# Machine Learning for Toxicogenetics and Drug Response Prediction

Jean-Philippe Vert
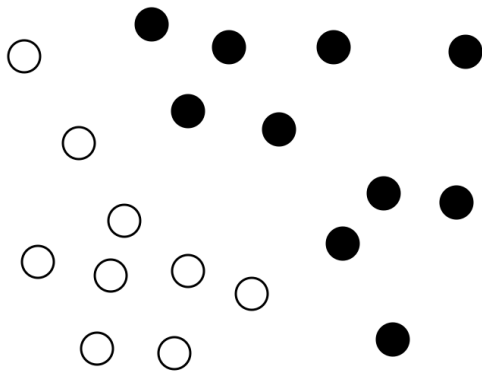
MINES ParisTech

institutCurie
Together, let's beat cancer.

UNIVERSITY OF CALIFORNIA BERKELEY · 1868 ·

Festival of Genomics, San Mateo, Nov 6, 2015

# Molecular stratification



Patients with same condition

DNA Profiling

Good responders

No Responders

Bad side effects

Diagnosis, prognosis, drug response prediction, ...

$n(=19) >> p(=2)$ : easy

# Machine learning formulation



$n(= 19) >> p(= 2)$ : easy

$n(= 19) >> p(= 2)$ : easy

$n(= 19) >> p(= 2)$ : easy

- $n = 10^2 \sim 10^4$ (patients)
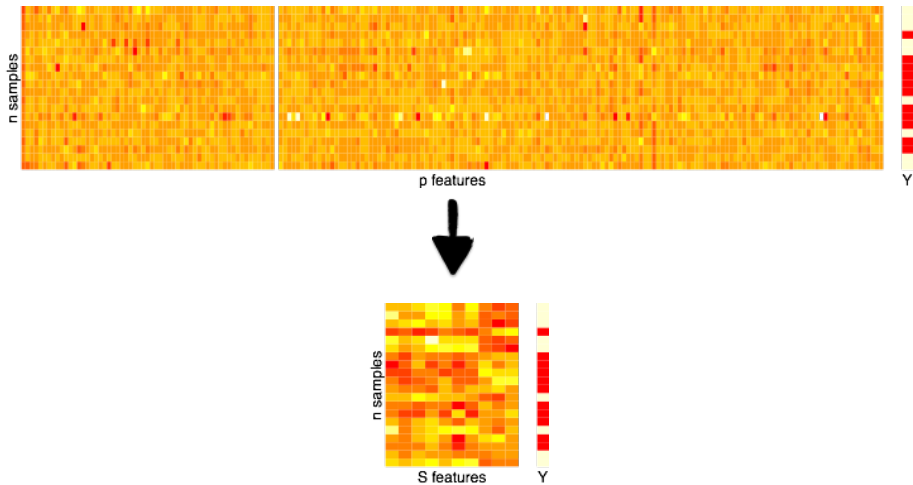- $p = 10^4 \sim 10^7$ (genes, mutations, copy number, ...)

Accuracy drops,

# Outline

# Outline

# Feature selection (a.k.a. *molecular signature*)

# Example: Breast cancer prognostic signature

# But...

**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)

## 3 genes in common

# 3 genes is the best you can expect given *n* and *p*
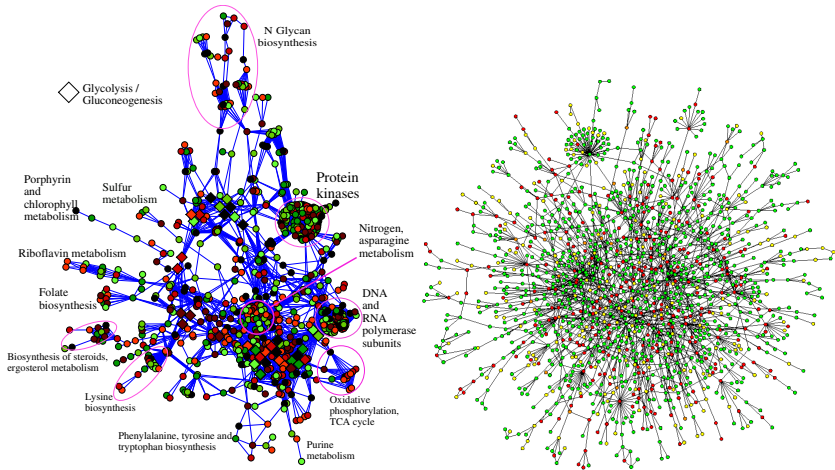
## The Influence of Feature Selection Methods on Accuracy, Stability and Interpretability of Molecular Signatures

Anne-Claire Haury[1,2,3,*], Pierre Gestraud[1,2,3], Jean-Philippe Vert[1,2,3]

1 Mines ParisTech, Centre for Computational Biology, Fontainebleau, France, 2 Institut Curie, Paris, France, 3 Institut National de la Santé et de la Recherche Médicale, Paris, France
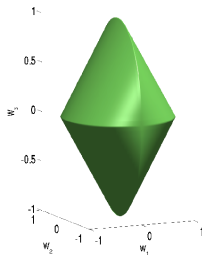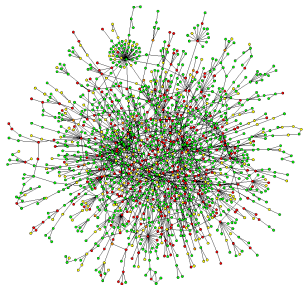
# Gene networks as prior knowledge



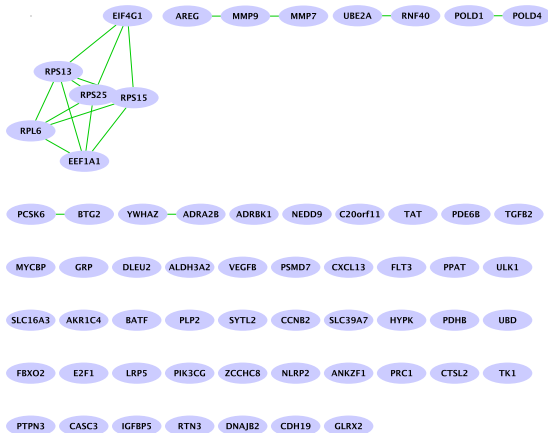Can we force the signatures to be "coherent" with a known gene network?

1. Using the network, define a non-smooth and convex subset of "candidate" signatures compatible with it



$$\Omega(\beta) = \sup_{\alpha \in \mathbb{R}^p : \forall i \sim j, \|\alpha_i^2 + \alpha_j^2\| \leq 1} \alpha^\top \beta.$$
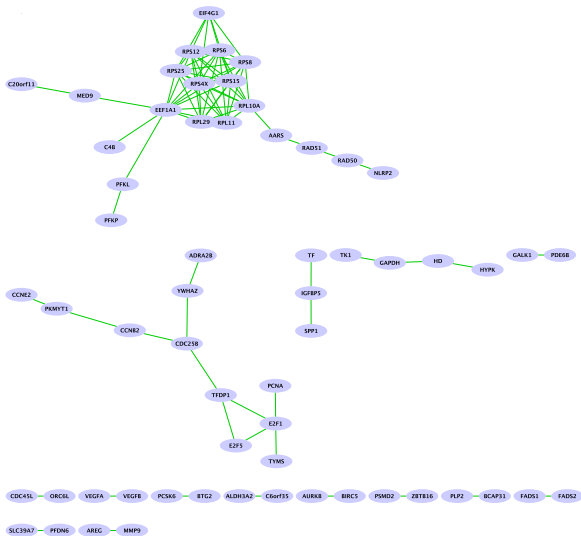
2. Among the candidates, find the best signature that explains the data (efficient optimization through convex programming)

# Lasso signature (accuracy 0.61)



*Breast cancer prognosis*
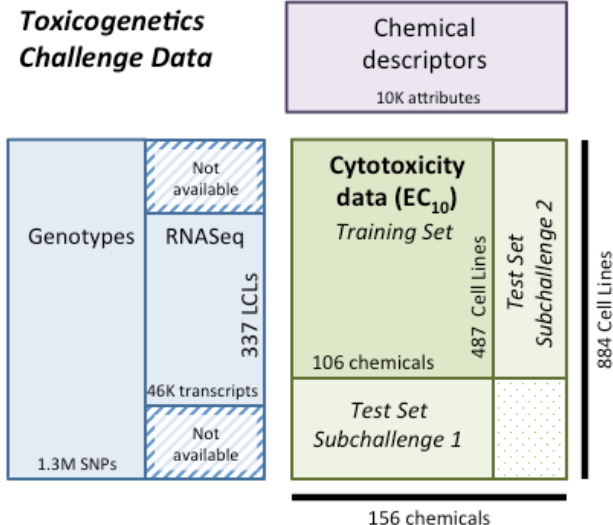
# Graph Lasso signature (accuracy 0.64)



*Breast cancer prognosis*

# Outline

Genotypes from the 1000 genome project
RNASeq from the Geuvadis project

Toxicogenetics Challenge Data

Chemical descriptors
10K attributes

Genotypes

RNASeq

Not available

337 LCLs

46K transcripts

1.3M SNPs

Not available

Cytotoxicity data (EC$_{10}$)
Training Set

106 chemicals

487 Cell Lines

Test Set Subchallenge 2

Test Set Subchallenge 1

884 Cell Lines

156 chemicals

n = 5E4

p = 1E10

# Crowd-sourcing: the DREAM8 Toxicogenetics challenge

## Bilinear regression

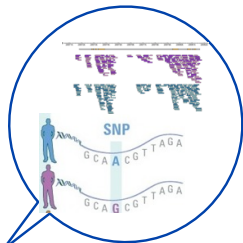- Cell line $X$, chemical $Y$, toxicity $Z$.
- Bilinear regression model:

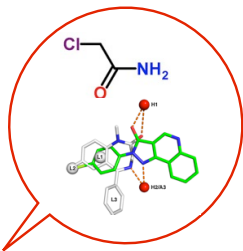$$Z = f(X, Y) + b(Y) + \epsilon,$$

- Estimation by kernel ridge regression:

$$\min_{f \in \mathcal{H}, b \in \mathbb{R}^m} \sum_{i=1}^{n} \sum_{j=1}^{m} \left( f(x_i, y_j) + b_j - z_{ij} \right)^2 + \lambda \|f\|^2,$$
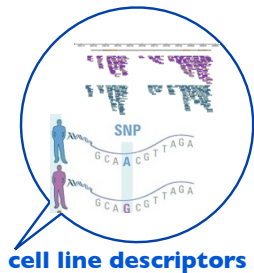
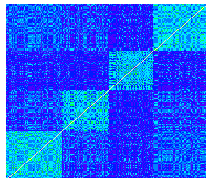- Solved in $O(max(n, p)^3)$

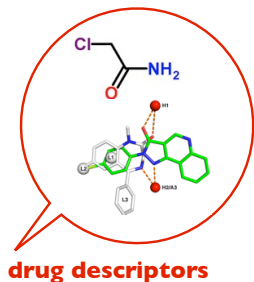**cell line descriptors**



**drug descriptors**

cell line descriptors

kernelized →

Kcell

drug descriptors

Kdrug

# Kernel Trick
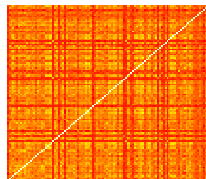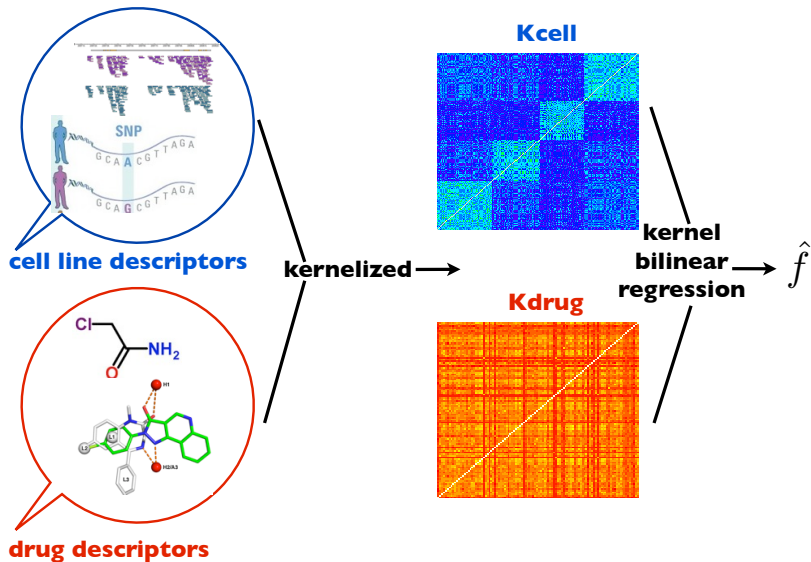
**Kcell**

**Kdrug**

cell line descriptors

drug descriptors

kernelized

**Kernel choice?**
. descriptors
. data integration
. missing data

kernel
bilinear
regression

$\hat{f}$

SNP

GCAACGTTAGA

GCAGCGTTAGA

Cl

NH₂

O

1. $K_{cell}$ :
   $\implies$ 29 cell line kernels tested
   $\implies$ 1 kernel that *integrate all information*
   $\implies$ deal with missing data

2. $K_{drug}$ :
   $\implies$ 48 drug kernels tested
   $\implies$ multi-task kernels

1. **$K_{cell}$** :
   - $\implies$ 29 cell line kernels tested
   - $\implies$ 1 kernel that *integrate all information*
   - $\implies$ deal with missing data
2. **$K_{drug}$** :
   - $\implies$ 48 drug kernels tested
   - $\implies$ multi-task kernels

**Covariates**
. linear kernel

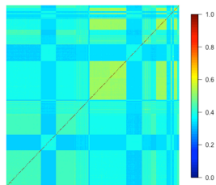**SNPs**
. 10 gaussian
kernels

**RNA-seq**
. 10 gaussian
kernels

**Covariates**
. linear kernel

**SNPs**
. 10 gaussian kernels

**RNA-seq**
. 10 gaussian kernels

**Integrated kernel**

1. **Dirac**
2. Multi-Task
3. Feature-based
4. Empirical
5. Integrated



independent regression for each drug

1. Dirac
2. **Multi-Task**
3. Feature-based
4. Empirical
5. Integrated



sharing information across drugs
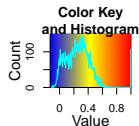
1. Dirac
2. Multi-Task
3. **Feature-based**
4. Empirical
5. Integrated

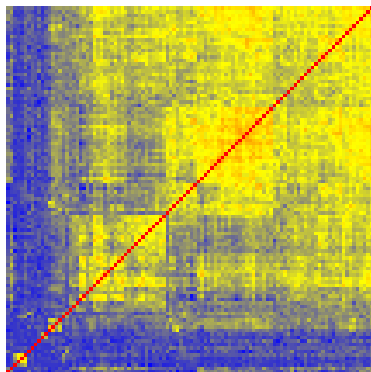Linear kernel and 10 gaussian kernels based on features:

- CDK (160 descriptors) and SIRMS (9272 descriptors)
- Graph kernel for molecules (2D walk kernel)
- Fingerprint of 2D substructures (881 descriptors)
- Ability to bind human proteins (1554 descriptors)

**Color Key
and Histogram**

Count

**Empirical correlation**

1. Dirac
2. Multi-Task
3. Feature-based
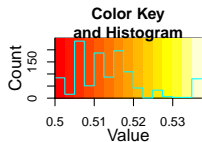4. **Empirical**
5. Integrated

# Multi-task drug kernels

1. Dirac
2. Multi-Task
3. Feature-based
4. Empirical
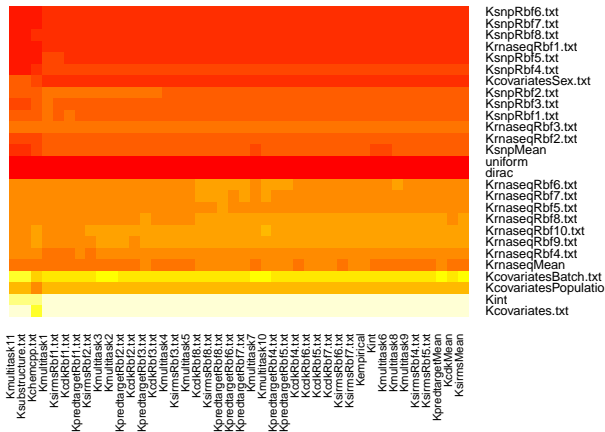5. **Integrated**

$$K_{int} = \sum_i K_i$$
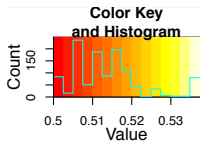
Integrated kernel:

- Combine all information on drugs
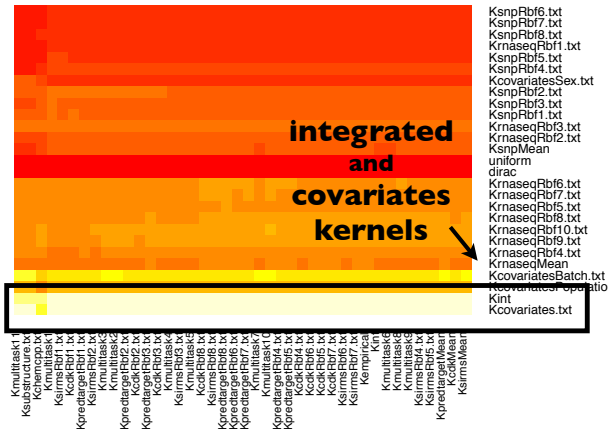
# 29x48 kernel combinations: CV results

**integrated kernel**

**Mean CI for cell line kernels**

Kint
Kcovariates.txt
KcovariatesBatch.txt
KcovariatesPopulation.txt
KrnaseqRbf10.txt
KrnaseqRbf9.txt
KrnaseqRbf8.txt
KrnaseqRbf6.txt
KrnaseqRbf7.txt
KrnaseqRbf5.txt
KrnaseqRbf4.txt
KrnaseqMean
KrnaseqRbf3.txt
KsnpRbf2.txt
KsnpRbf3.txt
KsnpRbf1.txt
KrnaseqRbf2.txt
KsnpMean
KsnpRbf4.txt
KcovariatesSex.txt
KrnaseqRbf1.txt
KsnpRbf5.txt
KsnpRbf8.txt
KsnpRbf6.txt
KsnpRbf7.txt
uniform
dirac

**batch effect**

0.50    0.51    0.52    0.53    0.54

Mean CI for chemicals kernels
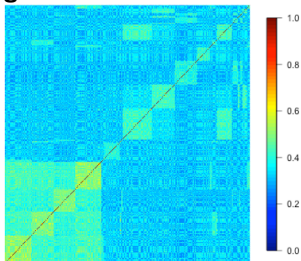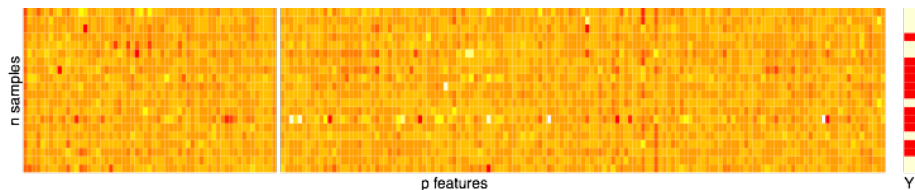
Mean CI for chemicals kernels

**Empirical kernel on drugs**

**Integrated kernel on cell lines**

- Small $n$ large $p$ $\implies$ regularized models with prior knowledge
- Heterogeneous data integration $\implies$ kernel methods
- Performance remains often disappointing!
- Progress arise by small steps

# Thanks

cbio@mines-paristech.fr , u900@curie.fr, sandrine@stat.berkeley.edu