# Biostatistics:
## *from genome to diseases prevention? – polygenic risk scores*

Turid Frahnow

Institute of Computational Biology

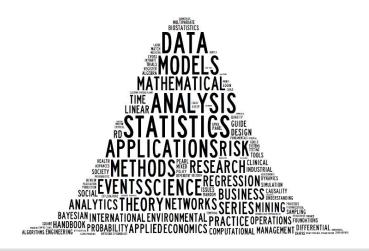Helmholtz Zentrum München, Germany

20th March, 2017

# Schedule

- Biostatistics in general
    - What is it?
    - Why we use it?
    - Applications

- Biostatistics in the context of polygenic risk scores
    - Biological background
    - What is it?
    - Why we use it?
    - Problems

# First of all…



Please interrupt me (immediately),
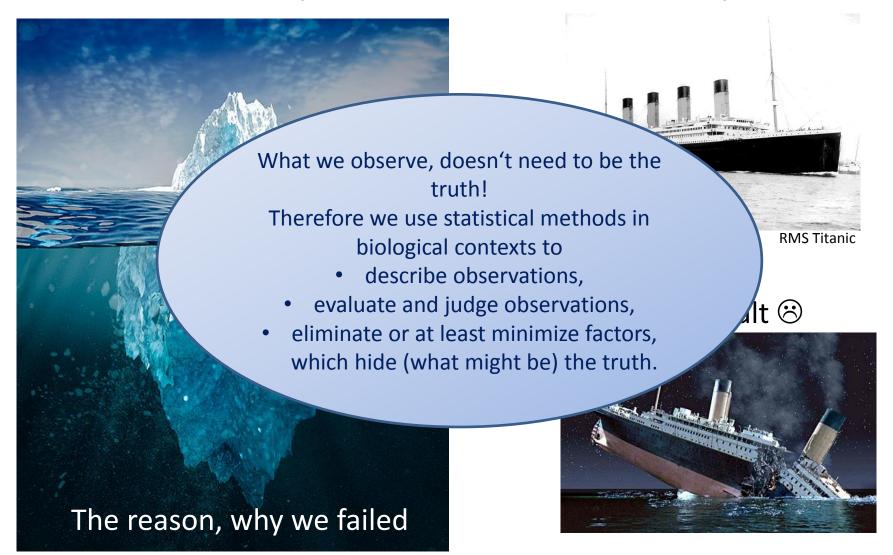if you don't understand something!

# **Biostatistics in general**

# Motivation

## Our observation / data

## Our decision/conclusion

RMS Titanic

What we observe, doesn't need to be the truth!
Therefore we use statistical methods in biological contexts to
- describe observations,
- evaluate and judge observations,
- eliminate or at least minimize factors, which hide (what might be) the truth.

The reason, why we failed

# Biostatistics - „measurement of life"

= medical statistics, biometry …

- Definition:

    „Biostatistics is the branch of statistics responsible for the proper interpretation of scientific data generated in biology, medicine, public health and other health or natural sciences."

- Or shorter: The **application** of statistics to biological/natural sciences.

- And important: It is much more, than using a t-test.

- (Should be) always interdisciplinary **!**

# It's about understanding….

**Biological Problem** >

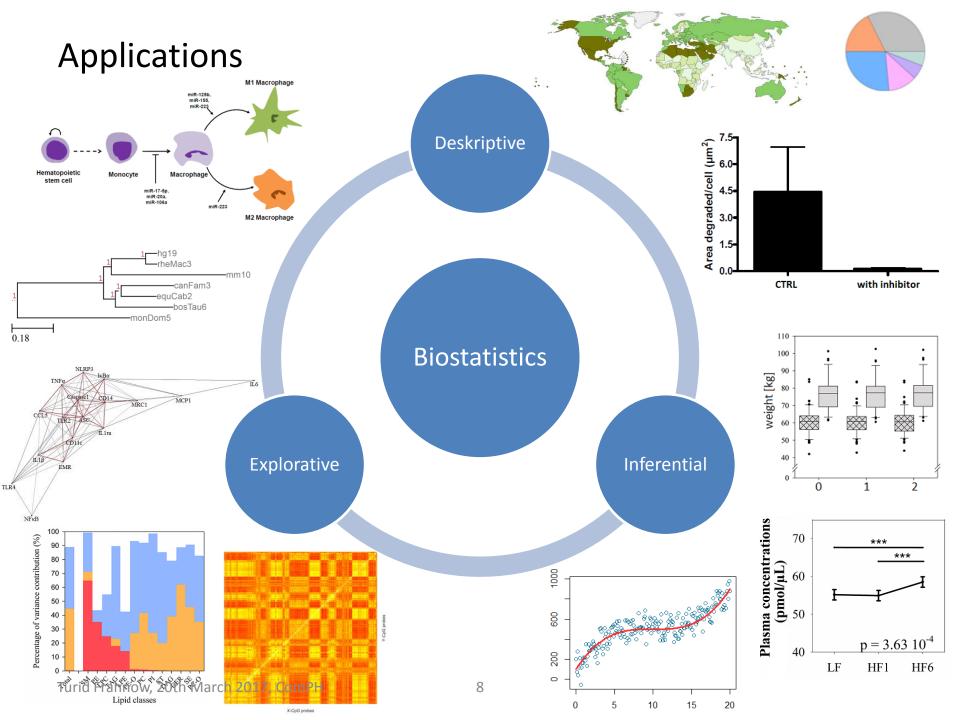   Experimental design >

      Experimental techniques>

         Data structure>

            (Transformation >)
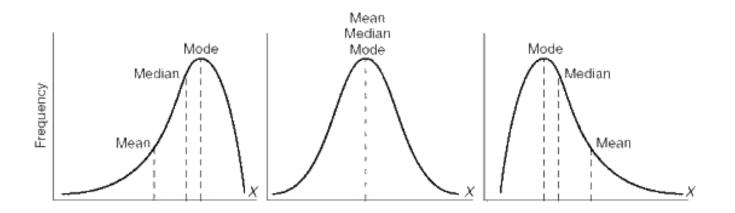
               Results >

               **Interpretation**



Statistics shows that teenage pregnancy drops dramatically after 20

# Applications



Deskriptive

Biostatistics

Explorative

Inferential

Pierre-Charles-Alexandre Louis (1787-1872) was a pioneer of the "numerical method" in medicine.

# Descriptive statistics

# The 3 m's of central tendency



- **Mean** = the sum of observations divided by the number of observations

- **Median** = the observation in the center when all observations are ordered from smallest to largest

- **Mode** = the most frequently occuring value among all observations

# Measures of Spread/Variation

- **Range** = Max – Min

- **Standard deviation** = amount of variability of the observations (range of 68% of the observations around the mean)

- **Standard error (of mean)** = amount of variability in the population mean

- **Interquartile range** = range between 25% and 75% quartile (range of 50% of the observations around the median)

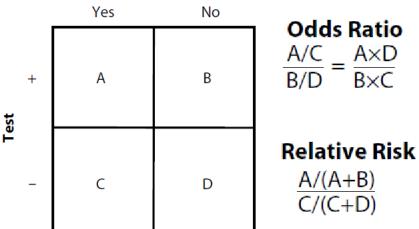- **Quartile coefficient of dispersion, coefficient of variation,** …

# Associations



deterministic            statistical

- **Correlation** = measurement of association

- **Odds ratios** = the odds of exposure in the group with disease divided by the odds of exposure in the control group (Case-Control design)

- **Relative risk** = the ratio of the incidence of disease in the treated group divided by the corresponding incidence of disease in the placebo group (Cohort design)
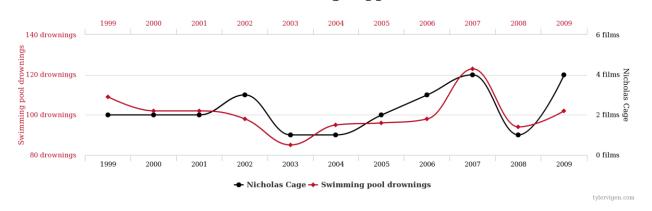


**Disease**

|      | Yes | No |
|------|-----|-----|
| + | A | B |
| − | C | D |

**Odds Ratio**

$$\frac{A/C}{B/D} = \frac{A \times D}{B \times C}$$

**Relative Risk**

$$\frac{A/(A+B)}{C/(C+D)}$$

# Correlation doesn't imply causation!

**Number of people who drowned by falling into a pool**
correlates with
**Films Nicolas Cage appeared in**

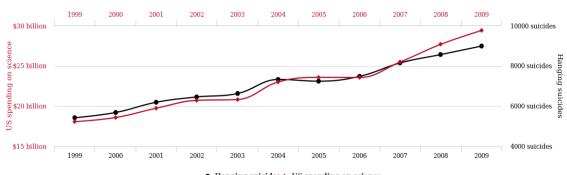r ≈ 0.6660

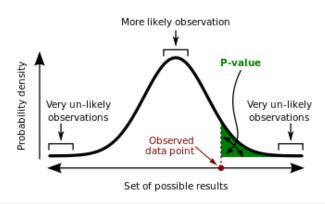**US spending on science, space, and technology**
correlates with
**Suicides by hanging, strangulation and suffocation**

r ≈ 0.9978

# Inferential and explorative statistics

# Measures of Quality

**Estimates**

- **Confidence intervals** = a range of values within which there is a high probability (95% by convention) that the true population value can be found.
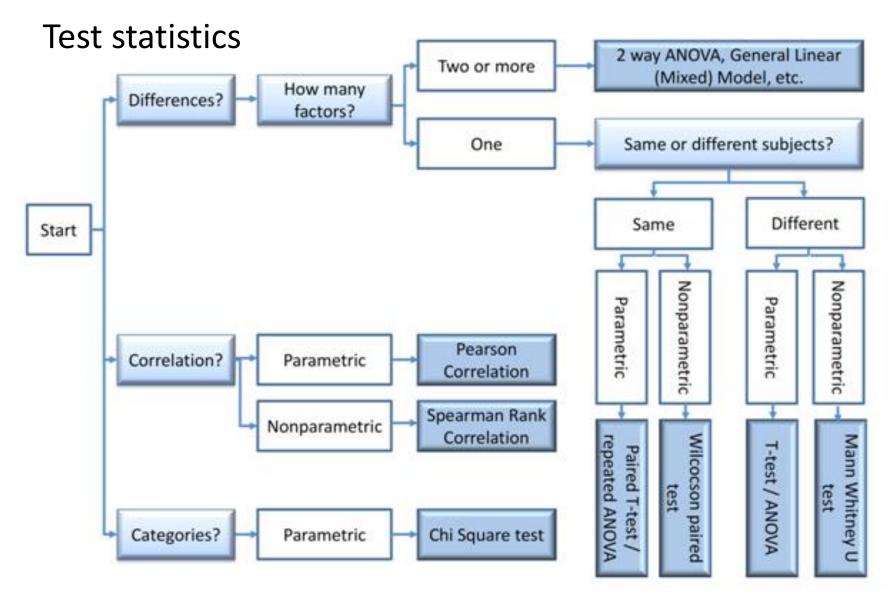
**Tests**

- **Sensitivity** = the ability of the test to identify correctly those who have the disease (true-positives).

- **Specificity** = the ability of the test to identify correctly those who do not have the disease (true-negatives).
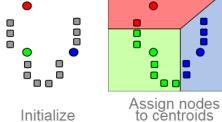
# Test statistics



The P-value is not the probability that the null hypothesis is true given the data.
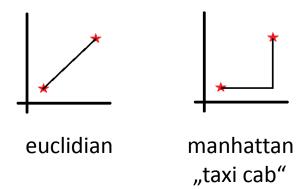Absence of evidence is not evidence of absence.

# Classification – Cluster analysis

- Grouping of similar observations by distance in p-dimensional system ( p = number of variables)

- Different strategies
  - K-means: partition of n observations in k cluster by choosing centroids



Initialize

Assign nodes to centroids

  - Hierarchical
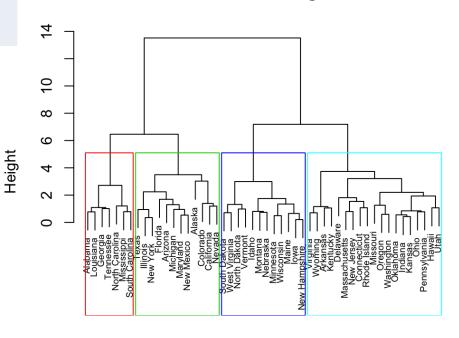    e.g. „Bottom up":     each observation starts in its own cluster, the pair of "closest" clusters are merges



euclidian      manhattan „taxi cab"

# Cluster analysis - example

| | Murder | Assault | UrbanPop | Rape | |
|---|---|---|---|---|---|
| Alabama | 1,24256408 | 0,7828393 | -0,5209066 | -0,00341647 | |
| Alaska | 0,50786248 | 1,1068225 | -1,2117642 | 2,48420294 | |
| Arizona | 0,07163341 | 1,4788032 | 0,9989801 | 1,04287839 | ... |
| Arkansas | 0,23234938 | 0,230868 | -1,0735927 | -0,1849166 | |
| California | 0,27826823 | 1,2628144 | 1,7589234 | 2,06782029 | |
| Colorado | 0,02571456 | 0,3988593 | 0,8608085 | 1,86496721 | |
| | | ... | | | |

e.g. euclidian distance
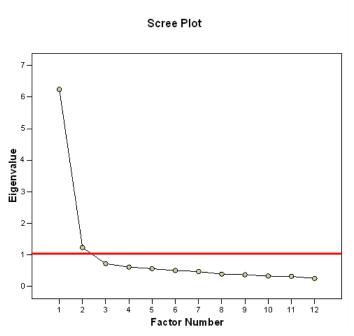


**Cluster Dendrogram**

# Classification – Principal component analysis (PCA)

- Dimension reduction

- Linear orthogonal transformation:
  n possibly correlated variables -> m linearly uncorrelated factors(n > m)
  = orthogonal basis set

- PCA can be thought of as fitting an n-dimensional ellipsoid to the data.

- Largest possible variance explained by first component/factor

- Different strategies to decide how many
  components are needed:
  - Eigenvalue > 1 (Kaiser rule)
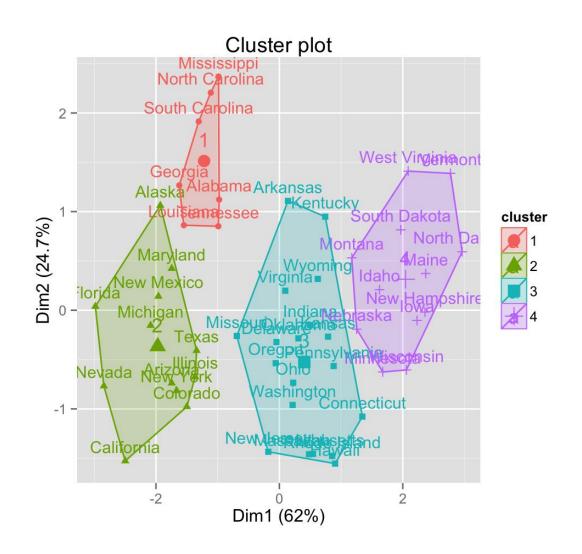  - variance explained > 90 %
  - Scree-plot „elbow" (Cattell rule)

# PCA - Example



**Scree Plot**

The factor with an eigenvalue of 1 accounts for as much variance as a single variable.
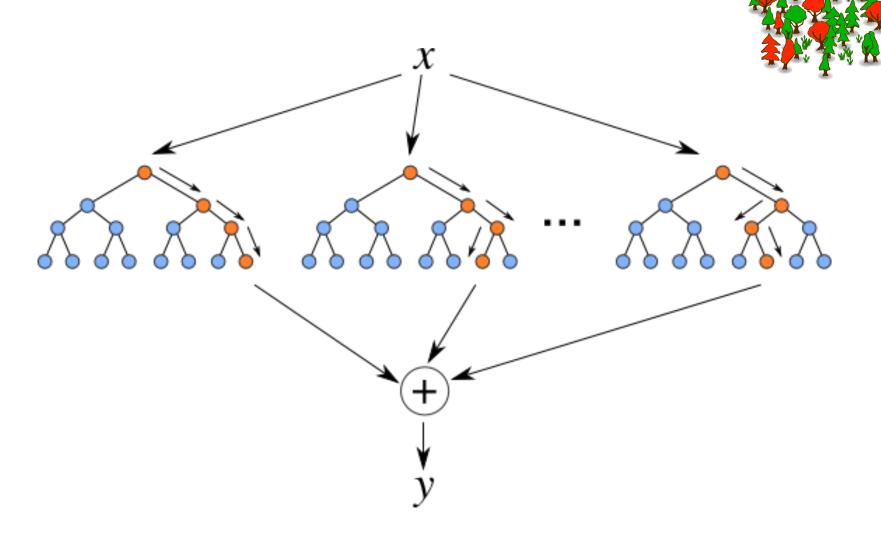
# A combined approach – HCPC Analysis
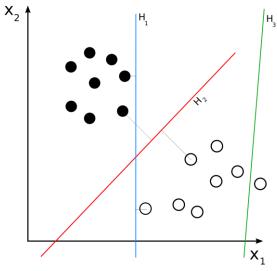
# Classification - Random Forrest

- Construction of multiple decision trees (up to several thousands)

- To rank the importance of variables

- „Tree bagging" - repeatedly fitting trees to random subset

- Output: Mode of classes or mean of prediction

- While the predictions of a single tree are highly sensitive to noise, the average of many (uncorrelated) trees is not.

# Classification - Random Forrest

# Support vector machines

- Large margin classifier

- Strategy: Find a hyperplane (p-1 dimensional space) which
    1) Separates the two classes,
    2) Maximizes the margin (distance to closest data points).

- linearly separable or up-transformation
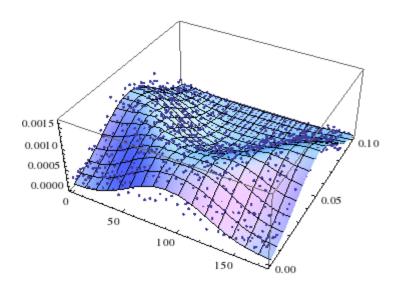


$H_2$ is the correct solution

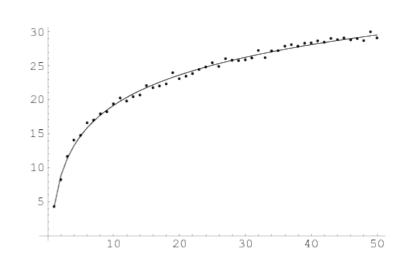# Regression / modeling – generalized linear model

- Many subtypes – here just an introduction – to study relationships.

- $g(y) = X\beta + \varepsilon$          observation y and predictor(s) X

- Goal: Prediction of future observations, assessment of effects of predictors on observation or general description of data.

- (Pseudo) coefficient of determination to judge, how well the model fit the data

- In the generalized case, y is continous. What if y is binary (health/disease)?
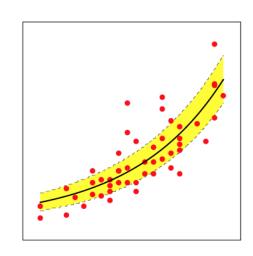  - (binary) logistic regression
  - Death/suvival processes

# Regression / modeling – overfitting

- We want to explain the data in the simplest way

- Unnecessary predictors will add noise to the estimation of other quantities that we are interested in

- Collinearity is caused by having too many variables trying to do the same job

- Procedures: Backward Elimination/Forward Selection

- Decision to keep / drop variables based on
  - Hypothesis tests
  - Information criteria (AIC, BIC)

- All procedures are heuristics, so try out!

# Regression / modeling – different examples

# Regression / modeling – Lasso regrularization

- lasso (least absolute shrinkage and selection operator)

$$\min_{\alpha,\beta} \frac{1}{N} \sum_{i=1}^{N} f(x_i, y_i, \alpha, \beta) \text{ to subject } \|\beta\|_1 \leq t$$

- Performs both variable selection (can set $\beta_i$ = 0 ) and regularization

- Enhance the prediction accuracy and interpretability

# Known problem - Multiple Testing

- Consider a case where you have 1000 hypotheses to test, and a significance level of $\alpha = 0.05$

50 significant hypotheses just by chance

- Different strategies to adjust $\alpha$ (for n hypotheses)
  - Bonferroni: $\alpha_{new} = \alpha/n$
  - Benjamini-Hochberg: control of the false discovery rate
  - Holm
  - …

# Biostatistics in the context of polygenic risk scores

# Different languages…



Nevertheless, it is not all about significance – especially for big data.

# 21st Century – The Century of big data

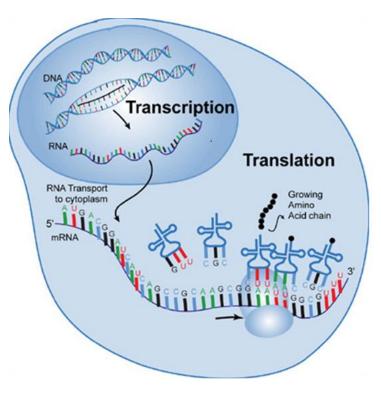- 90% of the data in the world today has been created in the last two years.

- Challenge:
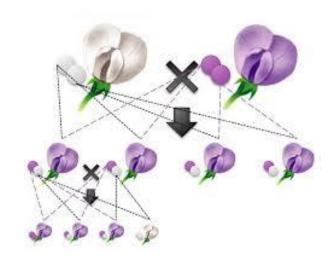  Traditional methods (and data processing software) are inadequate to deal with them

- Where do we find big data?
  e.g. in (epi)genetics…
  - Single nucleotide polymorphisms (SNPs)
  - Next generation sequencing
  - Histone modifications/DNA methylation



The Power of Healthcare Data
The Body as a Source of Big Data

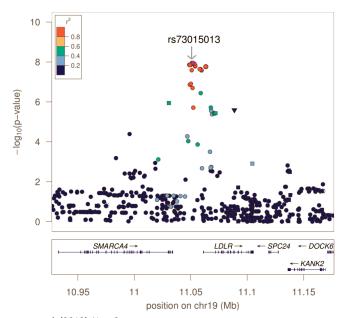# The central concept of gene-environment interactions

- It all starts with…peas and a monk (G. Mendel)
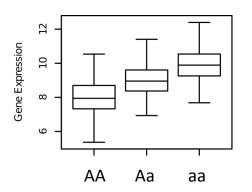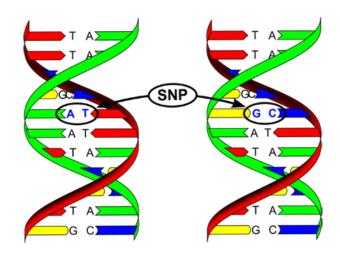
- Phenotype = Genotype + Environment





- 100 years later:
  Transcription + Translation and their regulation as key players of gene-environment interactions
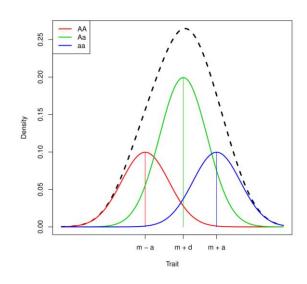
- DNA -> mRNA -> Protein
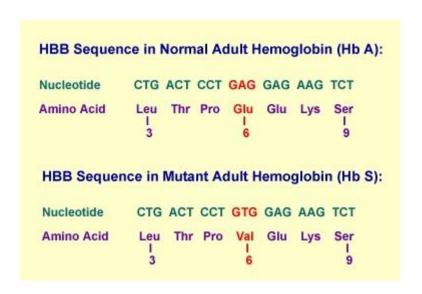
# SNPs – Single nucleotide polymorphisms
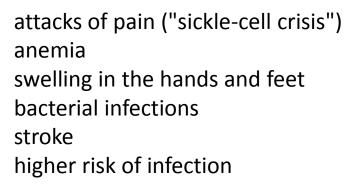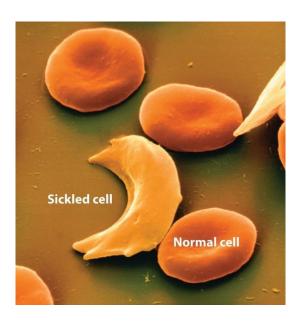


Kettunen et al. (2012) Nat. Genet.

# SNPs – Example (Sickle-cell anaemia)

**HBB Sequence in Normal Adult Hemoglobin (Hb A):**

| Nucleotide | CTG | ACT | CCT | GAG | GAG | AAG | TCT |
|---|---|---|---|---|---|---|---|
| Amino Acid | Leu (3) | Thr | Pro | Glu (6) | Glu | Lys | Ser (9) |

**HBB Sequence in Mutant Adult Hemoglobin (Hb S):**

| Nucleotide | CTG | ACT | CCT | GTG | GAG | AAG | TCT |
|---|---|---|---|---|---|---|---|
| Amino Acid | Leu (3) | Thr | Pro | Val (6) | Glu | Lys | Ser (9) |



Sickled cell

Normal cell

attacks of pain ("sickle-cell crisis")
anemia
swelling in the hands and feet
bacterial infections
stroke
higher risk of infection

# Each Disease, another SNP…



Each color is another disease/disorder

# What is a risk score?

- In mathematical terms: A RS is a classifier!

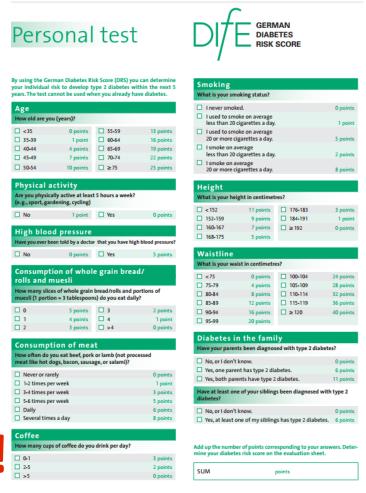  In <u>statistical learning</u>, we want to infer to unknown relationship f(·) in:

  $$\mathbf{y} = f(\mathbf{X}) + \boldsymbol{\epsilon}$$

  where **y** [sometimes g(y)] is the outcome, **X** is your data matrix, and **ε** is the error.

  In <u>classification</u>, y is binary ➜ yes/no, 1/0, true/false, healthy/sick

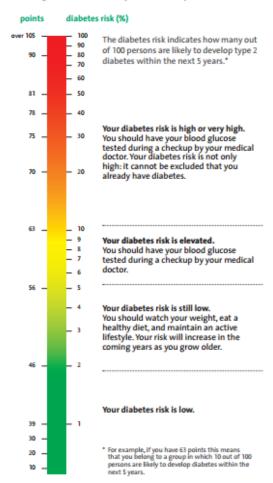  f(X) can have many forms, as we will see in the following

# The easiest way of risk prediction



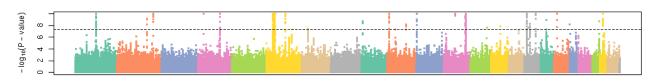## This is your estimated diabetes risk

By using the scale shown below, you can correlate your number of points with your diabetes risk. Please note that individuals with a low risk of diabetes may also develop this disease. On the other hand, high-risk individuals may remain healthy.
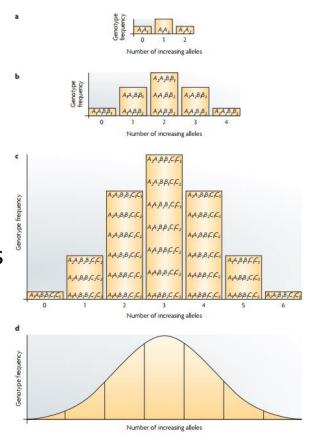
If it's not diabetes, then it's dementia…

# What is a polygenic risk score?

- Adding genetic informations to a risk score. "polygenic" = "more than one gene"

- Diseases/disorders (especially their risk) are quantitative traits!

- What does this mean? → additive genetic effects

- Identifying the variants, that are causing the disease.

- GWAS = genome-wide association studies.
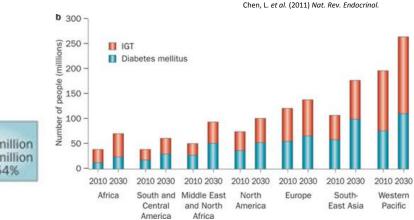
# Aim of polygenic risk scores

We want to predict patients at high risk for a disease based on clinical **and** genetic data

# What is the benefit?

If we are able to identify patients at high risk, we are able to…

- individualize the therapy,
- start therapy quiet early or (in best case) prevent the disease outbreak,
- (maybe) understand the mechanisms behind the disease/disorder,
- reduce the economic burden, e.g

*$174 billion in 2007*
*$245 billion in 2012*

only in the U.S. and only Type II Diabetes
(source: ADA)

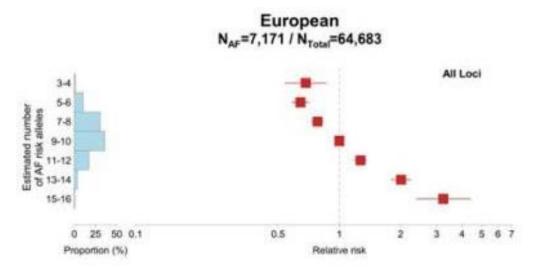Chen, L. *et al.* (2011) *Nat. Rev. Endocrinol.*

# Incidents vs. prevalence

- **Incidence** The number of new events (e.g. death or a particular disease) that occur during a specified period of time in a population at risk for developing the events.

- **Prevalence** The number of persons in the population affected by a disease at a specific time divided by the number of persons in the population at the time.
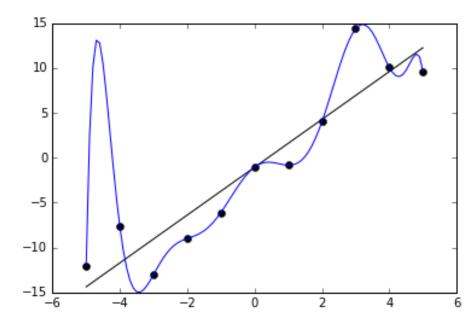
# The easiest way of a PRS

- Sum up the risk allells

- Example: Atrial Fibrillation



- Problem: Quiet different for each outcome/disease and subpopulation (sex/ethnicity/...).
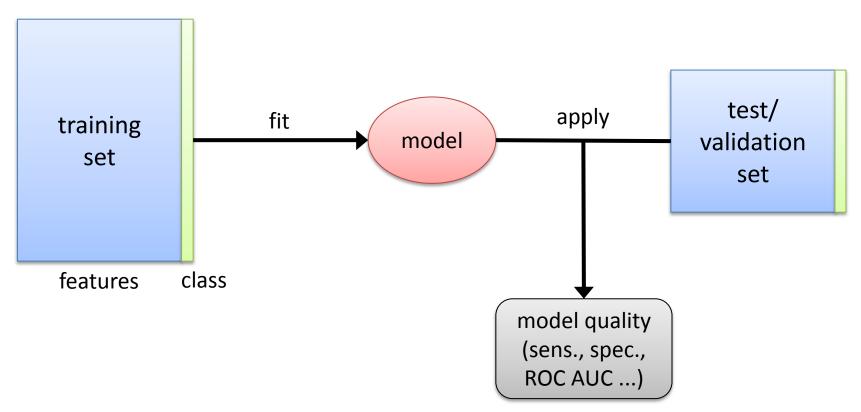
# Problems – Overfitting

- Multicolinearity = high intercorrelation between the predictors
  -> the information of one predictor might be contained in another one

- These classical statistical techniques were developed for low dimensional data (n >> p).

- The blue curve fits the data points too well.

# Validation

Ideal case:



if the model learned a pattern that does not generalize,
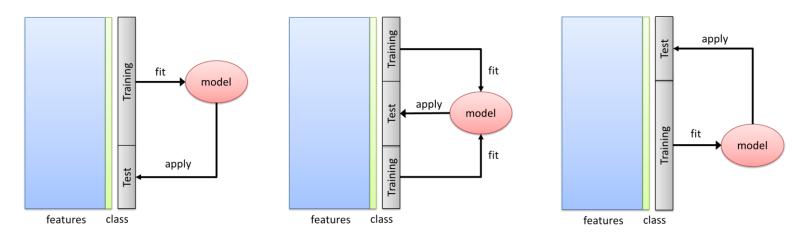the quality on the test set will be low

# Cross validation

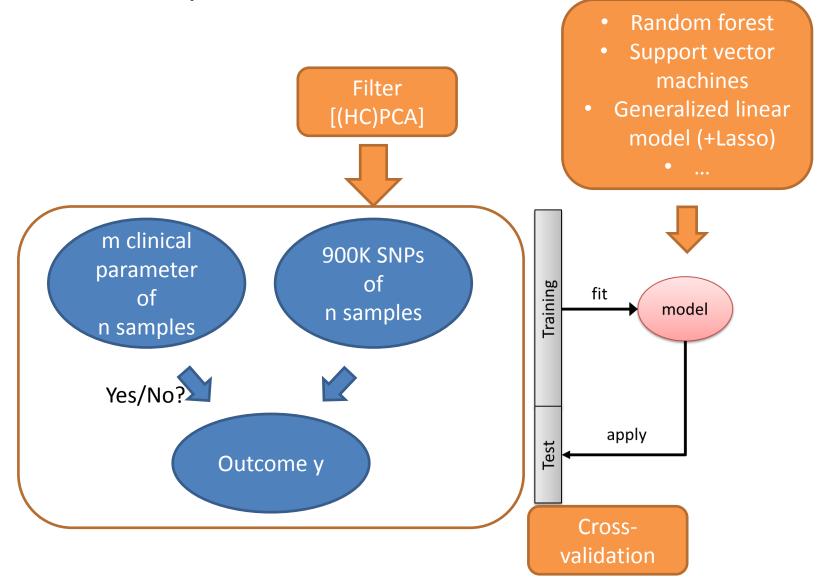Main motivation for cross-validation:

In reality, we can often not afford to leave out a substantial part of the data for validation.

Can we perform validation just using the training data?

Example: 3-fold cross-validation

# Scheme of risk prediction

# Thank you for your attention!