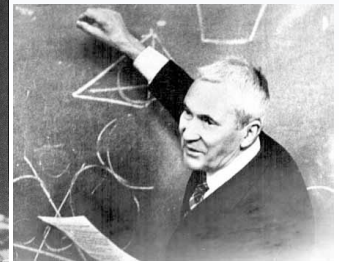
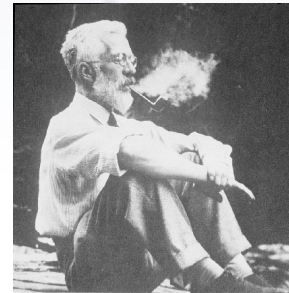
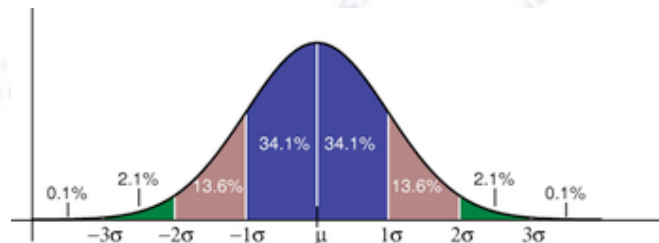


# Applied Statistics

## Course information 2017



Troels C. Petersen (NBI)



*"Statistics is merely a quantisation of common sense!"*

# Applied Statistics 2017

...all the technical stuff!

Technicals:

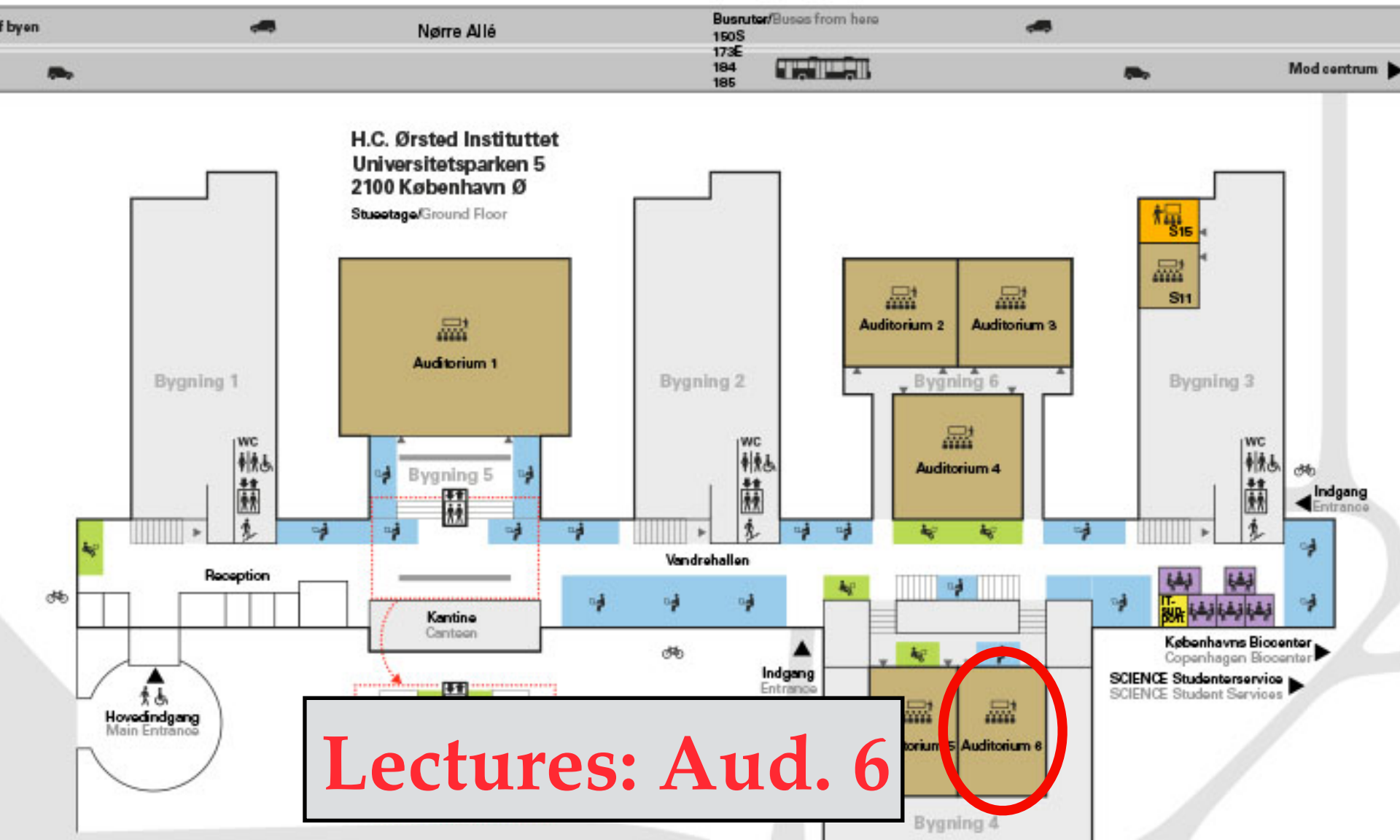
- Rooms and hours.
- Course structure and dates.
- Computers and software.
- Data sets.
- Literature.
- Curriculum.
- Problem set.
- Projects.
- Exam.
- Expectations.
- Goals.



The course webpage (central source of course information, bookmark or fail!):

<http://www.nbi.dk/~petersen/Teaching/AppliedStatistics2017.html>

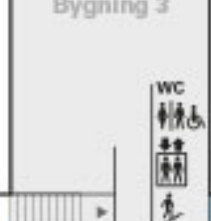
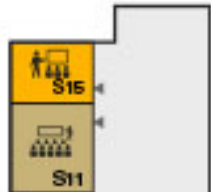
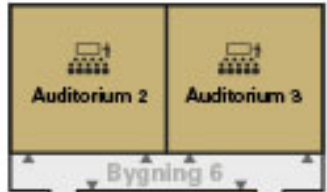
# Lectures & Exercises



# Lectures & Exercises

f byen      Nørre Allé      Busruter/Buses from here  
160S  
173E  
184  
185      Mod centrum ▶

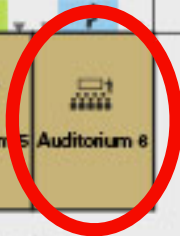
H.C. Ørsted Institutet  
Universitetsparken 5  
2100 København Ø  
Stueetage/Ground Floor



**Exercises: A111+A112**



**Lectures: Aud. 6**



Københavns Biccønter  
Copenhagen Biccønter  
SCIENCE Studenterservice  
SCIENCE Student Services



**My office**  
(building M, top floor)

**First Lab**

**Entrance to Auditorium A**

Blegdamsvej

# Hours & Rooms

## Hours:

Following block B, but using the morning hours 8:15 - 9:00 Monday and Friday for “self-studying”, *except first two Mondays!!!*

## Monday:

9:15 - 10:00 Lectures  
10:15 - 12:00 Exercises

## Tuesday:

13:15 - 14:00 Lectures  
14:15 - 17:00 Exercises

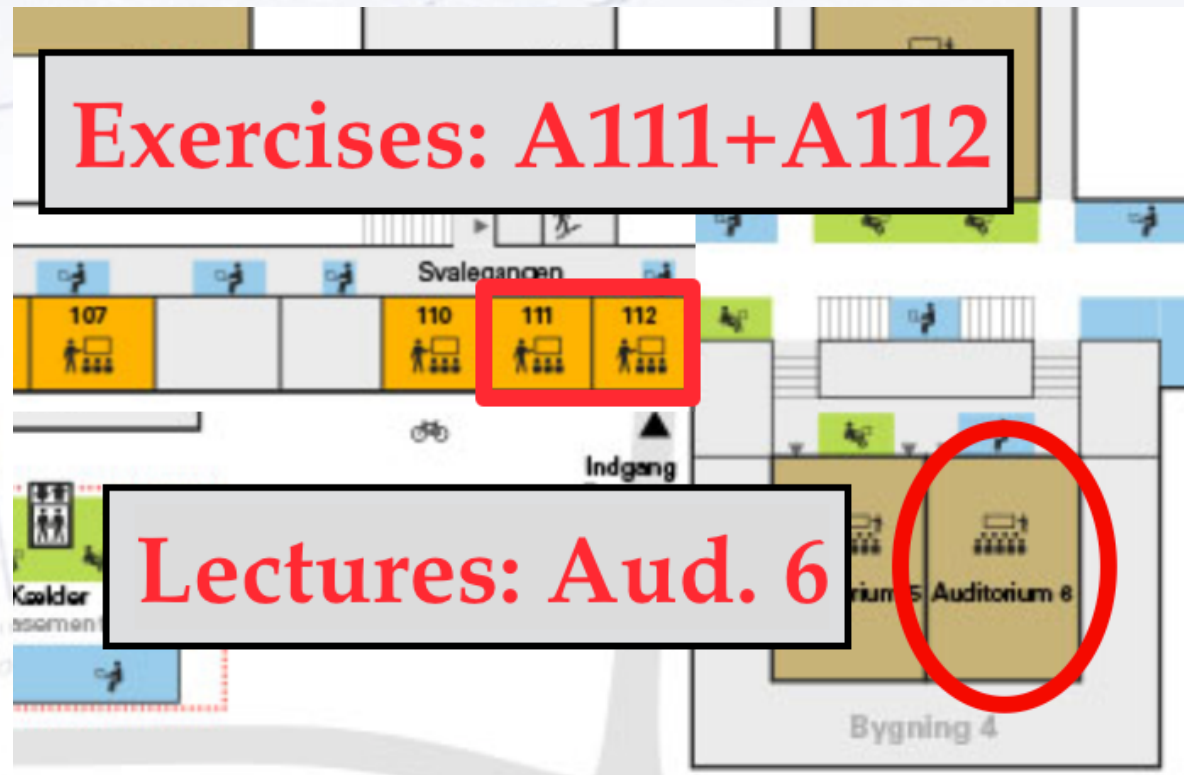
## Friday:

9:15 - 10:00 Lectures  
10:15 - 12:00 Exercises

## Rooms:

**Lectures: HCØ Auditorium 6 (math)**

**Exercises: HCØ A111 + A112 (I hope!)**



First week: Additional Python introduction Tuesday 13:15-13:45

# Hours & Re

Hours:

Following block B, but  
morning hours  
and

**Only exceptions:  
First day of course  
20th of November 8:15  
in Auditorium A  
...and the second Monday, starting 8:15 in First Lab**

F  
9  
10:

Additional Python introduction Tuesday 13:15-13:45

# Hours & Re

Hours:

Following block B, but  
morning hours  
and

**Only exceptions:  
...of course**

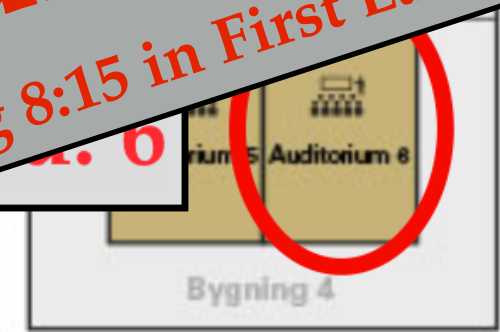
**Note for 2017:**

**There will be no teaching 22nd of December!**

(i.e. count it as a day for homework/preparation/project/etc.)

**20th of Nov  
in Auditorium 8  
...and the second Monday, starting 8:15 in First Lab**

F  
9  
10:

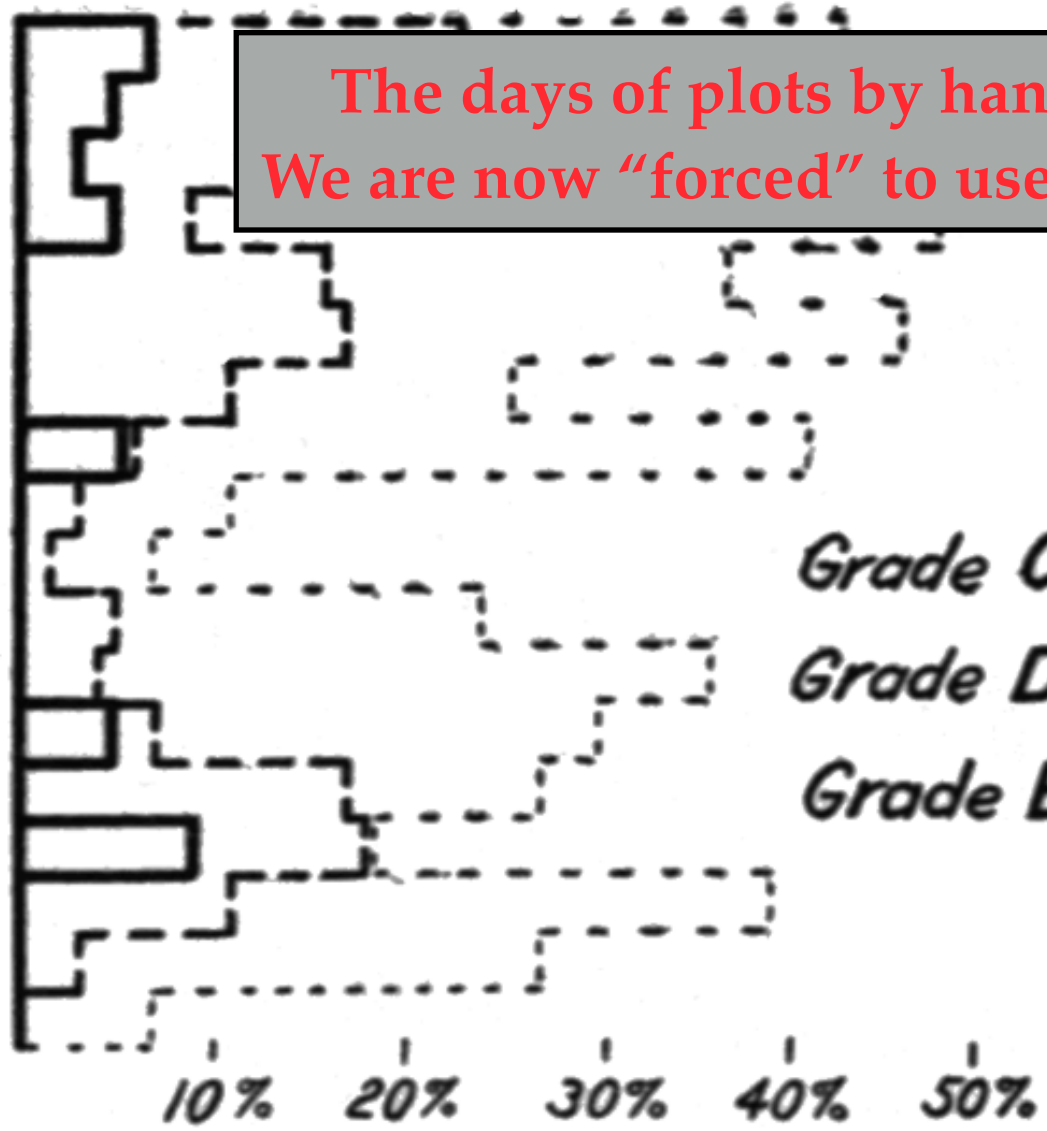


Additional Python introduction Tuesday 13:15-13:45



# Computers and software

*Fine Arts*  
*Economics*  
*English Lit.*  
*Slavic*  
*Metallurgy*  
*French*  
*History*  
*Semitic*  
*Philosophy*  
*English Comp.*  
*Latin*  
*Government*  
*Education*  
*Physics*  
*Mathematics*  
*Zoology*  
*Chemistry*  
*Greek*



The days of plots by hand are over!  
We are now "forced" to use computers!!!

Grade C = - - - - -  
Grade D = - . - - -  
Grade E = ———

# Computers and software

The times are *way past* pencil and/or calculator stage!

**Fast computers** is the *only* answer to do (any serious) data analysis.

Operating system:      **Linux/MAC OS/Windows**

Programming:          **Python** - version 2.7.X (or 3.5)

Editor:                  **Spyder, PyCharm, Sublime, Emacs** (or own favorit!)

## Installation:

<http://www.nbi.dk/~petersen/Teaching/Stat2017/installation.html>

Before course start, we will give an introduction to this (“Week 0”):

Wednesday 15th 11:15-14:00 in Aud. A: Help with **installation**.

Thursday 16th 11:15-14:00 in Aud. A: Introduction to **programming**.

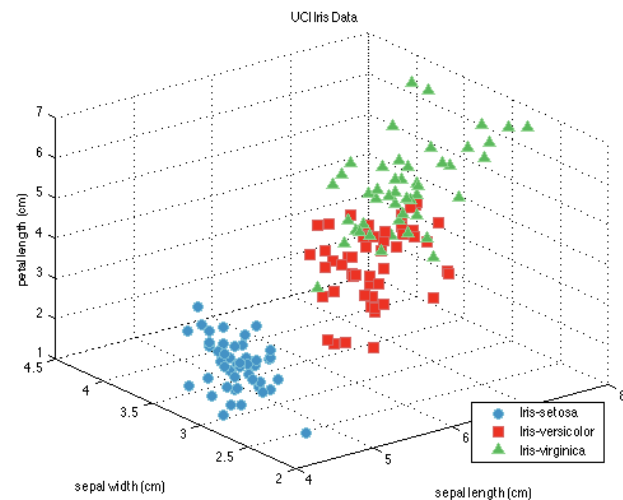
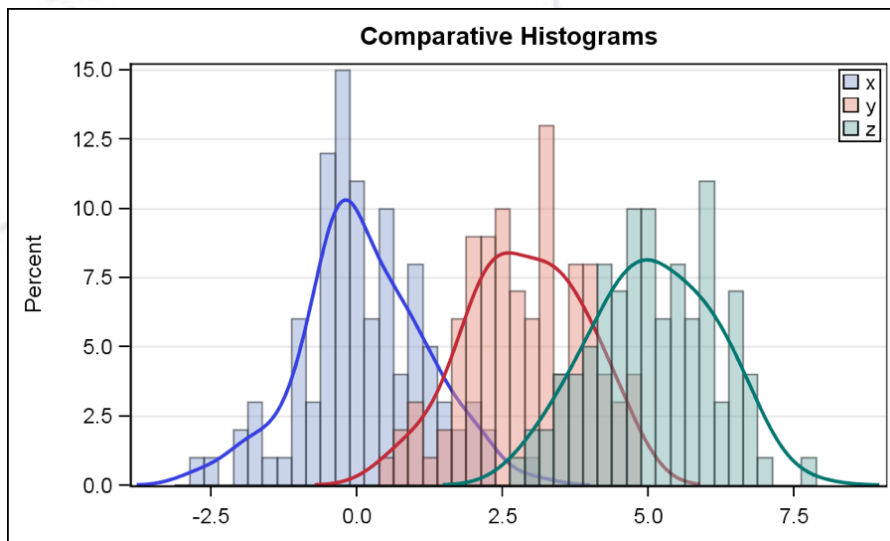
Also, during first week, first hour of Tuesday and possibly Friday after class, we will further introduce and train programming.

# Data sets

In general, any data set can be used for this course! If you happen to have an interesting and illustrative one, bring it to me/class!

I've tried my best to search for a large variety of data sets, but this is not always easy. Publicly available data sets are often old/small/biased/etc.

As a result, some data sets are from my own field (particle physics). This is both due to my access to data here, but also because particle physics is one of the fields providing *billions of measurements*.



# Literature

We use Roger J. Barlow's "Statistics", as it is an accessible introduction to statistics with many examples, and the best overall book (I think).

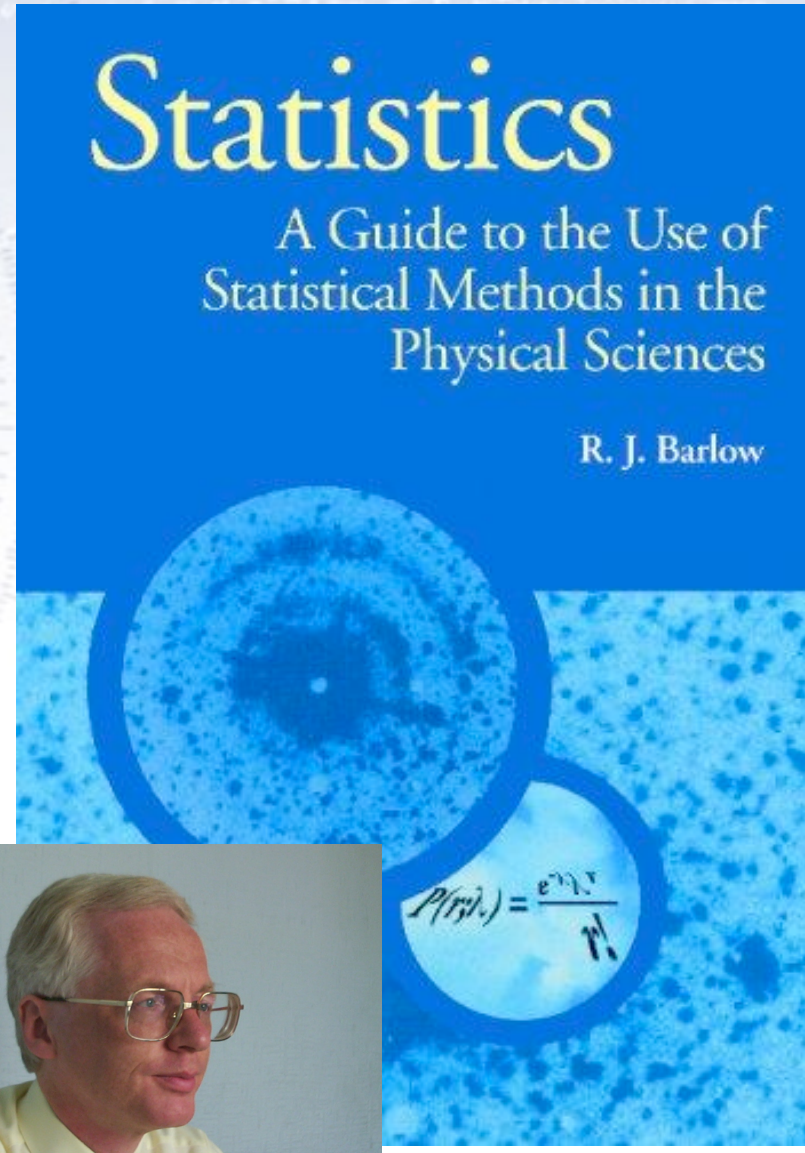
If anything, it is lacking a bit on how to generate random numbers according to a specific PDF and on categorising events.

I will occasionally also refer to:

- Bevington: Data Reduction & Error Analysis
- Cowan: Introduction to Statistics

...and notes from Particle Data Group!

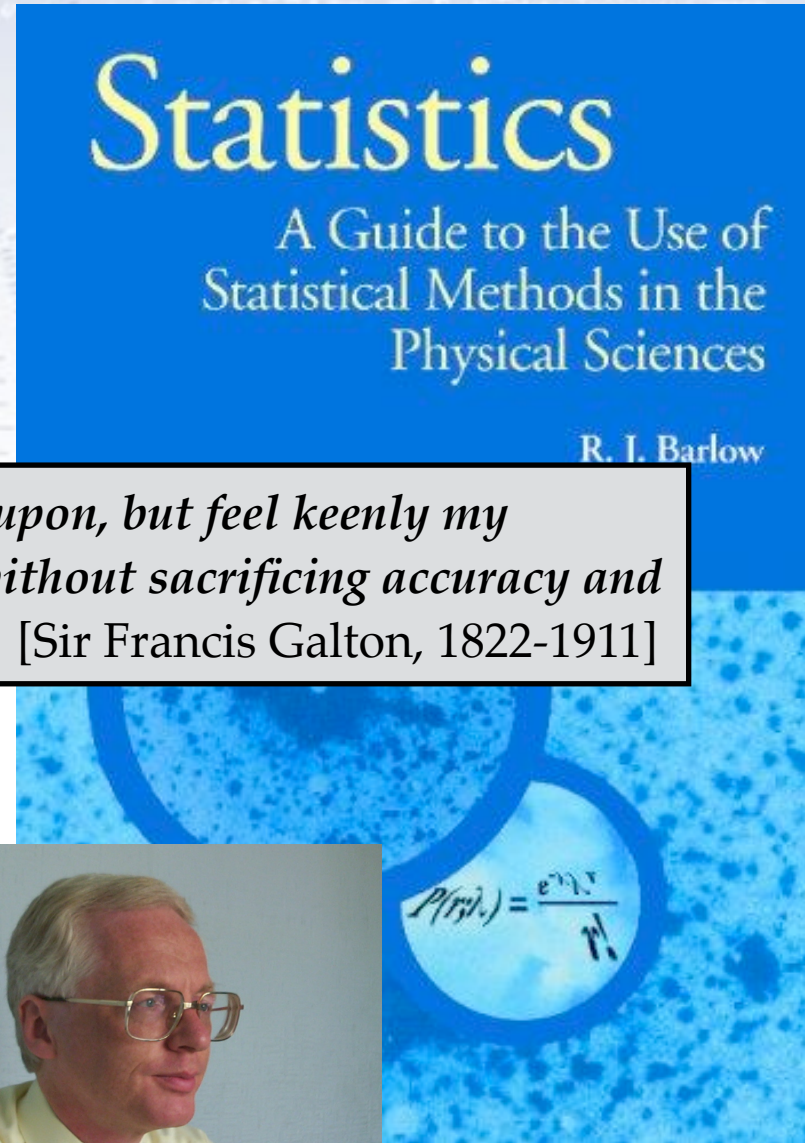
NOTE: There is a great abundance of notes, fora, etc. on both statistics but especially also Python on the web, which I encourage you to use (with a proper critical mind).



# Literature

We use Roger J. Barlow's "Statistics", as it is an accessible introduction to statistics with many examples, and the best overall book (I think).

If anything, it is lacking a bit on how to generate random numbers according to a specific PDF and on categorising events.



*"I have a great subject [statistics] to write upon, but feel keenly my literary incapacity to make it intelligible without sacrificing accuracy and thoroughness"*

[Sir Francis Galton, 1822-1911]

- B
  - Cowan: Introduction to Statistics
- ...and notes from Particle Data Group!

NOTE: There is a great abundance of notes, fora, etc. on both statistics but especially also Python on the web, which I encourage you to use (with a proper critical mind).



# Curriculum

The course will cover the following chapters in R. Barlow:

- Chapter 1 (All)
- Chapter 2 (All)  
Exercises: All, except 2.5 and 2.9.
- Chapter 3 (Except 3.2.2, 3.3.2, 3.4.2, 3.5.2)  
Exercises: All, except 3.7.
- Chapter 4 (All)  
Exercises: All, except 4.10.
- Chapter 5 (Except 5.1.3, 5.3.2, 5.3.3 (formal part), 5.3.4, 5.5)  
Exercises: 5.2
- Chapter 6 (Except 6.4.1, 6.7)  
Exercises: All
- Chapter 7 (Except 7.3.1)  
Exercises: All, except 7.1, 7.3, and 7.7.
- Chapter 8 (Except 8.4.4, 8.4.5, 8.5.1, and 8.5.2)  
Exercises: All, except 8.6.
- Chapter 10 (All)

# Core of Curriculum

The course will **focus mostly on** the following chapters in R. Barlow:

- Chapter 2: 2.1, 2.2, 2.3, 2.4.1, 2.4.2, 2.6
- Chapter 3: 3.1, 3.2, 3.2.1, 3.3, 3.3.1, 3.4.1, 3.4.7, 3.5.1
- Chapter 4: 4.1, 4.2, 4.3, 4.3.1, 4.3.2, 4.3.3
- Chapter 5: 5.1, 5.1.1, 5.1.2, 5.2, 5.6
- Chapter 6; 6.1, 6.2, 6.2.1, 6.2.2, 6.2.3, 6.2.4, 6.3, 6.4
- Chapter 8: 8.1, 8.2, 8.3, 8.4, 8.4.1, 8.4.2, 8.4.3

This is less than 80 pages, but... they do not only require reading!

**They request understanding!!!**

The plan is to go through this curriculum in 4-5 weeks, spending the rest of the time on applying it.

**It is through application that statistics is really understood.**

# Check list

In order for me to consider you inscribed in this course, you should fulfil the following check list (*preferably first day, and within first week!*):

- **Have read the course information** (these slides, on course webpage).  
Otherwise, you don't know what is going to happen.
- **Have your picture ("mug shot") taken.**  
Otherwise, I don't know who you are.
- **Have filled in the questionnaire** (on course webpage).  
Otherwise, I don't know what you know and don't.
- **Have measured the length of the lecture table in Auditorium A.**  
Otherwise, you haven't contributed to a common course dataset.
- **Have Python running on your laptop.**  
Otherwise, you can't follow the exercises or solve problems.
- **Have registered for exam!**  
Otherwise, the administration will kill us!

In order not to continuously be doing the above, we will be doing all of the steps today and possibly Friday after class, **only!**



# Problem set

During the course (week 3-4), I will give a larger problem set to be solved and handed in.

This will cover most of the curriculum covered at this point, and it *will count 15% in your final grade!!!*

It will require the use of computers and modifications of some of the code you have been running.

You are welcome to work in groups, but **each student must hand in their own solution**, and you should **state your collaboration**.

The final exam will somewhat resemble this problem set!



# Projects

During the course (week 2-3 and week 6-8) you will be working on a larger data analysis project for about 1-2 week(s).

Each of these is your chance to play with real data and gain experience of what planning an experiment and detailed data analysis requires!

These *will count 25% in your final grade!!!*

They will require the use of computers and modifications of some of the code you have been running.

You are encouraged to work in groups, and only one report (2-4 pages) is required from each group.

Real life problems will resemble these projects!



# Projects

## Project 1:

Attempt at precision measurement of the Earth's gravitation locally at NBI, using only "simple" means (OK - a little bit of cheating there).

You will be doing two separate experiments (both seen before):

- Simple pendulum.
- Ball rolling down an incline.

## Project 2:

*Whatever you would like to do*, as long as it involves (advanced) data analysis! However, be warned that getting data is hard! So start today, or ask me for data known to work (reasonably well), such as:

- Gravity measurements and influence of Moon and Sun.
- UFO sightings (yes, fun data!).
- Rutherford experiment (with modern setup, Ian Bearden).
- Identification of "V0" particle with ATLAS 2015 data.
- Data from a company that needs analysing.
- Fun data from others...



# Projects

## Project 1:

Attempt at precision measurement of the Earth's gravitation locally at NBI, using only "simple" means (OK - a little bit of cheating there).

You will be doing two separate experiments (both seen before):

- Simple pendulum.
- Ball rolling down an incline.

## Project 2:

*Whatever you would like to do*, as long as it involves (advanced) data analysis! However, be warned that getting data is hard! So start today, or ask me for data known to work (reasonably well), such as

- Gravity measurements and influence of Moon and Sun.
- UFO sightings (yes, fun data!).
- Rutherford experiment (with modern setup, Ian Bearden).
- Identification of "V0" particle with ATLAS 2015 data.
- Data from a company that needs analysing.
- Fun data from others...



# Exam

Exam will be a **28 hour take-home exam** with a problem set, which resembles the one previously given.

It will cover most of the curriculum, and it *will count 60% in your final grade!!!*

It will require the use of computers and modifications of some of the code you have been running.

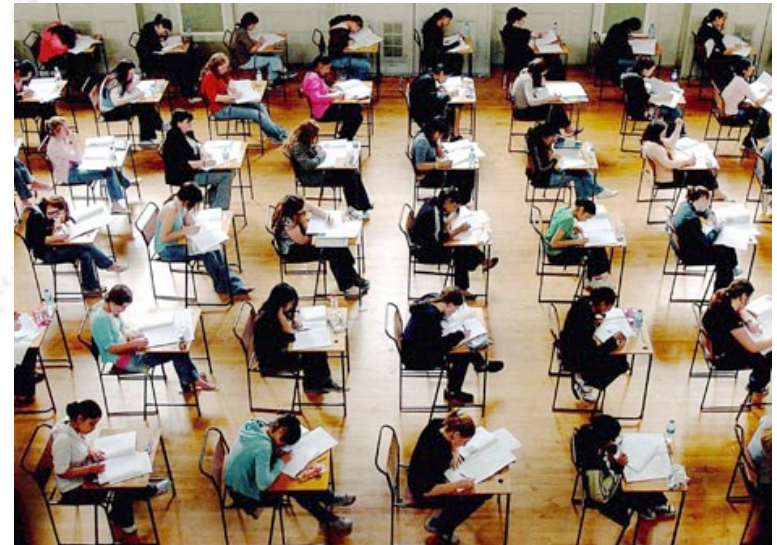
**You must work on your own!**

I will provide this exam on:

**Thursday the 18th of January 8:15am.**

It will then naturally have to be handed in:

**Friday the 19th of January before 12:00!**



# Challenges

During the course, there will be a few challenges:

- Best table measurement analysis.
- Most precise measure of  $g$  (to better than  $1/10000$ ?).
- A problem on the problem set.
- Project 2
- ???

They are meant as advanced exercises to those, who are not already challenged significantly by the course! They do not give credit, but will of course earn you advanced experience and impress me (who gives grades and might be writing you a letter of recommendation).

Don't stress over this - you can of course still earn the grade 12 without ever touching upon them.

# Expectations

I want (read: insist) this course to be useful to all of you!

Therefore, please give me feedback (during the course, thanks!) if you have anything to add / suggest / criticise / alter.

However, it is also through your active participation that you have this privilege (i.e. that I'll listen most).

This also means, that I will require much from you - as much as I can without spoiling the social life of your youth!

In return, I'll try to make statistics as interesting as possible (and not deprive you of your early mornings).

# Power for computers/you

There are no individual power sockets in Auditorium 6 (lectures) nor in A110 and A111 (exercises).

We will try our best to bring a few extension cords, but please charge your laptops before class.

Regarding your own energy, we will try to supply tea/coffee and a few cookies.

If that is of any interest, then bring a cup yourself.

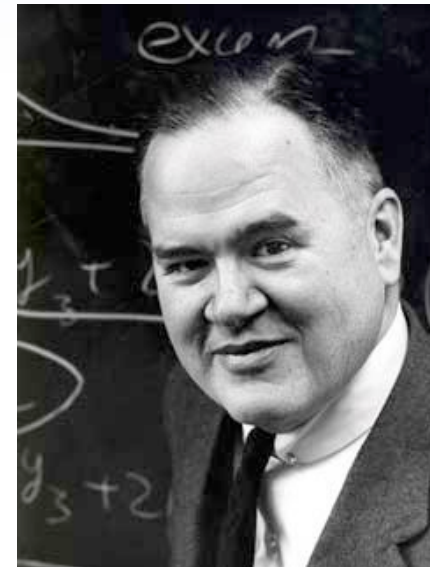




# Statistical practices

The famous statistician John Tukey (1915-2000) was quoted for wanting to teach:

- The **usefulness and limitation of statistics**.
- The importance of having methods of statistical analysis that are robust to violations of the assumptions underlying their use.
- The need to amass experience of the behavior of specific methods of analysis in order to provide guidance on their use.
- The importance of allowing the possibility of data's influencing the choice of method by which they are analysed.
- The need for statisticians to reject the role of “guardian of proven truth”, and to resist attempts to provide once-for-all solutions and tidy over-unifications of the subject.
- **The iterative nature of data analysis**.
- Implications of the increasing power, availability and cheapness of **computing facilities**.
- The training of statisticians.



*"Far better an approximate answer to the right question, which is often vague, than an exact answer to the wrong question, which can always be made precise." J. W. Tukey*

# References

*Roger J. Barlow: Statistics (course book!)*

*(A guide to the use of statistics methods in the physical sciences)*

Very good introduction, which goes further than Bevington. Very much to the point.

*Philip R. Bevington: Data reduction and error analysis.*

Classic introduction with very good examples - a standard reference in all of experimental physics [and so essentially all of physics!].

*Glen Cowan: Statistical Data Analysis*

A bit brief, but once you got the hang of statistics, this book contains much of what you will ever need, written in a useful or precise way.

# Top 10

## Most important things in applied statistics

1. Errors decrease with the **square root of N**
2. **ChiSquare** is simple, powerful, robust and provides a **fit quality** measure
3. **Binomial** distribution → **Poisson** distribution → **Gaussian** distribution
4. **Error propagation** is **craftsmanship** - **fitting** is an **art**
5. Error on a (Poisson) number,  $N$ :  $\sqrt{N}$  on a fraction,  $f=n/N$ :  $\sqrt{f(1-f)/N}$ .
6. **Correlations** are important and needs consideration
7. Hypothesis testing of  $H_0$  (null) and  $H_1$  (alt.) is done with a test statistic  $t$
8. The **likelihood** (ratio) is generally the optimal estimator (test)
9. Low statistics is terrible – needs special attention
10. Prior probabilities needs attention, i.e. Bayes' Theorem