



FORSCHUNGS  
**CAMPUS**

öffentlich-private Partnerschaft  
für Innovationen

SPONSORED BY THE



Federal Ministry  
of Education  
and Research

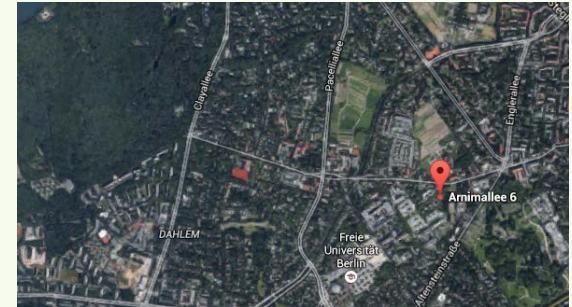
## *Omics-based Cancer Diagnostics*

Tim Conrad, Freie Universität Berlin

Forschungs Campus MODAL, **MedLAB**

# About me...

- Studied BioInformatics & Computer Science
- Ph.D. Mathematics
- Professor for Medical Bioinformatics  
@ Freie Universität Berlin
- Current research topics:
  - Network & compressed-sensing based analysis methods  
([Einstein Center for Mathematics](#))
  - Frameworks for mass-data processing ([Berlin Big Data Center](#))
  - Analysis of bio-medical mass-data ([Forschungscampus MODAL](#))



# After this talk you should know...

- ... what -omics data is
- ... that -omics data can be very large
- ... why it is useful to analyze -omics data
- ... a method how to analyze large -omics data
- ... a framework implementing this method

*Spoiler: we will do this during the hands-on session*

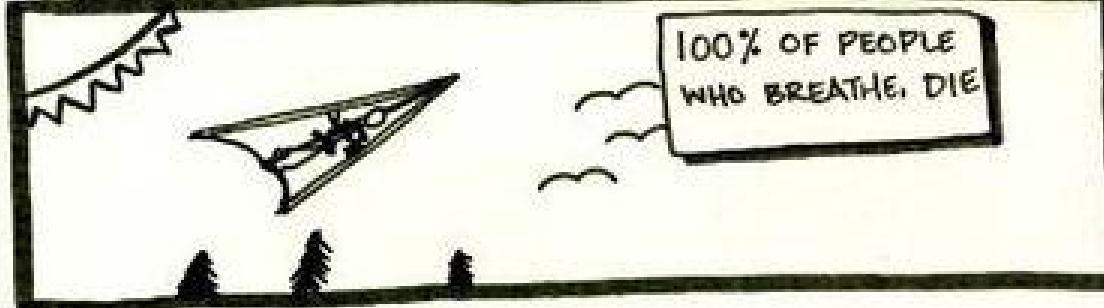
**ACHTUNG**



# Correlation is not Causation!

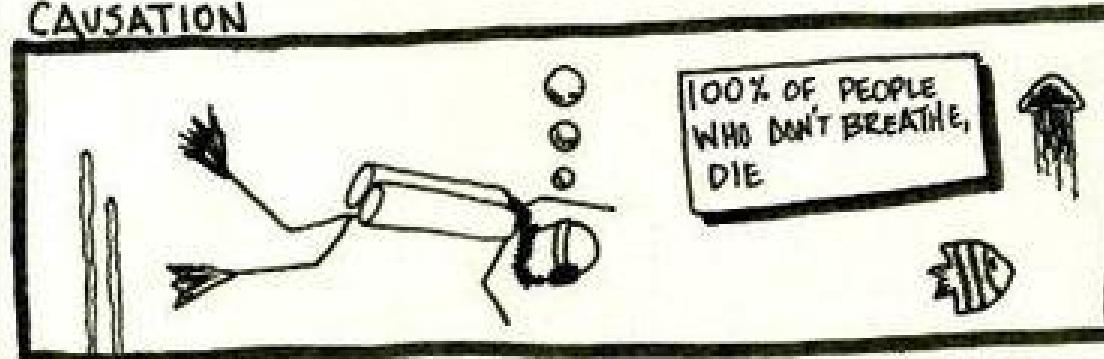
CORRELATION

[www.asandford.com](http://www.asandford.com)



CAUSATION

100% OF PEOPLE  
WHO DON'T BREATHE,  
DIE





**example**

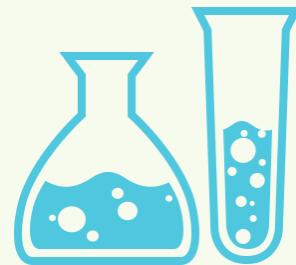
## TRAINING



## APPLICATION



**YES OR NO?**



## Introduction

Background & Goal

## Omics Data Analysis

How can it be done?

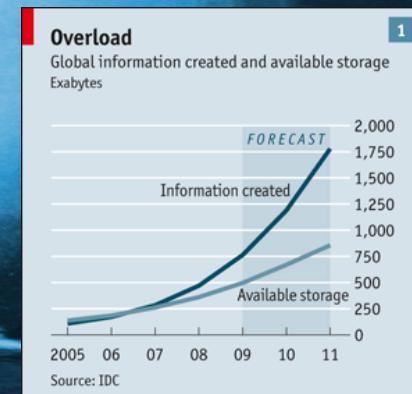
## Application

Reaching out  
to the real world





Healthcare is challenged by large amounts of data in motion that is **diverse**, **unstructured** and growing exponentially.



Source: iHT^2

Image source: Flickr



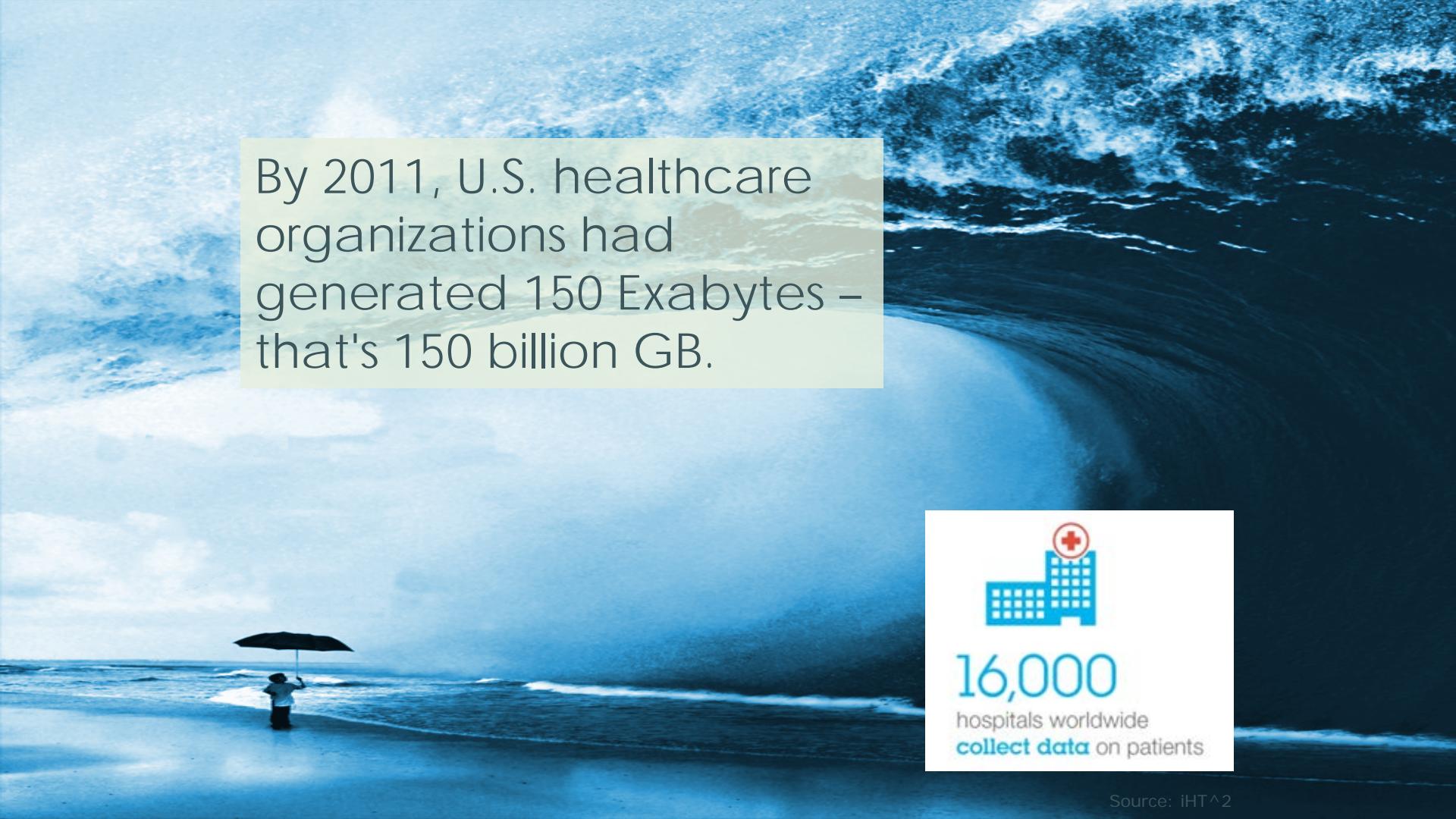
Data constantly streams in through interconnected sensors, monitors and instruments in real-time faster than a physician or nurse can keep up.



Patient monitoring equipment pumps out an average of

**1,000**

readings per second or  
86,400 readings in a day



By 2011, U.S. healthcare organizations had generated 150 Exabytes – that's 150 billion GB.



# 100 trillion human cells

3.2 billion base pairs of DNA

30,000 genes encoding proteins

10 million total distinct proteins in a person

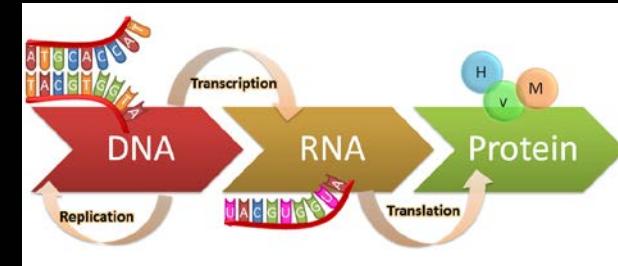
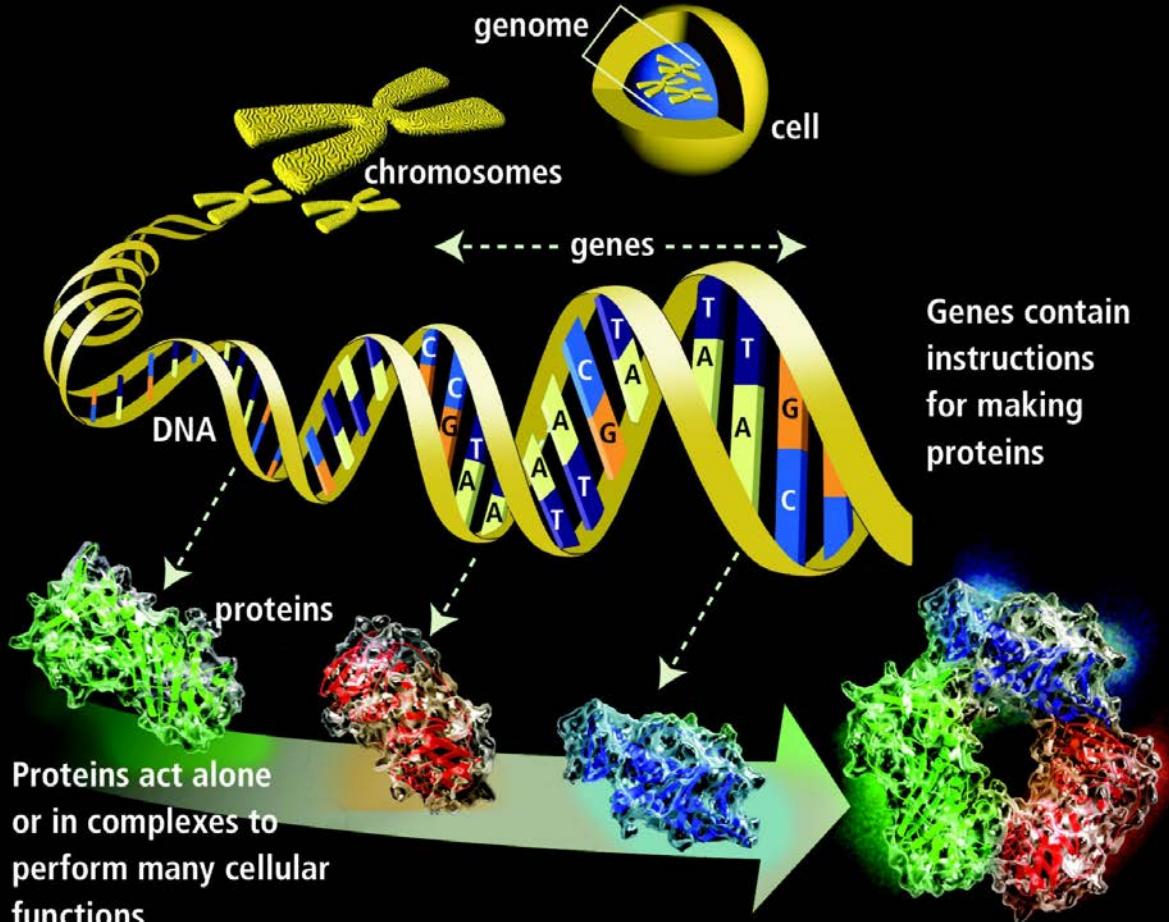
2000 distinct proteins functioning in a cell

60,000 reactions/cell/minute

100,000's of molecular events

50 or so organs and organ systems

## Intro - Central Dogma



# WHY IS THIS INTERESTING?

# 100 trillion human cells

3.2 billion base pairs of DNA

30,000 genes encoding proteins

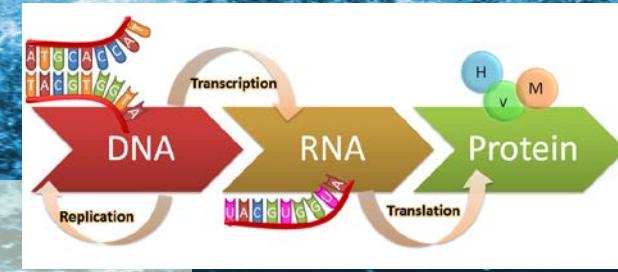
10 million total distinct proteins in a person

2000 distinct proteins functioning in a cell

60,000 reactions/cell/minute

100,000's of molecular events

50 or so organs and organ systems



# 100 trillion human cells

3.2 billion base pairs of DNA

30,000 genes encoding proteins

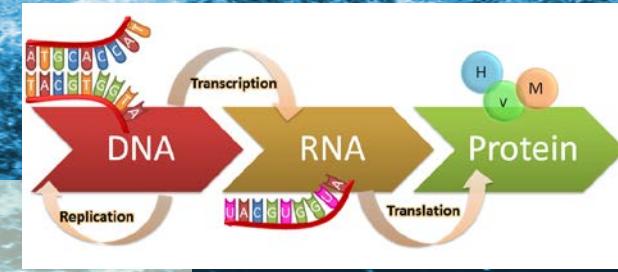
10 million total distinct proteins in a person

2000 distinct proteins functioning in a cell

60,000 reactions/cell/minute

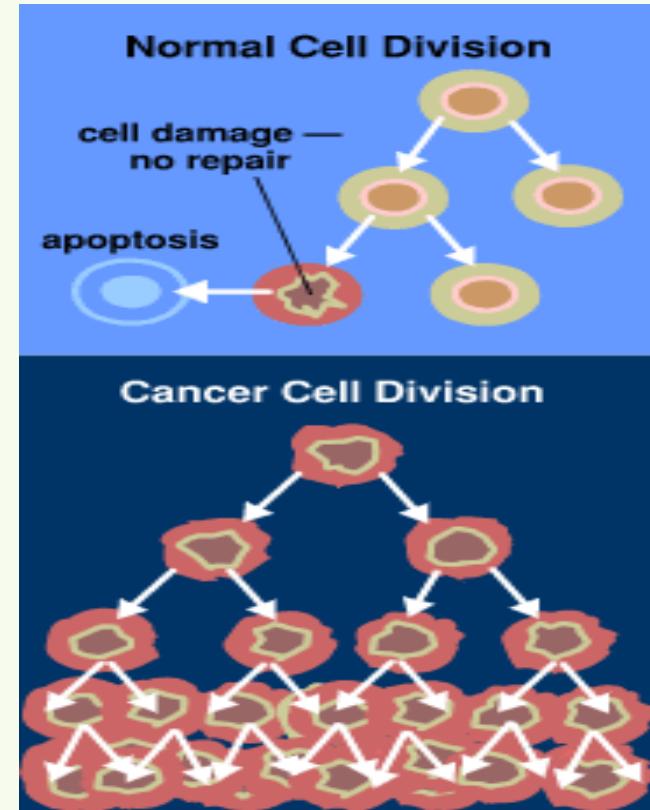
100,000's of molecular events

50 or so organs and organ systems

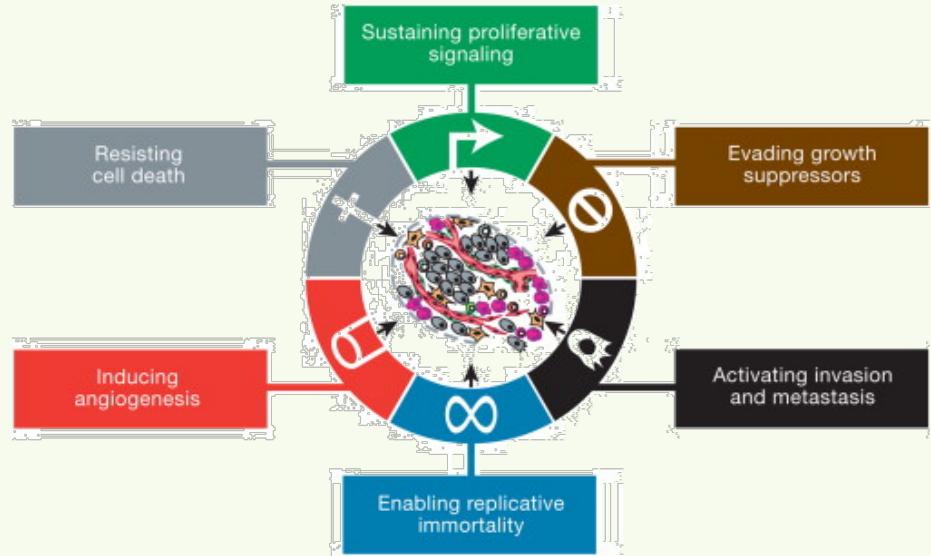
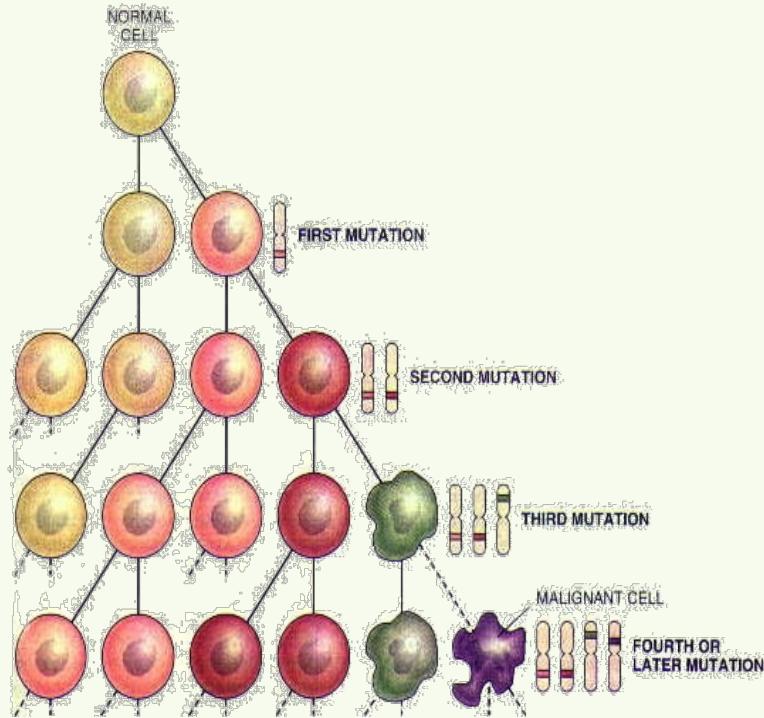


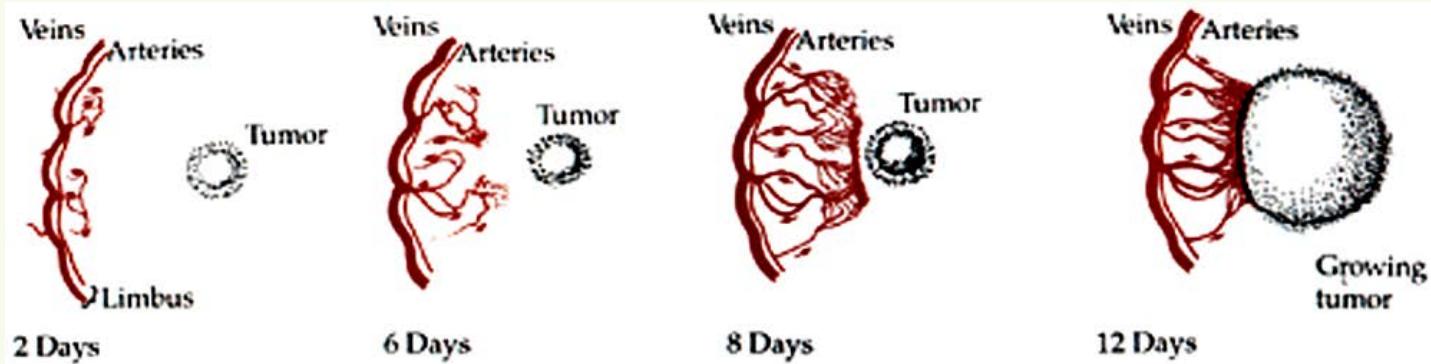
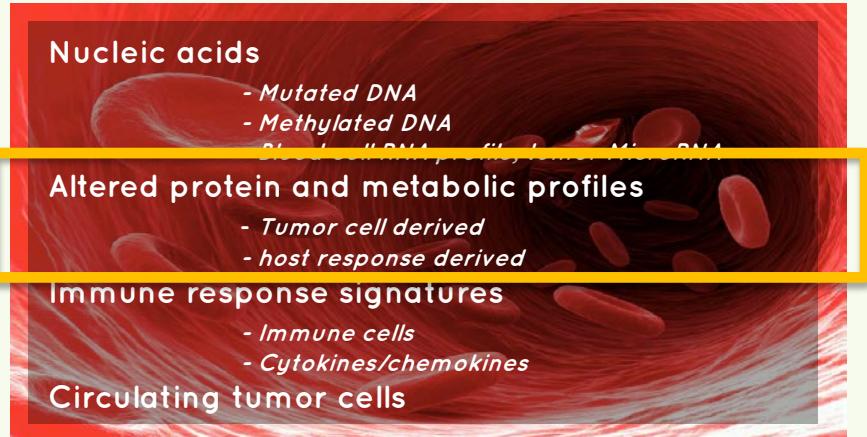
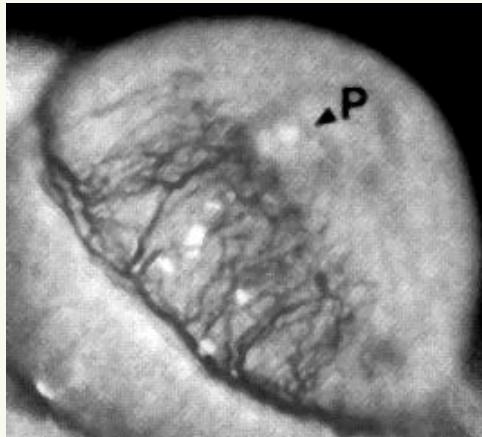
# CANCER MOLECULAR BIOLOGY IN A NUTSHELL

- All cancers derive from single cells that have **acquired the characteristics of continually dividing in an unrestrained manner** and invading surrounding tissues.
- Cancer cells behave in this abnormal manner because of **mutations in the DNA sequence of key genes**, which are known as cancer genes. Therefore all cancers are genetic diseases.



# Mutations in **multiple** cancer genes are required for the development and progression of a single cancer.







FACT: Diseases leave fingerprints in the human body (e.g. blood proteins)



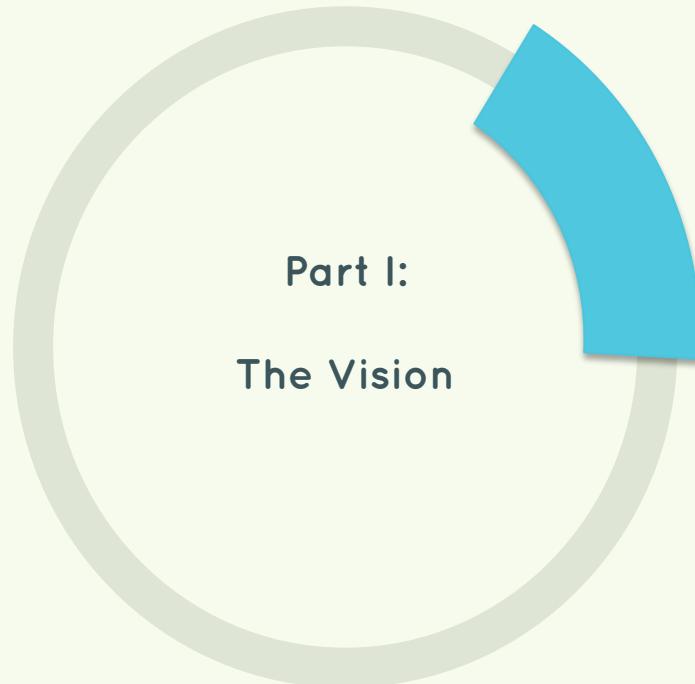
Diseases leave fingerprints in the human body (e.g. blood proteins)



BUT: 10 million total distinct proteins in a person



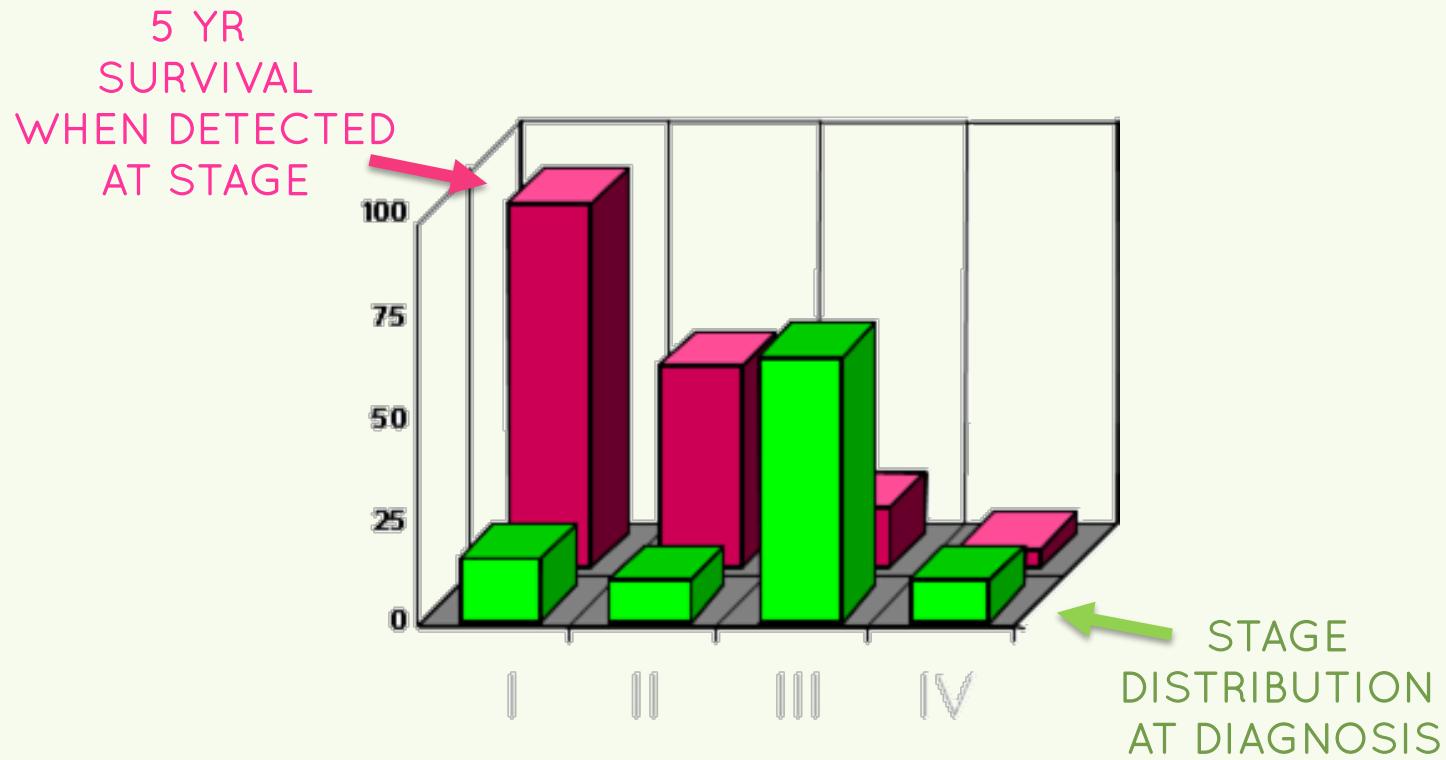
- All cancers result from changes in the DNA sequence of our genome.
- Occasionally, a series of these mutations alters the function of a set of critical genes, resulting in tumour formation.
- This data can be measured with -omics technology,  
e.g. from blood samples



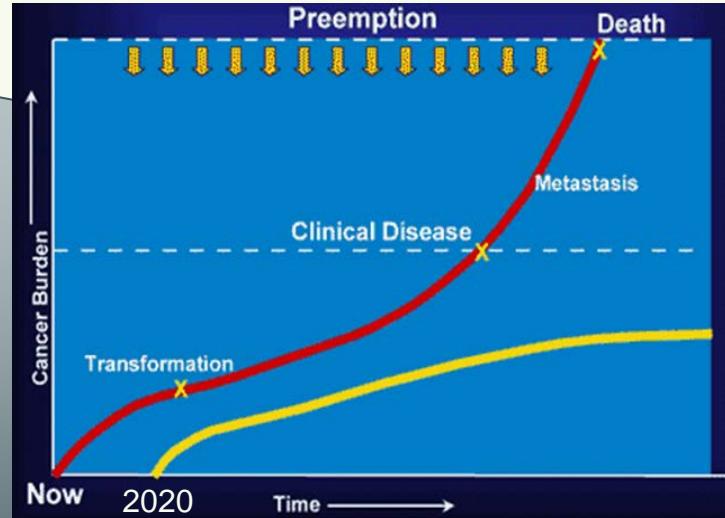
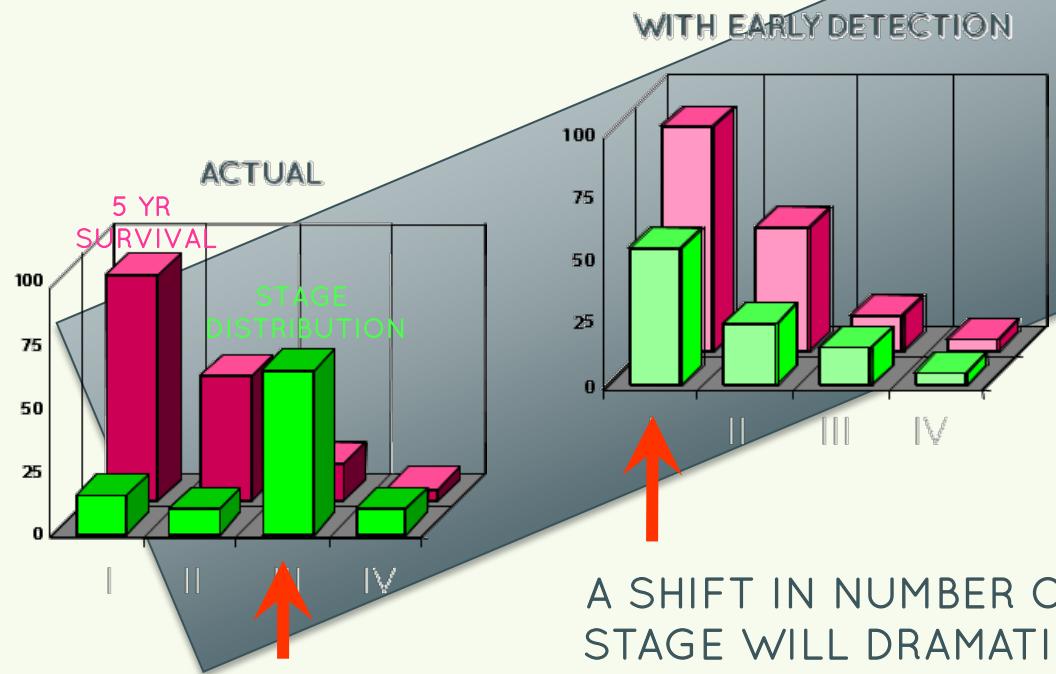


- There are approximately 200 types of cancer, each with different causes, symptoms and treatments
- In 2007, ~470.000 people were newly diagnosed with cancer in Germany
- An individual's risk of developing cancer depends on many factors, including age, lifestyle and genetic make-up
- One in three people in the Western world develop cancer and one in five die of the disease

# The main problem



Develop a minimal-invasive method for early cancer detection allowing medical doctors to start treatment earlier!

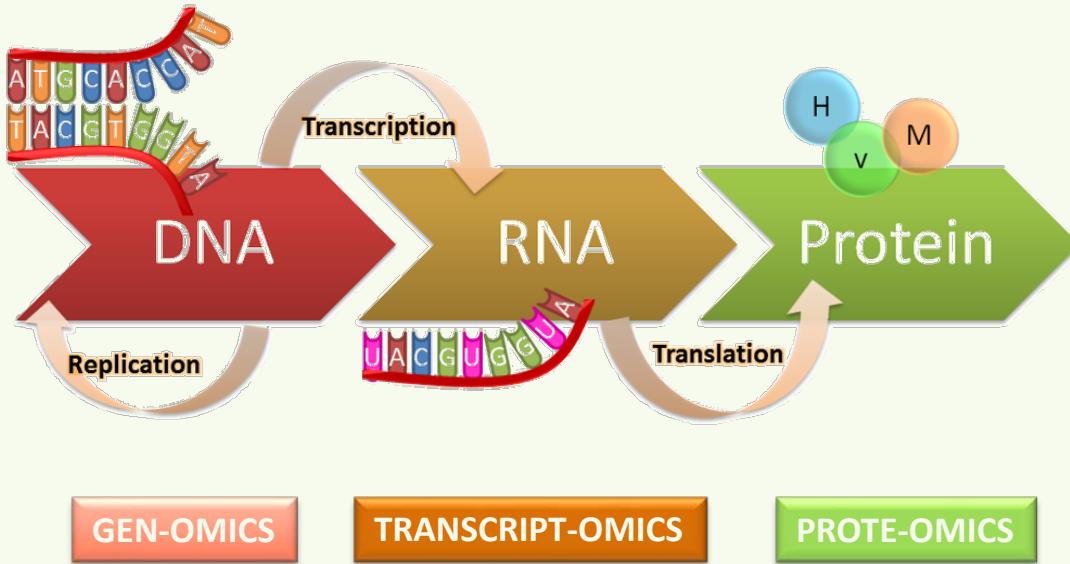


# The Strategy

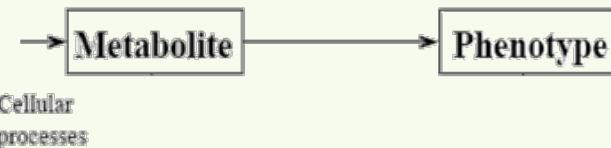
- Identify changes (signals) in the human body that only occur if a disease is present
- Use this information to test a patient whether the signal(s) are present
- Use (easily accessible) data, e.g. derived from blood, that allows to track very early changes



# THE DATA



Manifestation of interactions between “-omics” and the environment



What might happen  
What can happen

What happens/happened



# Size of OMICS data

- Genome : **tens of GBs** raw data
  - Transcriptome : **tens of GBs** raw data
  - Proteome : **several GBs** raw data
  - Metabolome : **hundreds of MBs** raw data
  
  - Phenome : potentially TBs
- Epigenome, lipidome, glycome, interactome, spliceome, mechanome, exposome, etc...





- All cancers result from changes in the DNA sequence of our genome.
- The changes are “transported” through the -omics “layers”
- Consequence: cancer cells leave traces (fingerprints), e.g. in the blood proteome.  
This can be measured by -omics technologies.

# THE (BASIC) METHODS

# Problem Statement

In biomedical mass-data analysis we are often interested in three main things:

- (1) **CLASSIFICATION** - Given the data: **which class**  
(e.g. “healthy” vs. “diseased”) does the patient belong to?
- (2) **REGRESSION** - Given the data: **which score**  
(e.g. “survival estimation”) does a patient get?
- (3) **FEATURE SELECTION** - Given the data: **which features**  
(e.g. expression level of genes A and D) are most relevant  
(e.g. for classification or regression)?

Find significantly changed protein concentrations to predict class

# CLASSIFICATION & PROTEOMICS



Diseases leave fingerprints in the human body (e.g. blood proteins)

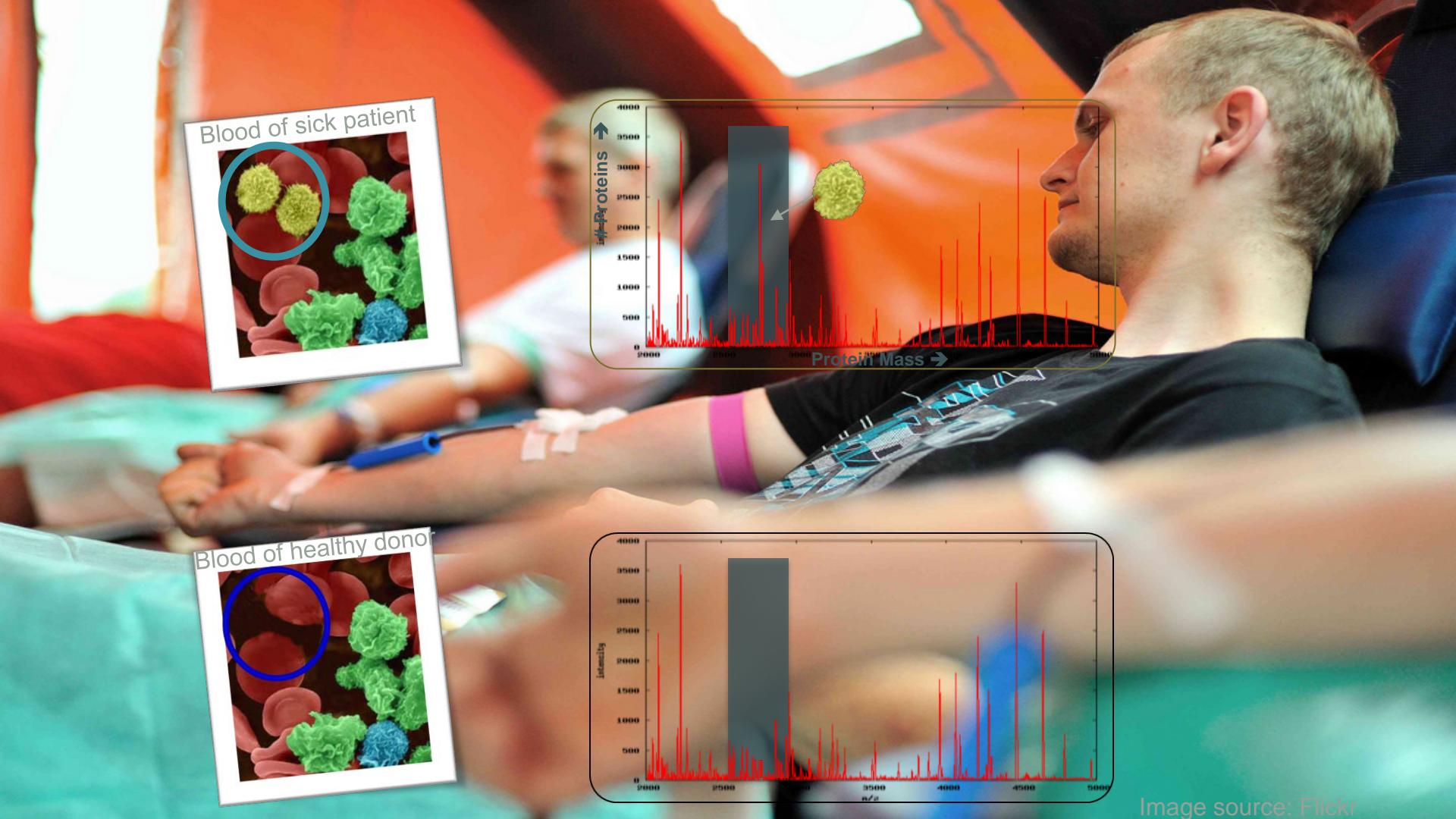


Image source: Flickr

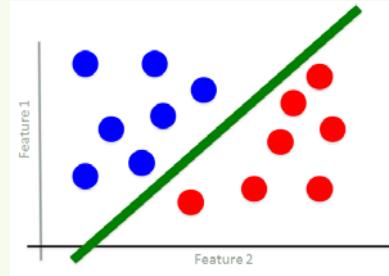
- It is possible to identify cancer specific fingerprints for early diagnosis based on mass spectrometry-based proteomics



1. Leichtle, A. and Nuoffer, J.-M. and Ceglarek, U. and Kase, J. and Conrad, T. O. F. and Witzigmann, H. and Thiery, J. and Fiedler, G. M. (2011) **Serum amino acid profiles** and their alterations in **colorectal cancer**. Metabolomics . ISSN 1573-3890 (In Press)
2. Leichtle, A. and Ceglarek, U. and Kase, J. and Conrad, T. O. F. and Hauss, J. and Witzigmann, H. and Thiery, J. and Fiedler, G. M. (2010) **Metabolome Alterations in Pancreatic Cancer**. Clinical Cancer Research . ISSN 1078-0432 (Submitted)
3. Strenziok, R. and Hinz, S. and Wolf, C. and Conrad, T. O. F. and Krause, H. and Miller, K. and Schrader, M. (2009) Surface-enhanced laser desorption/ionization time-of-flight mass spectrometry: **serum protein profiling** in **seminoma** patients. World J of Urology, 28 (2). pp. 193-197.
4. Fiedler, G. M. and Leichtle, A. and Kase, J. and Baumann, S. and Ceglarek, U. and Felix, K. and Conrad, T. O. F. and Witzigmann, H. and Weimann, A. and Schütte, Ch. and Hauss, J. and Büchler, M. and Thiery, J. (2009) **Serum Peptidome Profiling Revealed Platelet Factor 4 as a Potential Discriminating Peptide Associated With Pancreatic Cancer**. Clinical Cancer Research, 15 (11). pp. 3812-3819. ISSN 1078-0432
5. Strenziok, R. and Hinz, S. and Wolf, C. and Conrad, T. O. F. and Krause, H. and Lingnau, A. and Lein, M. and Miller, K. and Schrader, M. (2008) **Serum proteomic profiling by surface-enhanced laser desorption/ionization time-of-flight mass spectrometry in testicular germ cell cancer** patients. European Urology Supplements, 7 (3). p. 83.
6. Unpublished: lung cancer (2012)

Preliminary results: **fingerprints for 6 different cancer types** two of which have proven successful in first clinical studies.

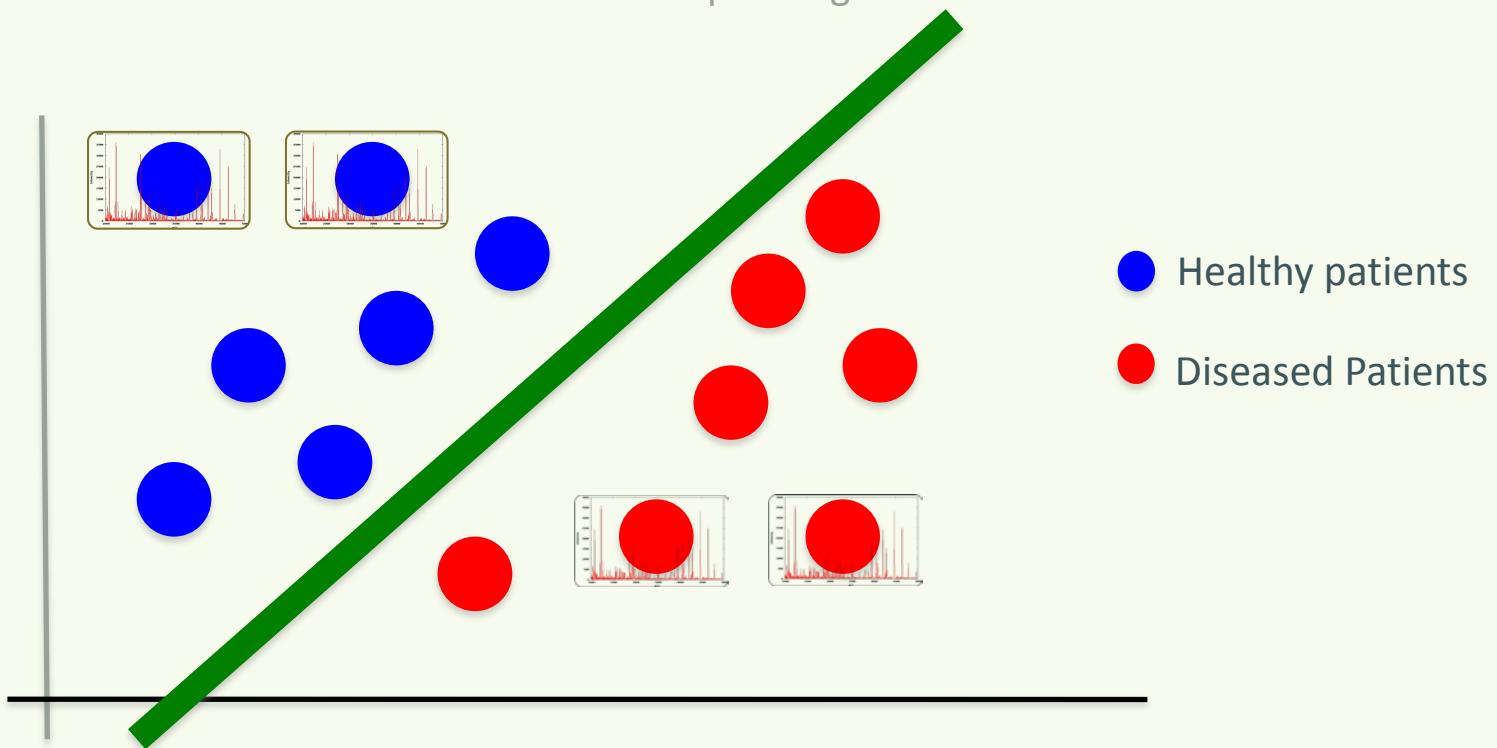
# Classification in General



- Classification is a supervised learning problem
- Preliminary task is to construct classification rule (some functional form) from the training data
- For  $p \ll n$ , many methods are available in classical statistics,
  - ◆ Linear (LDA, LR)
  - ◆ Non-Linear (QDA, KLR)
- However when  $n \ll p$  we face estimability problem (“Curse of dimensionality”)
- Some kind of data transformation is inevitable
- Well known techniques for  $n \ll p$ : PCA , SVM etc.

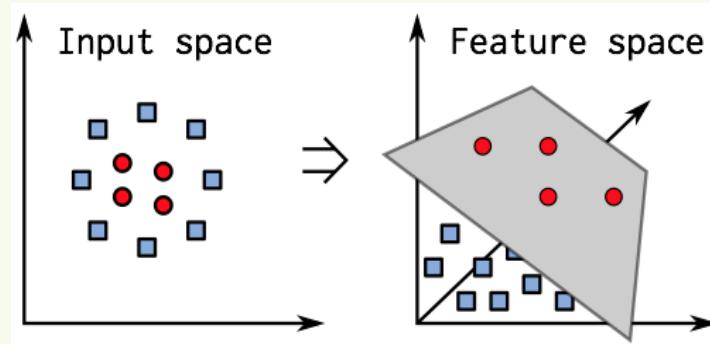
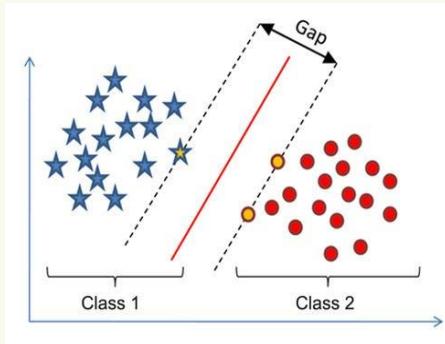
# Classification (e.g. SVM)

Determine separating line

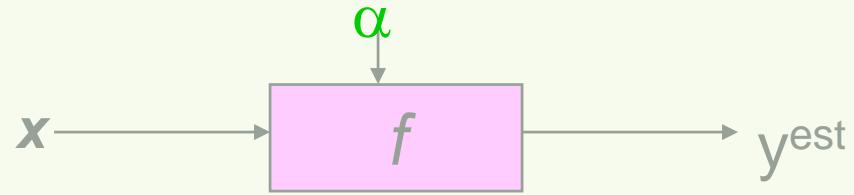


# Support Vector Machines (SVM)

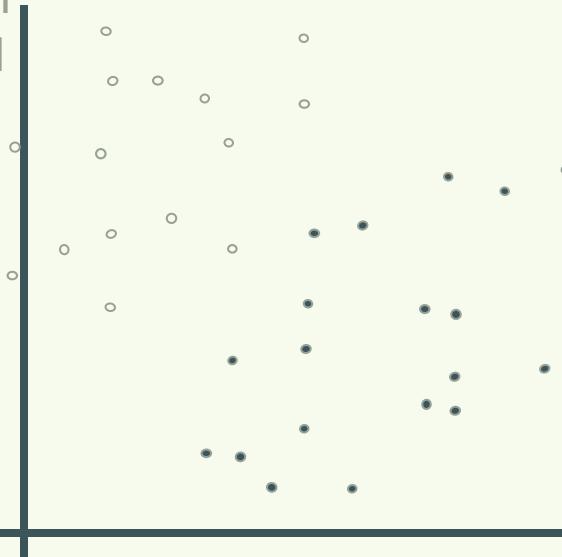
- Goal of SVM method: create a hyperplane separating data from 2 classes that maximizes margin between the hyperplane and the classes
- Formulation:  $\min \|\beta\|$  subject to  $\begin{cases} y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i \quad \forall i, \\ \xi_i \geq 0, \sum \xi_i \leq \text{constant}. \end{cases}$
- Kernels
- Extensions: Post-processing with ROC variable importance to reduce the number of features



# Linear Classifiers



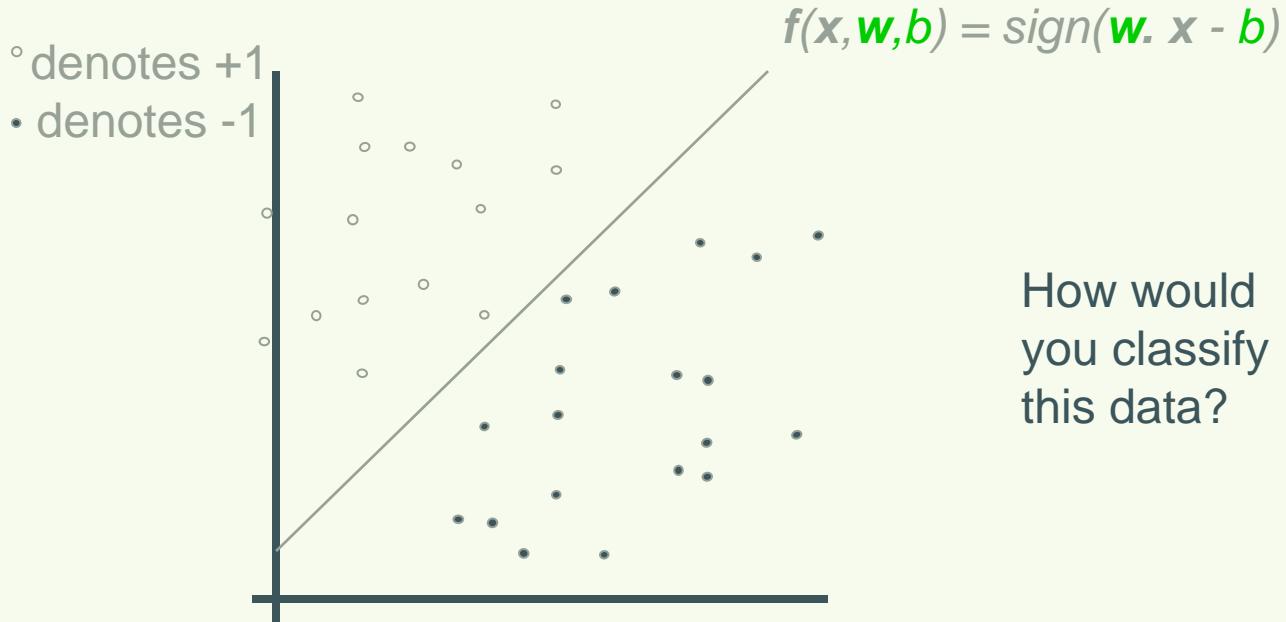
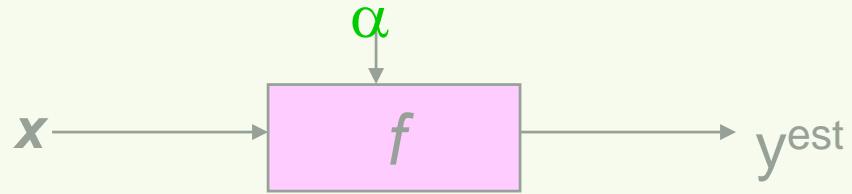
◦ denotes +1  
• denotes -1



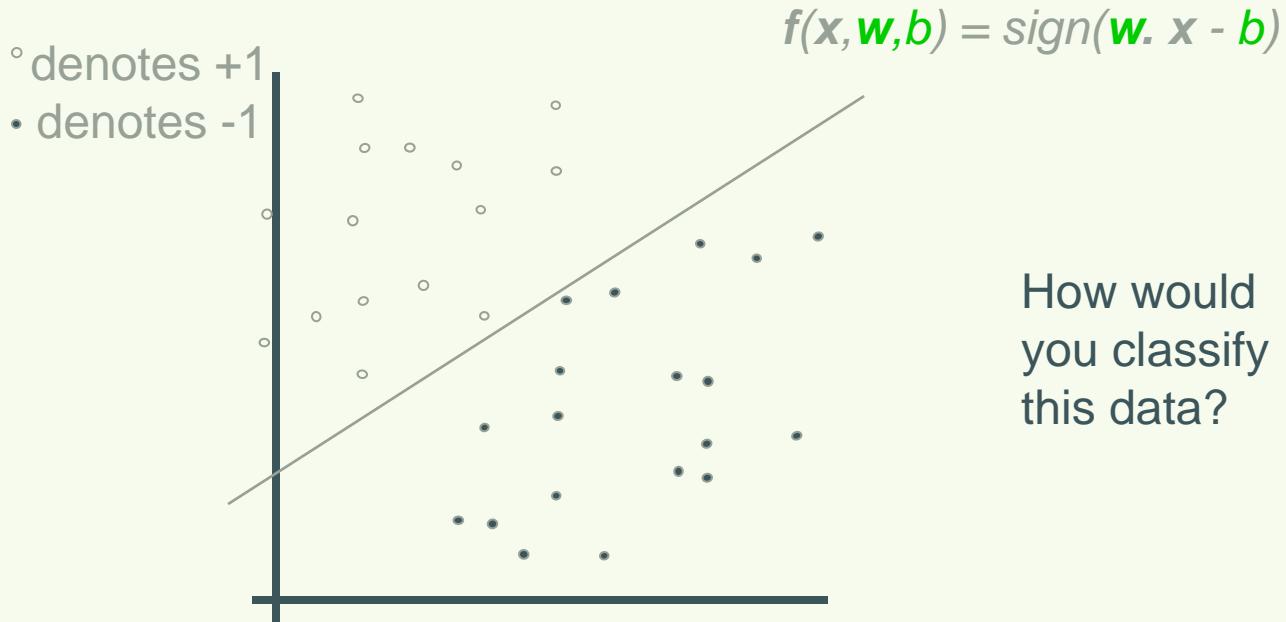
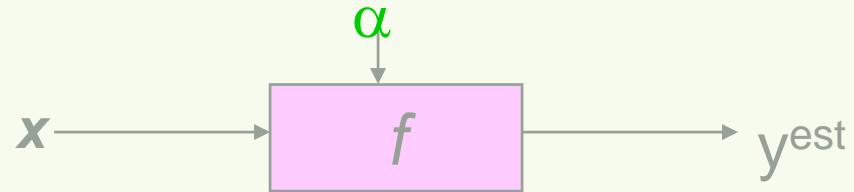
$$f(x, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

How would you classify this data?

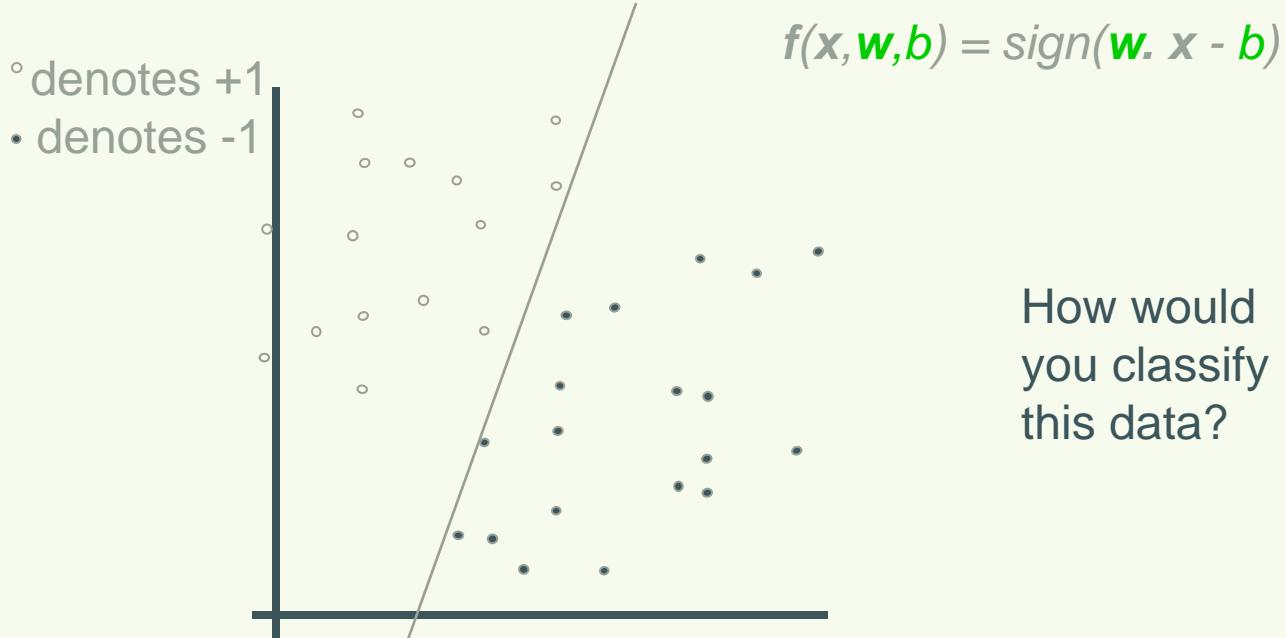
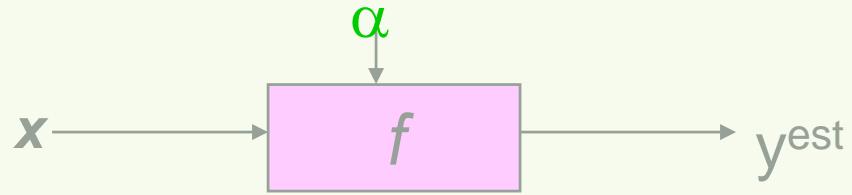
# Linear Classifiers



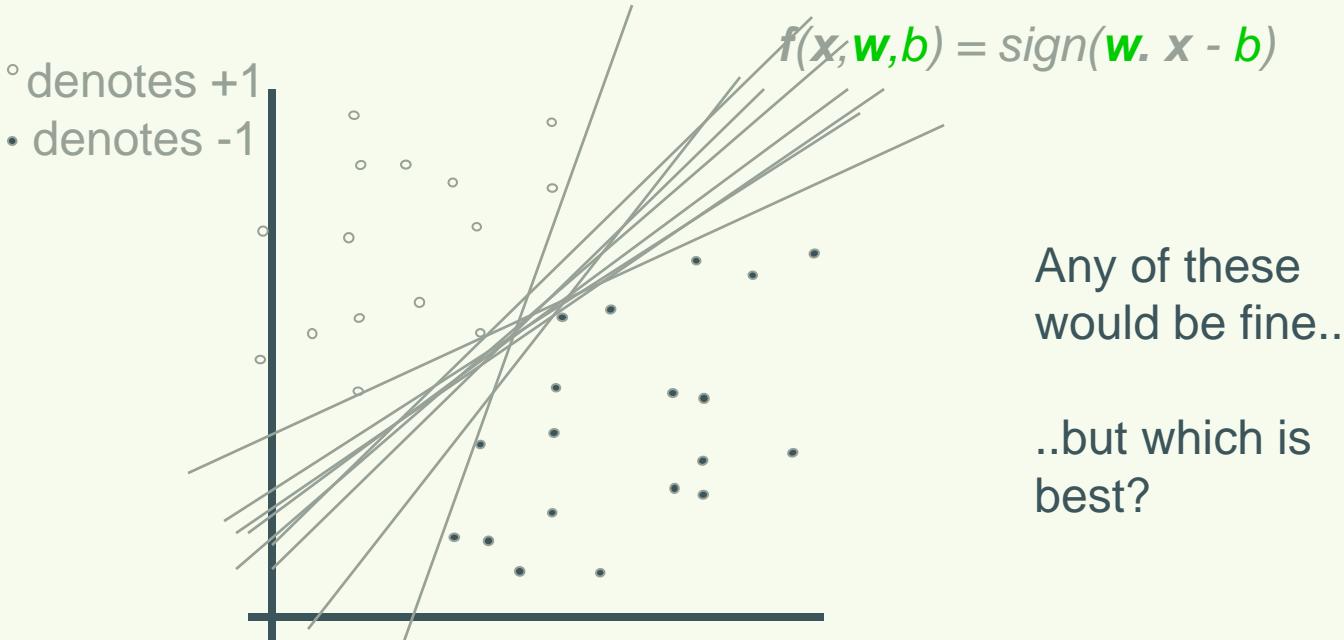
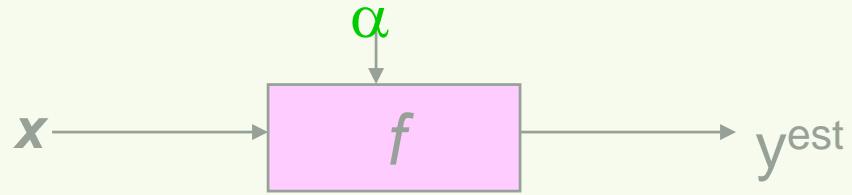
# Linear Classifiers



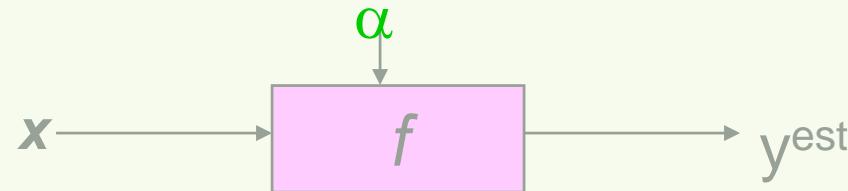
# Linear Classifiers



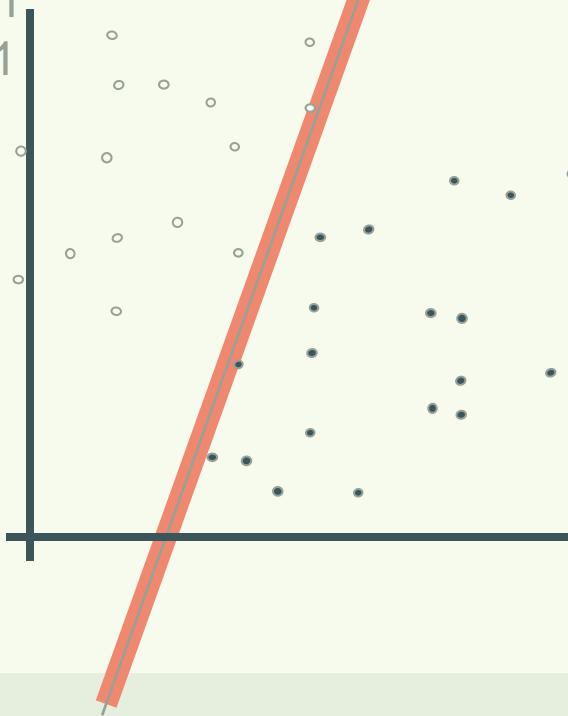
# Linear Classifiers



# Classifier Margin



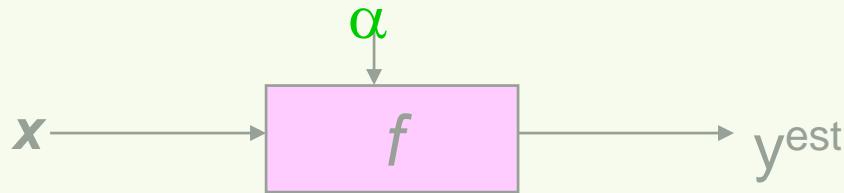
◦ denotes +1  
• denotes -1



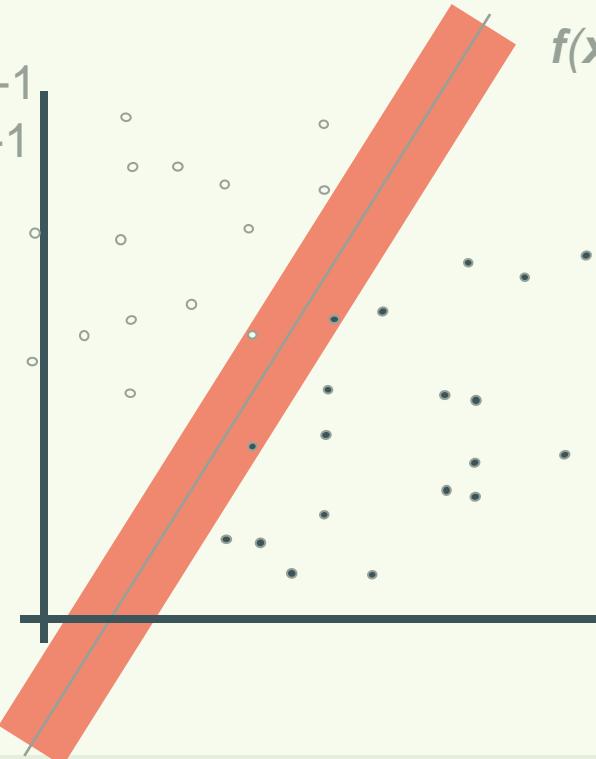
$$f(x, w, b) = \text{sign}(w \cdot x - b)$$

Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

# Maximum Margin

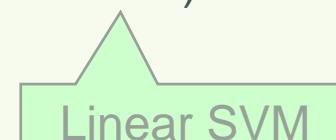


◦ denotes +1  
• denotes -1

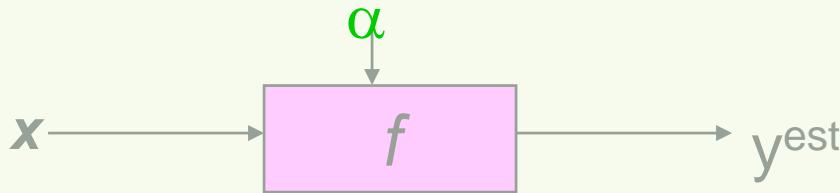


$$f(x, w, b) = \text{sign}(w \cdot x - b)$$

The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin. This is the simplest kind of SVM (Called an LSVM)

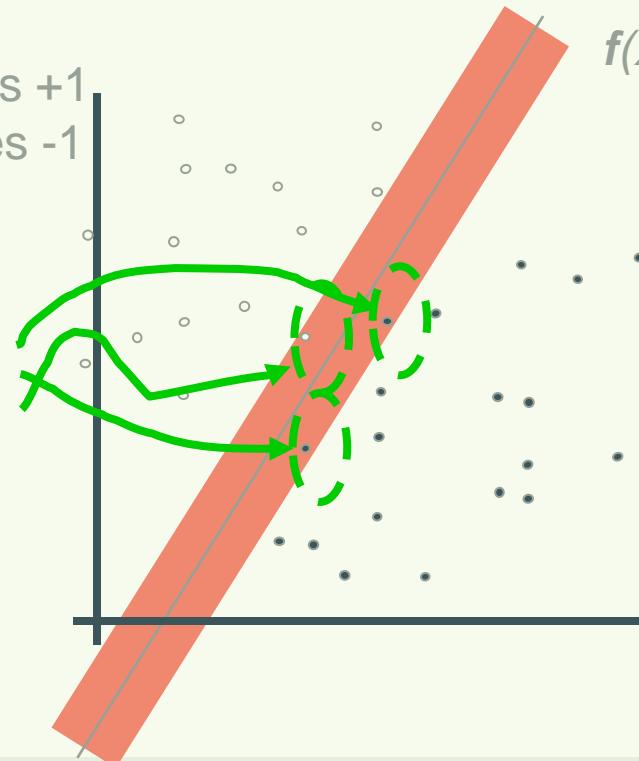


# Maximum Margin



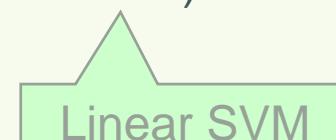
◦ denotes +1  
• denotes -1

**Support Vectors**  
are those  
datapoints that  
the margin  
pushes up  
against



$$f(x, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} - b)$$

The **maximum margin linear classifier** is the linear classifier with the, um, maximum margin. This is the simplest kind of SVM (Called an LSVM)



# Linear SVMs Mathematically

- We can formulate the *quadratic optimization problem*:

Find  $\mathbf{w}$  and  $b$  such that

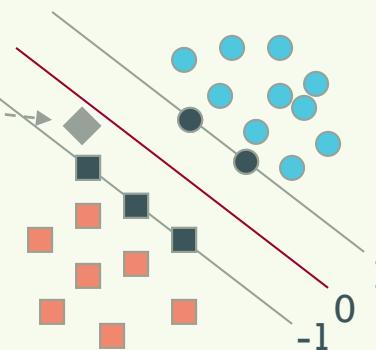
$$\rho = \frac{2}{\|\mathbf{w}\|} \quad \text{is maximized; and for all } \{(\mathbf{x}_i, y_i)\}$$

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \text{ if } y_i=1; \quad \mathbf{w}^T \mathbf{x}_i + b \leq -1 \quad \text{if } y_i= -1$$

# Classification with SVMs

- Given a new point  $\mathbf{x}$ , we can score its projection onto the hyperplane normal:
  - I.e., compute score:  $\mathbf{w}^T \mathbf{x} + b = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$ 
    - Decide class based on whether  $<$  or  $> 0$

- Can set confidence threshold  $t$ .
  - Score  $> t$ : yes
  - Score  $< -t$ : no
  - Else: don't know



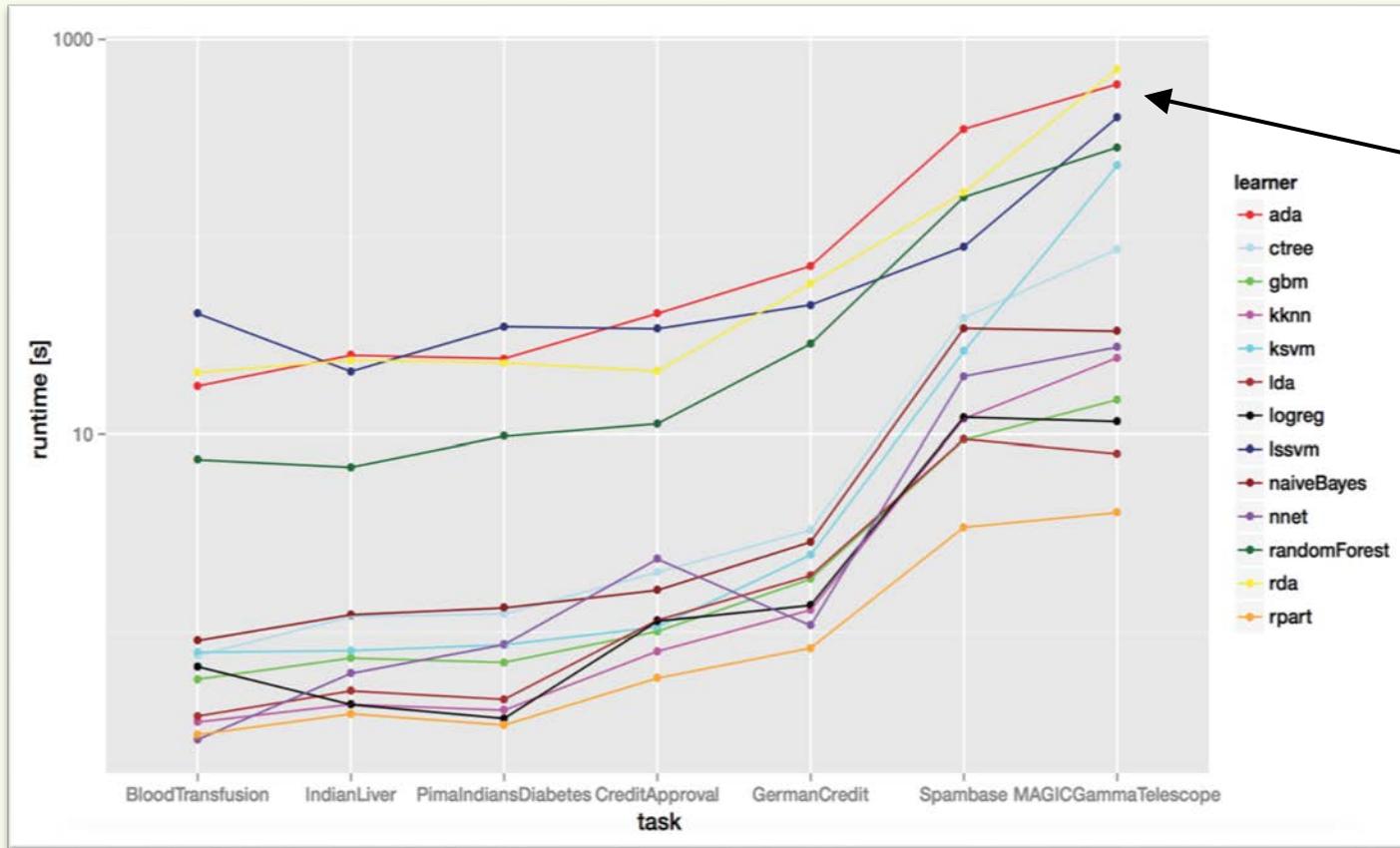
# Linear Regularized Loss Minimization

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

- ▶ support vector machines
- ▶ logistic regression
- ▶ lasso regression
- ▶ ridge regression
- ▶ etc...

# PRACTICAL ISSUE(S)

# Runtime analyses for machine learning R programs



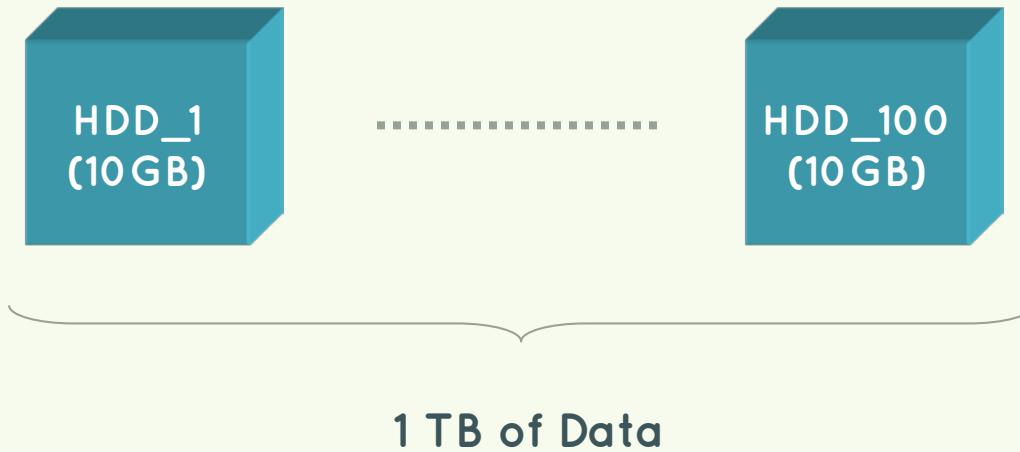
ML on a  
19,020 x 10  
matrix:

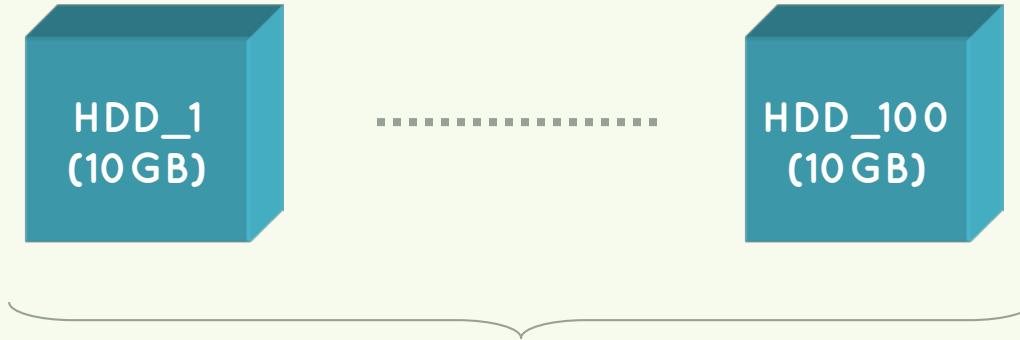
takes up to  
15 minutes

# Time & Size Problem

- What if data does not fit into memory any more?
- Goal: efficient solving of objective when data is distributed

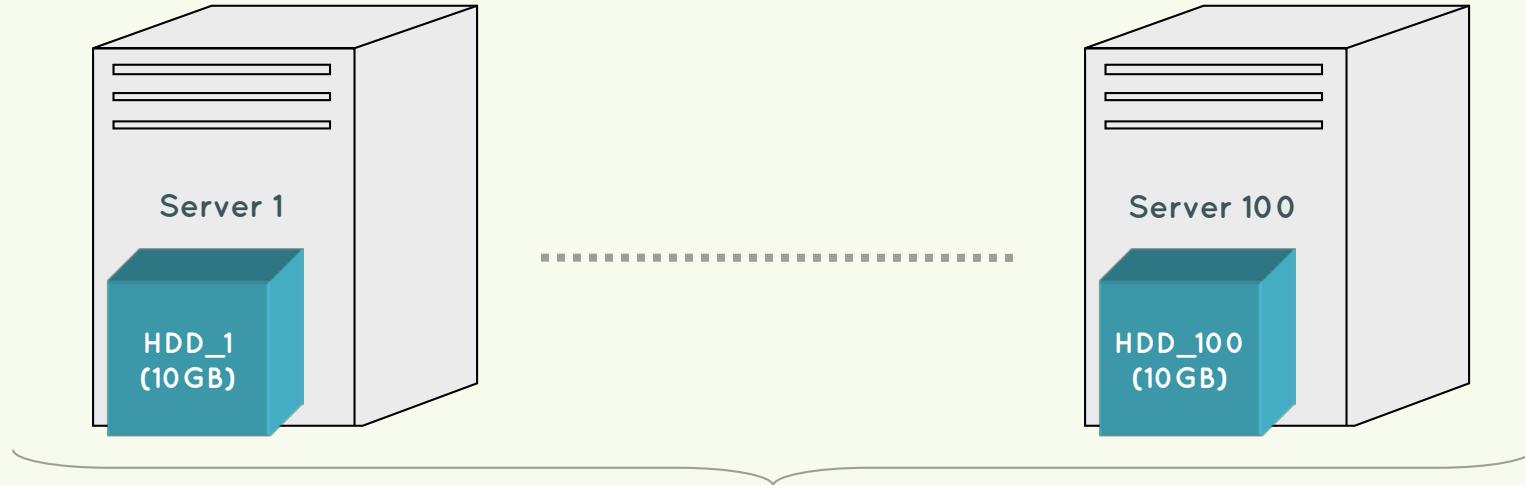
# Data Distribution





1 TB of Data in 100 HDDs, but how many computers should there be?  
When N=1, reading 1 TB requires 2.5 HOURS.

What should N be in order to give us appreciable speed-up on reads?



Given:

- 10 GB per drive
- $100 \times 10 \text{ GB drives} = 1 \text{ TB}$
- Read rate is 100 MB/second

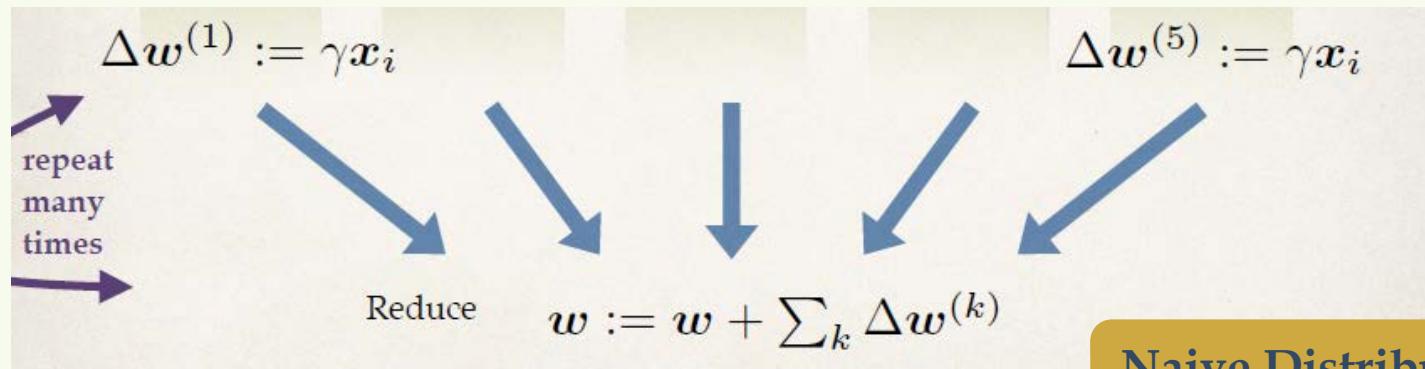
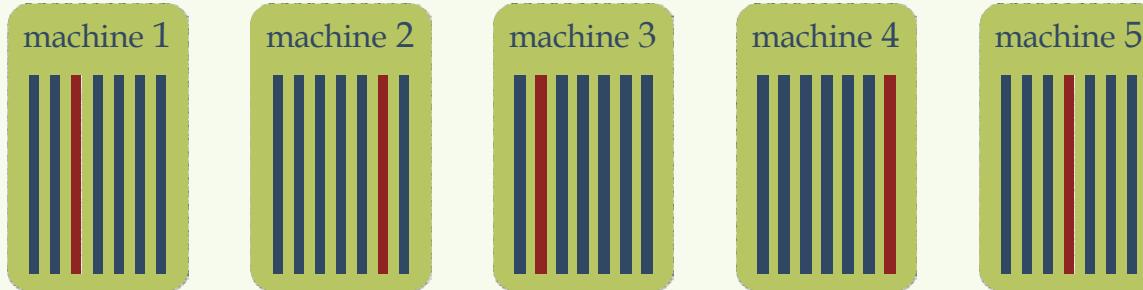
Full 1 TB of data can be read in 100 seconds :

$$10 \text{ GB} / 100 \text{ MB per second} = 100 \text{ seconds to read one drive.}$$

We read all 100 drives in parallel, and the computers can process the data read in parallel.

This is the architecture in which distributed computing frameworks shine, because not only is the data read in parallel, it is processed in parallel as well.

# Distributed Stochastic Optimization



**Naive Distributed SGD**

# The Cost of Communication

- ⊕ Reading  $v$  from Memory (RAM)

$100\text{ ns}$

$$v \in \mathbb{R}^{100}$$

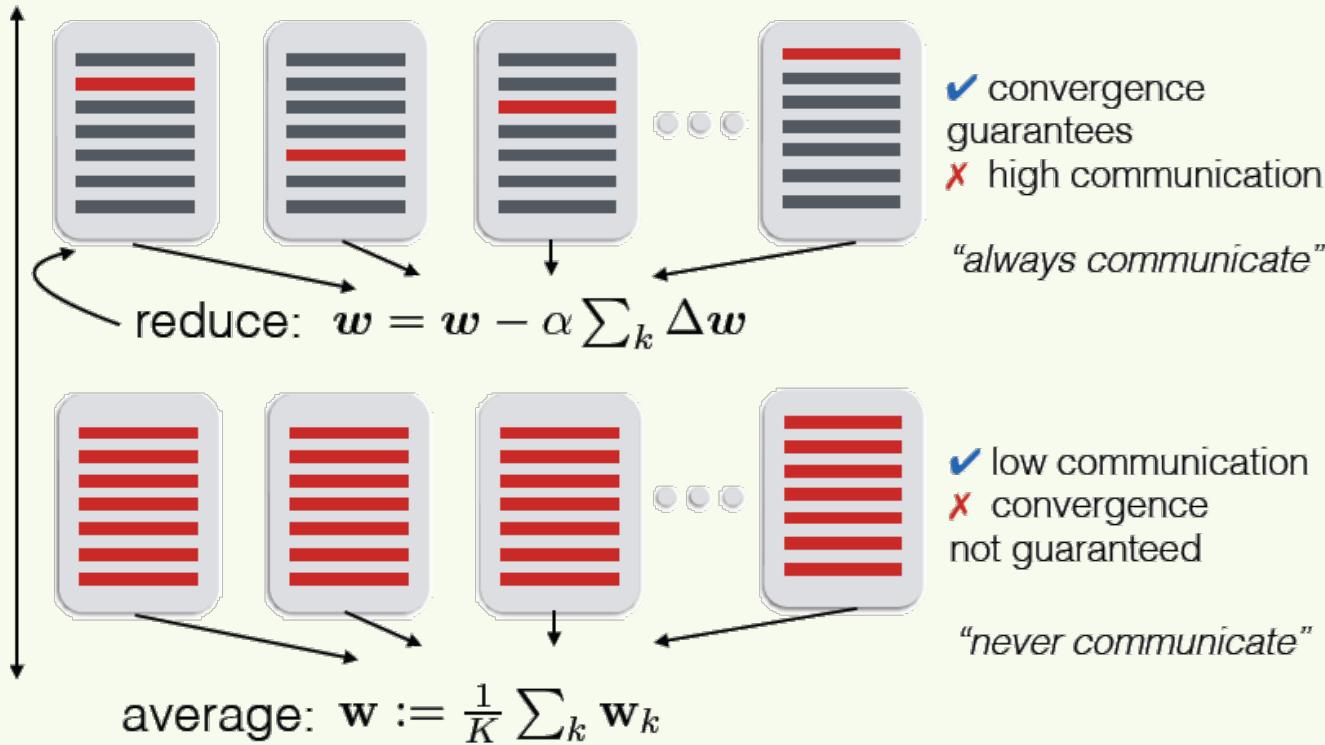
- ⊕ Sending  $v$  to another Machine

$500'000\text{ ns}$

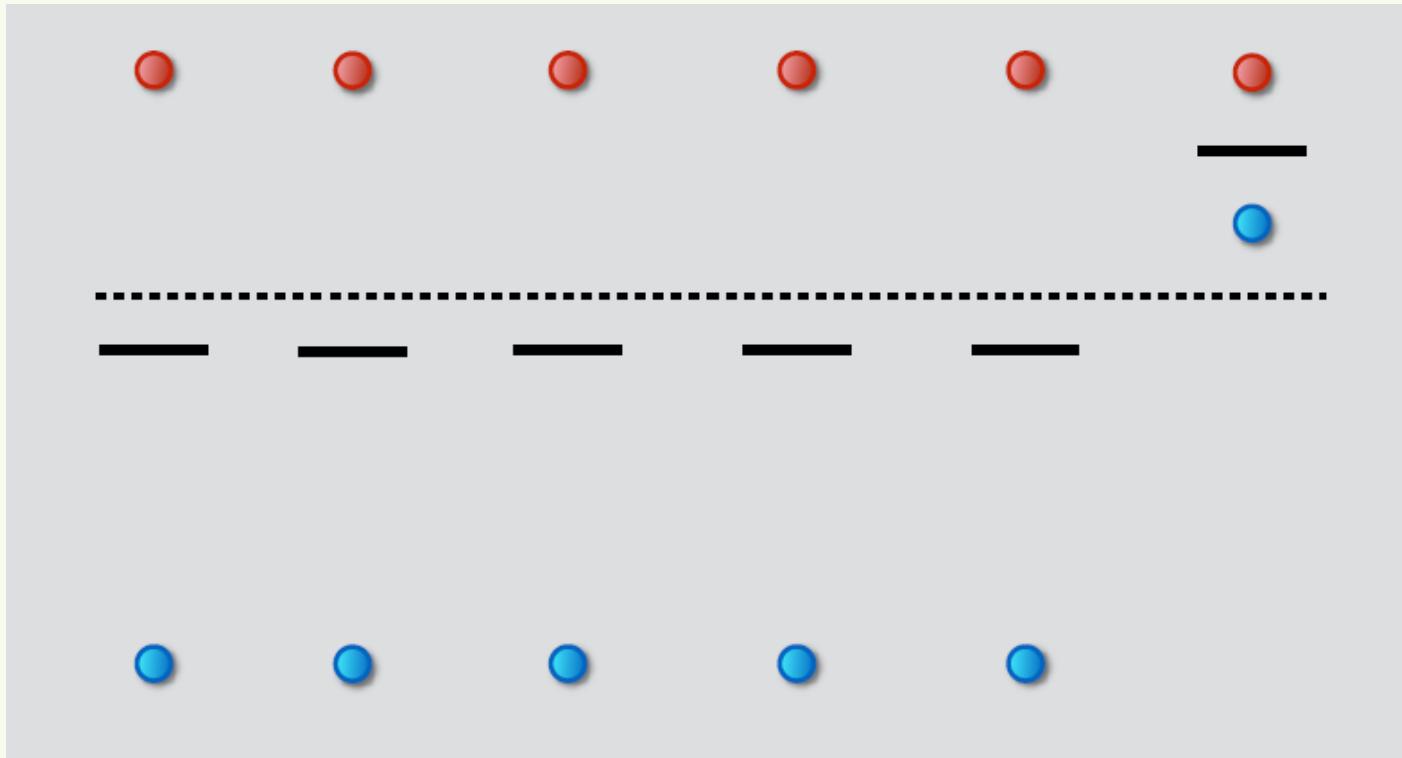
- ⊕ One Typical Map-Reduce Iteration (*Hadoop*)

$10'000'000'000\text{ ns}$

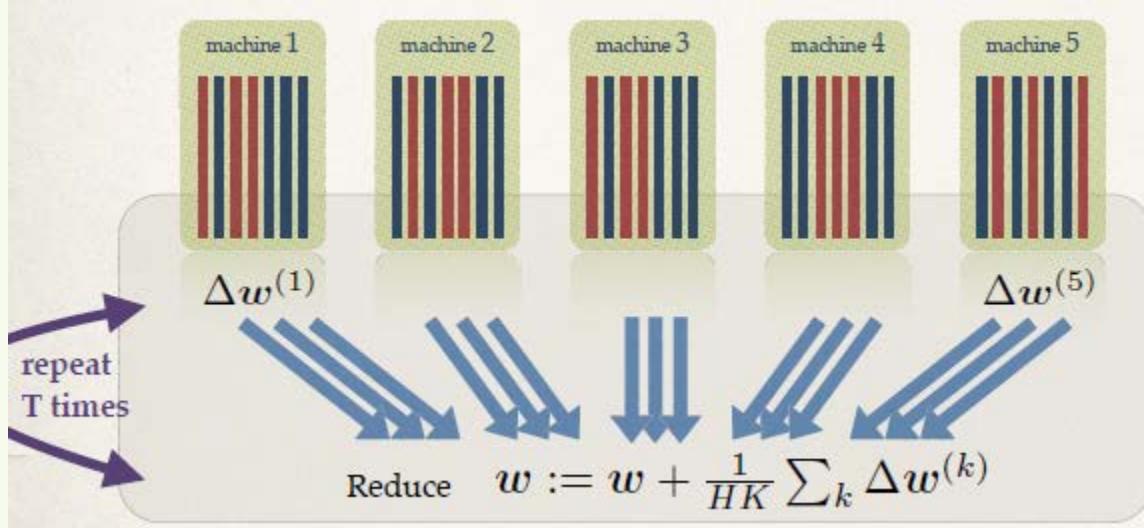
# Distributed Optimization



# One-shot strategy



# Something in the middle?



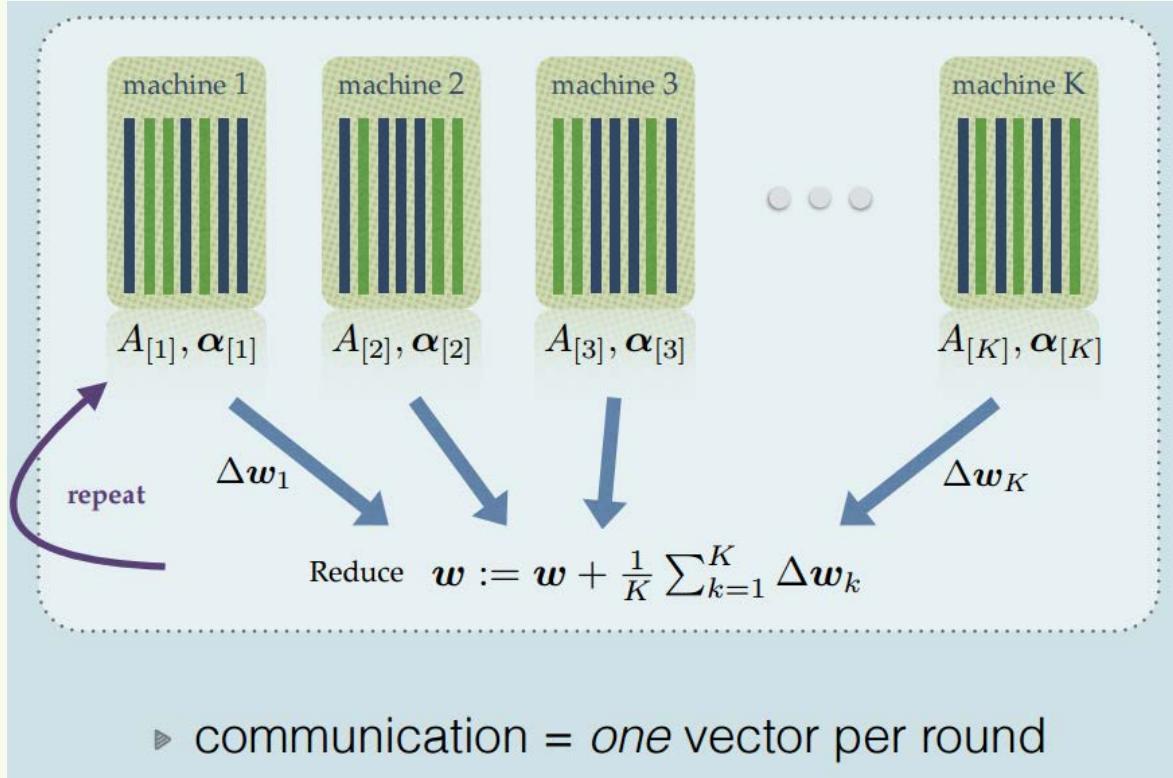
Mini-Batch SGD / CD

# local datapoints read: TH

# communications: T

convergence: ✓

# Something in the middle?



## setup

Primal problem formulation

$$\min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|^2 + \frac{1}{n} \sum_{i=1}^n \ell_i(\mathbf{w}^T \mathbf{x}_i)$$

data partitioned by examples  $A_i := \frac{1}{\lambda n} \mathbf{x}_i$

primal-dual correspondence  
 $w = A\alpha$

Dual problem

$$\max_{\alpha \in \mathbb{R}^n} -\frac{\lambda}{2} \|A\alpha\|^2 - \frac{1}{n} \sum_{i=1}^n \ell_i^*(-\alpha_i)$$

$$A_{\text{loc}} \alpha'_{\text{loc}} + w$$

$$\ell_i^*(s) := \sup_{t \in \mathbb{R}} \{st - \ell_i(t)\}$$

Information: local shared

Source: <http://m8j.net/data/List/Files-172/cocoa-NIPS-poster-final.pdf>

# Communication Efficient Distributed Dual Coordinate Ascent

---

## Algorithm 1: CoCoA

---

**Input:**  $T \geq 1$ , scaling parameter  $1 \leq \beta_K \leq K$  (default:  $\beta_K := 1$ ).

**Data:**  $\{(x_i, y_i)\}_{i=1}^n$  distributed over  $K$  machines

**Initialize:**  $\alpha_{[k]}^{(0)} \leftarrow 0$  for all machines  $k$ , and  $w^{(0)} \leftarrow 0$

**for**  $t = 1, 2, \dots, T$

**for all machines**  $k = 1, 2, \dots, K$  **in parallel**

$(\Delta\alpha_{[k]}, \Delta w_k) \leftarrow \text{LOCALDUALMETHOD}(\alpha_{[k]}^{(t-1)}, w^{(t-1)})$

$\alpha_{[k]}^{(t)} \leftarrow \alpha_{[k]}^{(t-1)} + \frac{\beta_K}{K} \Delta\alpha_{[k]}$

**end**

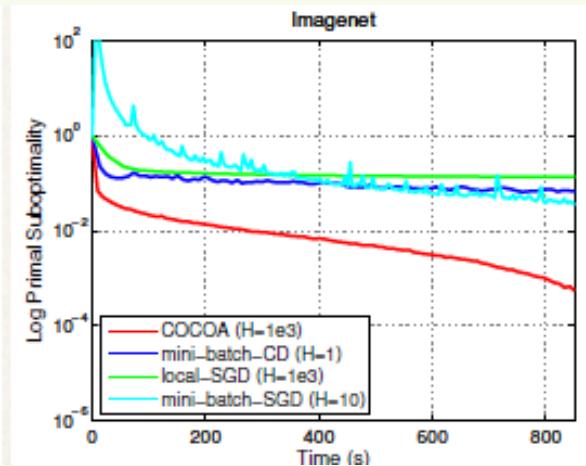
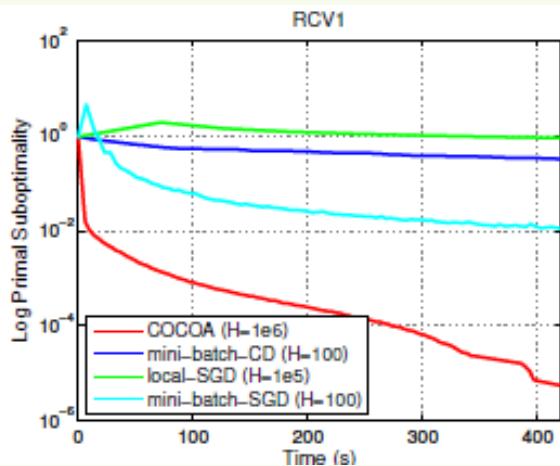
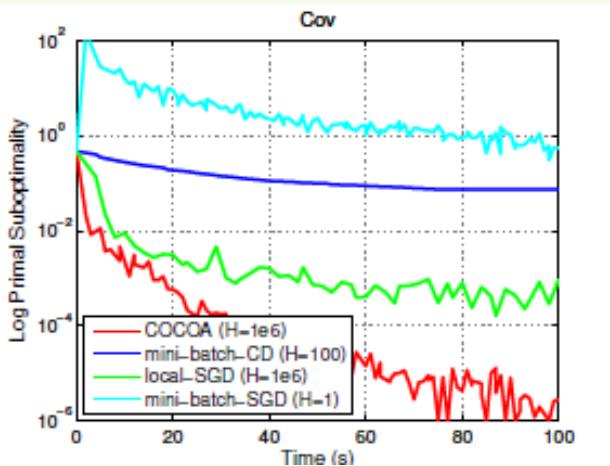
**reduce**  $w^{(t)} \leftarrow w^{(t-1)} + \frac{\beta_K}{K} \sum_{k=1}^K \Delta w_k$

**end**

---

# Benchmarks

Dataset	Training $n$	Features $d$	Sparsity	$\lambda$	Workers $K$
cov	522,911	54	22.22%	$1e-6$	4
rcv1	677,399	47,236	0.16%	$1e-6$	8
imagenet	32,751	160,000	100%	$1e-5$	32





**Part II:**  
**Real World  
Application**

# Big Data Frameworks

## Applications (Verticals)

<b>Personal Devices</b>	<b>Lifestyle</b>	<b>Connected Home</b>	<b>Industries</b>	<b>Industrial Internet</b>
Wearable Computing Fitness Health Family	Sports Cooking Pet Toys Music/Art Garden	Automotive Monitoring Security Tracker Hub	Retail Payment Loyalty Healthcare Automotive Infrastructure Agriculture Asset Tracking	Robotics Drones/Aerospace Green-tech 3D Scan/Print Smart Grid Corporates
pebble cookoo i'm @ striiiv APX MOTA GARMIN amigoo iFIT! JAWBONE MYFIT VALENCE BASIS fitbit TOMTOM LifeBEAM LUMO HAPIfork QUANTUM Lively SPIRE Withings QuantumLively HelloAliveCar Health remee mimo FILIP Sprouting ovuline Antenova greatcall Securicam mimo	SportS Smart Chef Scale THE ORANGE CHEF CO. Pantry Whistle PetPace pintofeed PetHub Petcube PetNet PetNet PetNet KAROTZ MAKIES ROLI CATCH GoPro plantlink Greenbox Koubachi	Quirky Radiator Labs netatmo LEVITON SmartThings Ubi nest LIFX X Geckoboard smarthome ULTRON INSTEON LIGHTING somfy lapka sense SUPERMECHANICAL tado <sup>®</sup> canary butterflye ring Lockitron RAVEN Kwikset GOJI scout smartalarm <sup>®</sup> Chipolo Linqueta locca! TrackR Homey Control revolv NINJABLOCKS NEXIA muzzley zonoff waze OpenXC	ignisight euclid instabot Boni payables ecobee Advancedcommerce vivint SAVANT Vera INSTEON Chamberlain PROTECH LIGHTING somfy birdi leeo ambient HomeMonitor dropcam august SCHLAEGER Genie KEY OP KEY lockitron Klikit Kwikset GOJI scout smartalarm <sup>®</sup> Chipolo Linqueta locca! TrackR Homey Control revolv NINJABLOCKS NEXIA muzzley zonoff waze OpenXC	Double Robotics ALDEBARAN iRobot ABB KUKA Double Robotics ROBOTIX EMPIRE Bluebeam Parrot Skybotix 3DR SKYRANGER SPIDEY DJI V-Smile STANLEY VITALITY MedMinder MediTracing AchieveTech GESUS CENTRAK redbeam Solera Zobe flic INRIX novdyDELPHI Airbridge dash waze OpenXC waveLink Beam KISI Johnson Controls Robin Schneider Electric eversys WIND RIVER MOTOROLA belkin DELL BOSCH NATIONAL INSTRUMENTS ARM LogMeIn Microsoft Honeywell PAN AMERICA SONY Atmel SIEMENS Qualcomm CISCO TOSHIBA SHARP

## Platforms & Enablement (Horizontals)

Connectivity/Dev Platforms	Software/Data Platforms	Open Source Platforms	Personal Interfaces	Security	Corporates
spark kynetx pinoccio ioBridge evthing! Ayila Networks EUPOWER resin.io Sympli TEBBEL blucity	iControl thingsquare carriots Keen IO SeeControl hings ConnectIO NewAer BERGO Axeda Yaler.net RacoWireless SpeedCart IFTTT greenwirc wot.io CyberSinge alliux Yo ThingWorx DN2K Xprespower iOTek thinkRF CANDI+ iboone bugsworm TempoIQ evercam.io covisint Jasper GigaThreadz ETHERIOS PubNub Bluebeam SensorCloud	webinos openHAB AllJoyn openINTERCONNECT KAR ThingSpeak GRID2HOME	NeuroSky Rayline LEAP gestigion apical ThingSpeak	inside SafeNet utimaco escrypt gemalto BASTILLE NETWORKS MOCANA	amazon hp LG intel HTC PHILIPS IBM SUN MICROSYSTEMS WIND RIVER MOTOROLA belkin DELL BOSCH NATIONAL INSTRUMENTS ARM LogMeIn Microsoft Honeywell PAN AMERICA SONY Atmel SIEMENS Qualcomm CISCO TOSHIBA SHARP

## Building Blocks

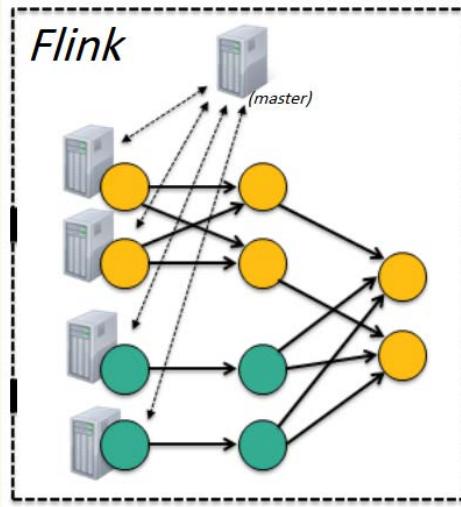
Protocols	M2M Networks	Portable WiFi	Telecom	M2M
Bluetooth WiFi 2G 3G 4G LTE CoAP 6LoWPAN LWM2M BTXML	Helium SIGFOX KORE stream aeris MACH36 MEMO	Open Garden GOOD SPEED BRICK L karma	at&t boostmobile verizon T-Mobile U.S. Cellular VimpelCom Vodafone airtel	Qualcomm FICOM Laird WEBOID QuSense Wireless seed giga
Cloud	Processors Sensors Parts/Kits Services	Services	Incubators	Funding
Google Cloud Platform Amazon Web Services redhat OMA iWAVE enModus HART MiWi M-Bus	iOS Windows Phone BlackBerry Processors Sensors TI TinkerForge	Makey Makey SAM littleBits TIE CIRCUIT Sculpteo	Highway1 LEARNERS LABS shapeways droon makeXYZ	KICKSTARTER indiegogo MedStarLabs
Microsoft Azure		3D HUBS	WEARABLE WORLD R/GA Accelerator TechShop	angelcam

© Matt Turck (@mattturck), David Rogg (@davidjrogg) & FirstMark Capital (@firstmarkcap) FIRSTMARK

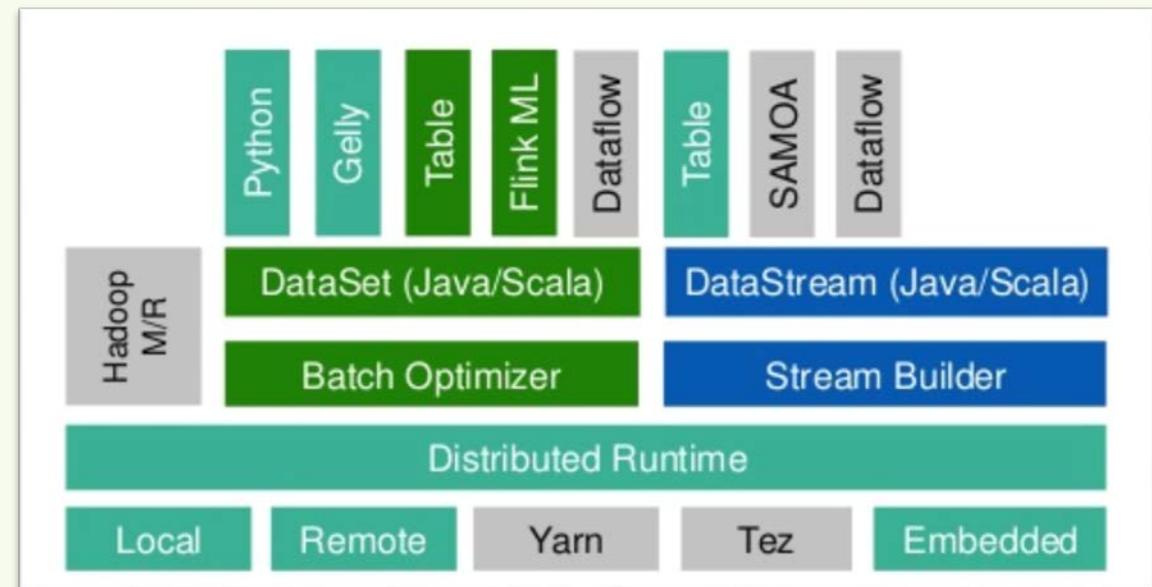
(2014)

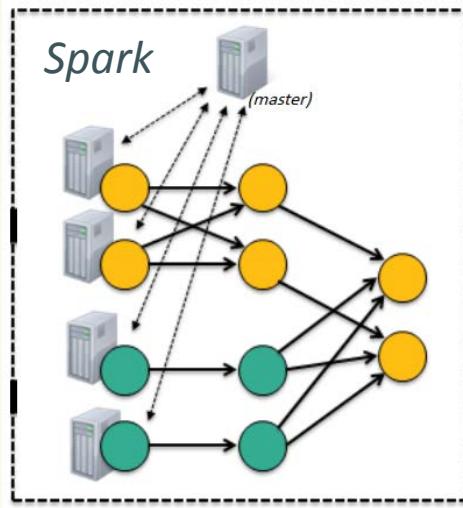


# Apache FLINK



- Allows to execute same code locally and on a compute cluster (distributed computing)
- Automatically distributes data over cluster nodes and allows real-time queries

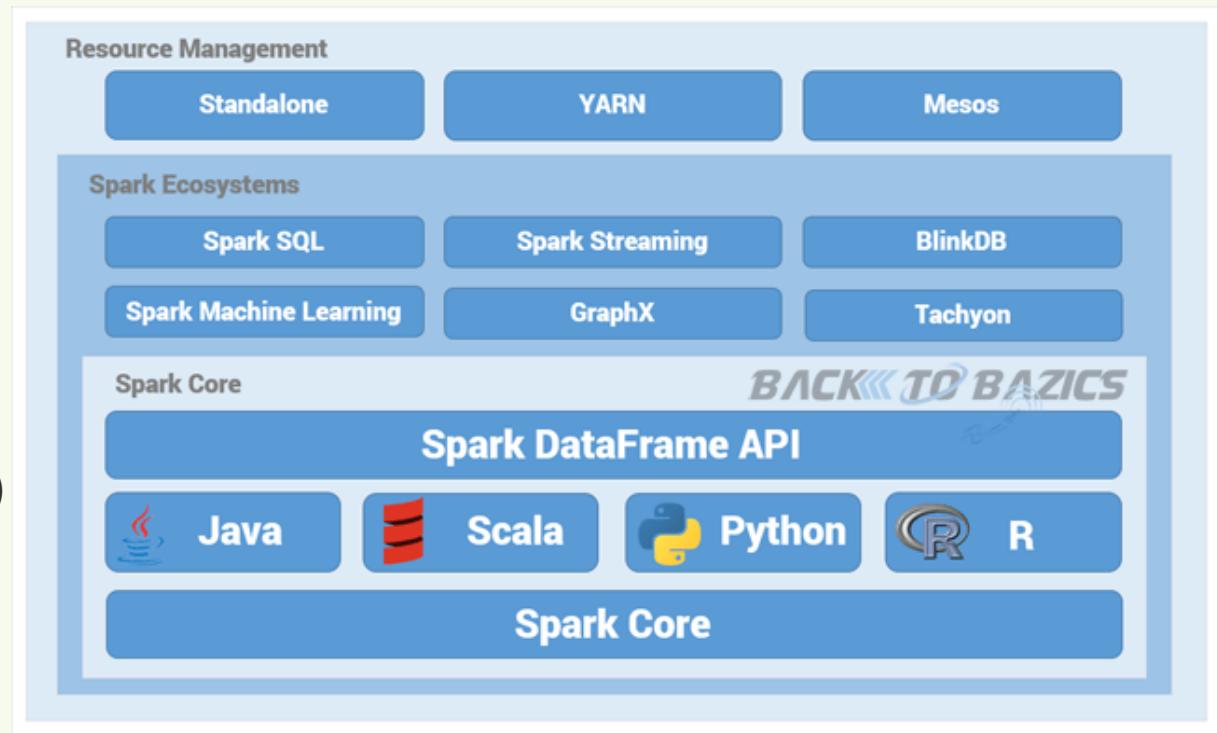




# Apache SPARK



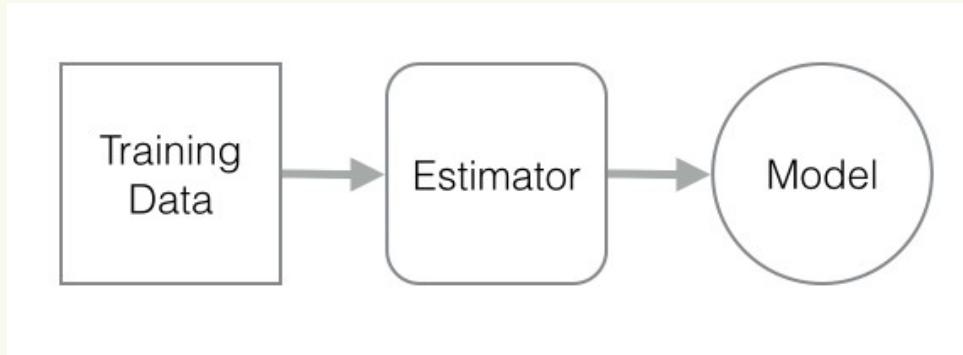
- Allows to execute same code locally and on a compute cluster (distributed computing)
- Automatically distributes data over cluster nodes and allows real-time queries



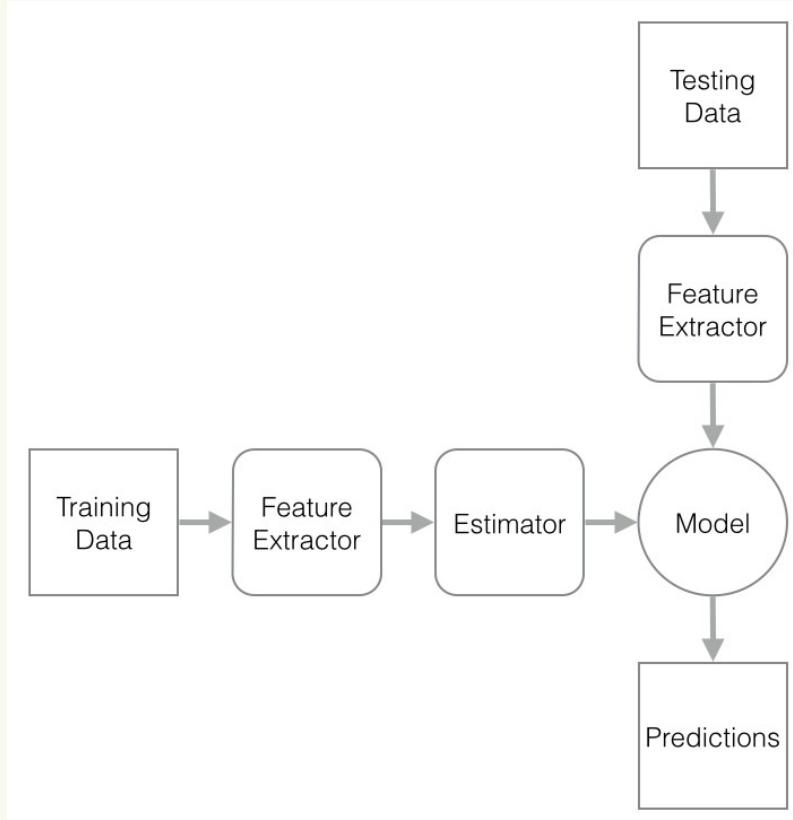


Machine learning using a framework

# Does ML work like that?



# More realistic scenario!



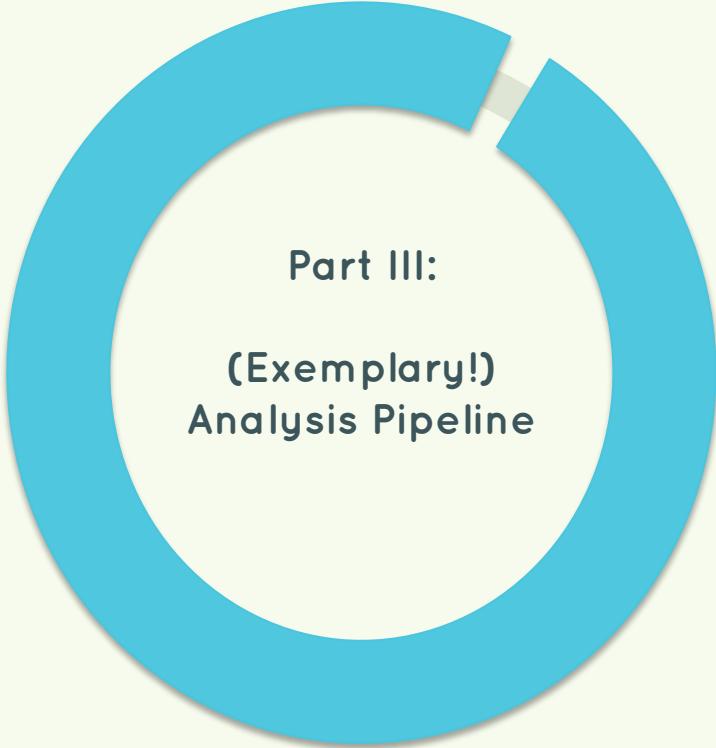
```
// Read the training data set, from a LibSVM formatted file
val trainingDS: DataSet[LabeledVector] = env.readLibSVM(pathToTrainingFile)

// Create the SVM Learner
val svm = SVM()
.setBlocks(10)
.setIterations(10)
.setLocalIterations(10)
.setRegularization(0.5)
.setStepsize(0.5)

// Learn the SVM model
svm.fit(trainingDS)

// Read the testing data set
val testingDS: DataSet[Vector] = env.readVectorFile(pathToTestingFile)

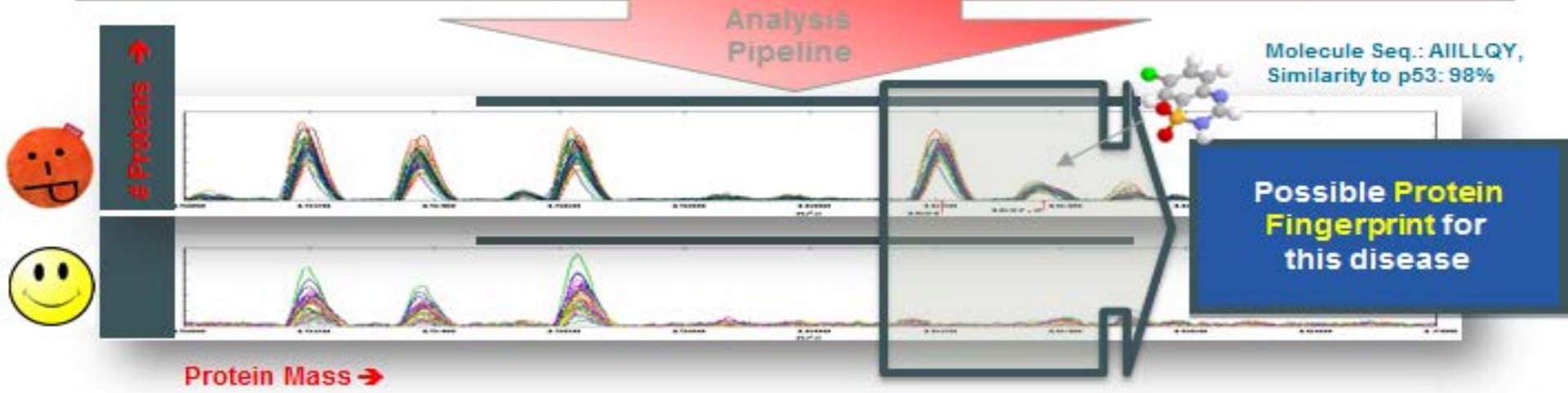
// Calculate the predictions for the testing data set
val predictionDS: DataSet[LabeledVector] = svm.predict(testingDS)
```



**Part III:**  
**(Exemplary!)**  
**Analysis Pipeline**

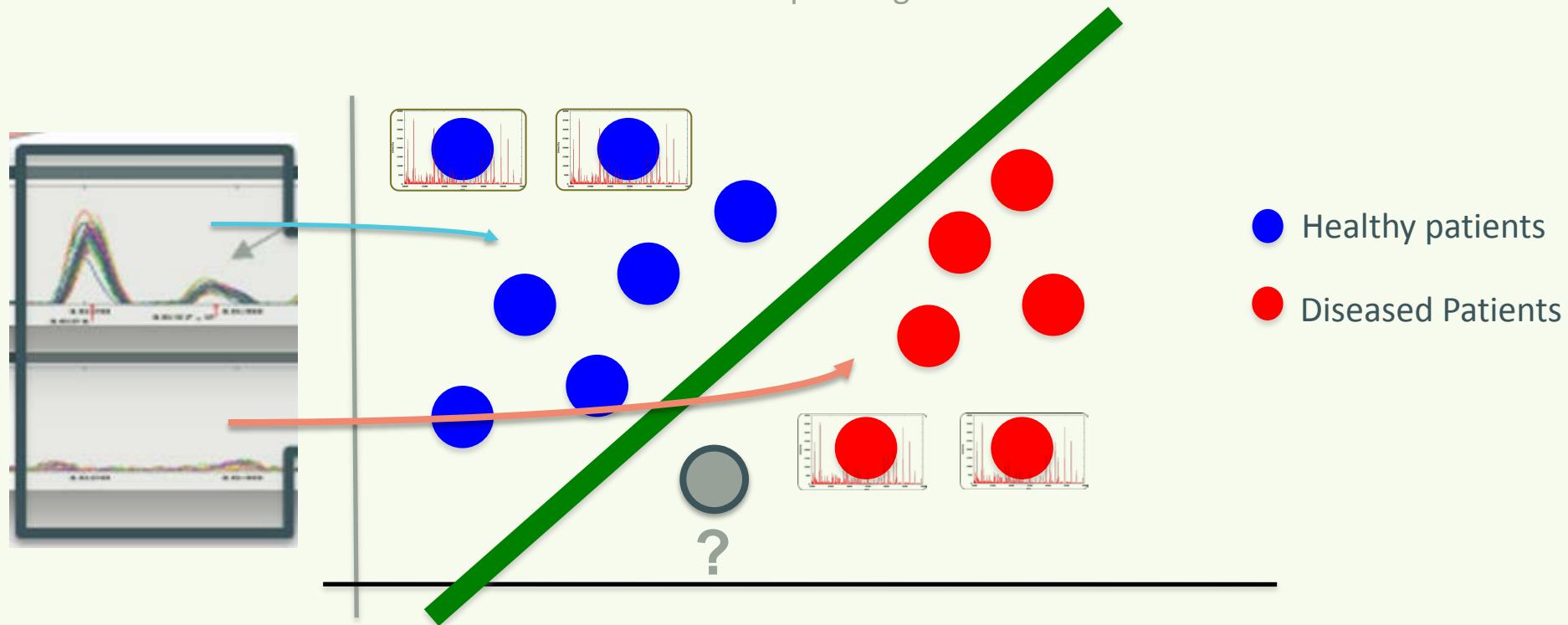
# GOAL



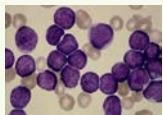
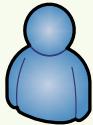


# Classification (e.g. SVM)

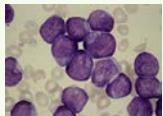
Determine separating line



# Biomarkers Discovery Workflow



Disease



Normal



GENE-DB

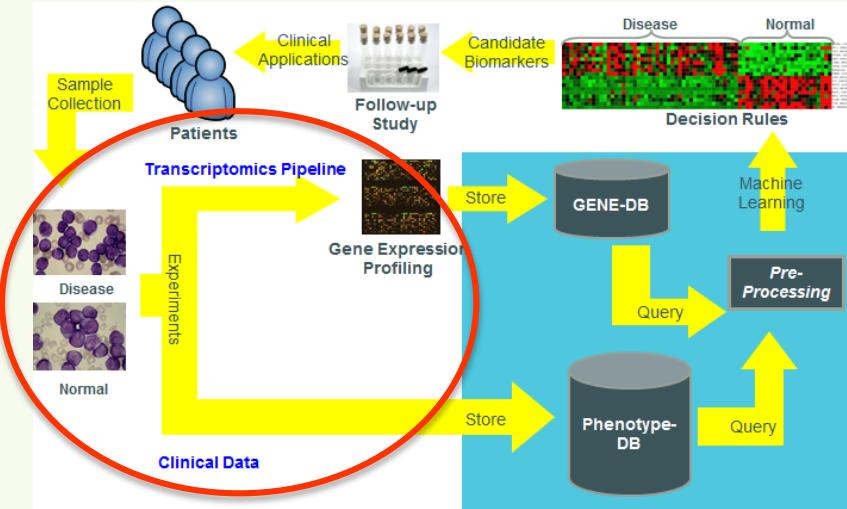


Phenotype-  
DB

# Workflow „Parts“

## Data

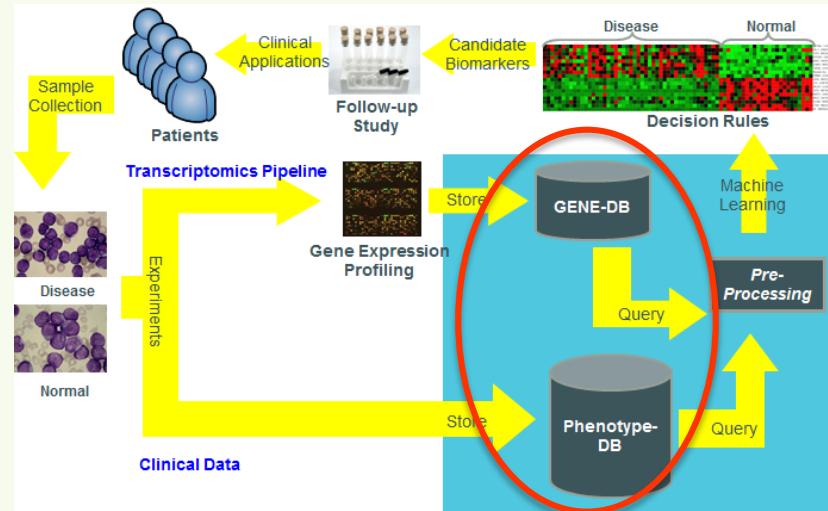
- The goal is to find differences between patients with different “properties”
- Thus: we need data from patients having these properties, e.g.
  - a „healthy“ group and a „diseased“ groups or
  - a group with different survival estimations



# Workflow „Parts“

## Storage

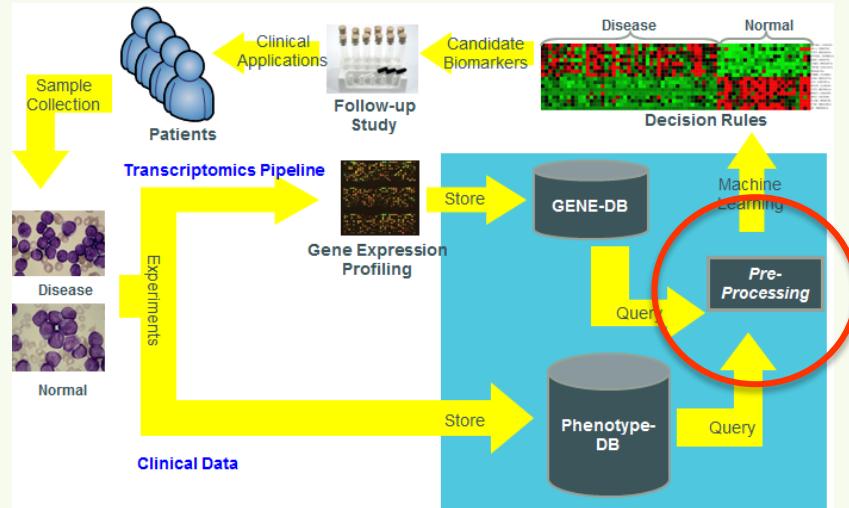
- The data is large, i.e. tens of Gigabytes per patient
- A storage system is needed that allows for fast access
- Hadoop's DFS (HDFS) distributes data on available machines



# Workflow „Parts“

## Pre-Processing

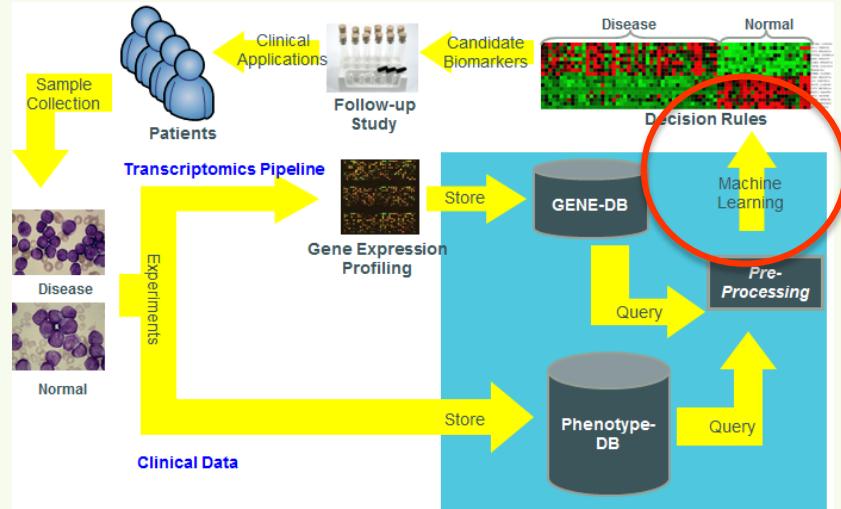
- Depending on the data-type, pre-processing might be needed
- E.g. noise reduction, filtering of „bad“ data, data-imputation, etc.



# Workflow „Parts“

## Machine Learning

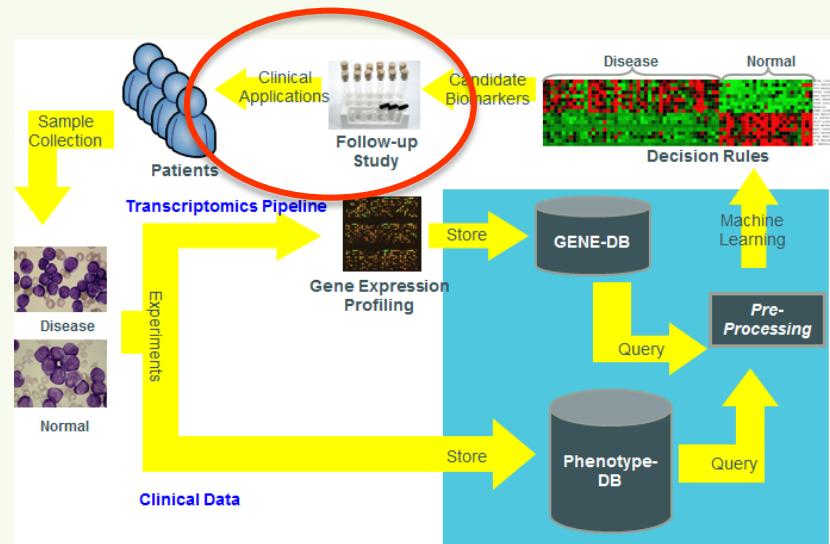
- Build model for
  - classification or
  - Regression
- Also interesting:
  - which features are important?
  - Causality?
  - Model?
  - Prediction?
  - ...



# Workflow „Parts“

## Validation

- Results should always be validated,  
e.g. in an independent lab



# After this talk you now should know...

- ... what -omics data is
- ... that -omics data can be very large
- ... why it is useful to analyze -omics data
- ... a method how to analyze large -omics data
- ... a framework implementing this method

*Spoiler: we will do this during the hands-on session*

Contact – Just write us!

Thank you for your attention!

BYE!

Questions?

*Tim Conrad*

Freie Universität Berlin  
conrad@math.fu-berlin.de

*Link*

[www.forschungscampus-modal.de](http://www.forschungscampus-modal.de)