# B4M36SAN BE4M36SAN Seminar/Tutorial

Alikhan Anuarbekov

Ali

anuarali@fel.cvut.cz

## cw.fel.cvut.cz

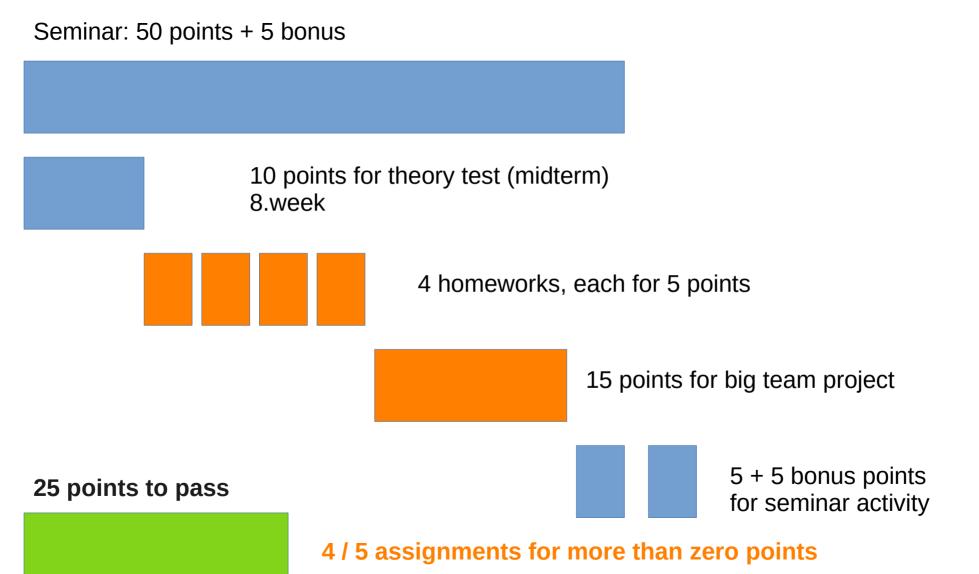
Total: 100 points

Seminar: 50 points + 5 bonus

Exam: 50 points

During 13 weeks of semester, programming, theory, project

At the end of the course, 1 big exam theory from lectures only



Every seminar (a total of 10 seminars):

15-25 minutes theoretical review -----> not covers all of lecture!

5-15 minutes assignment presentation

50-70 minutes for activity assignment solution

Attendance is obligatory, max. 2 missing tutorials

BRUTE -> solved and submitted till the end of day, earn 1 point (0.5 point if partially)

#### Big team project (Final assignment section in courseware)

- Start looking for 4 people team NOW
- Try to find a question/dataset that is interesting for you
- Look for previous year best projects for inspiration
- Ask question before the deadline for team creation (3.11)
- If no team on 3.11 (midterm week), then I will forcefully assign a team for you

#### **Tips and recommendations:**

- Do not focus on single question, try to combine two or more questions
  - Example: Not only look for medical question, but try to add socio-economy
- Try to find more than one dataset
  - Example: UNICEF + World Bank
- Start working early, at least to formulate a question



#### **Rstudio vs Python**

- In this tutorials I will go with R/Rstudio programming language
  - Has all math and formulas inside, just call functions
  - Syntax is the similar to Matlab/Javascript/....
- If you want, you could go with Python
  - I will not check you homework, Jan Blaha will, later feedback
  - I cannot help you with python questions
  - But python is much more easier and familiar :)
- R is used only by statisticians, python by everyone

Type in google: Rstudio download and R download

https://cran.r-project.org/

R (programming lang) e.g. same as Python

https://posit.co/download/rstudio-desktop/

RStudio (editor, IDE) e.g. same as Intellij/PyCharm File Edit Code View Plots Session Build Debug Profile Tools Help - Addins -R version -Go to file/function pp pvalue.cpp × Q list.in\$`circInte\>>> \_\_\_\_ workflow.Rmd × dfs.in\$circInter × **Environment History Connections Tutorial** Import Dataset ▼ | 1.38 GiB ▼ | 🍼 List → | (C) → Knit on Save Global Environment -■ Outline a Visual Source ---**Files Packages** Help Presentation Viewer return(list(pvalue=list.pval.out, 512 tpr=pvals1.tpr, fpr=pvals1.fpr, 513 514 precision=pvals1.prec, recall=pvals1.recall, Home > Desktop > work > R version > source > GPACDA 515 f1=pvals1.f1)) Name Size Modified 516 - 7 517 downloadMiRNAGO.R 801 B Nov 10, 202 518 - list.to.pr <- function(pvals.list.in){ downloadMultiMiR.R 1.3 KB Nov 10, 202 519 gene2disease.csv 118.4 KB Nov 10, 202 pvals.list.in\$pvalue[is.na(pvals.list.in\$pvalue)] <- 1</pre> 520 521 as.data.frame(pvals.list.in) %>% goterms.R 906 B Nov 10, 202 522 mutate(correct=as.factor(correct), score=-log(pvalue)/50) %>% goterms.RData 4.3 MB Nov 10, 202 523 pr\_curve(correct, score) mirna-go-terms.csv 84.3 KB Nov 10, 202 524 - } mirna-mrna-interactions.csv 15 MB Nov 10, 202 R Markdown \$ Chunk 9 \$ 575:1 mirnaGO.R 1.4 KB Nov 10, 202 Terminal × Background Jobs × Console pvalue.cpp 5.9 KB Nov 10, 202 R 4.1.2 · ~/Desktop/work/R version/ **README.txt** 5.2 KB Nov 10, 202 + } > rcPlot(list( 8 run.R 1.2 KB Nov 10, 202 "GEO"=all.geo.pvals, run.sh 207 B Nov 10, 202 "circInter"=all.inter.pvals, runDisease.sh "rand"=all.rand.pvals 154 B Nov 10, 202 + + )) diseasedata Called from: rcPlot(list(GEO = all.geo.pvals, circInter = all.inter.pvals, diseases-small rand = all.rand.pvals)) Browse[1]> c **GPACDA-paper-experiments** > install.packages("BSDA") result2024

R version - RStudio

### Rstudio/R basic programming

LOGICAL -> INTEGER -> NUMERIC -> CHARACTER

"boolean" "integer" "double" "string"

a = 1
class(a) = type of variable
vector = array of same type

https://www.w3schools.com/r/r\_get\_started.asp

https://www.datacamp.com/courses/free-introduction-to-r

Download today's tutorial script (just click, no programming):

#### **Tutorials**

. . . .

Т	Date	Teacher	Contents	Materials
1	23.9.	JB, AA, JK	Statistical testing, t- test, significance, power of the test.	san_intro.zip

### Rstudio/R library installation + problems

> I want to install a library that will compute z-test



```
Using PKG CFLAGS=
Using PKG LIBS=-lcurl
                ----- [ANTICONF] ------
Configuration failed because libcurl was not found. Try installing:
 * deb: libcurl4-openssl-dev (Debian, Ubuntu, etc)
 * rpm: libcurl-devel (Fedora, CentOS, RHEL)
If libcurl is already installed, check that 'pkg-config' is in your
PATH and PKG CONFIG PATH contains a libcurl.pc file. If pkg-config
is unavailable you can set INCLUDE DIR and LIB DIR manually via:
R CMD INSTALL --configure-vars='INCLUDE_DIR=... LIB_DIR=...'
                ----- [ERROR MESSAGE] -
<stdin>:1:10: fatal error: curl/curl.h: No such file or directory
compilation terminated.
ERROR: configuration failed for package 'curl'
```

\* removing '/usr/local/lib/R/site-library/curl'

'/tmp/RtmpNNghX1/downloaded packages'

installation of package 'curl' had non-zero exit status

The downloaded source packages are in

In install.packages("curl") :

Warning message:

## Rstudio/R library installation solution:



R library install fatal error: curl/curl.h: No such file or directory







### **Statistics 101 (basics)**

Prerequisites are in courseware (**Exam** section):

6. Sample questions pertaining to prerequisites:

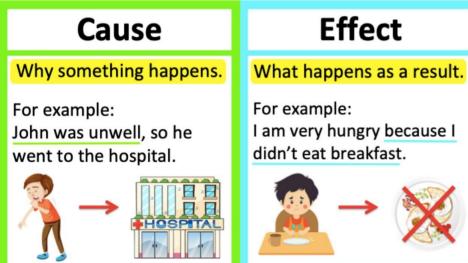
mstat\_min.pdf.

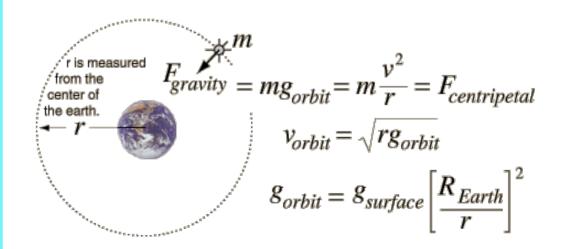
Today we will start with basics and go over it fast

If you don't understand anything from PDF, do a self study!

#### What is statistics?

- A typical scientific approach was to find the cause with the effect
- If we find all causes, we can 100% predict the effect
- A topic for math, algebra, physics, chemistry, ....





#### What is statistics?

In complex systems/objects, we cannot physically predict the effect with 100% accuracy

Laplace's demon principle
"If you were to measure all atom's position and speed,
you can predict everything in future and
know everything in past"

> but no one can measure all atoms in universe



Practical example:

Can you with 100% accuracy predict whether the drug will work on patient?

Human body: 3.2 billion DNA bases, 30 trillion cells -> any DNA problem, cancer cells, ... Surroundings: around 20 – 30 people interacts with you daily, you measure their bodies as well + Human's diet, behavior (smoking, sport,...), location (tropical disease?),.....

### **Bayesian statistics**

Instead of 100% accurate answer give the most probable one

> yes, if you have a 1 / 1,000,00 chance rare genetic disease, our drug will harm you, but in 999,999 / 1,000,000 cases it will help

Based on that principle:

Reality: complicated mechanism/equation that is impossible to calculate

$$\rho_t C_{p_t} \frac{\delta T_t}{\delta t}$$

$$= k_t \left[ \left( \frac{\delta^2 T_t}{\delta x^2} \right) + \left( \frac{\delta^2 T_t}{\delta y^2} \right) + \left( \frac{\delta^2 T_t}{\delta z^2} \right) \right]$$

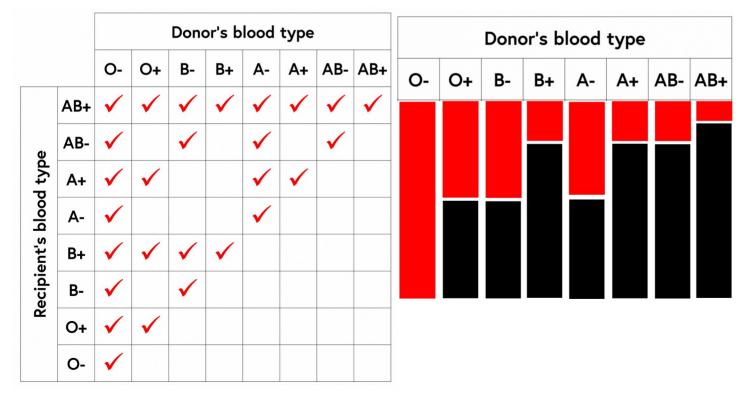
$$+ \rho_b C_{p_b} \dot{V} (T_A - T_V) + \dot{Q}_m$$



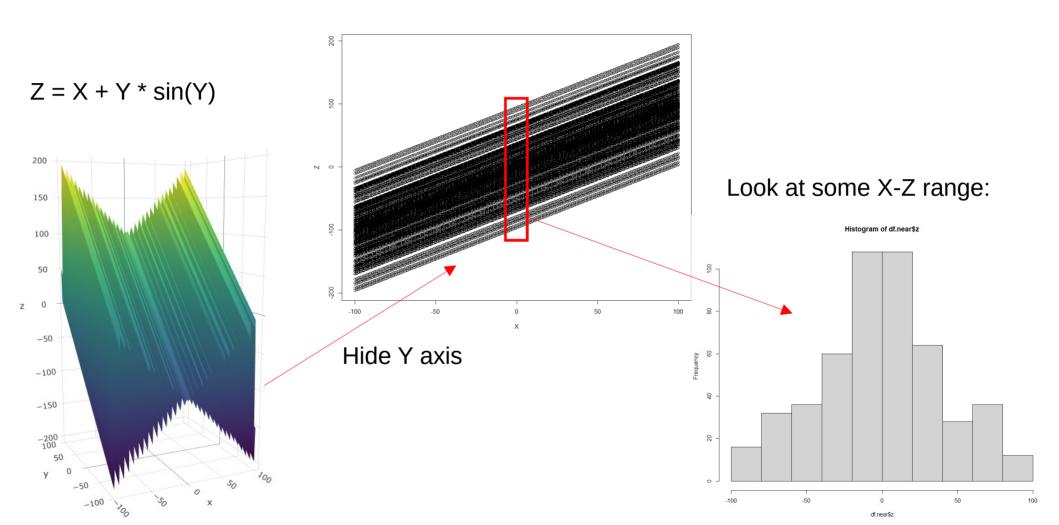
disease chance  $\approx$  age \*0.5disease chance = age  $*0.5 + \Delta_{error}$ 

### **Example (2D discrete space -> 1D distribution):**

You need to performed a blood transfusion, but you don't have time to measure patient blood type. You only know donor's blood type. Then you have an uncertainty about transfusion success



## **Example (2D continuous space -> 1D distribution):**

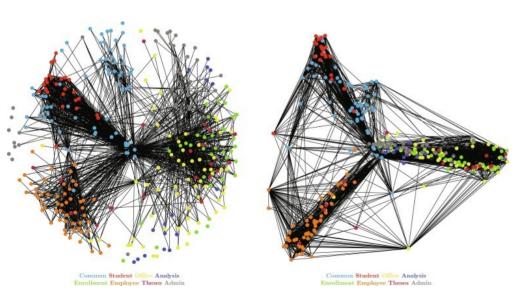


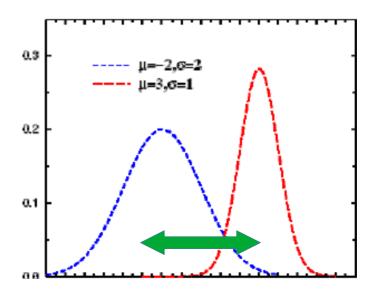
### **Bayesian statistics**

**Distribution** -> if we cut parts of complicated reality equation for subset of causes, lets say we want to predict if the drug will work based on age only, then **sometimes it works, sometimes not**.

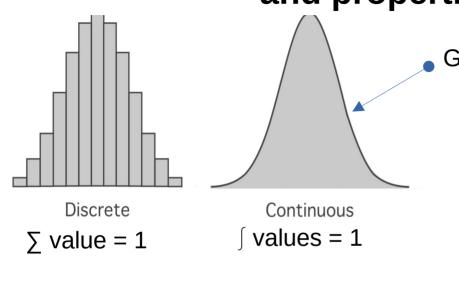
100% accurate complicated formula is some complicated curve in 1,000,000D (or even infinitely complex, intractable)

Statistical distribution if only age is considered (blue are age values when drug works, red are when drug does not work)
72,5% accurate





# Distribution can have different shapes, formulas and properties



$$f(x)=rac{1}{\sigma\sqrt{2\pi}}e^{-rac{1}{2}\left(rac{x+\mu}{\sigma}
ight)^2}$$

$$f(x) = \begin{cases} 0, & x < 0 \\ \lambda & \lambda^{k}, & x \ge 0 \end{cases}$$

Poisson:

Exponential:

$$f(x) = \frac{\lambda^r}{x!} e^{-\lambda}$$

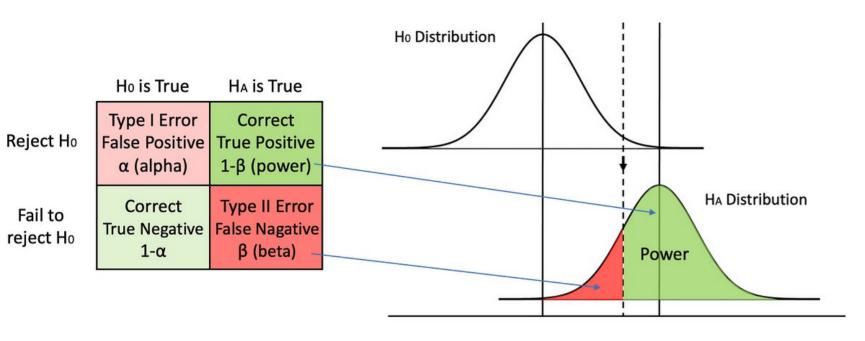
All of them have some parameters that we need to MEASURE / ESTIMATE

$$\mu = \frac{\sum \chi}{N}$$

What do these formulas say about distributions?

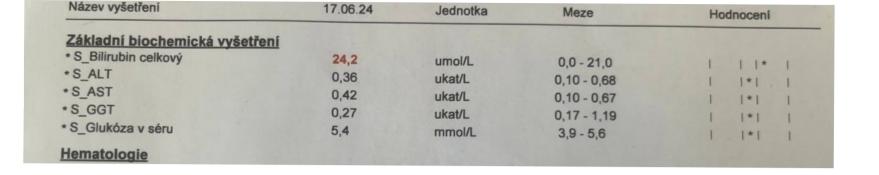
$$\sigma = \sqrt{rac{\sum (x_i - \mu)^2}{N}}$$

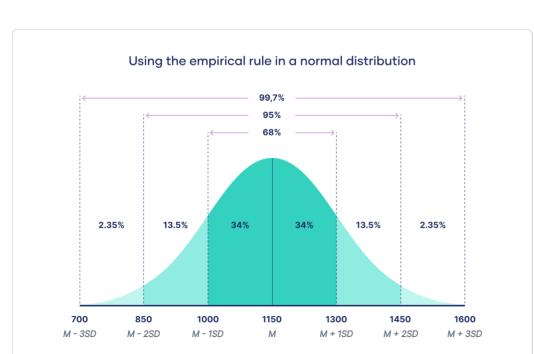
### **Bayesian statistics and Hypothesis testing**



Also – when choosing between two distributions (as for drugs), it **DOES** matter what is the distribution of interest.

Alternative Hypothesis ( $H_A$ ) = some new fact that I want to proof (discovery) Null Hypothesis ( $H_0$ ) = a typical assumption/belief that is held before your discovery





- I went to a doctor and they took a sample of my blood
- They measured multiple biomolecules
- If you have too many or too low of any molecules, something is wrong with you
- Scientists decided that **alpha=5%** is the border between "its okay" and "it is not okay"
- This is called Z-test
- It needs scientists to **measure** a lot of people to **estimate** (construct) a distribution e.g. → measure mean and variance:

$$\mu = rac{\sum \chi}{N}$$
  $\sigma = \sqrt{rac{\sum (x_i - \mu)^2}{N}}$ 

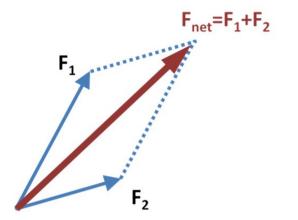
- > Do you know what is parametric and non-parametric test?
- > Do you know what is p-value?
- > What is confidence interval?
- > What is quantile? One-sided vs two-sided?

Show them in figure on previous slide

- 6. Sample questions pertaining to prerequisites:
  - mstat\_min.pdf.

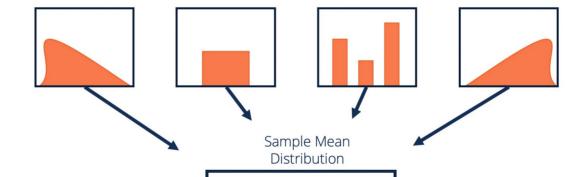
#### Why Gaussian normal distribution?

- Why we use a Gaussian normal? What if there is another distribution?
- Maybe, but typically when you have multiple factors/causes, they are typically summed



Example: Physics sum of multiple forces

 Then, when we compute their mean (sum / number), it is mathematically Gaussian! (Central limit theorem)



Normal Distribution

 Consequence – to be more precise, even in previous test they didn't just measured my biomolecules level, they measured it 3-5 times from big blood sample and took an average

### **Problem with Gaussian for small sample number**

$$\mu = \frac{\sum \chi}{N}$$

Mean is okay, it is calculated same way

$$=\sqrt{rac{\sum (x_i-\mu)^2}{N}}$$

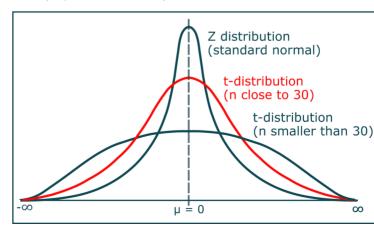
 $\sigma = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$  But if scientists do not have large enough number of people, variance will be very different from real one -> very big impact on correctness

#### t-student test (t-test)

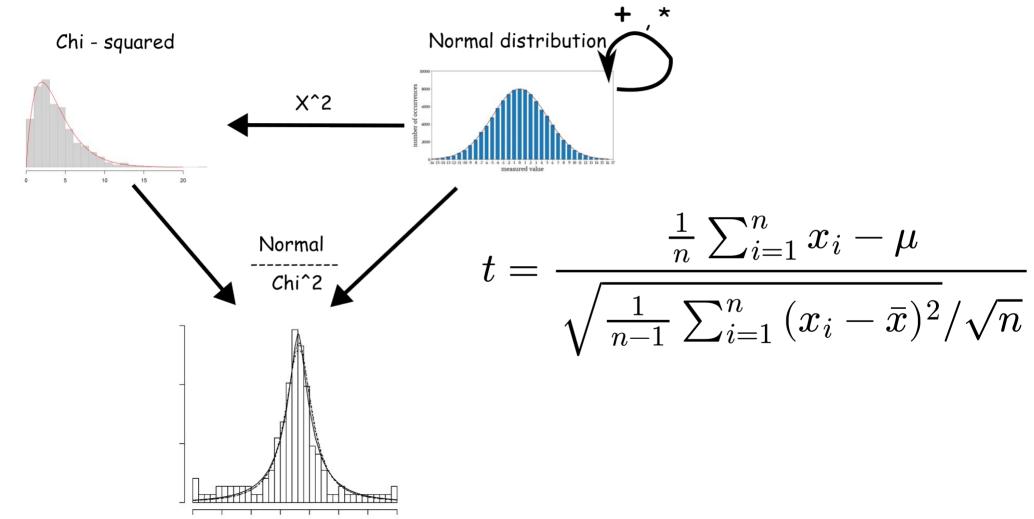
$$\frac{X-\mu}{\frac{1}{\sqrt{N}}S_H} \in not \ normal \ anymore \in t-student(df=N-1)$$

$$\frac{1}{\sqrt{N}}S_H$$

Instead of 
$$\sigma^2$$
 use  $S_H^2 = \frac{\sum (X_i - \mu)^2}{N - 1}$ 



# Where does t-distribution and other test distributions come from?



### **Statistics 101 (basics)**

Prerequisites are in courseware (**Exam** section):

6. Sample questions pertaining to prerequisites:

mstat\_min.pdf.

Today we will start with basics and go over it fast

If you don't understand anything from PDF, do a self study!

### Today's homework (voluntary) for 1 point

#### **Tutorials**

Т	Date	Teacher	Contents	Materials	
1	22.9.	JB, AA, JK	Statistical testing, t-test, significance, power of the test.	san_intro.zip, r_setup.zip, pres-1.pdf	

#### Further questions and tasks:

- 1. the power of the test that evaluates whether "the mean blood pressure after treparameters above to increase this power, report when you reach 80%,
- 2. see what happens with the test when violating assumptions with outliers, inject
- 3. violate other assumptions to see what happens (a distribution different from no
- 4. show the setting in which the central limit theorem can be applied, gradually ch

#### Goals:

- \*) Install Rstudio/R and run some code by Ctrl + Enter
- \*) Answer some theoretical questions from basic statistics
- \*) Or you could look into HTML file with plots/figures if you do not want to run R code
- \*) Submit the modified Rmd or plain text file with answers

# You want to ask an anonymous question? You have some problem with the seminar?

https://askany.me/session/acb96d8c-86d5-4763-86da-a2c88724

← Back to Home Live Q&A Session (Owner) - Mode: normal						
Mode: Normal Copy Session Link Session Owner						
Your Question						
I do not like the way you teach. I want to solve more examples, not more theory!						
Your Name (optional)						
Enter your name						
✓ Ask Anonymously						
Submit Question						

## End of statistics 101 Good luck!