# Machine learning for patient stratification from genomic information
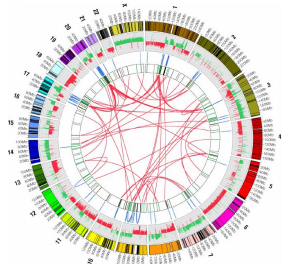
Jean-Philippe Vert

`jean-philippe.vert@ens.fr`
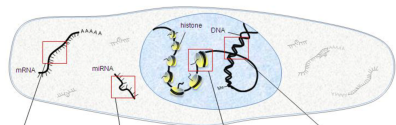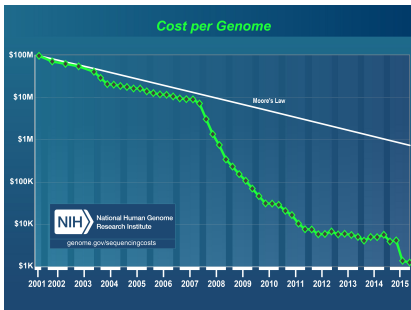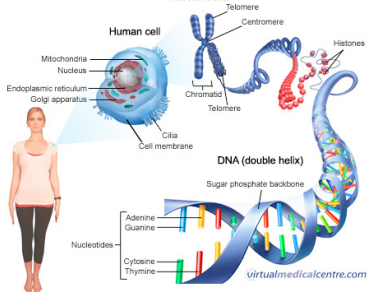
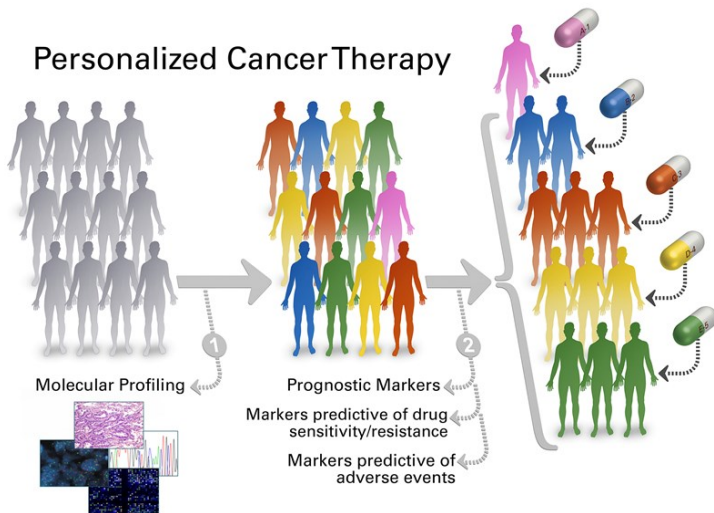Paris Sciences and Data, December 1, 2016

# Molecular data

Personalized Cancer Therapy

https://pct.mdanderson.org

# Learning from data: supervised classification/regression

- Patients with VS without relapse in 5 years
- Case where *n* (=19) patients $>>$ *p* (=2) markers

# Learning from data: supervised classification/regression

- Patients with VS without relapse in 5 years
- Case where $n$ (=19) patients $>> p$ (=2) markers

# Learning from data: supervised classification/regression

- Patients with VS without relapse in 5 years
- Case where $n$ (=19) patients $>> p$ (=2) markers

# Learning from data: supervised classification/regression

- Patients with VS without relapse in 5 years
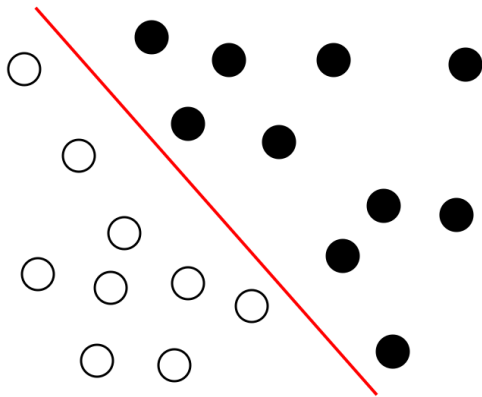- Case where $n$ (=19) patients $>> p$ (=2) markers

# Real data: $n \ll p$

- Gene expression



- Somatic mutations



- $n = 10^2 \sim 10^4$ (patients)
- $p = 10^4 \sim 10^7$ (genes, mutations, copy number, ...)
- Data of various nature (continuous, discrete, structured, ...)
- Data of variable quality (technical/batch variations, noise, ...)

## Consequence: limited accuracy

Breast cancer prognosis competition, $n = 2000$, Bilal et al (2013)



- C: 16 standard clinical data (age, tumor size, ...)
- M: 80k molecular features (gene expression, DNA copy number)
- P: incorporate prior knowledge

# Example: survival prediction from somatic mutations



- Data from TCGA (3.3k samples, 8 cancer types, >10k genes)
- Survival SVM on raw binary data, or processed by NSQN (Hofree et al., 2013) or NetNorm (Le Morvan et al., 2016).

# Consequence: unstable biomarker selection



**Gene expression profiling predicts clinical outcome of breast cancer**

Laura J. van 't Veer*†, Hongyue Dai‡‡, Marc J. van de Vijver*†, Yudong D. He‡, Augustinus A. M. Hart*, Mao Mao‡, Hans L. Peterse*, Karin van der Kooy*, Matthew J. Marton‡, Anke T. Witteveen*, George J. Schreiber‡, Ron M. Kerkhoven*, Chris Roberts‡, Peter S. Linsley‡, René Bernards* & Stephen H. Friend‡

*Divisions of Diagnostic Oncology, Radiotherapy and Molecular Carcinogenesis and Center for Biomedical Genetics, The Netherlands Cancer Institute, 121 Plesmanlaan, 1066 CX Amsterdam, The Netherlands
‡Rosetta Inpharmatics, 12040 115th Avenue NE, Kirkland, Washington 98034.

**Gene-expression profiles to predict distant metastasis of lymph-node-negative primary breast cancer**
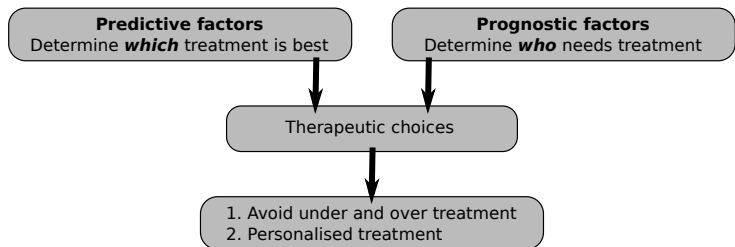
Yixin Wang, Jan G M Klijn, Yi Zhang, Anieta M Sieuwerts, Maxime P Look, Fei Yang, Dmitri Talantov, Mieke Timmermans, Marion E Meijer-van Gelder, Jack Yu, Tim Jatkoe, Els M J J Berns, David Atkins, John A Foekens

70 genes (Nature, 2002)                76 genes (Lancet, 2005)

<span style="color:red">3 genes in common</span>

van 't Veer et al. (2002); Wang et al. (2005)

# Some research directions

- Find a better representation



**One sample x**
**p features**

**Mapping f(x)**
**p(p-1)/2 bits**

- Incorporate prior knowledge

# From prognostic to predictive models
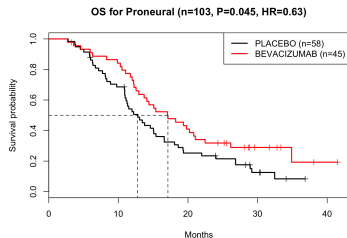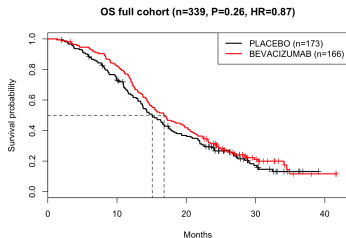


- Prognostic:
  - Predict outcome $Y$ of a disease on an untreated individual $X$
  - Standard supervised learning: model $Y = f(X)$ from observations of $(X_i, Y_i)$ pairs
- Predictive:
  - Predict the benefit in outcome $Y$ of a treatment $A$ on an individual $X$
  - We observe $(X_i, A_i, Y_i)$ but want to model $Y = f(X, A_1) - f(X, A_2)$
  - For each $X$ we only observe the outcome $Y$ under one treatment $A$ (cf e-marketing)

# Clinical trials for precision medicine?

1. Meta-analysis of clinical trials (typically *A/B testing*) to estimate predictive models



2. Dynamic trial to jointly optimize the predictive model and its performance (contextual multi-armed bandit problem)

- Lots of data
- $n << p$ is the rule, more and more...
- Limited impact so far for patients
- Active research
  - new representations $x \rightarrow \Phi(x)$
  - new learning techniques (structured sparsity, regularization, ...)
  - new experimental design strategies (contextual bandit)

# References

M. Hofree, J. P. Shen, H. Carter, A. Gross, and T. Ideker. Network-based stratification of tumor mutations. *Nat Methods*, 10(11):1108–1115, Nov 2013. doi: 10.1038/nmeth.2651. URL http://dx.doi.org/10.1038/nmeth.2651.

L. J. van 't Veer, H. Dai, M. J. van de Vijver, Y. D. He, A. A. M. Hart, M. Mao, H. L. Peterse, K. van der Kooy, M. J. Marton, A. T. Witteveen, G. J. Schreiber, R. M. Kerkhoven, C. Roberts, P. S. Linsley, R. Bernards, and S. H. Friend. Gene expression profiling predicts clinical outcome of breast cancers. *Nature*, 415(6871):530–536, Jan 2002. doi: 10.1038/415530a. URL http://dx.doi.org/10.1038/415530a.

Y. Wang, J. Klijn, Y. Zhang, A. Sieuwerts, M. Look, F. Yang, D. Talantov, M. Timmermans, M. Meijer-van Gelder, J. Yu, T. Jatkoe, E. Berns, D. Atkins, and J. Foekens. Gene-expression profiles to predict distant metastasis of lymph-node-negative primary breast cancers. *Lancet*, 365(9460):671–679, 2005. doi: 10.1016/S0140-6736(05)17947-1. URL http://dx.doi.org/10.1016/S0140-6736(05)17947-1.