

# Machine Learning in Biomedicine: From Basic Research to Clinical Translation

Lars Kaderali

Lübeck, November 11, 2019



Institute of Bioinformatics



Institut für Bioinformatik

Before we start a few words about myself...



- Computer Science @ University of Cologne
- PhD on a high dimensional regression problem: Predicting Cancer Patient Survival from Gene Expression Data
- Subsequently worked at Los Alamos National Laboratory (USA), German Cancer Research Center dkfz, University of Heidelberg, University Hospital Dresden and now University Medicine Greifswald
- Professor for Bioinformatics at Medical Faculty
- Short stints in industry – CEL Corporation in Canada, ITK in Germany, Partner in ISL GmbH (IT Security Company), Consultant for Beiersdorf AG
- Worked on AI in Biomedicine since 2000, moving more and more into applications over the years



Institut für Bioinformatik

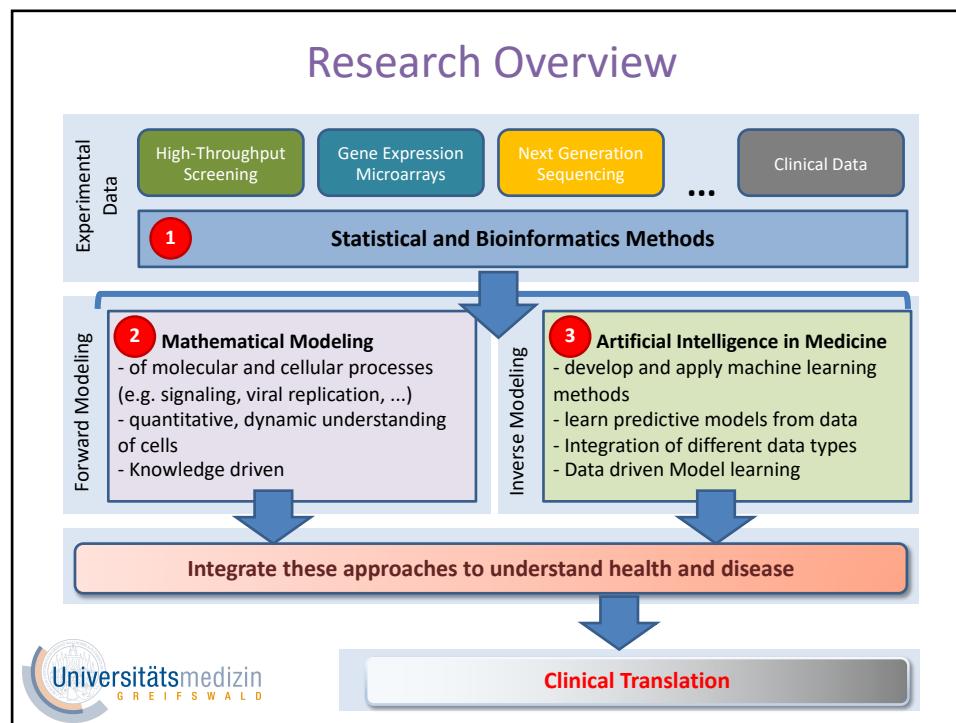
# Institute of Bioinformatics

- Founded 2015/2016
- Part of the University Hospital / Medical Faculty
- Two Professors, Staff approx. 20 Persons (Postdocs, PhD Students, Bachelor/Master Students, Administrative Staff), mostly funded through third-party grants, mostly with a maths / computer science background
- Former Institute of Biometry and Medical Informatics
- Teaching: Biometry/Biostatistics, Bioinformatics
  - Medicine
  - Human Biology BSc & MSc
  - Biomathematics BSc

 Universitätsmedizin  
GREIFSWALD




Institut für Bioinformatik

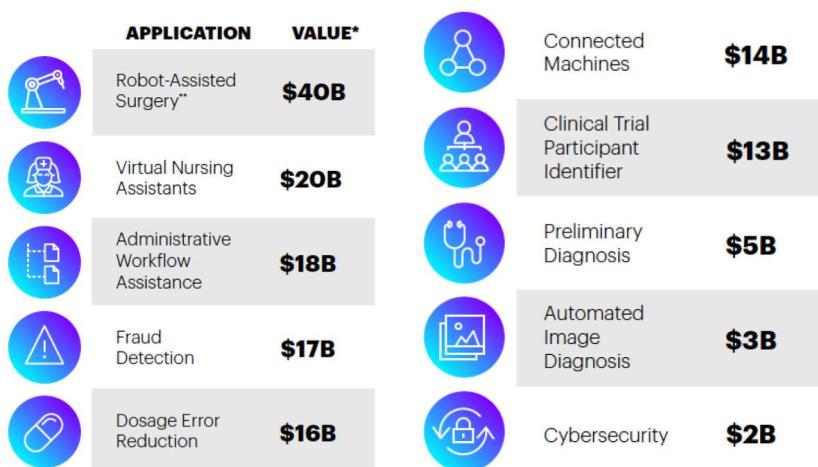


# AI IN MEDICINE: HOPE OR HYPE?



Institut für Bioinformatik

## Top 10 AI Health Applications



TOTAL = ~\$150B

Source: Accenture Analysis



## Potential of AI in Medicine

- AI as support / tool for medical doctors
  - Clinical Decision Making / Decision Support Systems
  - Complementation of human expertise, e.g. in Radiology
  - Access to current information (Scientific Literature, Medical Practice Guidelines etc)
  - Support for younger, less experienced doctors
  - 24x7 Availability
- Early Diagnosis
  - Prediction of Disease Outbreak and Progression, Prediction of Treatment Outcome
  - Feedback for patients and doctors, e.g. concerning patient's compliance
  - Reduction of diagnostic and therapeutic error
  - Increased patient safety and cost savings through AI
  - AI can extract information from huge patient numbers, more than any single doctor will ever see during his lifetime



Institut für Bioinformatik

## Successful Applications of Pattern Recognition in Medicine

- Detection of diabetic Retinopathy  
-> FDA approval for AI-based retina scan diagnostics
- Classification of Gendener using Retina Scans and AI:  
Human doctor 50%, AI 97% Acc
- Processing of Images from Radiology (MRI / CT Scans)
- Examples from Bioinformatics, e.g. Gene Signatures for Breast Cancer Prognosis
- Critical Examples: Prediction of Suicidality / Depression from Social Media Profiles



Institut für Bioinformatik

## MOLECULAR MECHANISMS OF AGEING

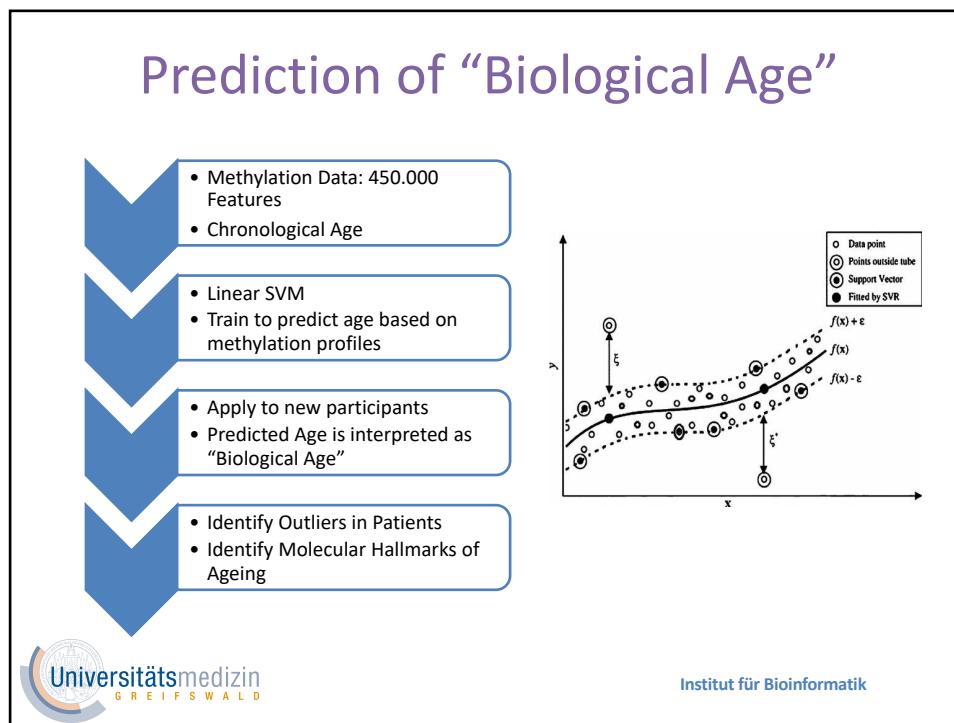
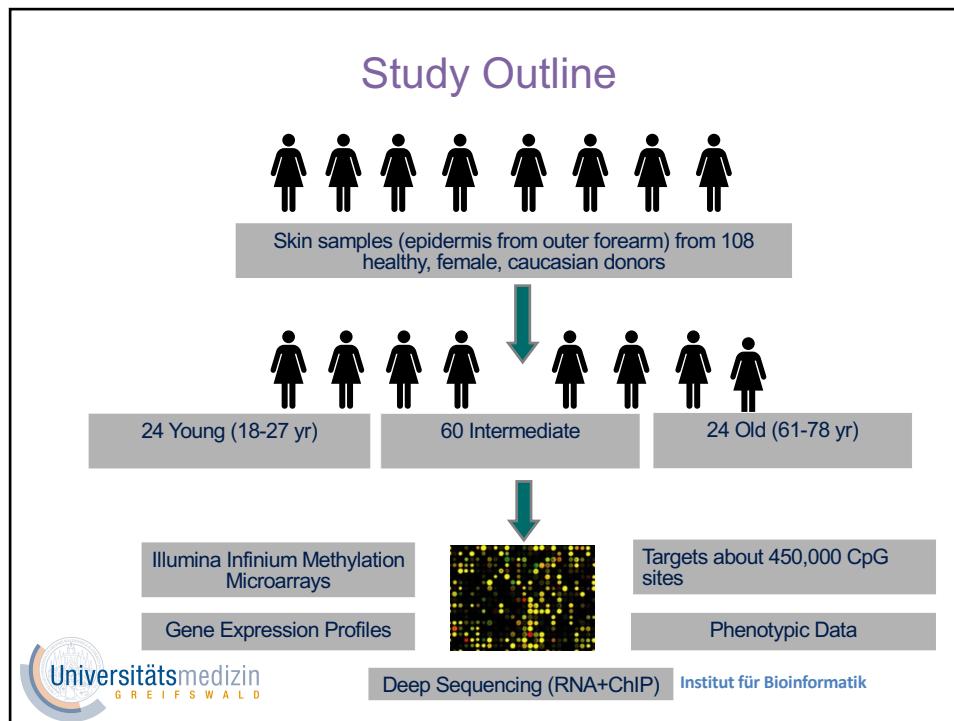


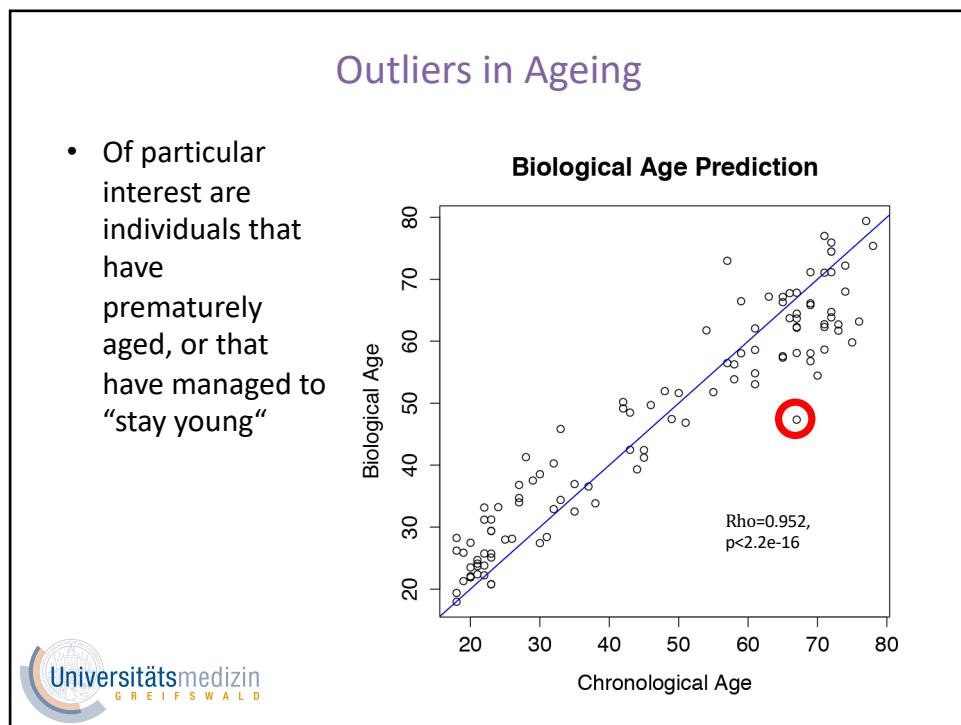
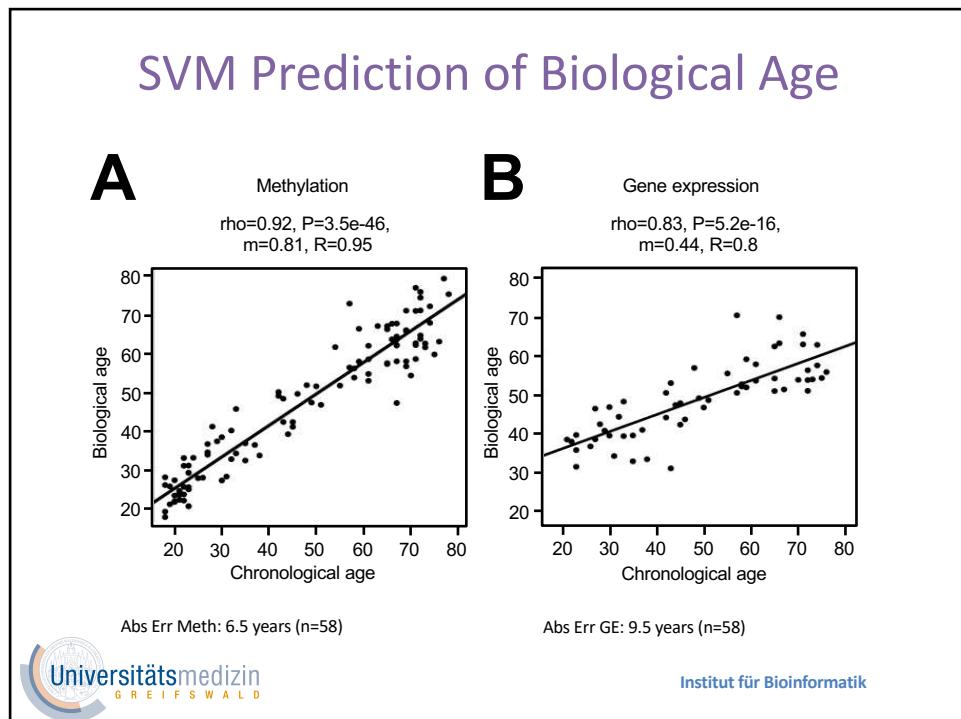
Institut für Bioinformatik

### The Quest for Eternal Youth

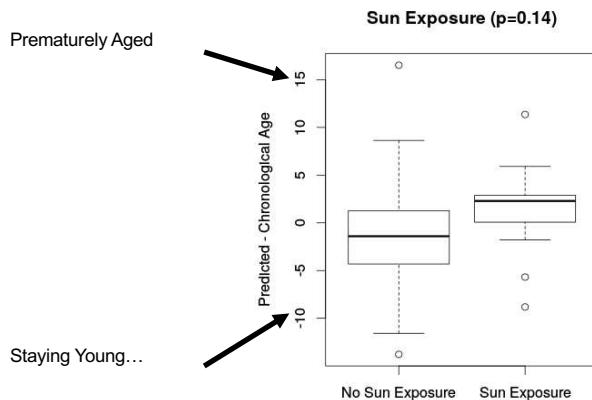


**Beiersdorf** **dkfz.** GERMAN CANCER RESEARCH CENTER IN THE HELMHOLTZ ASSOCIATION



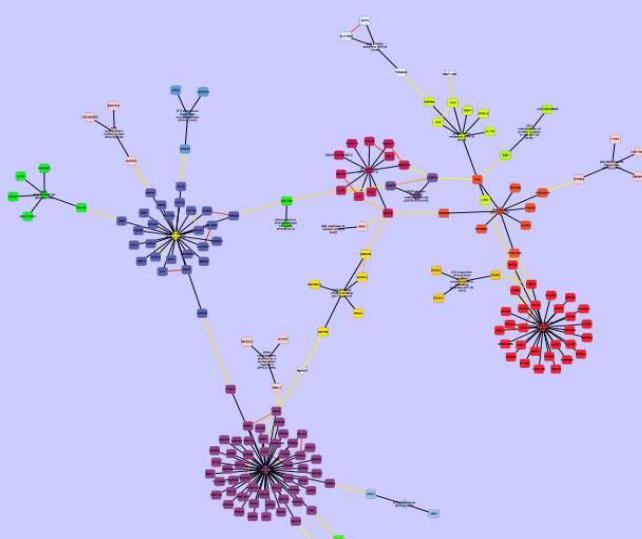


# Premature Ageing and Phenotype Correlations



Institut für Bioinformatik

## Multi-OMICs Signatures of „Oldies“ and „Youngsters“

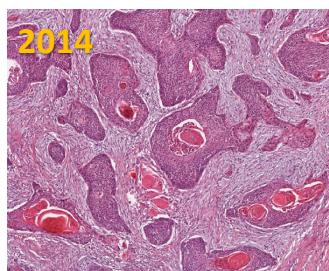


# TREATMENT OF LUNG CANCER IN SMOKERS



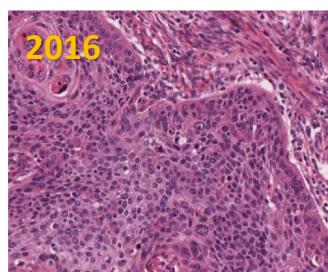
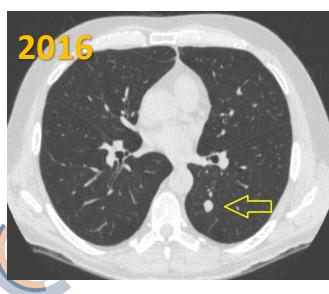
Institut für Bioinformatik

## Treatment of Lung Cancer in Smokers



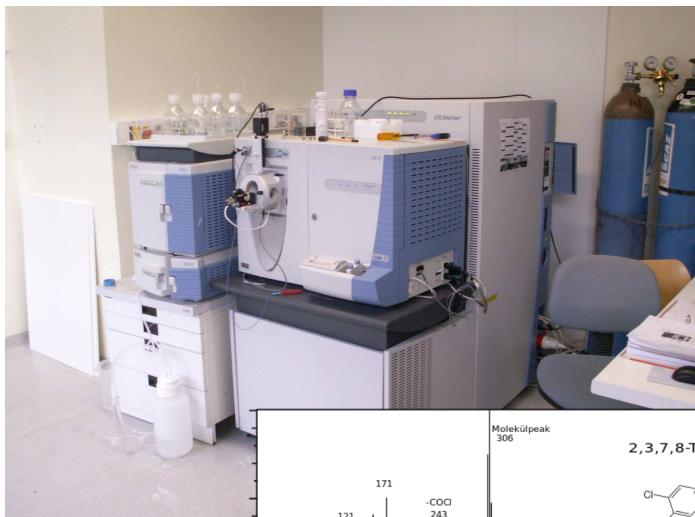
72 year old male, 50 packyears  
2014 - Squamous cell carcinoma of the oral cavity  
2016 - 1 cm solitary nodule in the lung with squamous cell carcinoma histology

- 1) Primary lung cancer? -> T1aN0M0?  
-> curative surgery?
- 2) Metastasis? -> Stadium IV?  
-> palliative chemotherapy?

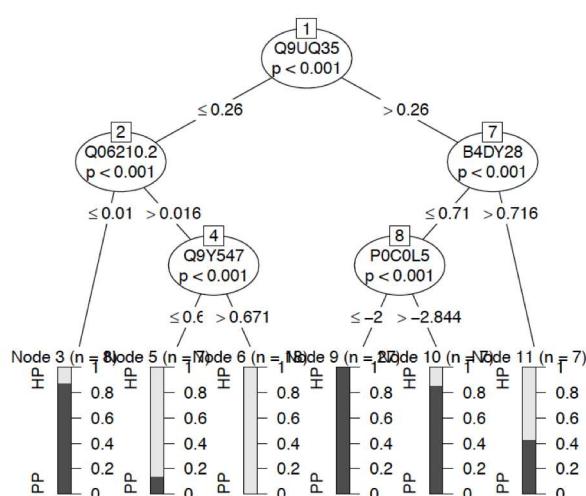


Institut für Bioinformatik

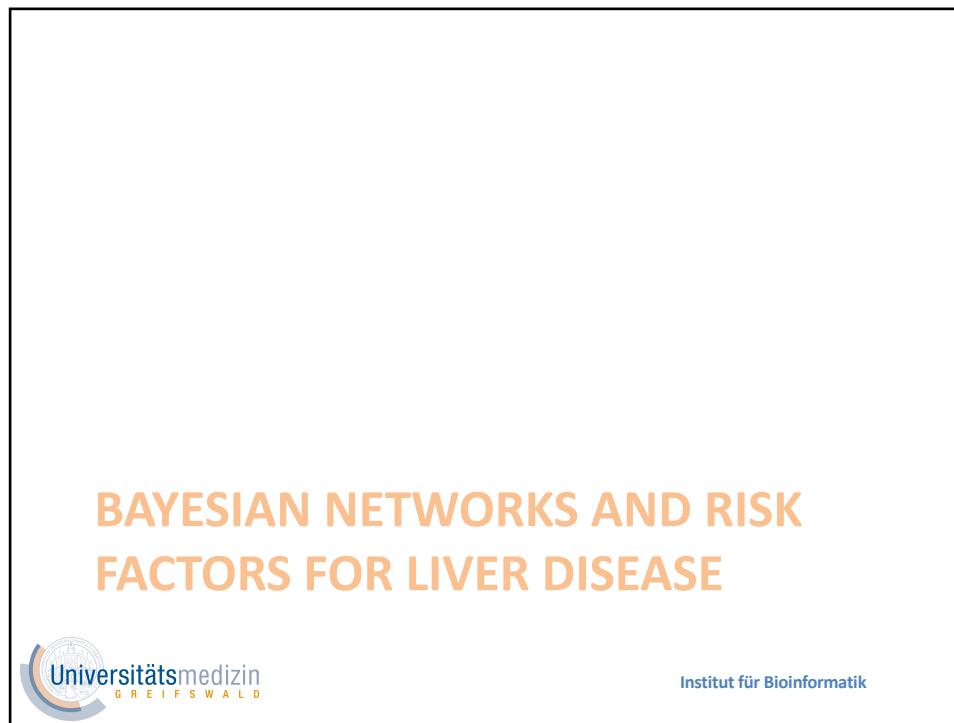
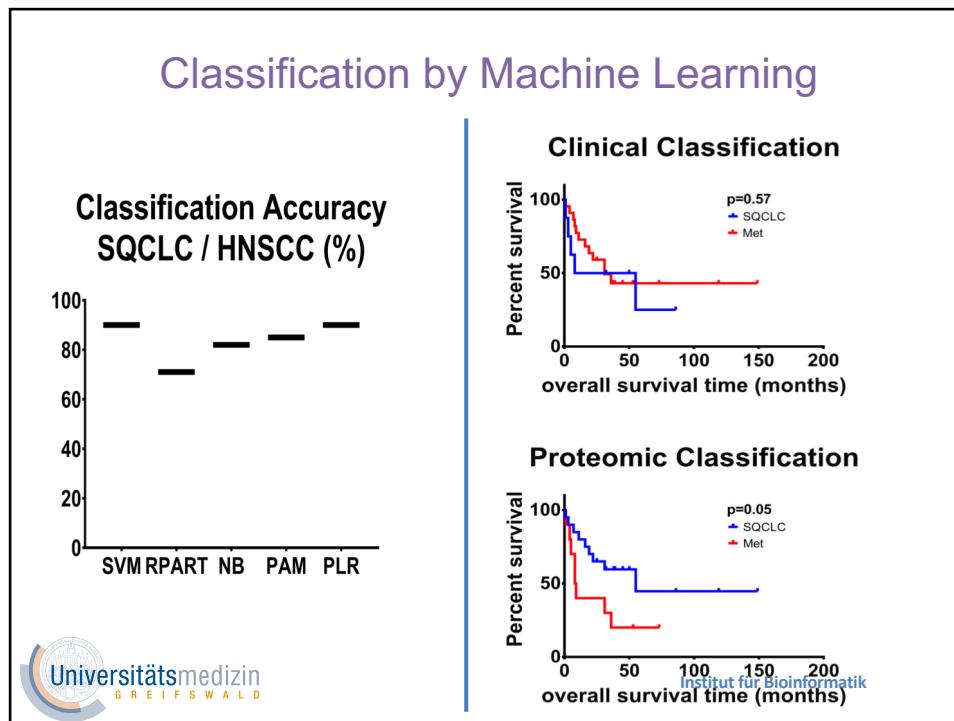
## SILAC Protein Quantification



## Classification of lung cancers

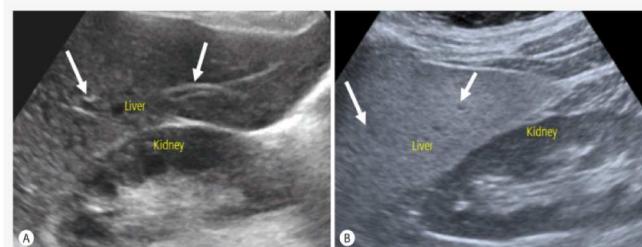
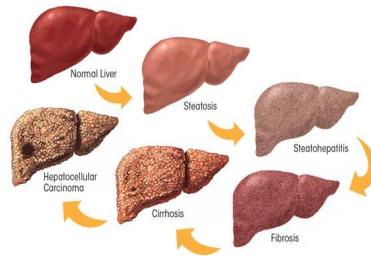


Institut für Bioinformatik



## Non-Alcoholic Fatty Liver Disease

- starts with **accumulation of fat** in the liver, progresses in 10-30% of the cases
- usually **no signs and symptoms** in early stages of NAFLD (sometimes fatigue or discomfort in the upper right abdomen)
- Leads to severe liver disease, cirrhosis & cancer
- Disease is underdiagnosed
- No noninvasive diagnostic method suitable for large-scale screening, diagnosis requires liver biopsy



**Universitätsmedizin  
GREIFSWALD**

Institut für Bioinformatik

## SHIP – the Study of Health in Pomerania



[www.shipstudy.de](http://www.shipstudy.de) 18.4.2005

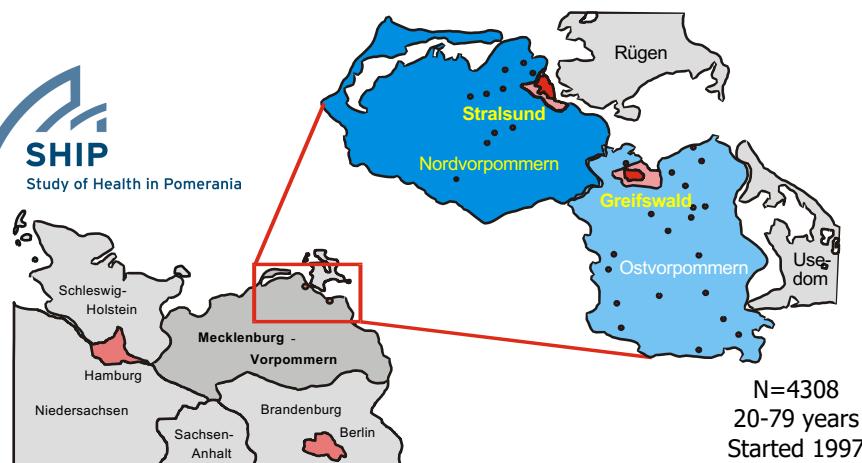
**Universitätsmedizin  
GREIFSWALD**

Institut für Bioinformatik

## SHIP – the Study of Health in Pomerania



**SHIP**  
Study of Health in Pomerania



Institut für Bioinformatik

## SHIP – the Study of Health in Pomerania

**SHIP-0**  
Baseline

n= 4308  
(68.8%)

**SHIP-1**  
5y Follow-up

n= 3300  
(84.5%)

**SHIP-2**  
10y Follow-up

n= 2333  
(63.5%)

**SHIP-3**  
15y Follow-up

n= 1718



Study of Health in Pomerania

1997-2001

2002-2006

n= 4420  
(50.1%)

2008-2012

2014-2018

Institut für Bioinformatik



**Universitätsmedizin**  
GREIFSWALD

## SHIP – the Study of Health in Pomerania

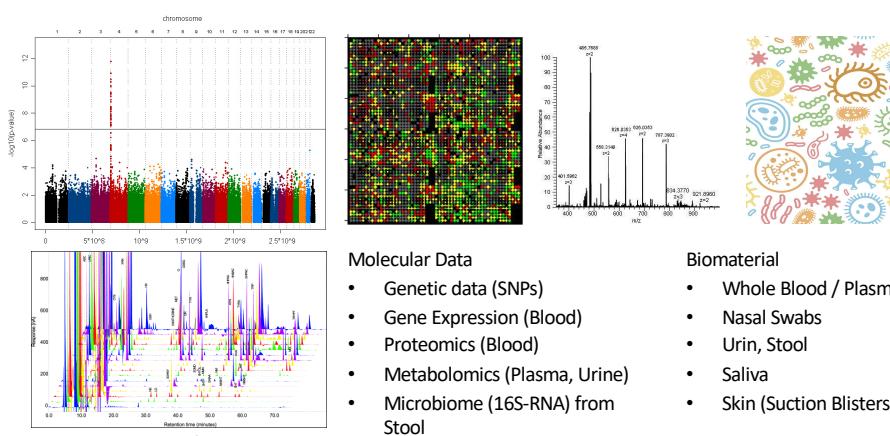


Völzke et al.; Int J Epidemiol 2011



Institut für Bioinformatik

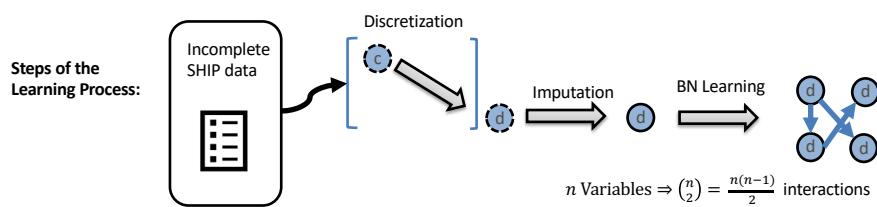
## Molecular Data in SHIP-Trend / Biomaterial



Institut für Bioinformatik

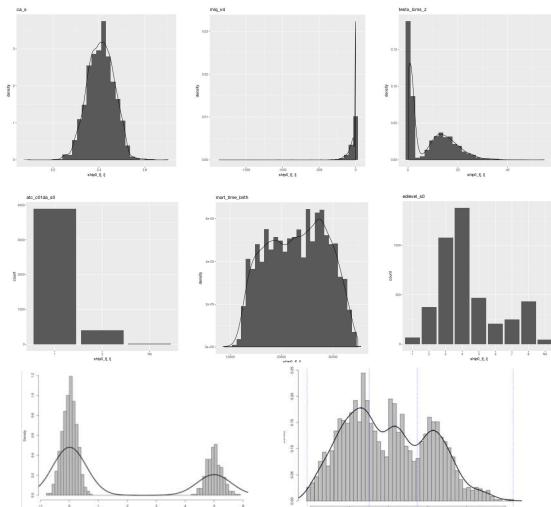
# Bayesian Network Idea

Idea: Learn a Bayesian Network from SHIP data to model all interactions of processes inside and outside of the liver, which possibly drive NAFLD



Institut für Bioinformatik

## What do our data look like?

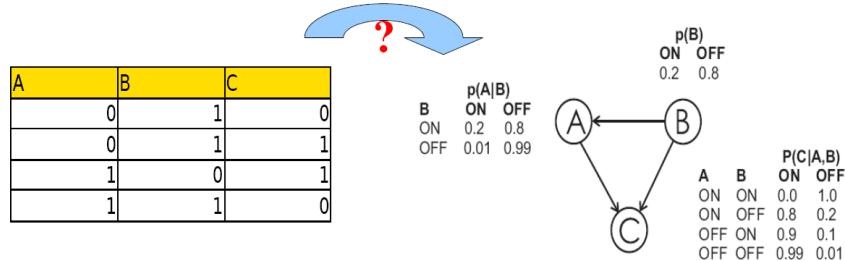


- Large number of variables, small number of patients ( $p > n$ )
  - Mix of discrete and continuously distributed variables
  - Variables at different scales
  - Large number of missing values
  - Highly correlated variables
  - High levels of noise



Institut für Bioinformatik

## Inferring a Bayesian Network from Data



- Two interdependent questions to solve: Find optimal model topology, and optimal model parameters
- Maximum Likelihood:  $\max p(D|M)$  / MCMC Sampling from the likelihood
- Due to low number of samples and noisy data usually does not work
- Instead, work with the posterior distribution

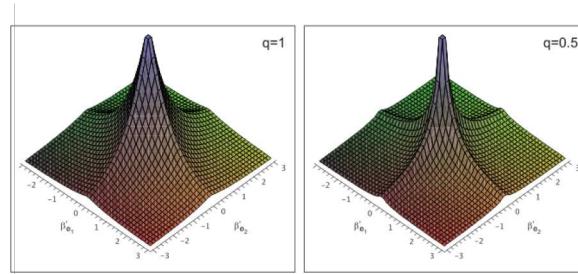
$$p(M|D) = \frac{p(D|M)p(M)}{\int_M p(D|M)p(M)dM}$$

- $P(M)$  can be used to include prior knowledge

## Prior Distribution

### Prior Distribution:

- If biological knowledge is available, can be included
- Simplest assumption is expectation of „sparse network“



$$p(M) = C e^{-\frac{1}{qs^q} |w|^q}$$

## Automatic Adaptive Grid Refinement

(Hydrodynamics)

= Method of adapting the accuracy of a solution within certain sensitive or turbulent regions of simulation, dynamically and during the time the solution is being calculated

Start computing a solution on a global grid

1. Apply an Error Estimator at every point of the grid
  2. Flag points with high estimate
  3. Create a finer grid at the parts of the grid where the estimate is high, so that every flagged point is contained in the finer grid
  4. The subgrid is nested in the coarse grid
- The procedure is iteratively performed until all error estimates fall below a fixed threshold

Application to Bayesian Networks:

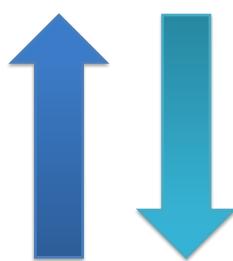
Start Learning a network between a small number of coarse clusters

1. Error Estimation for every cluster and flag clusters with high errors
2. Split clusters with high error
3. Add new clusters to the data and delete old ones
4. Refine network around new clusters



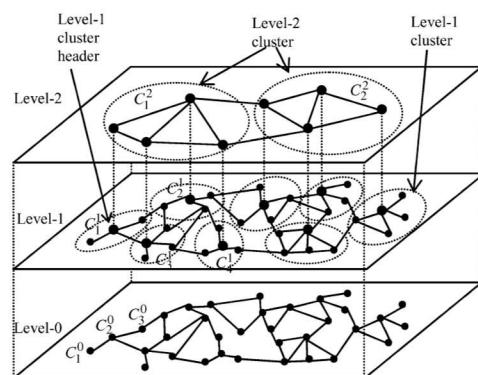
Institut für Bioinformatik

## Hierarchical Learning of Large Bayesian Networks



Idea: Learn hierarchy from data (via hierarchical clustering)

Learn networks between groups, not single variables

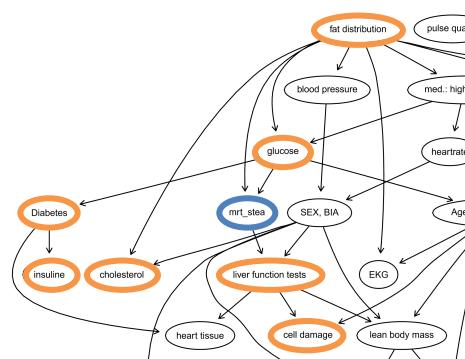


Refinement done taking relevant output variable (liver disease!) into account



Institut für Bioinformatik

## Results NAFLD



**Fat distribution:**  
BMI  
whtr  
Koerperfett\_in\_prozent  
etc.

**Liver Function Tests:**  
Alat, Asat, GGT

**Cell damage:**  
Lactatdehydrogenase  
Creatinkinase

**Cholesterol:**  
Chol\_hdl (HDL-Quotient)  
Triglycerides  
LDL Cholesterol  
Cholesterol (in total)  
HDL Cholesterol

Next steps: Refine Methods, Include Molecular Data



Institut für Bioinformatik

## Further Examples from our Work

- Intraoperative classification of tissue type during brain cancer operations
- Prediction of Mortality using Deep Learning on SHIP data
- Prediction of Mortality based on Electronic Health Records
- Using neuronal networks to predict cardiovascular disease on ECG data
- Predicting therapy response and survival in cancer patients using penalized survival models
- Risk factors for thyroid disease
- Prediction of therapy response for prediabetic patients
- Prediction of depression based on social media data



Institut für Bioinformatik

## Summer School “AI in Medicine”

- Zielgruppe: Naturwissenschaftliche Doktoranden und junge Ärzte
- 15 Mediziner + 15 Informatiker/Mathematiker
- **8. bis 12. Juni 2020 auf Hiddensee**
- Dozenten:

Tim Beissbarth, Unimedizin Göttingen	Harald Binder, Uniklinik Freiburg
Tim Friede, Unimedizin Göttingen	Nils Grabe, Uni Heidelberg
Christoph Lippert, HPI Potsdam	Andreas Stahl, Greifswald
Mario Stanke, Greifswald	Fabian Theis, Helmholtz München
Olaf Wolkenhauer, Rostock	Lars Kaderali, Greifswald
- Mehr Infos unter  
<http://www.kaderali.org>



Universitätsmedizin  
GREIFSWALD

JOACHIM  
HERZ  
STIFTUNG



## THANK YOU FOR YOUR ATTENTION



Universitätsmedizin  
GREIFSWALD

Institut für Bioinformatik