# Kernel methods
# for genomic data fusion

## Yves Moreau
### University of Leuven, Belgium

# Genetic diagnosis

- ## Main medical goals
  - End diagnostic odyssey
  - Estimate risk for next pregnancy
  - Predict disease progression, life expectancy, etc.
- ## Patient - deletion del(22)(q12.2)
  - Pulmonary valve stenosis
  - Cleft uvula
  - Mild dysmorphism
  - Mild learning difficulties
  - High myopia

# Deletion del(22)(q12.2)



- Deletion on Chromosome 22
  - ~0.8Mb
- Deletion contains NF2
  - NF2 ↔ acoustic neurinomas
  - Benign tumor, BUT
    - Hard to diagnose
    - Severe complications

# Candidate gene prioritization

# Data fusion

# Challenge of heterogeneous data

# Prioritization by example

- **Known/training genes**
  - Type 2 diabetes: 21 known genes in OMIM, 118 known genes in GAD
  - Manually curated gene set from Elbers et al., 2007
    - ACDC, ADRA2A, ADRA2B, ADRB1, ADRB2, ADRB3, LEP, LEPR, NR3C1, UCP1, UCP2, UCP3, PPARG, KCNJ11, TCF7L2
- **Candidate/test genes**
  - Prioritizations of a known region (from Elbers et al., 2007)
    - 12q24: 327 candidates

# Region 12q24: 327 candidates



TCF1 — Responsible for MODY, an uncommon monogenetic form of early onset T2D.

GPR109E
P2RX4
TBX3

NCOR2 — NCOR2 has an important role in the adipocyte by inhibiting adipocyte differentiation via repression of PPAR-g activity.

PTPN11

FZD10
ATP6V0A

SCARB1 — Key component in the reverse cholesterol transport pathway. Genetically associated with differences in insulin sensitivity in healthy subjects

McCarthy *et al.* (2006), Cohen *et al.* (2006), Perez-Martinez *et al.* (2005)

9

# Data fusion with order statistics



Candidate (test) genes

n data sources

Data source

Known (training) genes

n prioritizations

Overall prioritization

- *Aerts et al. Nature Biotech. 2006*

10

# Training of an attribute submodel



- A term is over-represented if its frequency inside the training set is significantly larger than its frequency over the genome
  - Gene Ontology, Interpro, KEGG & EST submodels

# Scoring of an attribute submodel



**Annotations**

| | p-value |
|---|---|
| Term *1* | 0.00054 |
| Term *4* | 0.00072 |
| … | … |
| Term *t* | 0.00457 |

Scoring derived from Fisher's omnibus statistic
- $S = -2 \, \bigcirc_i \log p_i$

| | Term 1 | | | Term 1 | … | | Term t |
|---|---|---|---|---|---|---|---|
| Candidate *1* | ✓ | □ | □ | ✓ | … □ | ✓ | □ |
| … | … | … | … | … | … … | … | … |
| Candidate *m* | □ | ✓ | ✓ | □ | … ✓ | □ | □ |

| | p-value |
|---|---|
| Candidate *1* | 0.0005 |
| … | … |
| Candidate *m* | 1.0 |

- Multiple species:
  - Human, mouse, rat, fly, worm
- Integration across species will soon be supported

http://www.esat.kuleuven.ac.be/endeavour

# Large-scale statistical validation

- Evaluation by an independent third party (pharma)

- MetaCore pathway and disease maps
  - 454 pathway maps with 10,053 pathway genes
  - 833 disease maps with 12,699 disease genes

- ROC curve for ranks

# A novel locus for congenital heart defect on chromosome 6q24-25



Patients with cardiac defects → Chromosome 6 → 105 candidate genes → Prioritization with endeavour ← 7 cardiac phenotypes

NHSL1
ECT2L
REPS1
...
PNLDC1
MAS1
IGF2R

c.622 C→T

G T C C A C

Mutation screen

★ Translocation t(2;6) carrier

| | Rank | | | | |
|---|---|---|---|---|---|
| | **1** | **2** | **3** | **4** | **5** |
| **Human CHDs** | TAB2 | GTF2H5 | GRM1 | HIVEP2 | HECA |
| **2nd Heartfield** | HIVEP2 | FBXO30 | PPP1R14 | TXLNB | OPRM1 |
| **1st Heartfield** | HIVEP2 | PPP1R14 | TXLNB | OPRM1 | GPR126 |
| **Left/right asymmetry** | CITED2 | IGF2R | GRM1 | MTRF1L | FBXO30 |
| **Valve formation** | TAB2 | UTRN | HIVEP2 | ZBTB2 | PEX3 |
| **Neural crest** | CITED2 | TAB2 | FBXO30 | GRM1 | HIVEP2 |
| **Vasculogenesis** | GPR126 | UTRN | IGF2R | CITED2 | TAB2 |
| **Combination** | TAB2 | HIVEP2 | CITED2 | GRM1 | FBXO30 |

Thienpont et al. *Am J Hum Genet.* 2010    15

# Kernel methods for genomic data fusion

# Kernel-based genomic data fusion



Kernel matrix
~ nonlinear extension of covariance/correlation matrix

Instead of using original data directly, use kernel matrix only
(Think of hierarchical clustering.)

Advantage 1: kernel matrices form a single type of object, regardless of the heterogeneity of the original data types

Advantage 2: all machine learning methods can be applied to kernels (classification, clustering, prioritization, ranking, etc.)

# Kernel data fusion (a.k.a. MKL)

# Prioritization by novelty detection

# One-class support vector machine

$$\boxed{\text{P:}} \min_{\vec{w}, \xi, \rho} \quad \frac{1}{2} \vec{w}^T \vec{w} - \frac{1}{\nu l} \sum_{k=1}^{l} \xi_k - \rho$$

$$\text{s.t.} \quad \vec{w}^T \phi(\vec{x}_k) \geq \rho - \xi_k, \ \ k = 1, \dots, N$$

$$\xi_k \geq 0, \ \ k = 1, ..., N.$$

$\vec{w}$: the norm vector of the separating hyperplane
$\vec{x}_k$: the training samples
$\nu$: a regularization term penalizing the outliers in the training samples
$\phi(\cdot)$: the feature map
$\rho$: the bias term
$\xi_k$: the slack variables
$N$: the number of training samples

$$\boxed{\text{D:}} \min_{\vec{\alpha}} \vec{\alpha}^T K \vec{\alpha}$$

$$\text{s.t.} \quad 0 \leq \alpha_k \leq \frac{1}{\nu N}, \ \ k = 1, ..., N$$

$$\sum_{k=1}^{N} \alpha_k = 1,$$

$\alpha_k$: the dual variables
$K$: the kernel matrix

# Kernel fusion for novelty detection



$$K = \mu_1 K_1 + \mu_2 K_2$$

# Kernel fusion in one-class SVM

- $L_\infty$-norm kernel fusion (De Bie et al., 2007)

$$\min_{\vec{\alpha}} \ t$$

$$\text{s.t.} \ \ t \geq \vec{\alpha}^T K_j \vec{\alpha}, \ j = 1, ..., p$$

$$0 \leq \alpha_k \leq \frac{1}{\nu N}, \ \ k = 1, ..., N$$

$$\sum_{k=1}^{N} \alpha_k = 1,$$

$p$: the number of kernel matrices
$K_j$: the $j$-th kernel matrix

- $L_2$-norm kernel fusion (Yu et al., 2009)

$$\min_{\vec{\alpha}} \ t$$

$$\text{s.t.} \ \ t \geq ||s_j||_2, \ j = 1, ..., p$$

$$s_j \geq \vec{\alpha}^T K_j \vec{\alpha}, \ j = 1, ..., p$$

$$0 \leq \alpha_k \leq \frac{1}{\nu N}, \ k = 1, ..., N$$

$$\sum_{k=1}^{N} \alpha_k = 1.$$

$s_j$: dummy variables

# $L_2$ vs. $L_\infty$ kernel fusion



Table 1: AUC values of LOO performance evaluated from 20 random repetitions. The paired Spearman correlation scores indicate the similarities of rankings obtained by different approaches compared with the target rankings (denoted as -).

|  | AUC | corr | corr | corr | corr |
|---|---|---|---|---|---|
| $L_\infty$ | 0.9045(0.0043) | - | 0.94 | 0.66 | 0.82 |
| $L_\infty(0.5)$ | 0.9176(0.0040) | 0.94 | - | 0.82 | 0.92 |
| $L_1$ | 0.9103(0.0035) | 0.66 | 0.82 | - | 0.90 |
| $L_2$ | **0.9219(0.0034)** | 0.82 | 0.92 | 0.90 | - |



| | |
|---|---|
| $L_\infty$: | 0.9018 |
| $L_\infty(0.5)$: | 0.9139 |
| $L_1$: | 0.9064 |
| $L_2$: | 0.9212 |

# A framework for kernel data fusion

# Kernel data fusion



Data source (multiple DBs, multiple organisms) — Formatted data — Contextualization — Kernel matrix

Data type

Kernel functions

Kernel combination

Prioritization

Clustering

Classification

# ETkL: Extract, Transform, Kernelize, Learn

- Systematic multi-tier framework for data integration
  - Resembles multi-tier architecture of complex IT systems and Extract-Transform-Load methodology of datawarehousing
    1. Database / web service sources
    2. Data reconciliation, cleaning, and warehousing, etc.
    3. Scaling, normalization, feature selection, etc.
    4. Computation and storage of kernels
    5. Learning
  - May require feedback loops  (e.g., feature selection)
- Scale up to large, heterogeneous databases
- 20,000 x 20,000 kernel matrices are ugly animals

# The No-Voodoo principle

- Given a data matrix D for a learning problem, the no voodoo principle states that, in the absence of prior knowledge or arbitrary assumptions, no information can be extracted about the problem except the information provided by the data matrix
  - In particular, no information can be created that wasn't initially present in the data
    - No amount of bagging, random projection, nonlinear high-dimensional feature map, etc. can extract information that was not present in the data (except through the implicit or explicit injection of constraints into the problem)
  - If two frameworks represent data in ways that are related in a one-to-one fashion, there is nothing that prevents the development of methods with identical accuracy (e.g., random projections vs. spectral methods)
  - If one method outperforms another on a given problem (remember the no free lunch theorem), it is because the methods are more or less efficient (in particular, in terms of generalization performance vs. retrospective accuracy) at capturing the available information or because the methods incorporate explicit or implicit constraints that are more or less relevant to the given learning task

# Handling large kernel matrices

- One way to handle large kernel matrices is via low-rank approximations
  - Store $r$ x $n$ instead of $n$ x $n$
- Cholesky decomposition
  - $K$ symmetric positive definite

$$\exists C (\text{lower triangular} \;\&\; \text{unique}) : K = CC'$$

# Incomplete Cholesky decomposition

- Incomplete Cholesky
  - $K$ symmetric positive semidefinite
  - Limit to rank $r \leq \mathrm{rank}(K)$
  - Add pivoting to capture more informative rows/columns first
  - Limit information loss to e.g. 5%

What if no or few genes known for a disease?

# Expression of candidate genes

- For positional cloning, checking expression of candidate genes is standard but has a low yield
    - No guarantee that disease gene itself is perturbed

- Existing prioritization methods (e.g., Endeavour) rely heavily on prior knowledge and hard to achieve "breakthroughs"

# Systems biology: network analysis

# Integrative protein network

- e.g., STRING

# Concept



**Candidate genes**

A

absolute diff. expr. level

distance between A and neighbors

B

absolute diff. expr. level

distance between A and neighbors

34

# Concept



**Candidate genes**

A strong candidate has many partners highly differentially expressed!

A

absolute diff. expr. level

distance between A and neighbors

B

absolute diff. expr. level

distance between A and neighbors

# Methods

- Machine-learning strategies
  - Naive ranking
    - Use only the differential expression of the candidate
  - Direct neighborhood ranking
    - Combine differential expression level of candidate with the average of the differential expression levels of the direct neighbors
  - Kernel ridge regression
    - Smooth a candidate's differential expression level by kernel ridge regression
  - Approximate heat kernel diffusion
    - Discrete low-accuracy approximation to the exponential diffusion kernel $\exp(\alpha L) = \exp(\alpha(D\text{-}A))$ takes direct and indirect association into account
  - Arnoldi diffusion
    - Memory-light high-accuracy approximation to exponential diffusion kernel using Krylov subspaces (Arnoldi algorithm)

# Methods & benchmark

- Benchmark: 40 KO experiments in mouse
  - Publicly available data sets from GEO - Affymetrix platform
  - Simple KO versus control
  - How well can we rank the KO gene?
  - Which algorithm and what combination of steps performs best?
- Preprocessing
  - RMA
  - GCRMA
  - MAS5
- Differential expression
  - Log2 ratio
  - Regularized t-statistic (CyberT)
  - Significant log2 ratio
- Different networks
  - STRING7, STRING8
  - PPI Network from BioGRID
  - PPI Network from I2D

| Database (mouse) | Number of genes | Number of interactions | Average node degree |
|---|---|---|---|
| STRING v7.1 | 16,566 | 820,177 | 49.5 |
| STRING v8.2 | 24,442 | 1,405,375 | 57.5 |
| BioGRID v2.0.61 | 1,417 | 2,026 | 2.5 |
| I2D v1.72 | 10,867 | 79,088 | 10.6 |

# Results

| Strategy | AUC | Error reduction relative to baseline |
|---|---|---|
| Simple expression ranking | 83.0% | baseline |
| Direct neighborhood ranking | 88.0% | 26.4% |
| Kernel ridge regression | 86.8% | 19.0% |
| Heat kernel | 92.3% | 52.8% |
| Arnoldi diffusion | 87.4% | 22.7% |

# A candidate gene for PCOS

- PolyCystic Ovary Syndrome (PCOS)
  - Major cause of infertility (chronic anovulation)
  - Hormonal dysfunction (hyperandrogenism)
  - Obesity (depending on diagnostic criteria)
  - Oligogenic disorder (no Mendelian inheritance)
- Two confirmed susceptibility loci
  - 19p13.2 -> FBN3 (Fibrillin 3)
  - 5q11.2
    - FST (follistatin) proposed, but infirmed in subsequent validation
- Expression data (GEO GDS2084)
  - Omental (belly) fat from patients vs. control
  - Affymetrix HG-U133A

# DDX4 as a PCOS candidate


ovary 3 wk old

- Prioritization of 5q11.2
    - FST ranks 2nd
    - DDX4 ranks 1st
        - Expression in ovary follicles (image = mouse Ddx4/Vasa)
        - A germline development gene (sperm and ovary only)
            - Plausible mechanism for infertility and hyperandrogenism
            - Mechanism not previously suggested for PCOS
            - Not a perfect candidate (male phenotype in mouse, not female)

# DDX4 expression neighborhood

# PINTA web tool

# Krylov subspace methods

- Function of a square matrix $f(A)$
  - Matrix inverse
  - Matrix exponential (e.g., exponential of graph Laplacian)
  - Computationally challenging for large matrices
- In applications, often no need for matrix function $f(A)$ directly, but only its evaluation $f(A).v$ at a point $v$
- Krylov methods – simplified argument
  - Cayley-Hamilton theorem for characteristic polynomial of square $n$ x $n$ matrix $A$

$$p_n(\lambda) = \det(\lambda I_n - A)$$

$$p_n(A) = 0$$

  - Any power of $A$ higher or equal to $n$ can be expressed in function of $A^{n-1}, A^{n-2}, ..., A^2, A, I$

- Any matrix function $f(A)$ can be expressed as a polynomial of degree $n-1$

- $f(A).v$ can be expressed as a linear combination of

$$\left\{ v, Av, A^2v, ..., A^{n-2}v, A^{n-1}v \right\}$$

- Krylov methods consists in projecting $f(A).v$ onto the subspace

$$S_m = \operatorname{span}\left\{ v, Av, A^2v, ..., A^{m-2}v, A^{m-1}v \right\}$$

- Only requires matrix-vector operations!

- The set of spanning vectors is kept orthogonal via QR orthogonalization

- In practice, often fast convergence $m \ll n-1$

# Networks vs. kernels

<Rant> A network is a matrix is a network

- And a symmetric similarity matrix is an undirected graph
- Implicitly, we mean more by a network
  - Sparse matrix
  - Edges have some underlying biological reality
  - e.g., KEGG metabolic network from one organism or regulonDB yeast transcriptional network
- Most predicted protein networks do not have such properties
  - Usually calculated much like similarity matrices
  - Why handle them as "biological" networks?
  - Networks useful in visualization, but should not be misleading
  - Network representation usually involves heavy thresholding and creates an information bottleneck </Rant>
- Gillis and Pavlidis suggest that networks and similarity matrices are almost equivalent under the no-voodoo principle

# What's wrong with network propagation?

- Kernel diffusion and network module biomarkers all seem to perform less strongly than expected.
- What is wrong?
  - Nothing? (Our expectation is unrealistic)



Best 50



Best 500

# Curse of the small world?

1. Data is not specific enough to highlight the right network neighborhoods
   -> improve experimental design (e.g., factorial design)
2. Our data is improperly scaled or normalized
   (propagating apples and oranges)
3. Our networks are bad
   (STRING > BioGrid, coverage more important than specificity)
4. Network is in fact thresholded propagation matrix already (e.g., STRING)
   (further propagation does not help much)
5. Our notion of neighborhood and diffusion is unsuitable (curse of the small world)
   (rough approximation to heat kernel works better than accurate one)
6. Our randomization procedures are unsuitable
   (propagation results are still apples and oranges)
7. Uncertainty propagation is the bottleneck
   (more complex model propagates more noise and thus destroys the advantage of added knowledge)

# Drug target prioritization

# Target prioritization

- **One drug, many targets**
  - Many targets are unknown
    - Side effects
    - Synergistic effects
  - Candidates identified by phenotypic screen
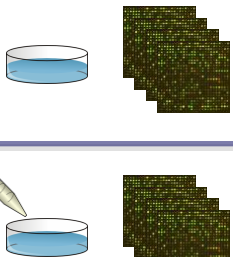  - Knowledge of a drug candidate's mode of action can help drug development
- **Predict targets based on gene expression following treatment**
  - Distinguish between genes targeted by the drug and indirectly regulated genes

# Network analysis of drug response

- Gene expression integrated with protein associations
- Neighborhood analysis
- Gene prioritization based on differential expression of functionally related network neighborhood



gene expression before and after drug administration → calculate differential expression value for each gene → map differential expression values to String protein association network → score genes based on differential expression of their neighborhood → set up ranking by correcting for neighborhood size

# Method

- **Filtering**
  - No filtering
- **Expression measure**
  - log ratio
- **Network**
  - STRING 8.2
- **Parameters**
  - N = 1, 2, 3, α = 0.9

$$p_\alpha = p_0 \left( I - \frac{\alpha}{N} L \right)^N$$

  - α = 0.9, β=0.1

$$p_\alpha = p_0 \left( I - \alpha L + \frac{\alpha^2}{4} L^2 \right) \Rightarrow p_\alpha = p_0 \left( I - \alpha L + \beta L^2 \right)$$

# Preliminary results: monoclonal antibodies

- Monoclonal antibodies specifically bind to one target

- 7 datasets from Gene Expression Omnibus:
    - tocilizumab: IL6
    - bevacizumab: VEGFA
    - rituximab: MS4A1
    - infliximab: TNF
    - h10H5: IGF1R
    - anti-CD25: IL2RA
    - LY2439821: IL17A

- Can we identify the target from expression response?

# Preliminary results: monoclonal antibodies

| | | tocilizumab | bevacizumab | rituximab | infliximab | h10H5 | anti-CD25 | LY2439821 | # in top 5% | # in top 10% |
|---|---|---|---|---|---|---|---|---|---|---|
| differential expression | | 1,057 | 6,896 | 176 | 4,281 | 12,279 | 522 | 1,992 | 2 | 3 |
| α=0.9 | N=1 | 982 | 99 | 342 | 142 | 9,109 | 48 | 227 | 5 | 6 |
| | N=2 | 597 | 279 | 268 | 2,254 | 5,055 | 102 | 766 | 5 | 5 |
| | N=3 | 720 | 446 | 151 | 2,628 | 4,645 | 148 | 840 | 4 | 5 |
| α=0.9, β=0.1 | | 763 | 93 | 186 | 758 | 7,780 | 38 | 454 | 6 | 6 |

- Why does the h10H5 target IGF1R rank this low?
  - Bad experiment?
  - Bad method?
  - No downstream transcriptional effect?
- Test other IGF1R inhibitor: BMS754807
  - ➔ For α=0.9 and β=0.1 IGF1R ranked at position 381

# Preliminary results: chemical drugs

- Chemical drugs can bind multiple targets

- 7 datasets from Gene Expression Omnibus
    - letrozole: CYP19A1
    - bicalutamide: AR
    - calcitriol: VDR
    - methylprednisolone: NR3C1
    - gefitinib: EGFR
    - methotrexate: DHFR
    - progesterone: PGR

- Can we identify the target from expression response?

# Preliminary results: chemical drugs

| | | letrozole | bicalutamide | calcitriol | methyl-prednisolone | gefitinib | methotrexate | progesterone | # in top 5% | # in top 10% |
|---|---|---|---|---|---|---|---|---|---|---|
| differential expression | | 14,055 | 700 | 4,262 | 4,316 | 8,612 | 871 | 79 | 2 | 3 |
| α=0.9 | N=1 | 2,460 | 23 | 887 | 5,109 | 1,848 | 650 | 919 | 2 | 4 |
| | N=2 | 1,658 | 7 | 387 | 6,402 | 915 | 683 | 118 | 4 | 5 |
| | N=3 | 2,030 | 11 | 595 | 5,283 | 990 | 500 | 112 | 4 | 5 |
| α=0.9, β=0.1 | | 1,349 | 1 | 55 | 8,246 | 1,076 | 819 | 325 | 4 | 6 |

- Why does the methylprednisolone target NR3C1 rank this low?
  - Bad experiment?
  - Bad method?
  - No downstream transcriptional effect?
- Test other NR3C1 agonist: fluticasone
  - ➜ For α=0.9 and β=0.1 NR3C1 ranked at position 40

**K.U.L. ESAT-SCD:** L. Tranchevent, Y. Shi, D. Nitsch, , R. Barriot, S. Leach, B. Coessens, S. Van Vooren
**K.U.L. CME-UZ:** J. Vermeesch, K. Devriendt, B. Thienpont, F. Hannes, J. Breckpot
**K.U.L. VIB:** D. Lambrechts, S. Maity, P. Carmeliet, S. Aerts, B. Hassan, P. Van Loo, P. Marynen
**U. Bristol:** T. De Bie
**INESC-ID, Lisbon:** J. Gonçalves, S.Madeira
**Novartis:** S. Schuierer, U. Dengler