# Day 1 - Part I
# Course introduction

1. Course Structure & Schedule, Course Materials, Exam
2. Motivations for the choice of Methods
3. What is OCBE?

Valeria Vitelli

Oslo Centre for Biostatistics and Epidemiology (OCBE)

Department of Biostatistics, UiO

valeria.vitelli@medisin.uio.no

MF9130E – Introductory Course in Statistics

08.04.2024

# Course Structure

## Structure Week 1: Flipped Classroom

- afternoon:
  - ▶ lecture on a topic + Q & A session
  - ▶ some free time for reflection / group work
- morning after: practical lab session on the same topic

## Structure Week 2: More Traditional Approach

- Half-day blocks on a specific topic
- Structure of each block:
  - ▶ lecture on a topic
  - ▶ practical lab session on the same topic

# Course Schedule

## Week 1

|    | Monday | Tuesday | Wednesday | Thursday | Friday |
|----|--------|---------|-----------|----------|--------|
| AM |        | Lab 1<br>Lab 2 | Distributions<br>Lab 4 | Lab 5<br>t-test | Lab 6<br>tables |
| PM | Intro Course<br>Data & Descript<br>Statistics | Intro Prob<br>Diagnostics<br>Lab 3 | Inference<br>part 1<br>t-test | Inference<br>part 2<br>table analysis | |

## Week 2

|    | Monday | Tuesday | Wednesday | Thursday |
|----|--------|---------|-----------|----------|
| AM | Non-parametrics<br>Lab | Study designs | Regression 2<br>Lab | Logistic<br>Lab |
| PM | Sample size<br>Power<br>Lab | Regression 1<br>Lab | Regression 3<br>lab<br>Course Summary | Survival<br>Lab |

# Overview for Day 1: "Data and Descriptive Statistics"

### This afternoon: Lectures in flipped classroom style

- **Introduction** to this course
- **Data and statistics** in medicine: Introduction and motivation
- **Descriptive statistics**
  - ▶ Data presentation: Central measures & Measures of variation
  - ▶ Graphical presentation of data
- **Self-study and Group-work**

### Tomorrow morning: Labs in flipped classroom style

- Introduction to statistical computing with **R**
- Descriptive statistics with **R**

### Course textbook chapters:

- Kirkwood and Sterne chapters 2-4
- Aalen chapters 1 and 2

# Links and Course Material

- **Course webpages:**
  `https://ocbe-uio.github.io/teaching_mf9130e/`
- We will mainly use the course webpages for all information and access to material. The webpages will be continuously updated throughout the course.

## Links and Course Material – Alternatives

- **Canvas room:** We will not use the Canvas room a lot, but Canvas is used for **emails** and general communication. Please let us know asap, if you do not have access to Canvas!

- **Official UiO course pages** with schedule, literature and details on admission rules, exam etc: `https://www.uio.no/studier/emner/medisin/med/MF9130E/`

# Computer exercises in R (starting tomorrow morning)

- You will need to have a laptop computer with access to R and RStudio for the labs.

- We advise that you install R/RStudio on your own laptop.

- Alternatively, you could register for a (free trial) account on a Posit Cloud server.

- See here for instructions: `https://ocbe-uio.github.io/teaching_mf9130e/get_started/get_started.html`

- Note: You can also access R/RStudio through the UiO Programkiosk: `https://www.uio.no/english/services/it/home-away/kiosk/`.

# Homework for tomorrow morning

- Go through the instructions above to get working access to R and RStudio. There will be a detailed introduction to R and RStudio tomorrow morning.

> 💡 **Note**
>
> It is recommended to have R and Rstudio installed on your laptop, this is because you have a better control of where you prefer to download data and course material. This is also useful when you want to analyse your own datasets. For example, you might have to upload datafiles to the server for Posit Cloud to work.
>
> However, if there is a problem with the installation, you can use Posit Cloud as an alternative.
>
> On Tuesday morning we will see if most people can successfully make R run on their laptop and make necessary adjustments.

# Exam

- Take-home exam.
- Will be published via Inspera at the end of the course.
- To be submitted within a specified deadline (4 weeks after the end of the course).
- A passed exam is required to get the course approved.

- More details on the third day of week 2.

Main course textbook: Kirkwood and Sterne (2003)



- Betty R. Kirkwood and Jonathan A. C. Sterne. **Essential Medical Statistics**. Second edition, Blackwell Science Ltd, 2003
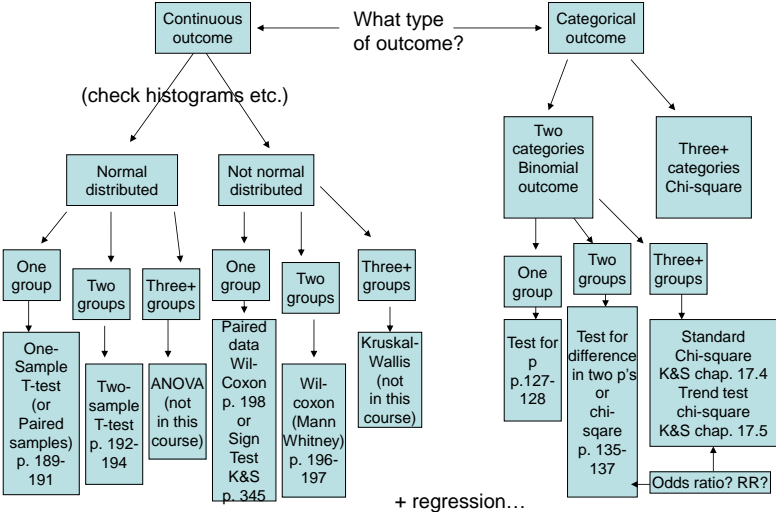- www.blackwellpublishing.com/essentialmedstats/

Norwegian alternative: Aalen (ed) *et al* (2006)

- Odd O. Aalen (red), Arnoldo Frigessi, Tron Anders Moger, Ida Scheel, Eva Skovlund, Marit B. Veierød. **Statistiske metoder i medisin og helsefag**. Gyldendal Akademisk 2006
- `www.med.uio.no/imb/studier/ressurser/statistikk/`
  `statistikkressurser-shs/aalen.html`

Page numbers in Aalen
Except K&S=Kirkwood &
Sterne

## Diagram for test choice

What type of outcome?

**Continuous outcome** ← → **Categorical outcome**

(check histograms etc.)

### Continuous outcome

**Normal distributed**

- One group
  - One-Sample T-test (or Paired samples) p. 189-191
- Two groups
  - Two-sample T-test p. 192-194
- Three+ groups
  - ANOVA (not in this course)

**Not normal distributed**

- One group
  - Paired data Wil-coxon p. 198 or Sign Test K&S p. 345
- Two groups
  - Wil-coxon (Mann Whitney) p. 196-197
- Three+ groups
  - Kruskal-Wallis (not in this course)

### Categorical outcome

**Two categories Binomial outcome**

- One group
  - Test for p p.127-128
- Two groups
  - Test for difference in two p's or chi-sqare p. 135-137
- Three+ groups
  - Standard Chi-square K&S chap. 17.4 Trend test chi-square K&S chap. 17.5

Odds ratio? RR?

**Three+ categories Chi-square**

+ regression…

# Many of the methods we cover can be seen as **linear models**.

- `https://lindeloev.github.io/tests-as-linear/`
- Regression models as well as most statistical tests:

## Common statistical tests are linear models

Last updated: 28 June, 2019. Also check out the *Python version*)

See worked examples and more details at the accompanying notebook: https://lindeloev.github.io/tests-as-linear

| | Common name | Built-in function in R | Equivalent linear model in R | Exact? | The linear model in words | Icon |
|---|---|---|---|---|---|---|
| **Simple regression: lm(y ~ 1 + x)** | **y is independent of x**<br>P: One-sample t-test<br>N: Wilcoxon signed-rank | t.test(y)<br>wilcox.test(y) | lm(y ~ 1)<br>lm(signed_rank(y) ~ 1) | ✓<br>for N >14 | One number (intercept, i.e., the mean) predicts **y**.<br>- (Same, but it predicts the *rank* of **y**.) | |
| | **P**: Paired-sample t-test<br>**N**: Wilcoxon matched pairs | t.test(y₁, y₂, paired=TRUE)<br>wilcox.test(y₁, y₂, paired=TRUE) | lm(y₂ - y₁ ~ 1)<br>lm(signed_rank(y₂ - y₁) ~ 1) | ✓<br>for N >14 | One intercept predicts the pairwise **y₂-y₁** differences.<br>- (Same, but it predicts the *signed rank* of **y₂-y₁**.) | |
| | **y ~ continuous x**<br>P: Pearson correlation<br>N: Spearman correlation | cor.test(x, y, method='Pearson')<br>cor.test(x, y, method='Spearman') | lm(y ~ 1 + x)<br>lm(rank(y) ~ 1 + rank(x)) | ✓<br>for N >10 | One intercept plus x multiplied by a number (slope) predicts **y**.<br>- (Same, but with *ranked* x and y) | |
| | **y ~ discrete x**<br>P: Two-sample t-test<br>P: Welch's t-test<br>N: Mann-Whitney U | t.test(y₁, y₂, var.equal=TRUE)<br>t.test(y₁, y₂, var.equal=FALSE)<br>wilcox.test(y₁, y₂) | lm(y ~ 1 + G₂)ᴬ<br>gls(y ~ 1 + G₂, weights=...)ᴮ ᴬ<br>lm(signed_rank(y) ~ 1 + G₂)ᴬ | ✓<br>for N >11 | An intercept for **group 1** (plus a difference if **group 2**) predicts **y**.<br>- (Same, but with one variance *per group* instead of one common.)<br>- (Same, but it predicts the *rank* of **y**.) | |
| **Multiple regression: lm(y ~ 1 + x₁ + x₂ + ...)** | **P**: One-way ANOVA<br>**N**: Kruskal-Wallis | aov(y ~ group)<br>kruskal.test(y ~ group) | lm(y ~ 1 + G₂ + G₃ + ... + Gₙ)ᴬ<br>lm(rank(y) ~ 1 + G₂ + G₃ + ... + Gₙ)ᴬ | ✓ | An intercept for **group 1** (plus a difference if group ≠ 1) predicts **y**.<br>- (Same, but it predicts the *rank* of **y**.) | |
| | **P**: One-way ANCOVA | aov(y ~ group + x) | lm(y ~ 1 + G₂ + G₃ + ... + Gₙ + x)ᴬ | ✓ | - (Same, but plus a slope on x.)<br>*Note: this is discrete AND continuous. ANCOVAs are ANOVAs with a continuous x.* | |
| | **P**: Two-way ANOVA | aov(y ~ group * sex) | lm(y ~ 1 + G₂ + G₃ + ... + Gₙ +<br>S₂ + S₃ + ... + Sₙ +<br>G₂*S₂ + G₃*S₃ + ... + Gₙ*Sₙ) | ✓ | Interaction term: changing **sex** changes the **y ~ group** parameters.<br>*Note: Run glm using the following arguments: glm(model., family="poisson").*<br>*As linear model, it uses an indicator (0 or 1) for each non-intercept levels of the group variable.*<br>*Similarly for Sₓ₌ₓ for sex. The first line (with G) is main effect of group, the second (with S) for sex and the third is the group × sex interaction. For two levels (e.g. male/female), line 2 would just be "S₂" and line 3 would be S₂ multiplied with each G.* | [Coming] |
| | **Counts ~ discrete x**<br>**N**: Chi-square test | chisq.test(groupXsex_table) | **Equivalent log-linear model**<br>glm(y ~ 1 + G₂ + G₃ + ... + Gₙ +<br>S₂ + S₃ + ... + Sₙ +<br>G₂*S₂ + G₃*S₃ + ... + Gₙ*Sₙ, family=...)ᴬ | ✓ | Interaction term: (Same as Two-way ANOVA.)<br>*Note: Now uses the multiplicative chi-square test is log(yᵢⱼ) = log(N) + log(βᵢ) + log(βⱼ) + log(αᵢβⱼ) where αᵢ and βⱼ are proportions. See more info in the accompanying notebook.* | Same as<br>Two-way<br>ANOVA |
| | **N**: Goodness of fit | chisq.test(y) | glm(y ~ 1 + G₂ + G₃ + ... + Gₙ, family=...)ᴬ | ✓ | (Same as One-way ANOVA and see Chi-Square note.) | 1W-ANOVA |

List of common parametric (P) non-parametric (N) tests and equivalent linear models. The notation y ~ 1 + x is R shorthand for y = 1 + a·x which most of us learned in school. Models in similar colors are highly similar, but really, notice how similar they *all* are across colors! For non-parametric models, the linear models are reasonable approximations for non-small sample sizes (see "Exact" column and click links to see simulations). Other less accurate approximations exist, e.g., Wilcoxon for the sign test and Goodness-of-fit for the binomial test. The signed rank function is signed_rank = function(x) sign(x) * rank(abs(x)). The variables G₂ and S₂ are "dummy coded" indicator variables (either 0 or 1) exploiting the fact that when Δx = 1 between categories the difference equals the slope. Subscripts (e.g., G₂ or yₓ) indicate different columns in data. In requires long-format data for all non-continuous models. All of this is exposed in greater detail and worked examples at https://lindeloev.github.io/tests-as-linear.

ᴬ See the note to the two-way ANOVA for explanation of the notation.
ᴮ Same model, but with one variance per group: gls(value ~ 1 + G₂, weights = varIdent(form = ~1|group), method="ML").

Why do we need statistics?

"Statistics is the science of collecting, summarizing, presenting
and interpreting data, and of using them to estimate the
magnitude of associations and test hypotheses"

Kirkwood and Sterne p. 1

The build-up of a research project

- **Planning**
- **Design**
- **Execution** (data collection)
- **Data analysis**
- **Presentation**
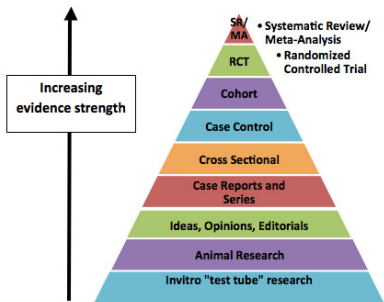- **Interpretation**
- **Publication**

**Statistics in <u>all</u> points**

Critical reading of publications

- Research **design**
- **Inclusion and exclusion criteria**
- **Sample size**
- **Exposure** (risk factor) and **confounding factors**
- **Outcome** (response)
- Statistical **analysis**
- **Bias**
- **Interpretation** of results

**Statistics in <u>all</u> points**

# Pyramid of evidence



Increasing evidence strength

- SR/MA • Systematic Review/Meta-Analysis
- RCT • Randomized Controlled Trial
- Cohort
- Case Control
- Cross Sectional
- Case Reports and Series
- Ideas, Opinions, Editorials
- Animal Research
- Invitro "test tube" research

- **Grading the evidence** for practice guidelines after susceptibility of threats to internal validity
- **Health literacy guide** designed to help students find and assess sources of quality health information: `https://libraryguides.unh.edu/c.php?g=326606&p=2191225`

# Oslo Centre for Biostatistics and Epidemiology (OCBE)

- ... is a joint centre of UiO (Department of Biostatistics, IMB) and OUS (Biostatistics and Epidemiology group at Forskningsstøtte). Approx. 80 people in total.

- **Research**: Methodological research in several areas, e.g.
  - ▶ Statistics for High-dimensional Data
  - ▶ Causal Inference
  - ▶ Clinical and Cancer Epidemiology
  - ▶ Hybrid Modeling / Mathematical Oncology
  - ▶ Infectious diseases
  - ▶ Clinical Trials Unit (CTU)

- **Statistical advising** for researchers at the Medical Faculty, OUS and Helse-Sør-Øst

- **Teaching** at MedFak: professional study programme for Medicine, Master's programmes in Clinical Nutrition and International Community Health, PhD courses

# OCBE Statistical Advising Service



```
https://www.med.uio.no/imb/english/research/centres/
ocbe/advising/
```