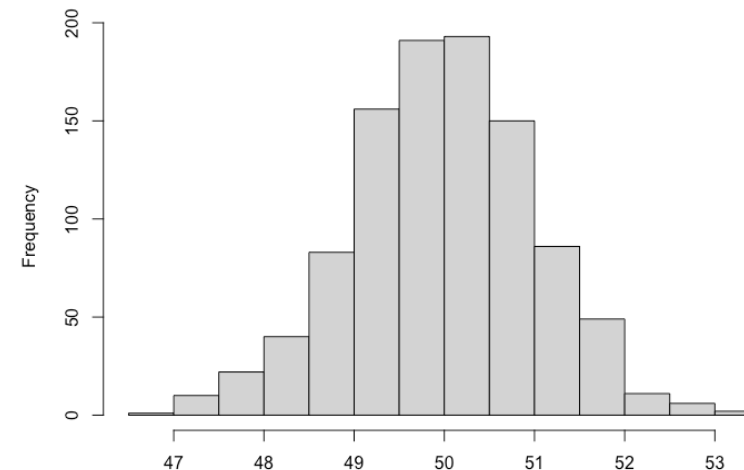variance (**s**), standard deviation ($\sigma$; **SD**; 68%; 95%)
interquartile range (**IQR=Q3-Q1**; middle 50%)

coefficient of variation [%] ($\sigma / \mu$)
quartile coefficient of dispersion [%]
(**IQR/(Q3+Q1)**)

# Median vs. Mean

**median**
- the middle value in distribution
- 50% left / 50% right
- less sensitive to outliers/extreme
- only numerical data

**mean ($\mu$; $\bar{y}$)**
- sum of all values divided by their number
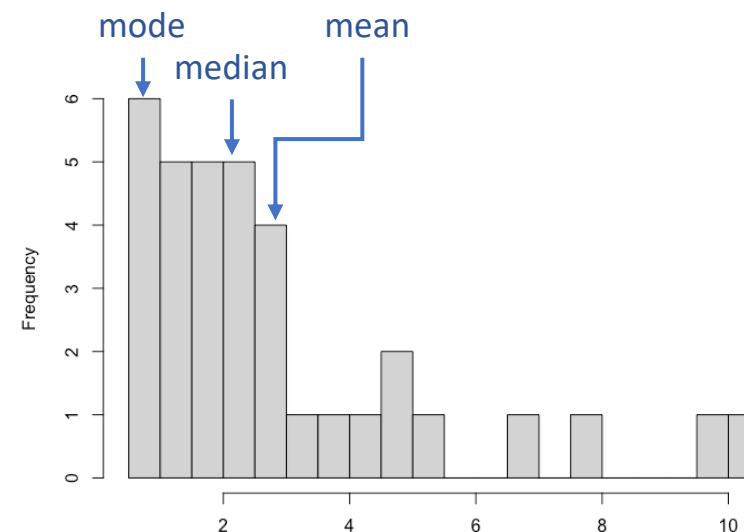- sensitive to outliers/extremes
- only numerical data

**Normal distribution**
mean & SD
median & IQR

**Non-normal distribution**
median & IQR (only!)

# Normal (Gaussian) distribution

- the "bell"-shape, symmetrical around mean
- mathematically **fully** characterized by mean ($\mu$) and standard deviation ($\sigma$)
- $\mu = 0$; $\sigma = 1$; AUC = 1;

**Significance of normal distribution**
- many biomedical parameters present normal distribution
- many statistical tests thus **assume normally distributed data** (if not met, the logic of analysis fails)

**Central limit theorem**

histogram of means from many non-gaussian samples will present normal distribution

# How to choose proper measures?

**Always check your data first!**

- visual inspection of histograms

- statistical tests for testing data normality
  - Shapiro-Wilk W test
  - Kolmogorov-Smirnov test
  - Lilliefors test

# Statistical inference

**Research hypothesis**
A statement specifying the existence of some relationship, difference, mechanism, process, etc.

**Statistical hypothesis**
Redefinition of research hypothesis into a measurable form.

**Hypothesis testing**
A sequence of steps allowing us to either accept or reject the hypothesis.

In order to do so, one has to follow **the rules**…

*One ring to rule them all…*
*…and in the darkness bind them.*

# Statistical inference

**Sir Karl Raimund Popper (1902 – 1994)**

**Criterion of falsifiability**
- the main scientific criterion
- In order to prove something, try to reject the negation thereof
- critical rationalism („popperism")
- philosophy of science
- historical context of Eastern Block 😵‍💫

**Rules:**
- study groups selection
- formulation of hypotheses
- hypothesis testing
- making decisions and drawing conclusion

Popper
The Logic of Scientific Discovery

# Statistical inference - workflow

```
┌─────────────────────────────────────────┐
│      Define the scientific problem        │  ←————  The idea
└─────────────────────────────────────────┘
                    ↓
┌─────────────────────────────────────────┐
│      Specify the research hypothesis      │  ←
└─────────────────────────────────────────┘
                    ↓
┌─────────────────────────────────────────┐
│    Choose appropriate statistical test    │  ←————  Planning
└─────────────────────────────────────────┘
                    ↓
┌─────────────────────────────────────────┐
│  Select properly the control and study groups │  ←
└─────────────────────────────────────────┘
                    ↓
┌─────────────────────────────────────────┐
│              Collect data                 │  ←————  „Experiment"
└─────────────────────────────────────────┘
                    ↓
┌─────────────────────────────────────────┐
│    Apply the appropriate statistical test │  ←————  Analysis
└─────────────────────────────────────────┘
                    ↓
┌─────────────────────────────────────────┐
│  Based on test result, accept or reject the│  ←————  Conclusions
│              hypothesis                   │
└─────────────────────────────────────────┘
```

# Scientific problem & Hypotheses

**Scientific problem**

"Does the XY disease affect the patients' IQ?"

**Statistical hypotheses**

we want to check whether our data allows us to reject this hypothesis as untrue

**The "null" hypothesis (H$_0$)**
The average IQ of people with the XY disease does not differ from the one among healthy people (without the XY disease).

$$H_0: \overline{IQ_{XY}} = \overline{IQ_{healthy}}$$

**The alternative hypothesis (H$_A$)**
The average IQ of people with the XY disease differs from the one among healthy people (without the XY disease).
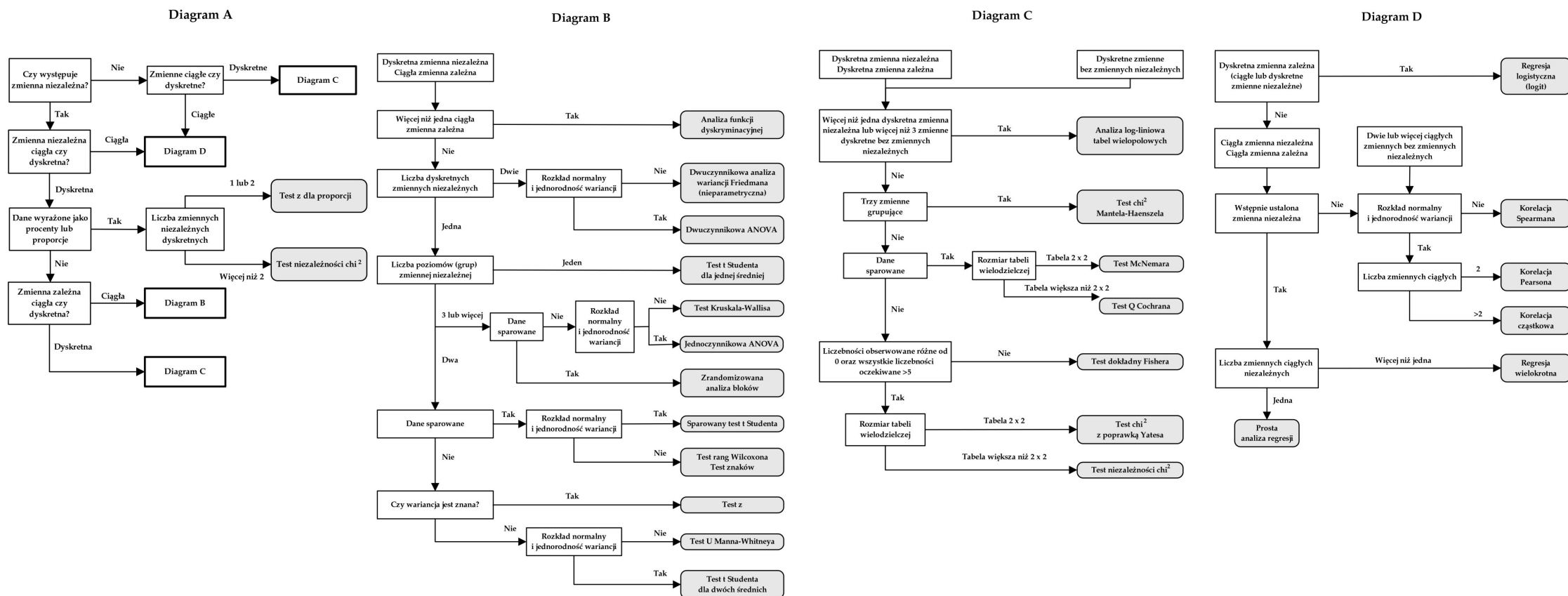
$$H_A: \overline{IQ_{XY}} \neq \overline{IQ_{healthy}}$$

# Choosing appropriate statistical test

Not a trivial problem, many aspects need to be considered

- type of data: numerical (continuous, discrete) / categorical / proportions
- normality of the data
- dependent/independent variable
- data pairing
- data censoring (right censoring – survival analysis)
- number of levels of categorical variables
- number of compared groups
- additional (confounding) factors

- **experimental design**
  - hierarchical design
  - balanced/unbalanced design
  - fixed/random effects
  - etc.

# Choosing appropriate statistical test



*Watała: Biostatystyka - wykorzystanie metod statystycznych w pracy badawczej w naukach biomedycznych. α-medica press, Bielsko-Biała, 2002.*

www.gumed.edu.pl

# Basic test characteristics

**Parametric tests**
- make several assumptions concerning the distribution of data
- if not met, their logic fails and provide unreliable results
- limited use but very "powerful" (general linear model, ANOVA, …)

**Nonparametric tests**
- less or no assumptions concerning the data distribution and other characteristics
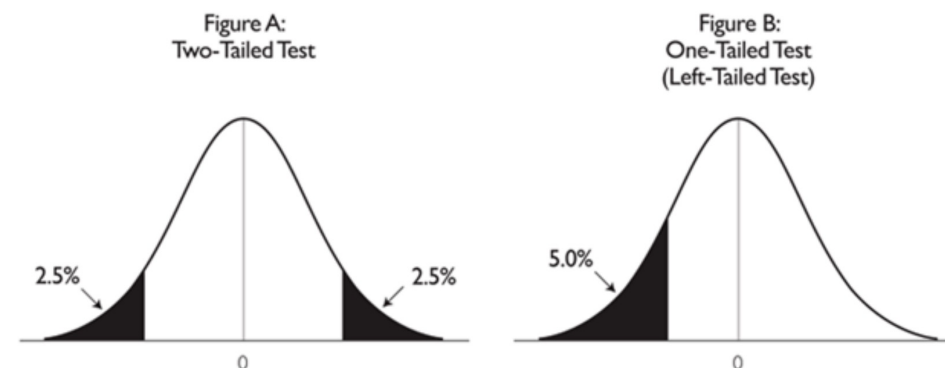- more "universal"

**One-tailed tests**
- assume one specific change (decrease | increase)

**Two-tailed tests**
- assume change in general (decrease & increase)
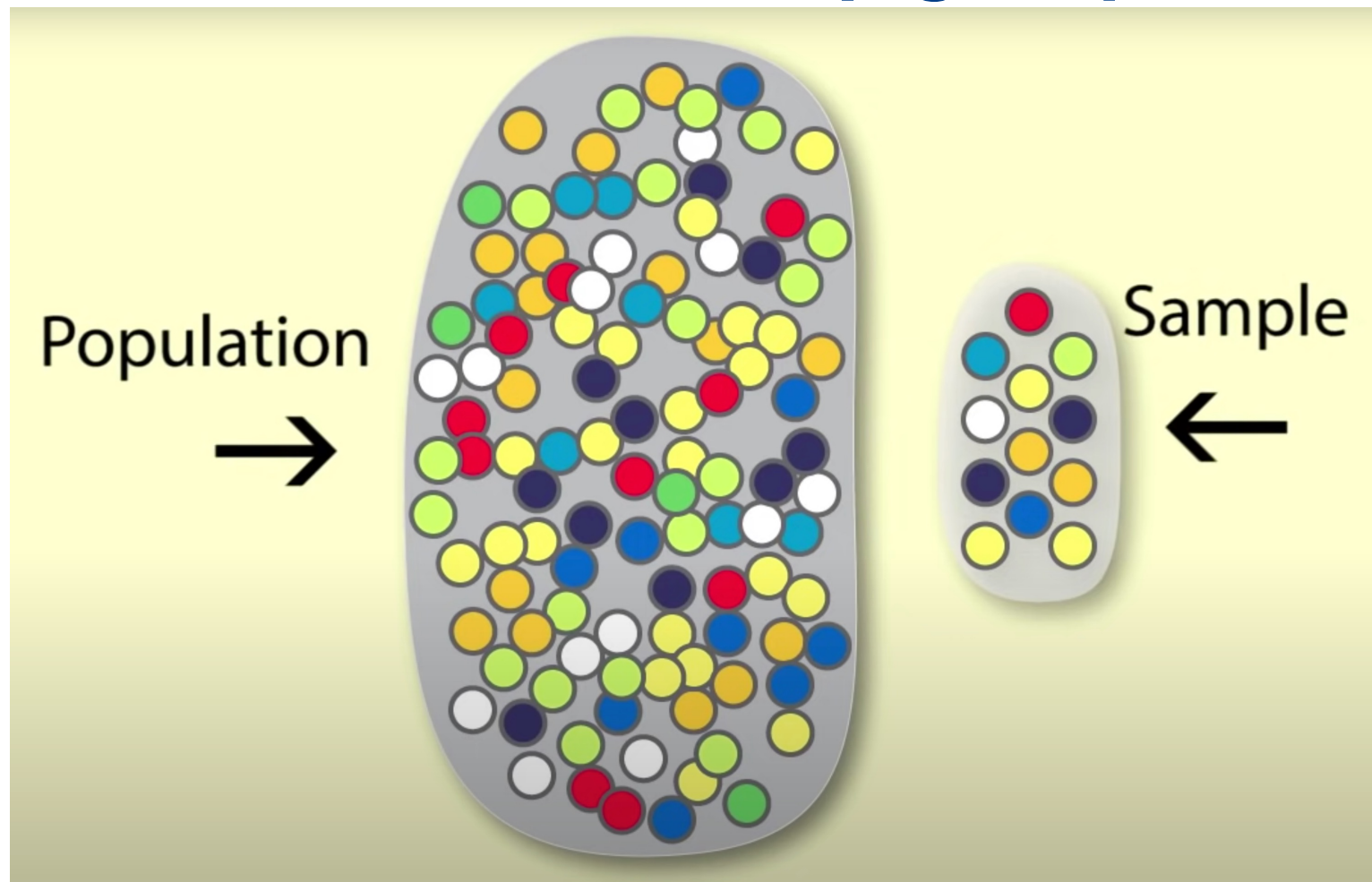


Two-Tailed Versus One-Tailed Hyphothesis Tests

Figure A: Two-Tailed Test — 2.5%  2.5%  0

Figure B: One-Tailed Test (Left-Tailed Test) — 5.0%  0

# Parametric vs. non-parametric tests

| Parametric test | Non-Parametric equivalent |
|---|---|
| Paired t-test | Wilcoxon Rank sum Test |
| Unpaired t-test | Mann-Whitney U test |
| Pearson correlation | Spearman correlation |
| One way Analysis of variance | Kruskal Wallis Test |

| Input Variable | Outcome variable | | | | | |
|---|---|---|---|---|---|---|
| | Nominal | Categorical (>2 Categories) | Ordinal | Quantitative Discrete | Quantitative Non-Normal | Quantitative Normal |
| Nominal | $X^2$ or Fisher's | $X^2$ | $X^2$-trend or Mann-Whitney | Mann-Whitney | Mann-Whitney or log-rank[a] | Student's $t$ test |
| Categorical (2>categories) | $X^2$ | $X^2$ | Kruskal-Wallis[b] | Kruskal-Wallis[b] | Kruskal-Wallis[b] | Analysis of variance[c] |
| Ordinal (Ordered categories) | $X^2$-trend or Mann-Whitney | e | Spearman rank | Spearman rank | Spearman rank | Spearman rank or linear regression[d] |
| Quantitative Discrete | Logistic regression | e | e | Spearman rank | Spearman rank | Spearman rank or linear regression[d] |
| Quantitative non-Normal | Logistic regression | e | e | e | Plot data and Pearson or Spearman rank | Plot data and Pearson or Spearman rank and linear regression |
| Quantitative Normal | Logistic regression | e | e | e | Linear regression[d] | Pearson and linear regression |

https://www.healthknowledge.org.uk

# Selection of study groups
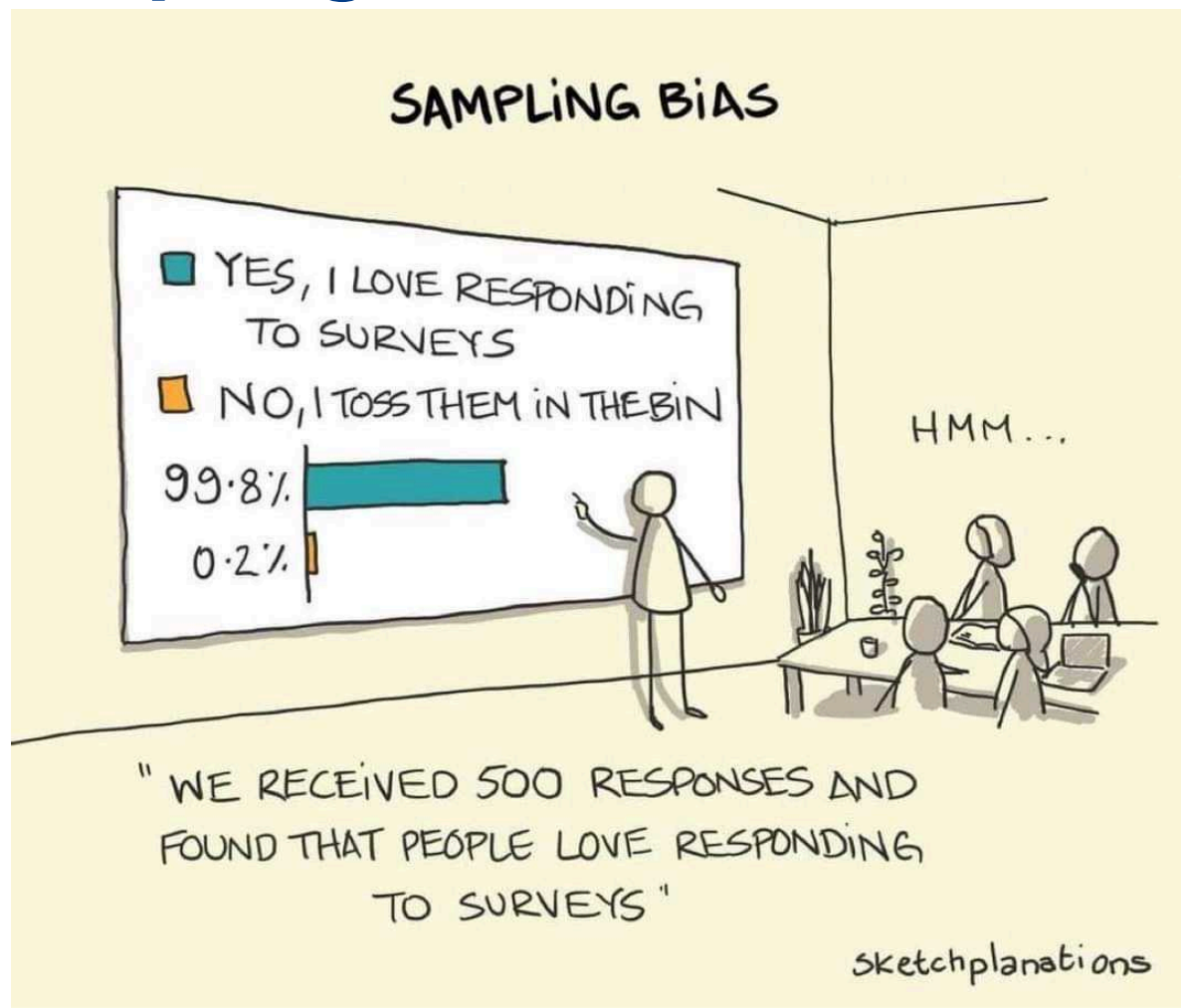
Population → Sample ←

- representative and appropriately selected groups
- inclusion/exclusion criteria (diseases, medications)
- certain degree of uncertainty

- **sampling bias**

**Study size estimation**
- how large should the study/control group be?
- ethical, financial, temporal, logistical aspects

*Dr Nic's Maths & Stats YT Channel (https://www.youtube.com/channel/UCG32MfGLit1pcqCRXyy9cAg)*

# Sampling bias



Source: Internet

- sampling strategy, in which a certain members of population have higher or lower sampling probability

- if not accounted for during data analysis, it's effect can be erroneously attributed to the phenomenon under study

- **healthy user bias** (overestimating health of general population)

- **Berkson's fallacy** (underestimating health of general population)

# Study groups selection precautions

- completely randomly selected volunteers; is it even feasible? (control group searched within specific occupation, "healthy" hospital visitors)

- how to match the study and control groups in terms of age, when the examined disease appears only in certain age groups? Will it then be possible to exclude other factors influencing the examined parameters in the compared control group?? (centenarians)

- how to perform randomization if subjects representing the study group are only rarely encountered?

- how to conduct drug effect testing (placebo, single/double blind study, the "noninferiority" problem)

# Data collection

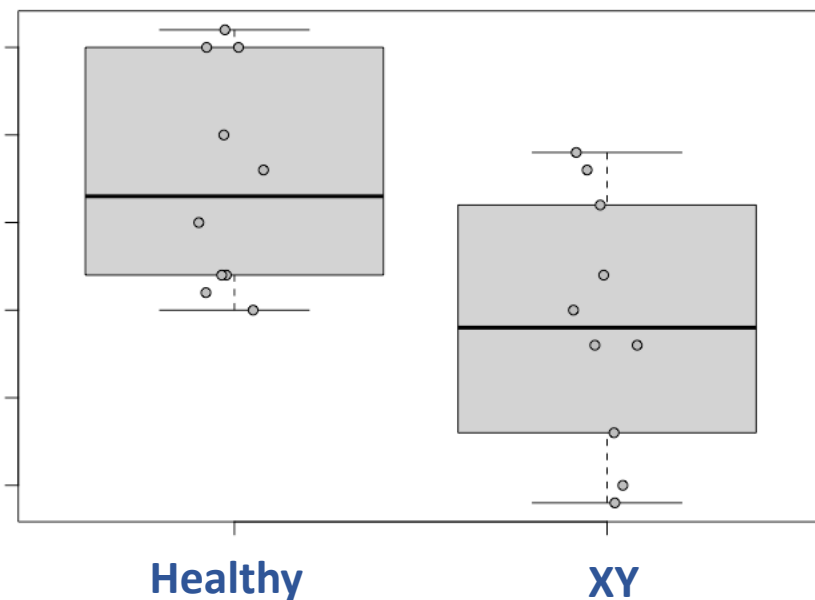**Various types of experiments:**
- basic research
- population studies
  - observational
    - cohort / case-control studies
    - cross-sectional / monitoring (*longitudinal*: *prospective* / *retrospective*)
  - interventional (experimental)
- clinical trials
- questionnaire studies


- ***independence, randomness of all observations***

# Hypothesis testing

**statistical tests**
- observed differences between the groups are just by chance or rather indicate a kind of regularity (pattern)
- transform observed differences into *statistics*

**remember to check the data distribution first!**



Healthy          XY

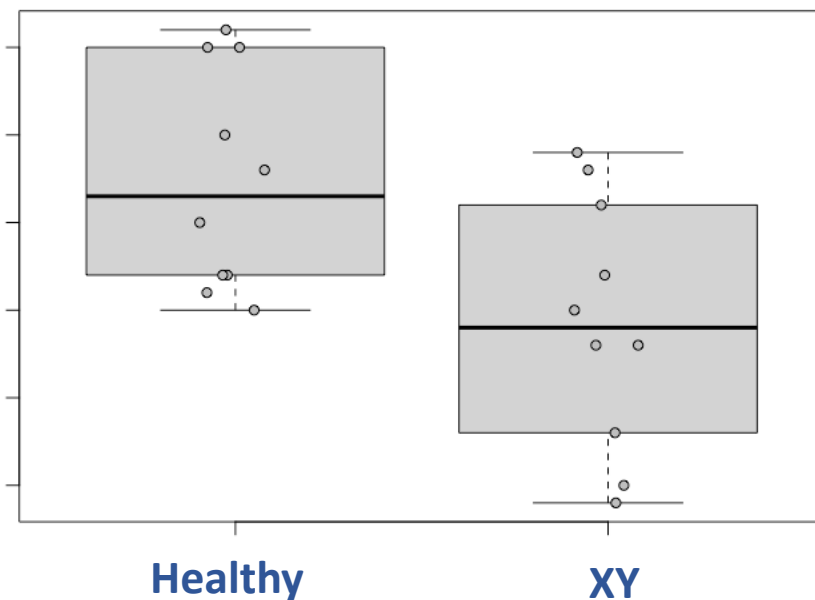$$t = \frac{\overline{x}_1 - \overline{x}_2}{s\sqrt{(1/n_1 + 1/n_2)}}$$

$$s = \sqrt{\left[\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)}\right]}$$

# Hypothesis testing

**Why** calculate *t statistic* out of real difference?

**Uncertainty**
- exact distribution of IQ is unknown
- the representativeness of groups is unknown

Distribution of $t$ under the $H_0$ hypothesis

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s\sqrt{(1/n_1 + 1/n_2)}}$$

$$s = \sqrt{\left[\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1+n_2-2)}\right]}$$

$f(t)$

observed t

$\alpha$

$t$

$0$   $t\alpha,\nu$

Healthy     XY

# Making decisions

Life Sciences
$\alpha$ = 0.05; $\alpha$ = 0.01

Physical sciences
$\alpha$ = 0.0000003 (Higg's boson)

|  | Real world | |
|---|---|---|
| **Test (experiment) result** | $x_1 = x_2$ ($H_0$ is true) | $x_1 \neq x_2$ ($H_A$ is true) |
| Reject $H_0$ $x_1 \neq x_2$ | **Type I error** (significance) $\alpha$, p | **Correct decision** (test power) |
| Do not reject $H_0$ $x_1 = x_2$ | **Correct decision** (1- significance) | **Type II error** (1- test power) |

The higher the test power, the better
The lower the significance level, the better (?!?)

**What is the result?**
- the effect size (i.e. difference between means)
- level of significance (the degree of our (un)certainty)

# Group size estimation

In order to generalize the study/control groups must be of required sizes.

- required study group size (**n**) depends on several parameters
    - statistical power (the ability to reject $H_0$ when it is false; 80%; $\mathbf{z_\beta}$)
    - level of significance (the probability of rejecting $H_0$ when it is *de facto* true; 5%; $\mathbf{z_\alpha}$)
    - effect size (i. e. assumed difference in mean values ($\boldsymbol{\delta}$))
    - assumed spread (e. g. variance; $\boldsymbol{\sigma^2}$)

$$z_\beta = \frac{\delta}{\sqrt{\dfrac{2\sigma^2}{n}}} - z_{\alpha/2}$$

$$n \approx \frac{\sigma^2\left(z_\beta + z_{\alpha/2}\right)^2}{\delta^2}$$

- larger the spread ➔ higher **n**
- smaller the difference ➔ higher **n**

# Group size estimation

$$n \approx \frac{\sigma^2\left(z_\beta + z_{\alpha/2}\right)^2}{\delta^2}$$

- independent observations
- numerical data
- trying to prove the difference between groups

| | musimy znać | | równanie |
|---|---|---|---|
| **(a) istotność różnic** | | | |
| 1. pojedyncza średnia | $u, v$ | jak poniżej | $\dfrac{(u+v)^2\,\sigma^2}{(\mu-\mu_0)^2}$ |
| | $\mu\text{-}\mu_0$ | różnica między średnią badaną $\mu$ i średnią teoretyczną $\mu_0$ ($H_0$) | |
| | $\sigma$ | odchylenie standardowe | |
| 2. pojedyncza częstość | $\mu$ | częstość | $\dfrac{(u+v)^2\,\mu}{(\mu-\mu_0)^2}$ |
| | $\mu_0$ | wartość dla $H_0$ | |
| | $u, v$ | jak poniżej | |
| 3. pojedyncza proporcja | $\pi$ | proporcja | $\dfrac{\{u\sqrt{[\pi(1-\pi)]}+v\sqrt{[\pi_0(1-\pi_0)]}\}^2}{(\pi-\pi_0)^2}$ |
| | $\pi_0$ | wartość dla $H_0$ | |
| | $u, v$ | jak poniżej | |
| 4. porównanie dwóch średnich *(liczebność każdej grupy)* | $u, v$ | jak poniżej | $\dfrac{(u+v)^2(\sigma_1^2+\sigma_2^2)}{(\mu_1-\mu_2)^2}$ |
| | $\mu_1\text{-}\mu_2$ | różnica między średnimi | |
| | $\sigma_1,\ \sigma_2$ | odchylenie standardowe | |
| 5. porównanie dwóch częstości *(liczebność każdej grupy)* | $u, v$ | jak poniżej; $\mu_1, \mu_2$ częstości | $\dfrac{(u+v)^2(\mu_1+\mu_2)}{(\mu_1-\mu_2)^2}$ |

| | musimy znać | | równanie |
|---|---|---|---|
| 6. porównanie dwóch proporcji *(liczebność każdej grupy)* | $u, v$ | jak poniżej; $\pi_1, \pi_2$ proporcja | $\dfrac{\{u\sqrt{[\pi_1(1-\pi_1)+\pi_2(1-\pi_2)]}+v\sqrt{[2\bar{\pi}(1-\bar{\pi})]}\}^2}{(\pi_2-\pi_1)^2}$ |
| | | | gdzie $\quad \bar{\pi}=\dfrac{\pi_1-\pi_2}{2}$ |
| 7. badanie typu *case-control* *(liczebność każdej grupy)* | $\pi_1$ | proporcja w grupie kontrolnej wystawiona na działanie czynnika | $\dfrac{\{u\sqrt{[\pi_1(1-\pi_1)+\pi_2(1-\pi_2)]}+v\sqrt{[2\bar{\pi}(1-\bar{\pi})]}\}^2}{(\pi_2-\pi_1)^2}$ |
| | | | gdzie $\quad \bar{\pi}=\dfrac{\pi_1+\pi_2}{2}$ |

GDAŃSKI UNIWERSYTET MEDYCZNY

www.gumed.edu.pl

# What should you take from this part?

1. **Plan** your experiments properly (in cooperation with a specialist)
2. **Check** your data prior to analysis (if it meets assumptions of tests)
3. Use **appropriate measures** of centrality and spread
4. Don't be afraid to use **other than parametric tests**
5. The **p-value is not a result** of experiment!
   It is the effect size (difference), while the p-value tells something about how certain you are when it comes to the effect size.

6. Contact **CABiB** if you need help

Idea → Planning → Funds → Data collection → Analysis → Conclusions → Publication

# Selected practical issues

# Group size estimation

- determine a number of factors *a priori*
  - the level of significance
  - statistical power of the test
  - the effect size
  - the dispersion measure (variance)

"to start from something is better than starting from nothing"

**How to formulate the question concerning the group size?**

❌ „ So how many subjects do I need to include in the study to obtain statistically significant difference?"

✅ „Assuming the variability of data at the level of 25%, the level of statistical significance below 0.01, statistical power above 80% and the planned use of one-tailed Student's t-test, how many patients should be included in the study to show a 15% higher level of biomarker expression in the group subjected to treatment?"

✅ " If the study and control groups include both 30 subjects per group, then assuming the variability data at the level of 25%, the significance level below 0.01, the statistical power above 80% and the planned use of one-tailed Student's t-test, how large an increase in the biomarker expression level will I be able to evaluate as statistically significant?"

**What if we expect a non-Gaussian distribution of the data?**

# Group size estimation: two means

**Required information:**

- expected means ($\mu_1$; $\mu_2$)
- expected spreads ($\sigma_1$; $\sigma_2$)
- u – a value related to level of significance
- v – a value related to statistical power of test

$$\frac{(u+v)^2(\sigma_1^2 + \sigma_2^2)}{(\mu_1 - \mu_2)^2}$$

| significance | u/$z_{alfa/2}$ | u/$z_{alfa}$ |
|---|---|---|
| 0,1 | 1,6449 | 1,2816 |
| 0,05 | 1,9600 | 1,6449 |
| 0,01 | 2,5758 | 2,3263 |
| 0,0125 | 2,5000 | 2,2400 |
| 0,005 | 2,8070 | 2,5758 |
| 0,0024 | 3,0357 | 2,8202 |
| 0,00125 | 3,2272 | 3,0233 |
| 0,001 | 3,2905 | 3,0902 |
| 0,0001 | 3,8906 | 3,7190 |
| 0,00001 | 4,4172 | 4,2649 |
| | * two-tailed | * one-tailed |

| power [%] | u/$z_{beta}$ |
|---|---|
| 60 | 0,2533 |
| 70 | 0,5244 |
| 75 | 0,6745 |
| 80 | 0,8416 |
| 85 | 1,0364 |
| 90 | 1,2816 |
| 95 | 1,6449 |
| * one-sided value | |

# Group size estimation: three means

**Required information:**

- expected means ($\mu_1$; $\mu_2$; $\mu_3$)
- expected spreads ($\sigma_1$; $\sigma_2$; $\sigma_3$)
- $u$ – a value related to level of significance
- $v$ – a value related to statistical power of test

$$\frac{(u+v)^2(\sigma_1^2 + \sigma_2^2)}{(\mu_1 - \mu_2)^2}$$

- Bonferroni correction for multiple testing (adjust p)

| significance | $u/z_{alfa/2}$ | $u/z_{alfa}$ |
|---|---|---|
| 0,1 | 1,6449 | 1,2816 |
| 0,05 | 1,9600 | 1,6449 |
| 0,01 | 2,5758 | 2,3263 |
| 0,0125 | 2,5000 | 2,2400 |
| 0,005 | 2,8070 | 2,5758 |
| 0,0024 | 3,0357 | 2,8202 |
| 0,00125 | 3,2272 | 3,0233 |
| 0,001 | 3,2905 | 3,0902 |
| 0,0001 | 3,8906 | 3,7190 |
| 0,00001 | 4,4172 | 4,2649 |
|  | * two-tailed | * one-tailed |

| power [%] | $u/z_{beta}$ |
|---|---|
| 60 | 0,2533 |
| 70 | 0,5244 |
| 75 | 0,6745 |
| 80 | 0,8416 |
| 85 | 1,0364 |
| 90 | 1,2816 |
| 95 | 1,6449 |
| * one-sided value | |

# Group size estimation: three groups; non-Gaussian distribution

**Required information:**

- expected means ($\mu_1$; $\mu_2$; $\mu_3$)
- expected spreads ($\sigma_1$; $\sigma_2$; $\sigma_3$)
- u – a value related to level of significance
- v – a value related to statistical power of test

$$\frac{(u+v)^2(\sigma_1^2 + \sigma_2^2)}{(\mu_1 - \mu_2)^2}$$

- Bonferroni correction for multiple testing (adjust p)

- Add +15% to each group's N

| significance | u/$z_{alfa/2}$ | u/$z_{alfa}$ |
|---|---|---|
| 0,1 | 1,6449 | 1,2816 |
| 0,05 | 1,9600 | 1,6449 |
| 0,01 | 2,5758 | 2,3263 |
| 0,0125 | 2,5000 | 2,2400 |
| 0,005 | 2,8070 | 2,5758 |
| 0,0024 | 3,0357 | 2,8202 |
| 0,00125 | 3,2272 | 3,0233 |
| 0,001 | 3,2905 | 3,0902 |
| 0,0001 | 3,8906 | 3,7190 |
| 0,00001 | 4,4172 | 4,2649 |
|  | * two-tailed | * one-tailed |

| power [%] | u/$z_{beta}$ |
|---|---|
| 60 | 0,2533 |
| 70 | 0,5244 |
| 75 | 0,6745 |
| 80 | 0,8416 |
| 85 | 1,0364 |
| 90 | 1,2816 |
| 95 | 1,6449 |
| * one-sided value | |

# Group size estimation: two proportions

**Required information:**

- expected proportions ($\pi_1$; $\pi_2$)

- OR, alternatively:
- expected proportion $\pi_1$ and effect size (OR, RR, …)

- u – a value related to level of significance
- v – a value related to statistical power of test

$$\frac{\{u\sqrt{[\pi_1(1-\pi_1)+\pi_2(1-\pi_2)]}+v\sqrt{[2\bar{\pi}(1-\bar{\pi})]}\}^2}{(\pi_2-\pi_1)^2}$$

|  | Females | Males | Total |
|---|---|---|---|
| Smoking | **10** | **90** | *100* |
| Non-smoking | **110** | **30** | *140* |
| *Total* | *120* | *120* | *240* |

| significance | u/$z_{alfa/2}$ | u/$z_{alfa}$ |
|---|---|---|
| 0,1 | 1,6449 | 1,2816 |
| 0,05 | 1,9600 | 1,6449 |
| 0,01 | 2,5758 | 2,3263 |
| 0,0125 | 2,5000 | 2,2400 |
| 0,005 | 2,8070 | 2,5758 |
| 0,0024 | 3,0357 | 2,8202 |
| 0,00125 | 3,2272 | 3,0233 |
| 0,001 | 3,2905 | 3,0902 |
| 0,0001 | 3,8906 | 3,7190 |
| 0,00001 | 4,4172 | 4,2649 |
|  | * two-tailed | * one-tailed |

| power [%] | u/$z_{beta}$ |
|---|---|
| 60 | 0,2533 |
| 70 | 0,5244 |
| 75 | 0,6745 |
| 80 | 0,8416 |
| 85 | 1,0364 |
| 90 | 1,2816 |
| 95 | 1,6449 |
| * one-sided value | |

# Group size estimation: more complex situation

## Calculating the required study size for testing the new diagnostic tool (e.g. AI-based classifier)

# Group size estimation: more complex situation

**Calculating the required study size for testing the new diagnostic tool (e.g. AI-based classifier)**

- **paired** data
- data as **proportions**
- trying to prove **noninferiority** of the new diagnostic tool (that it is <u>not worse</u> compared to standard test)

RTG/CT/MRI/Histo

AI

## Table 1

*Data structure of a matched pair $2 \times 2$ table*

| New test | Standard test | | Total |
|---|---|---|---|
| | Response (1) **+** | Nonresponse (2) **−** | |
| Response (1) | $a(\pi_{11})$ | $b(\pi_{12})$ | $a + b(\pi_N)$ |
| Nonresponse (2) | $c(\pi_{21})$ | $d\ (\pi_{22})$ | $c + d(1 - \pi_N)$ |
| Total | $a + c(\pi_S)$ | $b + d(1 - \pi_S)$ | $n(1.0)$ |

Predicted **+**

Predicted **−**

# Group size estimation: more complex situation

$$n_{TS} = \left\{ \frac{z_{(1-\alpha)}\sqrt{(1+\delta_0)\bar{\pi}_{21} + (\delta_1\pi_S + \pi_{21})(\delta_0-1)}}{(\delta_1-\delta_0)\pi_S} \right.$$

$$+ z_{(1-\beta)}\left\{ 2\delta_0\pi_{21} - \delta_0(1-\delta_0)\pi_S \right.$$

$$\left. + (\delta_1-\delta_0)\pi_S[1-(\delta_1-\delta)\pi_S] \right\}^{1/2}$$

$$\left. \div (\delta_1-\delta_0)\pi_S \right\}^2 .$$

$$\pi_{21}: \quad \begin{array}{l} \min[(2-\delta_1)\pi_S/2,\ (1+\pi_S)/2-\delta_1\pi_S] \text{ for } \delta_1 \leq 1 \\ \min[(1-\delta_1\pi_S)/2,\ \pi_S/2] \text{ for } 1 < \delta_1 \leq 1/\pi_S \end{array}$$

### Sample Size Determination for Establishing Equivalence/Noninferiority via Ratio of Two Proportions in Matched-Pair Design

Man-Lai Tang,[1,*] Nian-Sheng Tang,[2] Ivan Siu-Fung Chan,[3] and Ben Ping-Shing Chan[4]

# Group size estimation: more complex situation
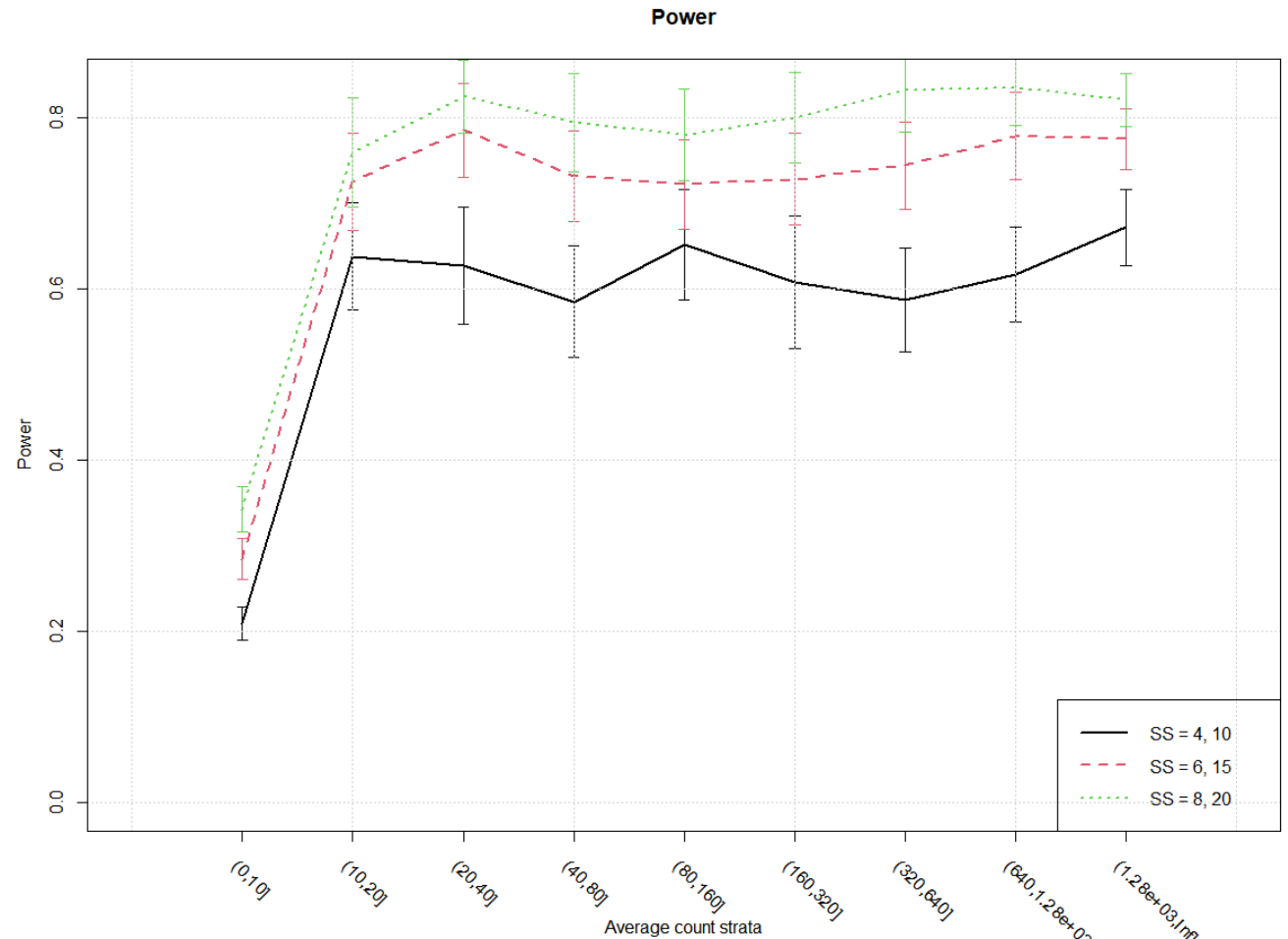
## Circulating tumour cells mRNA transcriptomics study

Q: how many reads should a differentially-expressed miRNA transcript have to provide reliable results?

*Assumptions*
- *N = 5000 transcripts;*
- *DEG = 0.1; FDR = 0.1*
- *log(FC) = 0.5 (~3.16x)*

- *4 and 10 subjects per group*
- *6 and 15 subjects per group*
- *8 and 20 subjects per group*

**What is the statistical power?**

*simulation: 50 iterations*
**(R;** *proper* package) (free)



Power

# Group size estimation: dynamic processes

**Clinical trial assessing the effectiveness of stem cell wound dressings against commercially available dressings**

_Literature_
- **Area = 6.84 * e$^{-0.124 * T}$**
- **k$_T$ = -0.124 week$^{-1}$**
- approx. 50% area reduction each 6 weeks

_Assumptions_
- **10% faster healing (k$_T$ = -0.138 week$^{-1}$)**
- 80% statistical power
- 0.05 level of significance
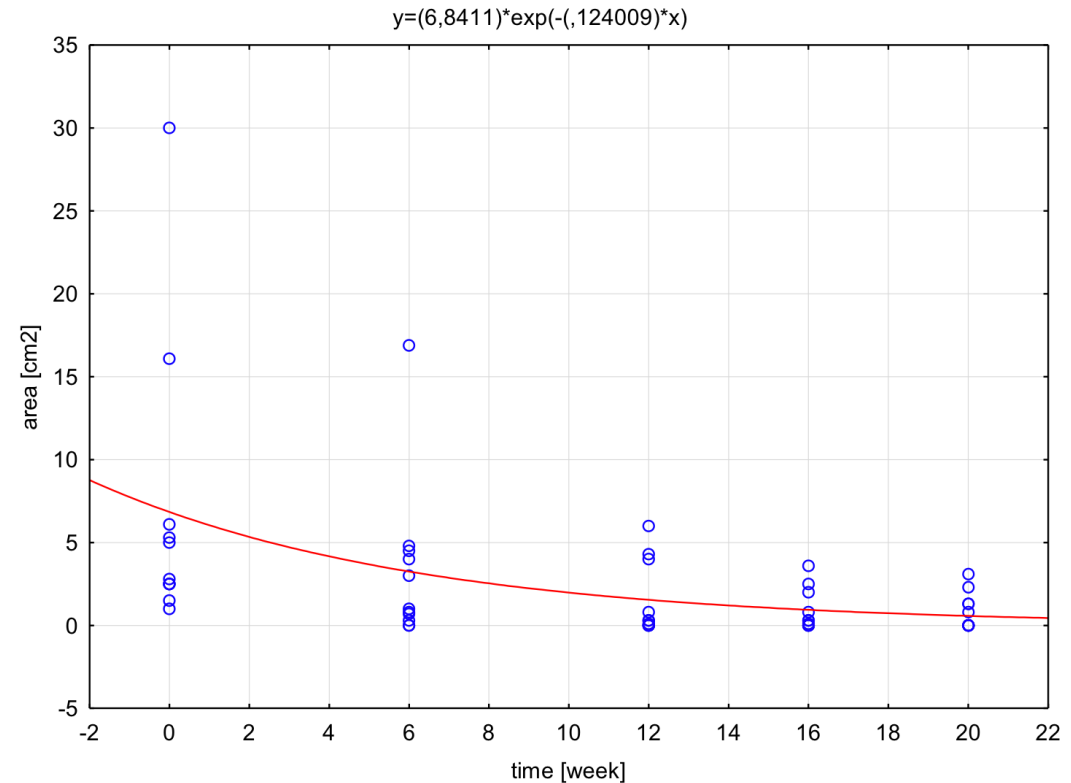- one-sided tests

**What is the required group size?**



Fig. 1. Wound surface area reduction following the use of commercially available dressing [1]

# Group size estimation: DIY solutions

**First spreadsheet:**

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | u = | 0,8416 | | numerator = | 228,69 |
| 2 | v = | 2,5000 | | denominator = | 4,00 |
| 3 | mi1 = | 3,000 | | N = | 57,17 |
| 4 | mi2 = | 1,000 | | N corretced to ro = | 33,27 |
| 5 | SD = | 3,2 | | * covariates | |
| 6 | ro = | 0,66 | | | |
| 7 | | | | | |
| 8 | | | | | |
| 9 | ratio | -66,7% | | | |

**Upper green box (two-tailed):**

| power [%]= | 80 | | $z_\beta$ = | 0,8416 | | | $z_\beta$ = | 0,8416 |
|---|---|---|---|---|---|---|---|---|
| significance = | 0,01 | | $z_\alpha$ = | 2,5758 | N = 103,0085 | | $z_\alpha$ = | 2,5758 | N = 103,0085 |
| | | | $\mu_1$ = | 2 | | | $\Delta\mu$ = | -1 |
| two-tailed | | | $\mu_2$ = | 1 | | | $\sigma$ = | 2,1 |
| | | | $\sigma$ = | 2,1 | | | | |

**Lower green box (one-tailed):**

| power [%]= | 80 | | $z_\beta$ = | 0,8416 | | | $z_\beta$ = | 0,8416 |
|---|---|---|---|---|---|---|---|---|
| significance = | 0,0125 | | $z_{\alpha/2}$ = | 2,2400 | N = 83,75814 | | $z_{\alpha/2}$ = | 2,2400 | N = 83,75814 |
| | | | $\mu_1$ = | 2 | | | $\Delta\mu$ = | -1 |
| one-tailed | | | $\mu_2$ = | 1 | | | $\sigma$ = | 2,1 |
| | | | $\sigma$ = | 2,1 | | | | |

**Second spreadsheet (Example table):**

| | A | B | | D | E | F | G | H | I | J | K | L |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | A = | 1,90 | | Example table | | Standard test (S) | | | | $\pi_S$ = | 0,790 | - czułość standardowego testu diagnostycznego zasotosowanego u osób chorych (CT, radiologia, itp.) |
| 2 | B = | -0,2501 | | | | + | - | total | | AI accuracy = | 0,900 | - accuracy narzędzia AI (względem testu standardowego) |
| 3 | C = | 0,0064069 | | New AI test (N) | + | 711 | 21 | 732 | $\pi_N$ | $\pi_{21}$ = | 0,0790 | |
| 4 | $\pi_{21est}$ = | 0,0968 | | | - | 79 | 189 | 268 | | $\pi_N$ = | 0,732 | - spodziewana RZECZYWISTA czułość nowego testu w oparciu o AI |
| 5 | | | | | total | 790 | 210 | 1000 | | $\delta_0$ = | 0,90000 | - najniższa tolerowana czułość WZGLĘDNA metody AI (czułość względna = pi(N)/pi(S); d1 <= d0) |
| 6 | człon$_1$ = | 0,5274 | | | | $\pi_S$ | | | | $\delta_1$ = | 0,92658 | - zakładana rzeczywista czułość WZGLĘDNA metody AI (czułość względna = pi(N)/pi(S)) |
| 7 | człon$_2$ = | 0,2554 | | | | | | | | $\delta$ = | 0,93 | |
| 8 | | | | | | | | | | $\pi_{N(lim)}$ = | 0,711 | |
| 9 | | | | | | | | | | $\alpha$ = | 0,05 | - pożądany poziom istotności |
| 10 | | | | | | | | | | $\beta$ = | 0,80 | - pożądana moc statystyczna |
| 11 | $N_{TS}$ = 1390 | | | | | | | | | $z_\alpha$ = | 1,6449 | - wymagana wielkość próby badanej |
| 12 | $\pi_N \geq$ 0,711 | | | | | | | | | $z_{1-\beta}$ = | 0,8416 | - dolna granica rzeczywistej czułość narzędzia AI |
| 13 | | | | | | | | | | | | |

# Group size estimation: available software solutions

- **G\*Power** (free)
- **R** (*pwr* package) (free)

- *a priori vs. a posteriori*

```
> library(pwr)
> pwr.t.test(n = NULL,
+            d=1.64399, sig.level = 0.05, power = 0.80,
+            alternative = "two.sided", type = "two.sample")

       Two-sample t test power calculation

              n = 6.91427
              d = 1.64399
      sig.level = 0.05
          power = 0.8
    alternative = two.sided

NOTE: n is number in *each* group
```

# Group size estimation: survival analysis

**GDAŃSKI UNIWERSYTET MEDYCZNY**

**UCSF Clinical & Translational Science Institute**
Explore the *Training in Clinical Research* Program at UCSF

*Required information*

- expected Hazard Ratio (HR)
- proportion of exposed/unexposed subjects

- $u$ – a value related to level of significance
- $v$ – a value related to statistical power of test

*Output*

- required number of events (!!!)

**Sample Size Calculators**
for designing clinical research

Home

**Calculators**
- CI for proportion
- CI for mean
- Means - effect size
- Means - sample size
- Proportions - effect size
- Proportions - sample size
- CI for proportion - sample size
- Survival analysis - sample size
- Prevalence
- CI for risk ratio
- More calculators...

**Calculator finder**

**About calculating sample size**

**About us**

## Sample size – Survival analysis

Two calculators for two-group survival analysis.

**Calculator 1: Number of events, given relative hazard.**

**Instructions:** Enter parameters in the green cells. Answers will appear in the blue box below.

| | | |
|---|---|---|
| α (two-tailed) = | 0.05 | Threshold probability for rejecting the null hypothesis. Type I error rate. |
| β = | 0.2 | Probability of failing to reject the null hypothesis under the alternative hypothesis. Type II error rate. |
| $q_1$ = | 0.5 | Proportion of subjects that are in Group 1 (exposed) |
| $q_0$ = | 0.5 | Proportion of subjects that are in Group 0 (unexposed); 1-$q_1$ |
| RH = | 1.6 | Relative hazard (Group 1/Group 0) |

**Calculate events**

The standard normal deviate for α = $Z_\alpha$ = 1.9600
The standard normal deviate for β = $Z_\beta$ = 0.8416

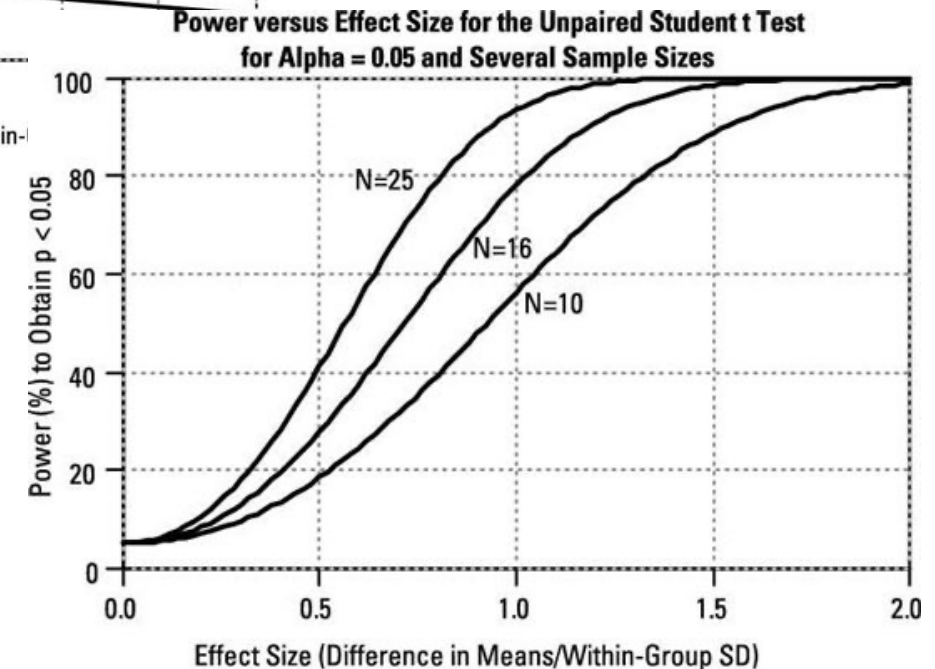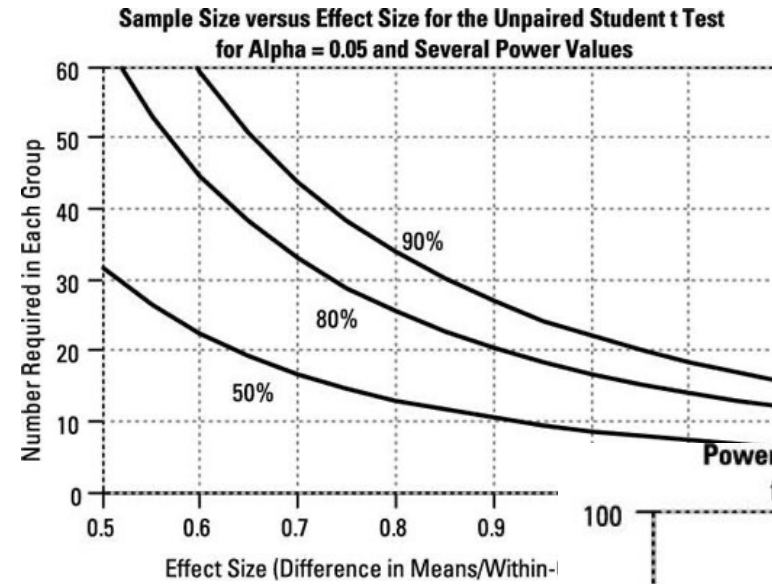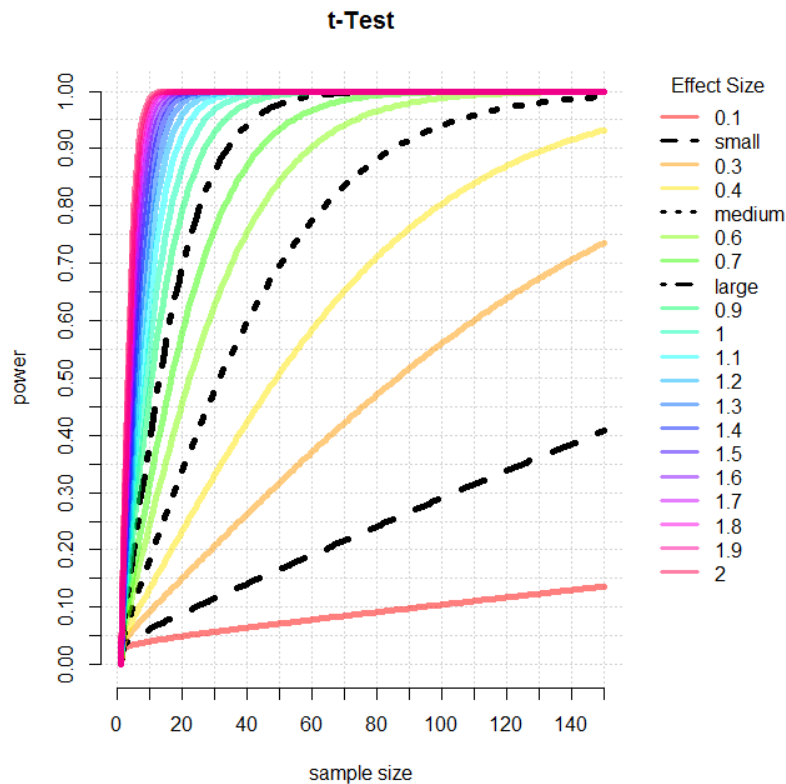$A = (Z_\alpha + Z_\beta)^2 = 7.8489$

$B = (\log(RH))^2 q_0 q_1 = 0.0552$

Total events needed = A/B = **142**

🔒 sample-size.net/sample-size-survival-analysis/

# Group size estimation: graphical outputs interpretation

- *p*: constant (0.01 or 0.05)
- *σ*: beyond our control

**N**, *δ*, **1-*β***

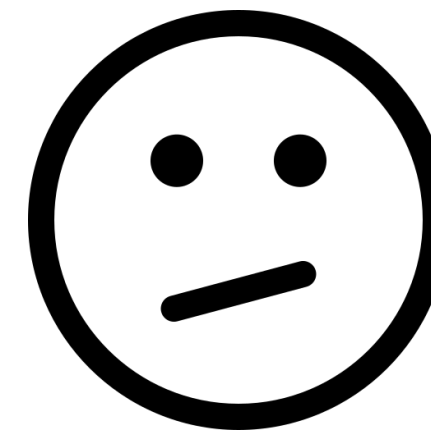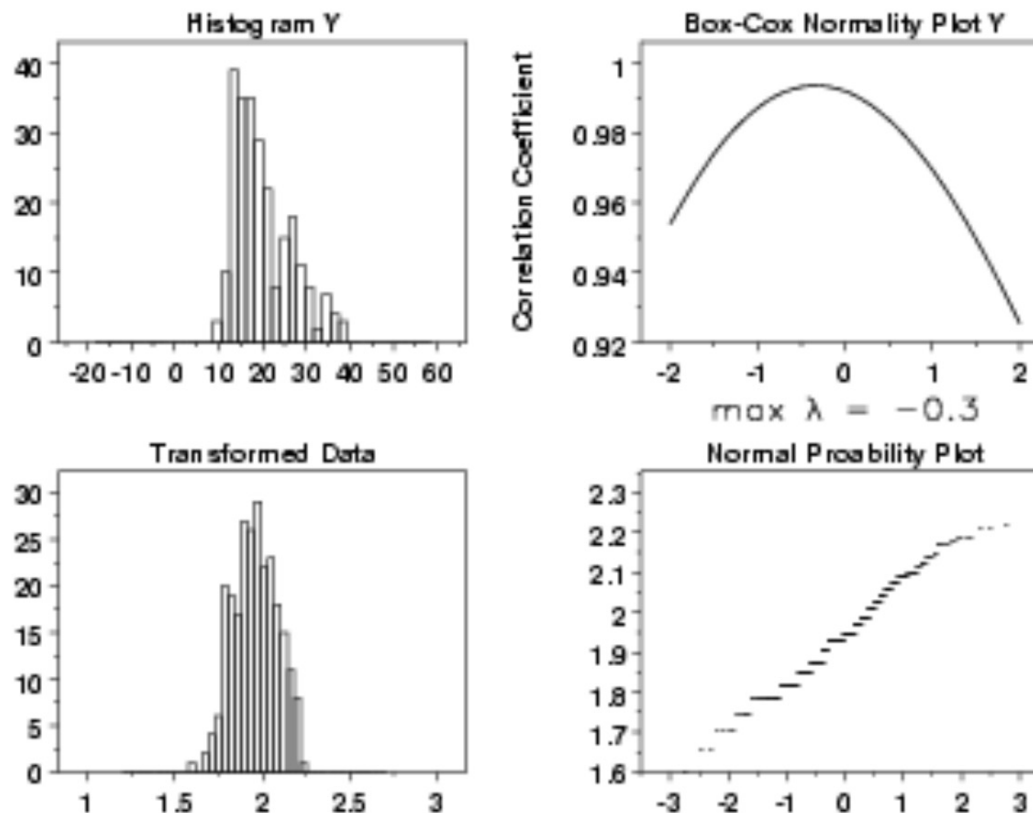# Transformed data or non-parametric?

| Situation | Transformation |
|---|---|
| *right-skewed distribution* | |
| lognormal | x' = log(x) |
| more skewed than lognormal | x' = 1/x |
| less skewed than lognormal | x' = sqrt(x) |
| *left-skewed distribution* | |
| moderately skewed | $x' = x^2$ |
| more skewed | $x' = x^3$ |
| *nonhomogeneous variances* | |
| SD proportional to means | x' = log(x) |
| SD proportional to means$^2$ | x' = 1/x |
| SD proportional to sqrt(means) | x' = sqrt(x) |
| *data as percentages (0-100%)* | p' = arcsin(sqrt(p)) |
| *proportions (p, X/n, 0-1)* | p' = arcsin(sqrt(p)) |

*Watała: Biostatystyka - wykorzystanie metod statystycznych w pracy badawczej w naukach biomedycznych. α-medica press, Bielsko-Biała, 2002.*

# Transformed data or non-parametric?



Box-Cox transformation
Rank-transformation

Box, G. and Cox, D. (1964) An Analysis of Transformations. Journal of the Royal Statistical Society. Series B (Methodological), 26, 211-252.

# Experimental design: selection of the control group

❌ **More than just 1 control group**

| healthy **(HC)** | healthy + comorbidities **(DC)** | disease **(sick?)** | disease + comorbidities **(sicker?)** |

✅

```
              control                                   COVID
             /       \                                 /      \
comorbidities -    comorbidities +        comorbidities -    comorbidities +
```

**Additional advantages of the model
(possible answers to various pre- or post-planned questions within one analysis)**

# Experimental design: multiple comparisons

## Bonferroni correction

$$p_0 = 1 - \sqrt[n]{1 - p} \qquad\qquad p_0 \sim p/n$$

| N = 3 | $p_0 \approx 0.016$ |
| N = 10 | $p_0 \approx 0.005$ |
| N = 30 | $p_0 \approx 0.0017$ |

Very conservative, restrictive

## FDR (Benjamini & Hochberg)
different logic
less restrictive

## Unjustified use of corrections
relevant answer to question of no interest



J. R. Statist. Soc. B (1995)
57, No. 1, pp. 289–300

**Controlling the False Discovery Rate: a Practical and Powerful Approach to Multiple Testing**

By YOAV BENJAMINI† and YOSEF HOCHBERG

Tel Aviv University, Israel

[Received January 1993. Revised March 1994]

Education and debate

### What's wrong with Bonferroni adjustments

Thomas V Perneger

Institute of Social and Preventive Medicine, University of Geneva, CH-1211 Geneva 4, Switzerland
Thomas V Perneger, *medical epidemiologist*

Correspondence to: Dr Perneger perneger@cmu.unige.ch

*BMJ* 1998;316:1236–8

When more than one statistical test is performed in analysing the data from a clinical study, some statisticians and journal editors demand that a more stringent criterion be used for "statistical significance" than the conventional P<0.05.[1] Many well meaning researchers, eager for methodological rigour, comply without fully grasping what is at stake. Recently, adjustments for multiple tests (or Bonferroni adjustments) have found their way into introductory texts on medical statistics, which has increased their apparent legitimacy.[2][3] This paper advances the view, widely held by epidemiologists, that Bonferroni adjustments are, at best, unnecessary and, at worst, deleterious to sound statistical inference.[4][5]

#### Adjustment for multiple tests

Bonferroni adjustments are based on the following reasoning.[1-3] If a null hypothesis is true (for instance, two treatment groups in a randomised trial do not differ in terms of cure rates), a significant difference (P<0.05) will be observed by chance once in 20 trials. This is the type I error, or α. When 20 independent tests are performed (for example, study groups are compared with regard to 20 unrelated variables) and the null hypothesis holds for all 20 comparisons, the chance of at least one test being significant is no longer 0.05, but 0.64. The formula for the error rate across the study is $1-(1-\alpha)^n$, where n is the number of tests performed. However, the Bonferroni adjustment deflates the α applied to each, so the study-wide error rate remains at 0.05. The adjusted significance level is $1-(1-\alpha)^{1/n}$ (in this case 0.00256), often approximated by α/n (here 0.0025). What is wrong with this statistical approach?

#### Problems

**Summary points**

Adjusting statistical significance for the number of tests that have been performed on study data—the Bonferroni method—creates more problems than it solves

The Bonferroni method is concerned with the general null hypothesis (that all null hypotheses are true simultaneously), which is rarely of interest or use to researchers

The main weakness is that the interpretation of a finding depends on the number of other tests performed
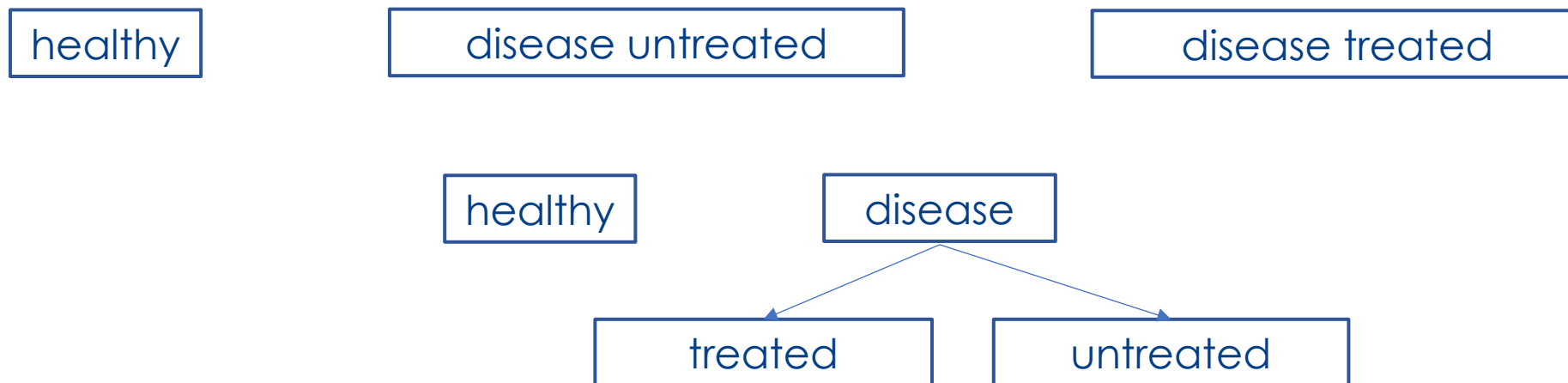
The likelihood of type II errors is also increased, so that truly important differences are deemed non-significant

Simply describing what tests of significance have been performed, and why, is generally the best way of dealing with multiple comparisons

difference in remission rates between two chemotherapeutic treatments could be interpreted as statistically significant or not depending on whether or not survival rates, quality of life scores, and complication rates were also tested. In a clinical setting, a patient's packed cell volume might be abnormally low, except if the doctor also ordered a platelet count, in which case it could be deemed normal. Surely this is absurd, at least within the current scientific paradigm. Evidence in data is what the data say—other considerations, such as how many other tests were performed, are irrelevant.

# Experimental design: several separate analyses or a model?



✓

❌ Breaking down analyses into smaller partial analyses
Multiple comparisons using the simplest possible tests (Student's *t* test)
Required corrections for multiple testing (Bonferroni; FWER; FDR)

✓ Hierarchical (nested) designs
ANOVA & post-hoc tests
General Linear Model, General Additive Model, …

# Lack of appropriate test?

## Permutation (randomization) tests

**TABLE 18.4** | **DRP scores for third-graders**

| Treatment group | | | | Control group | | | |
|---|---|---|---|---|---|---|---|
| 24 | 61 | 59 | 46 | 42 | 33 | 46 | 37 |
| 43 | 44 | 52 | 43 | 43 | 41 | 10 | 42 |
| 58 | 67 | 62 | 57 | 55 | 19 | 17 | 55 |
| 71 | 49 | 54 | | 26 | 54 | 60 | 28 |
| 43 | 53 | 57 | | 62 | 20 | 53 | 48 |
| 49 | 56 | 33 | | 37 | 85 | 42 | |

$$x_T - x_C = 51.472 - 41.522 = \mathbf{9.954}$$

| permutation | statistic |
|---|---|
| 1 | -0,7039 |
| 2 | -1,2505 |
| 3 | 3,1221 |
| 4 | -8,0828 |
| 5 | 2,1201 |
| 6 | 4,9441 |
| 7 | -1,8882 |
| 8 | -0,7039 |
| 9 | -0,2484 |
| 10 | -4,8033 |
| | |
| max = | 4,9441 |
| n (higher) = | 0 |
| **P {(n(higher) / [n(all)+1]} =** | **0,0000** |

# Lack of appropriate test?

## Permutation (randomization) tests

| no permuations | max stat. | n (higher) | p (n higher / n all) |
|---|---|---|---|
| 10 | 4,9441 | 0 | 0,0000 |
| 50 | 13,4162 | 1 | 0,0200 |
| 100 | 10,7743 | 2 | 0,0200 |
| 200 | 13,5073 | 3 | 0,0150 |
| 500 | 16,4224 | 9 | 0,0180 |
| 1000 | 12,8696 | 11 | 0,0110 |
| 10000 | 16,6957 | 129 | **0,0129** |
| 50000 | 18,6998 | 699 | **0,0140** |
| 100000 | 17,6977 | 1409 | **0,0141** |
| 500000 | 18,6087 | 6912 | **0,0138** |
| | | t-test | 0,0264 |
| | | rang-sum | 0,0127 |



**Permutation Distribution**

Hesterberg T & Monaghan, Shaun & S Moore, David & Clipson, Ashley & Epstein, Rachel & H Freeman, W & New York, Company. (2005). Bootstrap Methods and Permutation Tests. Introduction to the Practice of Statistics. 14.
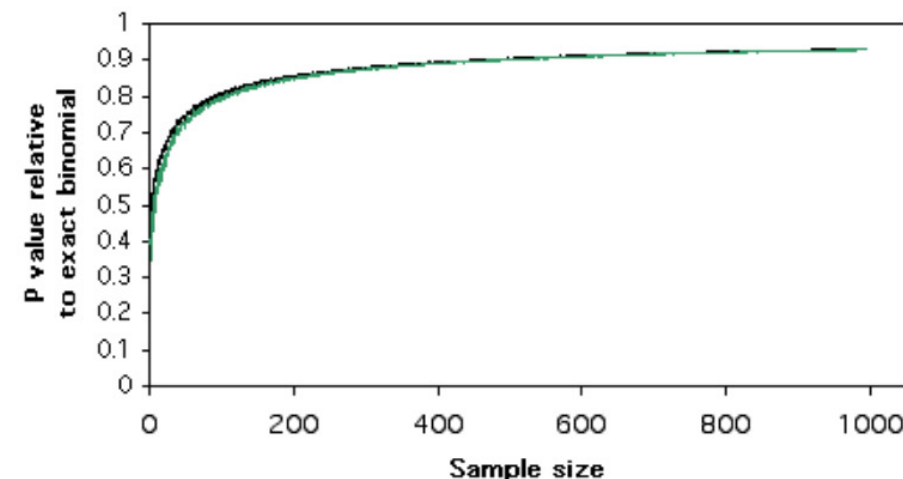
# A permutation test: Fisher's exact test

**Fisher's exact test**
- an example of commonly used permutation test
- alternative to Pearson's chi-squared test
- if the count in any cell is below 5
- but...

For every test in the case of which some statistic is being computed (t, z, F, chi2, …) there is also a permutation version thereof

Instead of comparing the calculated value of test statistic against the values in the corresponding tables, one can compared it against the distribution obtained empirically *en route* resampling

|  | F | M | Total |
|---|---|---|---|
| Smoking | 1 | 9 | 10 |
| Non-smoking | 11 | 3 | 14 |
| Total | 12 | 12 | 24 |



*P* values of chi-square and *G*–tests, as a proportion of the *P* value from the exact binomial test.

www.biostathandbook.com/small.html  (5/06/2019)

# What is your result?

p-value vs. effect size

**p-value is not your result**
*it only tells how reliable your observed effect size is*
*it only tells how often you're about to see the same effect size in repeated experiment*

**Effect size is your result**
*are the changes biologically relevant?*
*are the changes clinically relevant?*

**Publish or parish**
*striving for "p-value"*
*nonsignificant results are not being published*
*sufficient power must be shown in order to publish results with p>0.05*

# An example: Data Torturing

**StatSoft Polska**                                    **DaneWiedzaSukces.pl**

**„IF YOU TORTURE THE DATA LONG ENOUGH, IT WILL CONFESS" –
NAUKOWA DOCIEKLIWOŚĆ A *DATA TORTURING* NA PRZYKŁADZIE
PRACY KLINICZNEJ Z ZAKRESU RADIOTERAPII ONKOLOGICZNEJ**

*Bartłomiej Tomasik*
*Zakład Biostatystyki i Medycyny Translacyjnej, Uniwersytet Medyczny w Łodzi*

# Data torturing

**StatSoft Polska**

DaneWiedzaSukces.pl

**„IF YOU TORTURE THE DATA LONG ENOUGH, IT WILL CONFESS" – NAUKOWA DOCIEKLIWOŚĆ A *DATA TORTURING* NA PRZYKŁADZIE PRACY KLINICZNEJ Z ZAKRESU RADIOTERAPII ONKOLOGICZNEJ**

*Bartłomiej Tomasik*
*Zakład Biostatystyki i Medycyny Translacyjnej, Uniwersytet Medyczny w Łodzi*

**Aim:** Assessment of the impact of IMRT radiotherapy on the development of xerostomia in patients with Head and Neck Cancer

A questionnaire survey + salivary gland scintigraphy
Searching for risk factors of severe xerostomia (grade 3/4) one year after IMRT

**The power analysis:**
assumed statistical power of 80%
assumed level of significance of 5%

~ 100 subjects per group (control; study)
data collection should take 2 years

# Data torturing

**2 years of data collection**
- only 53 subjects
- 30 scintigraphic examinations
- 40 questionnaires

**a posteriori power analysis:**
< 20 % (max **!**)

**So what about it now?**
Analyse?
Continue collecting data?

Tabela 1. Charakterystyka grupy włączonej do badania.

| Zmienna | Kategoria | Liczba (%) |
|---|---|---|
| Płeć | Kobiety | 12 (22,64%) |
| | Mężczyźni | 41 (77,36%) |
| Cecha T | T1 | 1 (1,89%) |
| | T2 | 20 (37,74%) |
| | T3 | 24 (45,28%) |
| | T4 | 8 (15,09%) |
| Cecha N | N0 | 33 (62,26%) |
| | N1 | 6 (11,32%) |
| | N2 | 14 (26,42%) |
| Lokalizacja guza | Krtań/gardło dolne | 25 (47,17%) |
| | Jama ustna/gardło środkowe | 28 (52,83%) |
| Pełne badanie scyntygraficzne | Tak | 30 (56,60%) |
| | Nie | 23 (43,40%) |
| Wypełnione kwestionariusze QLQ H&N-35 | Tak | 40 (75,47%) |
| | Nie | 13 (24,53%) |

# Data torturing

**Xerostomia severity (degree) vs. risk group (location)**
- Fisher's exact test
- p=0.048 (**!**)


- one-sided test (**!**)

- trying to find an argument for usage of one-sided test *post factum*
- it was not taken into account while planning the experiment

Tabela 2. Podsumowująca tabela dwudzielcza z wynikami testu Fishera dla porównania występowania nasilonej kserostomii w zależności od grupy ryzyka.

| Lokalizacja | Podsumowująca tabela dwudzielcza | | |
|---|---|---|---|
| | Kserostomia ≥3 | Kserostomia <3 | Wiersz - razem |
| Krtań/gardło dolne | 4 | 16 | 20 |
| %kolumny | 28,57% | 61,54% | |
| %wiersza | 20,00% | 80,00% | |
| %ogółu | 10,00% | 40,00% | 50,00% |
| Jama ustna/gardło środkowe | 10 | 10 | 20 |
| %kolumny | 71,43% | 38,46% | |
| %wiersza | 50,00% | 50,00% | |
| %ogółu | 25,00% | 25,00% | 50,00% |
| Kolumna - razem | 14 | 26 | 40 |
| % z całej grupy | 35,00% | 65,00% | 100,00% |
| Dokładny jednostronny test Fishera | p=0,048 | | |

# Data torturing

**Some of the questionnaires were repeated after several years**

Q Cochran's test
- p=0.083
- observable trend (**!?**)

- initial lack of statistical power
- decreasing number of subjects in analysis

- one-sided test (**!**)
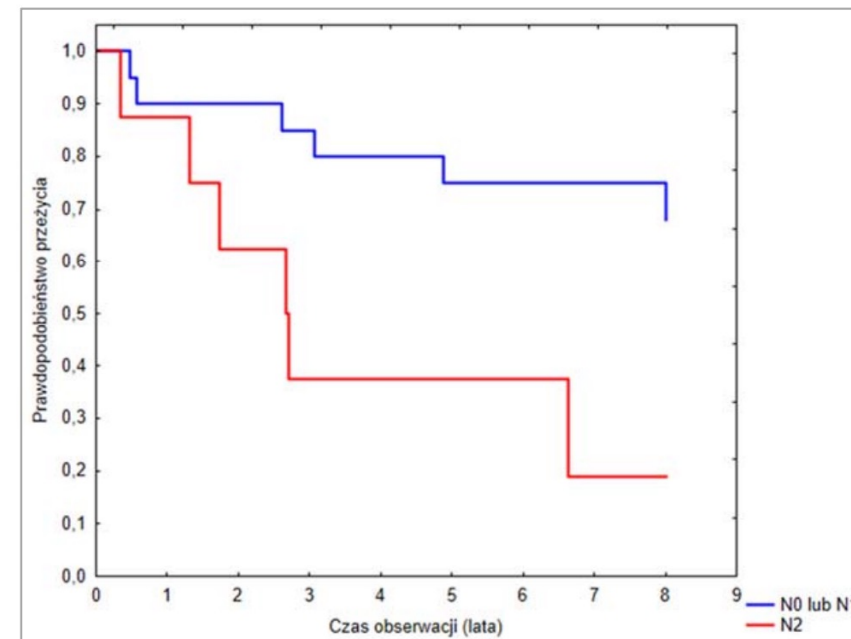  (trying to find an argument *post factum*)

Tabela 3. Wyniki testu Q Cochrana porównującego subiektywną ocenę kserostomii w pierwszym i drugim badaniu ankietowym.

| Lokalizacja | Podsumowująca tabela dwudzielcza | | |
|---|---|---|---|
| | Kserostomia ≥3 | Kserostomia <3 | Wiersz - razem |
| Pierwsze badanie ankietowe | 6 | 9 | 15 |
| %kolumny | 40,00% | 60,00% | |
| %wiersza | 40,00% | 60,00% | |
| %ogółu | 20,00% | 30,00% | 50,00% |
| Drugie badanie ankietowe | 9 | 6 | 15 |
| %kolumny | 60,00% | 40,00% | |
| %wiersza | 60,00% | 40,00% | |
| %ogółu | 30,00% | 20,00% | 50,00% |
| Kolumna - razem | 15 | 15 | 30 |
| % z całej grupy | 50,00% | 50,00% | 100,00% |
| Test Q Cochrana | Q=3,00 | | |
| | p=0,083 | | |

# Data torturing

**A survival analysis was planned in addition**

- the log-rank test

- after several trials and divisions of the whole study groups into subgroups, a statistically significant result was obtained (p=0.021)

- researchers were satisfied ☺

- have they planned the survival analysis?
- have they checked whether the study provides enough power for survival analysis?
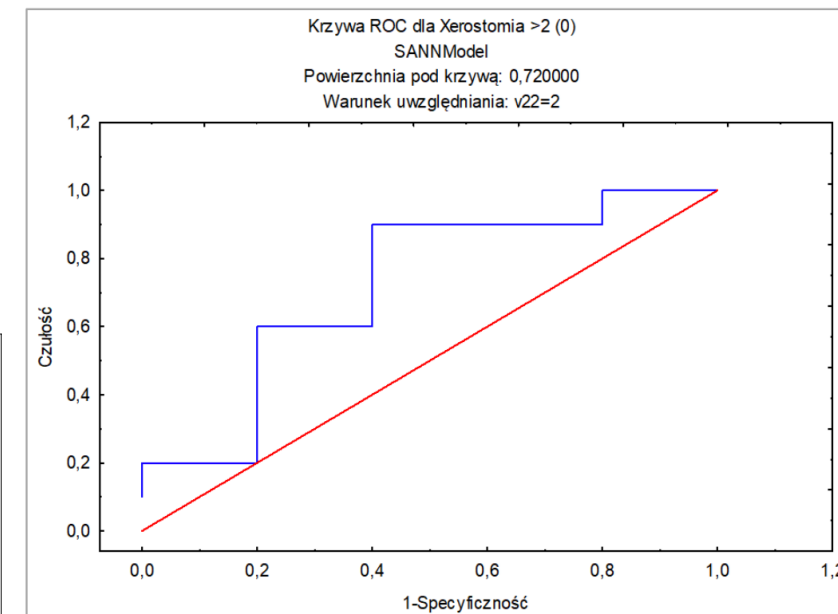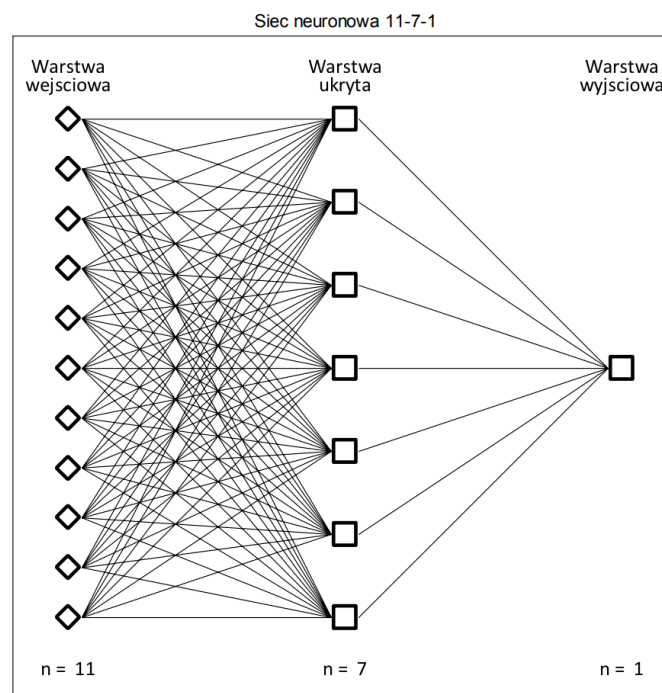


| | Oceny parametrów | | |
|---|---|---|---|
| | Poziom efektu | Poziom p | HR (95%CI) |
| Zajęcie węzłów chłonnych | N2 | 0,021 | 3,92 (1,24-12,45) |

Rys. 6. Krzywe Kaplana-Meiera przedstawiające różnice w przeżyciu w zależności od zajęcia węzłów chłonnych oraz tabela z wynikami analizy proporcjonalnego hazardu Coxa.

# Data torturing

**Searching for predictors of xerostomia worsening**

- several ROC curve analyses

- no success, therefore the authors used neural networks
  (11 $L_0$ neurones; 7 layers)

- obtained network provided accuracy of 80% and 0.72 area under the curve (AUC)

- researchers were fully satisfied ☺

- lack of interpretability?



Siec neuronowa 11-7-1

Warstwa wejsciowa | Warstwa ukryta | Warstwa wyjsciowa

n = 11    n = 7    n = 1



Krzywa ROC dla Xerostomia >2 (0)
SANNModel
Powierzchnia pod krzywą: 0,720000
Warunek uwzględniania: v22=2

# Data torturing

**Problems**

- excessive optimism in the planning phase regarding data collection (not taking into account that some patients may not meet the inclusion criteria)

- very low statistical power

- analyses within subgroup that were not previously planned (additional subdivisions of groups)

- conducting previously unplanned analyses (survival analysis, ROC analysis)

- multiple hypotheses testing without appropriate corrections

- abuse of more and more complicated analytical methods in order to prove the assumed thesis („hypothesis driven science"; Fisher's test → Neural Network)

- the blind pursuit of statistically significant results without considering their clinical significance

# Data torturing

**What should we do?**
- it cannot be avoided

But we should be cautious to:

- were the other (new?) hypotheses made before or during the experiment?

- were all the subgroup analyses planned in advance, prior to experiment?

- are all the obtained results related to the main hypothesis of the study?

- is it not necessary, at lease in the case of some statistical tests, that appropriate multiple testing corrections be applied?

**Thanks for your attention...**

# Literature

## CO MOŻNA WYCISNĄĆ Z TYCH DANYCH?

*Andrzej Stanisz, Collegium Medicum Uniwersytetu Jagiellońskiego w Krakowie, Zakład Biostatystyki i Informatyki Medycznej*

## JAK SKUTECZNIE WYKORZYSTYWAĆ METODY STATYSTYCZNE W PLANOWANIU I PRZEPROWADZANIU EKSPERYMENTU NAUKOWEGO?

*Cezary Watała, Uniwersytet Medyczny w Łodzi, Zakład Zaburzeń Krzepnięcia Krwi KDL; Uniwersytecki Szpital Kliniczny nr 2 im. WAM*

## JAK PLANOWAĆ DOŚWIADCZENIA NAUKOWE Z WYKORZYSTANIEM METOD STATYSTYCZNYCH?
## TESTOWANIE HIPOTEZ STATYSTYCZNYCH: MIĘDZY MOCĄ STATYSTYCZNĄ A NIEMOCĄ DECYZYJNĄ

*Cezary Watała, Uniwersytet Medyczny w Łodzi, Zakład Zaburzeń Krzepnięcia Krwi; Uniwersytecki Szpital Kliniczny nr 2 im. WAM*

5

## Sample size and significance – somewhere between statistical power and judgment prostration

Cezary Watała

Department of Haemostatic Disorders, Medical University of Lodz, Poland

Corresponding author:
Prof. Cezary Watała
Department of Haemostatic Disorders
Medical University of Lodz
Medical University Hospital No. 2
113 Zeromskiego Street
90-549 Lodz, Poland
Phone: +48 42 6393471
Fax: +48 42 6787567
E-mail: cwatala@csk.umed.lodz.pl

**Abstract**

When performing scientific research we are so "embraced" to use the tool of inductive logic in our reasoning that we often express more generalized opinions on the population of interest based on relatively small sample(s) of a general population. What we take care about in such situations is that chosen segments are representative for a whole set of elements in the general population. To cope with such a demand we always want to know how large our selected subpopulation should be to enable us to detect the experimental effect of interest not only at a certain level of significance, but also with the highest possible power of statistical reasoning. Thus, when designing our experiment, we have to compromise between a sample size not too small to ensure that our sample is sufficiently representative, and not too large to benefit from the sampling procedure at all. The tools for the estimation of minimum required sample size and the analysis of power, which help us to make quick decisions on how to compromise reasonably between significance, statistical power and sample size, are discussed in this paper.

# Literature