



# Structured feature selection in high dimension for precision medicine

Chloé-Agathe Azencott

Center for Computational Biology (CBIO)  
Mines ParisTech – Institut Curie – INSERM U900  
PSL Research University, Paris, France

February 15, 2019 – PASADENA Workshop

<http://cazencott.info>

[chloe-agathe.azencott@mines-paristech.fr](mailto:chloe-agathe.azencott@mines-paristech.fr)

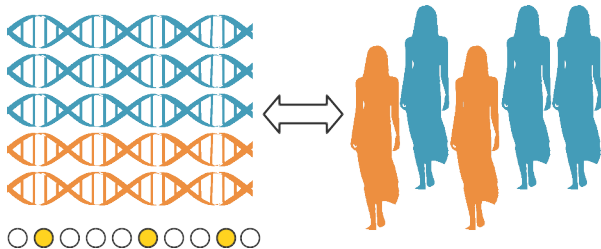
@cazencott

# Precision Medicine

- ▶ Treatment **adapted to the (genetic) features** of the patient.  
E.g. Trastuzumab for HER2+ breast cancer.
- ▶ Identify **similarities** between patients that exhibit **similar phenotypes**: susceptibilities, prognoses, responses to treatment.

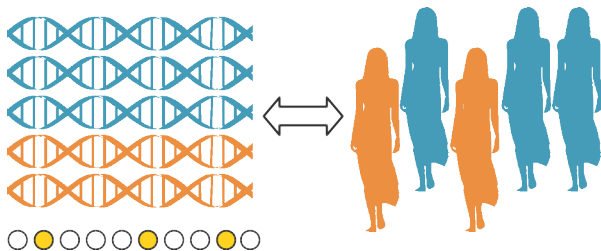


# From genotype to phenotype



**Which genomic features explain the phenotype?**

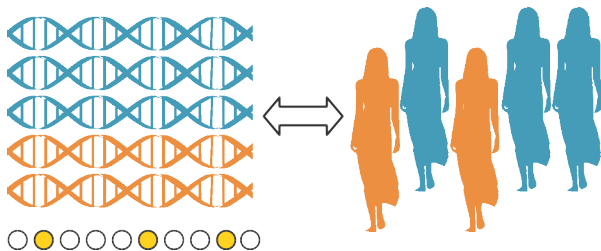
# From genotype to phenotype



## Which genomic features explain the phenotype?

- 80 000 proteins;
- 200 000 mRNA;
- 10 million SNPs;
- 28 million CpG islands.

# From genotype to phenotype



## Which genomic features explain the phenotype?

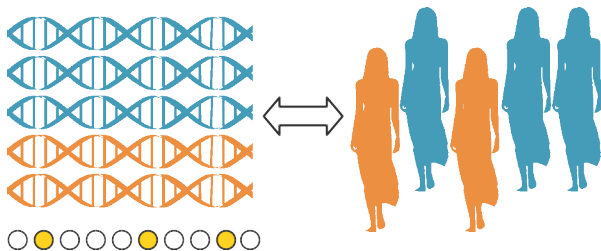
$p = 10^5 - 10^7$  **genomic features**

- 80 000 proteins;
- 200 000 mRNA;

$n = 10^3 - 10^5$  **samples.**

- 10 million SNPs;
- 28 million CpG islands.

# From genotype to phenotype



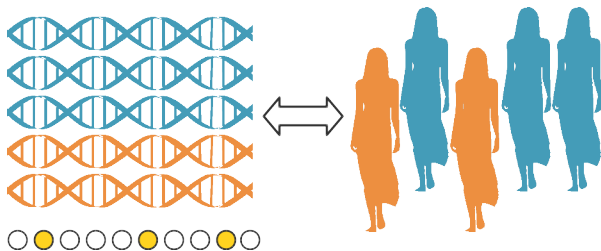
## Which genomic features explain the phenotype?

$p = 10^5 - 10^7$  **genomic features**       $n = 10^3 - 10^5$  **samples.**

- 80 000 proteins;
- 200 000 mRNA;
- 10 million SNPs;
- 28 million CpG islands.

**High-dimensional** (large  $p$ ), **low sample size** (small  $n$ ) data.

# From genotype to phenotype



Which genomic features explain the phenotype?

$p = 10^5 - 10^7$  **genomic features**       $n = 10^3 - 10^5$  **samples.**

- 10 million **S**ingle  
**N**ucleotide **P**olymorphisms.

**Genome-W**ide **A**ssociation **S**tudies.

# Missing heritability

GWAS **fail to explain** most of the **inheritable variability** of complex traits.

Many possible reasons:

- non-genetic / non-SNP factors
- heterogeneity of the phenotype
- rare SNPs
- weak effect sizes
- **few samples in high dimension ( $p \gg n$ )**
- joint effects of **multiple SNPs.**



# Integrating prior knowledge: Network-guided GWAS

Joint work with Dominik Grimm, Yoshinobu Kawahara, Karsten Borgwardt, and Héctor Climente González.

# Integrating prior knowledge

Use additional data and **prior knowledge** to **constrain** the feature selection procedure.

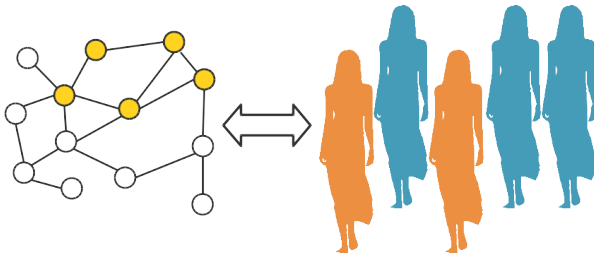
- **Consistent** with previously established knowledge;
- More easily **interpretable**;
- **Statistical power**.

Prior knowledge can be represented as **structure**:

- Linear structure of the genome;
- Groups: e.g. pathways;
- **Networks** (molecular, 3D structure).

# Network-guided biomarker discovery

- ▶ **Biological networks** help understanding disease.
- ▶ Goal: Find a **set of explanatory features** compatible with a **given network** structure.



C.-A. Azencott (2016). **Network-guided biomarker discovery**, LNCS.

# Integrating prior network knowledge

## ► Network-constrained lasso:

$$\arg \min_{\beta \in \mathbb{R}^p} \underbrace{\frac{1}{2} \sum_{i=1}^n \left( y^i - \sum_{j=1}^p \beta_j x_{ij} \right)^2}_{\text{loss}} + \underbrace{\lambda \sum_{j=1}^p |\beta_j|}_{\text{sparsity}} + \underbrace{\eta \sum_{j=1}^p \sum_{k=1}^p \beta_j L_{jk} \beta_k}_{\text{connectivity}}$$

## ► Graph Laplacian $L \rightarrow \beta$ varies **smoothly** on the network.

$$L_{jk} = \begin{cases} 1 & \text{if } j = k \\ -W_{jk} / \sqrt{d_j d_k} & \text{if } j \sim k \\ 0 & \text{otherwise.} \end{cases}$$

C. Li and H. Li (2008). **Network-constrained regularization and variable selection for analysis of genomic data**, *Bioinformatics*, 24, 1175–1182.

# Regularized relevance

Set  $\mathcal{V}$  of  $p$  variables.

- ▶ **Relevance score**  $R : 2^{\mathcal{V}} \rightarrow \mathbb{R}$

Quantifies the importance of any subset of variables for the question under consideration.

Ex : correlation, HSIC, statistical test of association.

- ▶ **Structured regularizer**  $\Omega : 2^{\mathcal{V}} \rightarrow \mathbb{R}$

Promotes a sparsity pattern that is compatible with the constraint on the feature space.

Ex : cardinality  $\Omega : \mathcal{S} \mapsto |\mathcal{S}|$ .

- ▶ **Regularized relevance**

$$\arg \max_{\mathcal{S} \subseteq \mathcal{V}} R(\mathcal{S}) - \lambda \Omega(\mathcal{S})$$

# Network-guided GWAS

- ▶ **Additive test of association** SKAT: [Wu et al. 2011]

$$R(\mathcal{S}) = \sum_{j \in \mathcal{S}} c_j \quad c_j = (\mathbf{X}^\top (\mathbf{y} - \mu))_j^2.$$

- ▶ **Sparse Laplacian regularization:**

$$\Omega : \mathcal{S} \mapsto \sum_{j \in \mathcal{S}} \sum_{k \notin \mathcal{S}} W_{jk} + \alpha |\mathcal{S}|.$$

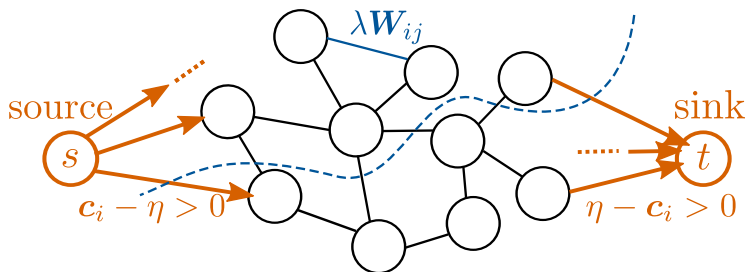
- ▶ **Regularized maximization of  $R$ :**

$$\arg \max_{\mathcal{S} \subseteq \mathcal{V}} \underbrace{\sum_{j \in \mathcal{S}} c_j}_{\text{association}} - \underbrace{\eta |\mathcal{S}|}_{\text{sparsity}} - \lambda \underbrace{\sum_{j \in \mathcal{S}} \sum_{k \notin \mathcal{S}} W_{jk}}_{\text{connectivity}}.$$

# Minimum cut reformulation

The graph-regularized maximization of score  $Q(*)$  is equivalent to a  $s/t$ -min-cut for a graph with adjacency matrix  $\mathbf{A}$  and two additional nodes  $s$  and  $t$ , where  $\mathbf{A}_{ij} = \lambda \mathbf{W}_{ij}$  for  $1 \leq i, j \leq p$  and the weights of the edges adjacent to nodes  $s$  and  $t$  are defined as

$$\mathbf{A}_{si} = \begin{cases} c_i - \eta & \text{if } c_i > \eta \\ 0 & \text{otherwise} \end{cases} \quad \text{and} \quad \mathbf{A}_{it} = \begin{cases} \eta - c_i & \text{if } c_i < \eta \\ 0 & \text{otherwise} \end{cases} .$$



**SConES: Selecting Connected Explanatory SNPs.**

# Comparison partners

- ▶ **Univariate linear regression**

$$\arg \min_{\beta_j \in \mathbb{R}} \frac{1}{2} \|\mathbf{y} - \beta_j \mathbf{x}_j\|_2^2.$$

- ▶ **Lasso**

$$\arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \eta \|\boldsymbol{\beta}\|_1.$$

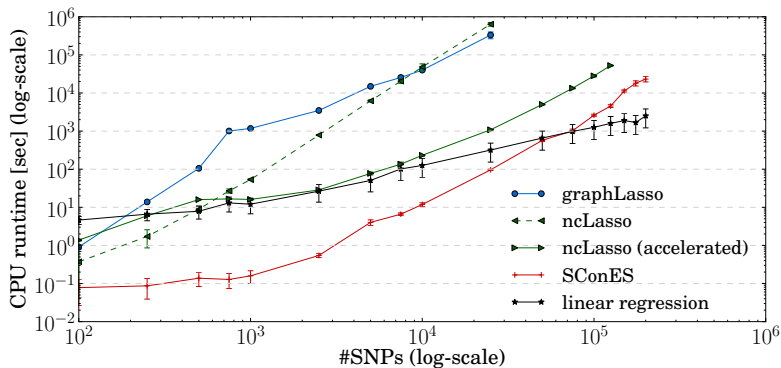
- ▶ **Feature selection with sparsity and connectivity constraints**

$$\arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \eta \|\boldsymbol{\beta}\|_1 + \lambda \Omega(\boldsymbol{\beta}).$$

- **nCLasso**: network connected Lasso [Li and Li, Bioinformatics 2008]
- Overlapping group Lasso [Jacob et al., ICML 2009]
  - **groupLasso**: E.g. SNPs near the same gene grouped together.
  - **graphLasso**: 1 edge = 1 group.



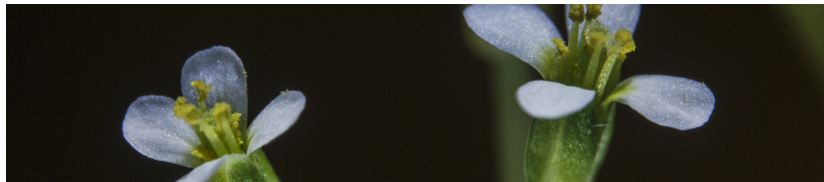
# Runtime



$n = 200$  exponential random network (2 % density)

# Experiments: Performance on simulated data

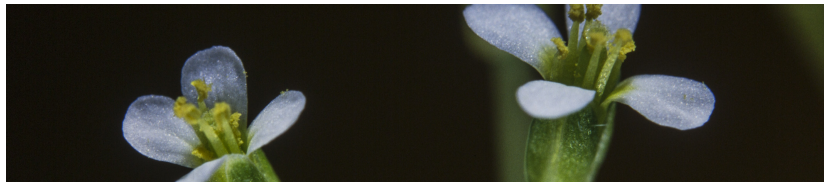
- ▶ *Arabidopsis thaliana* genotypes:  
n=500 samples, p=1 000 SNPs,  
TAIR Protein-Protein Interaction data  $\approx 50 \cdot 10^6$  edges.



- ▶ Higher **power** and lower **FDR** than comparison partners except for groupLasso when groups = causal structure.
- ▶ Systematically **better than relaxed** version (ncLasso).
- ▶ Fairly robust to **missing edges**.
- ▶ Fails if network is **random**.

# Experiments: Performance on real data

- ▶ *Arabidopsis thaliana* genotypes:  
n  $\approx$  150 samples, p  $\approx$  170 000 SNPs,  
165 **candidate genes** [Segura et al., Nat Genet 2012].



- ▶ SConES selects **about as many SNPs** as other network-guided approaches but **they tag more candidate genes**.
- ▶ **Predictivity** of the selected SNPs:
  - ▶ In half the cases, **lasso** outperforms all other approaches;
  - ▶ In the remaining cases, **SConES** outperforms all other approaches.

Image source: Jean Weber / INRA via Flickr.

## SConES: Selecting Connected Explanatory SNPs

- ▶ selects **connected**, **explanatory** SNPs;
  - ▶ incorporates **large networks** into GWAS;
  - ▶ is **efficient**, **effective** and **robust**.
- C.-A. Azencott, D. Grimm, M. Sugiyama, Y. Kawahara and K. Borgwardt (2013) **Efficient network-guided multi-locus association mapping with graph cuts**, *Bioinformatics* 29 (13), i171–i179 doi:10.1093/bioinformatics/btt238.  
<https://github.com/chagaz/sfan>
- H. Climente, C.-A. Azencott (2017). **martini: GWAS incorporating networks in R**, doi:10.18129/B9.bioc.martini.  
[Bioconductor/martini](https://bioconductor.org/packages/martini)

# Finding interactions between a target SNP and the rest of the genome.

Joint work with Lotfi Slim, Jean-Philippe Vert,  
and Clément Chatelain.

- ▶  **$p$  variables**  $X_1, X_2, \dots, X_p \in \{0, 1, 2\}$ ;
- ▶ one **target variable**  $A \in \{-1, 1\}$ ;
- ▶ **outcome**  $Y$ .

**Which of the  $p$  variables interact with  $A$  towards  $Y$ ?**

- ▶  $p$  **variables**  $X_1, X_2, \dots, X_p \in \{0, 1, 2\}$ ;
- ▶ one **target variable**  $A \in \{-1, 1\}$ ;
- ▶ **outcome**  $Y$ .

**Which of the  $p$  variables interact with  $A$  towards  $Y$ ?**

- ▶ **GBOOST:** For each  $j = 1, \dots, p$ , LRT between
  - a full logistic regression model on  $(X_j, A, A.X_j)$ ;
  - a main-effect logistic regression model on  $(X_j, A)$ .

- ▶  $p$  **variables**  $X_1, X_2, \dots, X_p \in \{0, 1, 2\}$ ;
- ▶ one **target variable**  $A \in \{-1, 1\}$ ;
- ▶ **outcome**  $Y$ .

**Which of the  $p$  variables interact with  $A$  towards  $Y$ ?**

- ▶ **GBOOST:** For each  $j = 1, \dots, p$ , LRT between
  - a full logistic regression model on  $(X_j, A, A.X_j)$ ;
  - a main-effect logistic regression model on  $(X_j, A)$ .
- ▶ **product Lasso:** Lasso on  $(X_1, X_2, \dots, X_p, A, A.X_1, A.X_2, \dots, A.X_p)$ .



# Modeling epistasis

- ▶  $Y = \mathbb{E}[Y|A = a, X] + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma^2).$
- ▶  $Y = \mu(X) + A.\delta(X) + \epsilon,$ 
  - $\mu(X) = \frac{1}{2} (\mathbb{E}[Y|A = 1, X] + \mathbb{E}[Y|A = -1, X])$
  - $\delta(X) = \frac{1}{2} (\mathbb{E}[Y|A = 1, X] - \mathbb{E}[Y|A = -1, X]).$

# Modeling epistasis

- ▶  $Y = \mathbb{E}[Y|A = a, X] + \epsilon \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$ .
- ▶  $Y = \mu(X) + A.\delta(X) + \epsilon$ ,
  - $\mu(X) = \frac{1}{2} (\mathbb{E}[Y|A = 1, X] + \mathbb{E}[Y|A = -1, X])$
  - $\delta(X) = \frac{1}{2} (\mathbb{E}[Y|A = 1, X] - \mathbb{E}[Y|A = -1, X])$ .
- ▶ SNPs in **epistasis** with  $A = \text{support}$  of  $\delta(X)$ .

# Modified outcome

$$\delta(X) = \frac{1}{2} (\mathbb{E}[Y|A = 1, X] - \mathbb{E}[Y|A = -1, X]).$$

- ▶ Introduce  $\tilde{A} = \frac{1}{2}(A + 1) \in \{0, 1\}$ .

$$\delta(X) = \frac{1}{2} \mathbb{E} \left[ Y \left( \frac{\tilde{A}}{\pi(\tilde{A} = 1|X)} - \frac{1 - \tilde{A}}{\pi(A = -1|X)} \right) \middle| X \right].$$

- ▶ **Modified outcome:**

$$\tilde{Y} = Y \left( \frac{\tilde{A}}{\pi(\tilde{A} = 1|X)} - \frac{1 - \tilde{A}}{\pi(A = -1|X)} \right).$$

# Propensity scores

$$Y = \mu(X) + A.\delta(X) + \epsilon \quad \delta(X) = \frac{1}{2} \mathbb{E} \left[ \tilde{Y} | X \right].$$

- ▶ Modified outcome was first proposed for **clinical trials**:

- $A$ : **treatment**;
- $X$ : **clinical covariates**;
- $Y$ : clinical trial **outcome**.

L. Tian et al. (2014). **A simple method for estimating interactions between a treatment and a large number of covariates**. JASA 109, 1517–1532.

- ▶ In GWAS,  $A$  and  $X$  are **not independent** because of **linkage disequilibrium**.

⇒ **propensity score**  $\pi(A|X)$ .

# Propensity scores

- ▶ Estimate **propensity scores**  $\pi(A|X)$
- ▶ Use **genomic structure**  $\Rightarrow$  **Hidden Markov Model**.
  - **Hidden states**: contiguous clusters of phased haplotypes;
  - **Emission states**: SNPs.
- ▶ Typically used for
  - ▶ imputing **missing values**;  
P. Scheet and M. Stephens (2006). **A fast and flexible statistical model for large-scale population genotype data**, AJHG 78, 629–44.
  - ▶ constructing **knockoffs** for **FDR control**.  
M. Sesia, C. Sabatti and E. J. Candès (2018). **Gene hunting with hidden markov model knockoffs**, Biometrika.

# Modified outcome variants

$$\tilde{Y} = Y \left( \frac{\tilde{A}}{\pi(\tilde{A} = 1|X)} - \frac{1 - \tilde{A}}{\pi(A = -1|X)} \right).$$

- ▶ Propensity scores tend to be close to 0.
- ▶ **Shifted modified outcome:**  $\pi(\tilde{A}|X) \leftarrow \pi(\tilde{A}|X) + \xi$ .
- ▶ **Robust modified outcome.**

J. M. Robins, A. Rotnitzky, and L. P. Zhao (1994). **Estimation of regression coefficients when some regressors are not always observed**, J. Am. Stat. Ass., 427 (89), 846–866.

# Evaluating the support of $\delta$

- ▶  $\delta(X) = \frac{1}{2} \mathbb{E}[\tilde{Y}|X]$ .
- ▶ Use an **elastic net** regression to relate  $\tilde{Y}$  and  $X$ :

$$\arg \min_{\beta \in \mathbb{R}^p} \frac{1}{n} \sum_{i=1}^n \left( \tilde{Y}_i - \beta^\top X_i \right)^2 + \lambda \left( (1 - \alpha) \|\beta\|_1 + \alpha \|\beta\|_2^2 \right).$$

$\alpha$  small  $\rightarrow$  **sparsity**.

- ▶ Add **stability selection**
  - ▶  $B$  bootstrap samples;
  - ▶ rank features based on the **area under the stability path**.

A.-C. Haury et al. (2012), **TIGRESS: Trustful Inference of Gene REGulation using Stability Selection**, BMC Sys. Bio. 6.

# Simulations

$$\pi_Y = \underbrace{\beta_{i,V}^\top X_V}_{\text{synergy with A}} + \underbrace{\beta_W^\top X_W}_{\text{marginal effects}} + \underbrace{X_{Z_1}^\top \text{diag}(\beta_{Z_1, Z_2} X_{Z_2})}_{\text{quadratic effects}}.$$

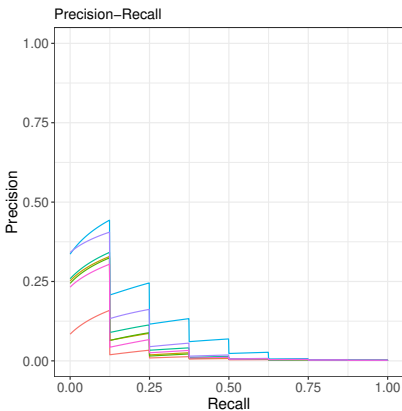
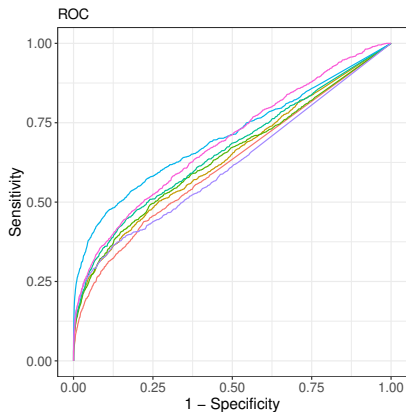
$$\pi_Y = \text{logit}(P(Y = 1 | \tilde{A} = i, X)).$$

- $p = 5\,000, n = 500.$
- $|V| = |W| = |Z_1| = |Z_2| = 8$
- $|V \cap W| = 2, |V \cap Z_1| = 2.$



# Simulations

$$\pi_Y = \underbrace{\beta_{i,V}^\top X_V}_{\text{synergy with A}} + \underbrace{\beta_W^\top X_W}_{\text{marginal effects}} + \underbrace{X_{Z_1}^\top \text{diag}(\beta_{Z_1, Z_2} X_{Z_2})}_{\text{quadratic effects}}.$$



— Outcome weighted learning — Normalized modified outcome — Robust modified outcome — GBOOST  
— Modified outcome — Shifted modified outcome — Product LASSO

## epiGWAS: Detecting epistasis with a target SNP.

- ▶ searches for a **sum of quadratic effects** with the **target SNP**;
- ▶ accounts for **main effects**;
- ▶ models **linkage disequilibrium**.

L. Slim, C. Chatelain, C.-A. Azencott, J.-P. Vert. (2018) **Novel methods for epistasis detection in genome-wide association studies**, BioRxiv.

[CRAN/epiGWAS](#)

# Looking ahead

## ▶ **Robustness/stability**

Stability selection is time consuming.

## ▶ **Complex interaction patterns**

epiGWAS is limited to a sum of quadratic interactions between one target SNP and the rest of the genome.

## ▶ **Statistical significance**

- **Significant pattern mining** [Llinares-López et al, Bioinformatics 2018].
- **Post-selection inference**
  - For the **lasso** [Lee et al., AoS 2016].
  - For **higher-order interactions** [Suzumura et al., ICML 2017].
  - Ongoing work with L. Slim on **kernel PSI**.
- Controlling **FDR** with **knockoffs** [Sesia et al., Biometrika 2018].

## CBIO:

Héctor Climente González, Lotfi Slim, Jean-Philippe Vert (Google Brain).

## Formerly MLCB Tübingen:

Karsten Borgwardt (ETH Zürich, Switzerland), Dominik Grimm (Weihenstephan, Germany), Mahito Sugiyama (National Institute of Informatics, Japan).

## Osaka University & RIKEN AIP:

Yoshinobu Kawahara.

## Sanofi:

Clément Chatelain.



source: <http://www.flickr.com/photos/wwworks/>



## Paris Women in Machine Learning and Data Science.

- ▶ March 12, 19:30
  - Human body extraction from images** – Gül Varol (INRIA Willow).
  - Data is beautiful, please don't ruin it** – Anne-Marie Tusch (Criteo Lab).
  - Salary negotiation workshop** – Natalie Cernecka.
- ▶ March 28, 19:00 – **Femmes, sciences et société**
  - Femmes, probabilités et finances** – Nicole El Karoui.
  - La féministe, l'économiste et la cité** – Hélène Périvier.
  - Discussion ouverte.