# Logic Programming for Big Data in Computational Biology

## Nicos Angelopoulos

Wellcome Sanger Institute
Hinxton, Cambridge

nicos.angelopoulos@sanger.ac.uk

18.9.18

# overview

- knowledge for Bayesian machine learning over model structure
- applied knowledge representation for biological data analytics

# Bayesian inference of model structure (Bims)

A Bayesian machine learning system that can model prior knowledge by means of a probabilistic logic programming.

Nonmeclature

- DLPs = Distributional logic programs
- Bims = Bayesian inference of model structure

Timeline

- Theory (York, 2000-5)
- Applications (Edinburgh, 2006-8, IAH 2009, NKI 2013)
- Bims library and theory paper 2015-2017

# Bims Overview

- ▶ syntax of DLPs

- ▶ a succinct classification tree prior program

- ▶ Bayesian learning of model structure

- ▶ learning classification and regression trees

- ▶ Bayesian learning of Bayesian networks

- ▶ the bims library

# DLPs- description

We extend LP's clausal syntax with probabilistic guards that associate a resolution step using a particular clause with a probability whose value is computed on-the-fly.

The intuition is that this value can be used as the probability with which the clause is selected for resolution.

Thus in addition to the logical relation, a clause defines over the objects that appear as arguments in its head, it also defines a probability distribution over aspects of this relation.

DLPs example

$$member(H, [H|\_T]).$$
$$member(El, [\_H|T]) :-$$
$$member(El, T).$$
$$(C_1)$$

$$L :: length(List, L) \sim El :: umember(El, List) \quad (G_1)$$

$$\frac{1}{L} :: L :: umember(El, [El|Tail]). \quad (C_2)$$

$$1 - \frac{1}{L} :: L :: umember(El, [H|Tail]) :- \quad (C_3)$$
$$umember(El, Tail).$$

# DLPs probabilistic goals

$$\frac{1}{L} :: L :: umember(El, [El|Tail]). \qquad (C_4)$$

$$1 - \frac{1}{L} :: L :: umember(El, [H|Tail]) \; :- \qquad (C_5)$$
$$K \text{ is } L - 1,$$
$$K :: umember(El, Tail).$$

$$? - umember(X, [a, b, c]).$$

$$X = a \quad (1/3 \text{ of the times } = 1/3);$$
$$X = b \quad (1/3 \text{ of the times } = 2/3 * 1/2);$$
$$X = c \quad (1/3 \text{ of the times } = 2/3 * 1/2 * 1).$$

# simple tree prior

M=nd(x2,1,nd(x1,0,lf,lf),lf)

$(C_0)$   $cart(\zeta, \xi, M, Cart) :-$
         $\psi_0$ is $\zeta$,
     $\psi_0$:   $split(0, \zeta, \xi, M, Cart)$.

$(C_1)$     $\psi_D$:   $split(D, \zeta, \xi, M_B, nd(F, Val, L, R)) :-$
         $\psi_{D+1}$ is $\zeta * (1 + D)^{-\xi}$,
         $D_1$ is $D + 1$,
         $r\_select(F, Val, M_B, L_B, R_B)$,
     $\psi_{D+1}$: $split(D_1, \zeta, \xi, L_B, L)$,
     $\psi_{D+1}$: $split(D_1, \zeta, \xi, R_B, R)$.

$(C_2)$ $1 - \psi_D$:   $split(D, \zeta, \xi, M_B, lf)$.
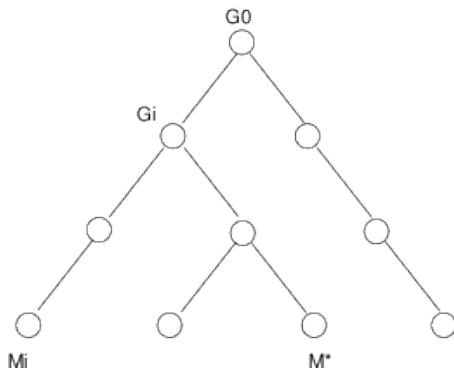
# Bims theory

Bayes' Theorem

$$p(M|D) = \frac{p(D|M)p(M)}{\sum_M p(D|M)p(M)}$$

Metropolis-Hastings

$$\alpha(M_i, M_*) = min\left\{\frac{q(M_*, M_i)P(D|M_*)P(M_*)}{q(M_i, M_*)P(D|M_i)P(M_i)}, 1\right\}$$

# DLP defined model space

From $M_i$ identify $G_i$ then sample forward to $M_\star$.
$q(M_i, M_\star)$ is the probability of proposing $M_\star$ when $M_i$ is the current model.

# Pyruvate kinase interactors

### objective
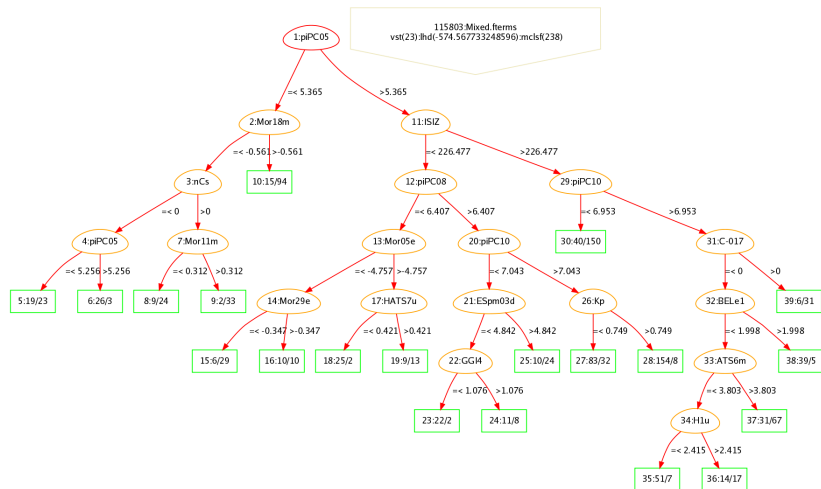improve chances of discovering binding molecules based on examples from screened chemical libraries.

### pyruvate kinase affinity data
582 Active and 582 Inactive. Dragon software produces 1500 property descriptors for each molecule, about 1100 were used.
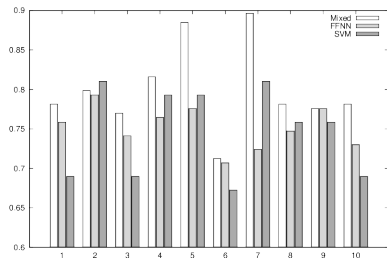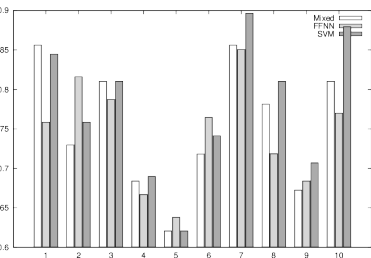
### ten-fold cross-validation
Compared to Feed Forward Neural Networks and Support Vector Machines by splitting the data into ten train/test segments.
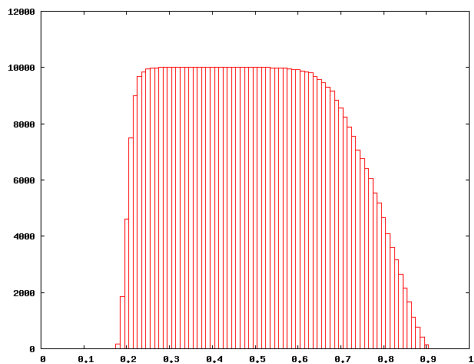
# best likelihood model

## ten-fold validation



$$Sensitivity = \frac{T^+}{T^+ + F^-}$$

$$Specificity = \frac{T^-}{T^- + F^+}$$

# molecules of Eduliss according to BCarts

# Bims: Bayesian inference of model structure

Released in 2016 as an easily installable SWI-Prolog library
Includes                      (IJAR paper in 2017)

- ▶ priors and likelihoods for: CARTs and Bayesian networks
- ▶ hooks for user defined models

Probabilistic logic programming

- ▶ thesis: probabilistic finite domains
- ▶ PLP workshop and IJAR associated issues (5th edition)

# knowledge-based computation biology

- graphical models
  (focal adhesion dynamics, NKI, 2011-3)
- proteomics functional analysis
  (TKSilac,KSR1,ATG9A, Imperial, 2014-5)
- mutational profiling
  (14MG, Sanger, 2016-8)

# Graphical models of FAD

Graphical models (aka Bayesian networks) can provide a network view of dependencies among variables, capturing much richer information than pairwise correlations.
In this project, microscopy based variables characterising focal adhesion in time are connected for a number of conditions in the HGF pathway.

untreated → HGF

prolonged HGF → HGF + 007

— negative dependence
— positive dependence

* thickness of line indicates strength of correlation

# tkSilac: tyrosine kinase screen

- ▶ MCF7 cell line
- ▶ 33 SILAC runs
- ▶ 65/66 expressed tyrosine kinases

- ▶ 4739 quantified in some experiment
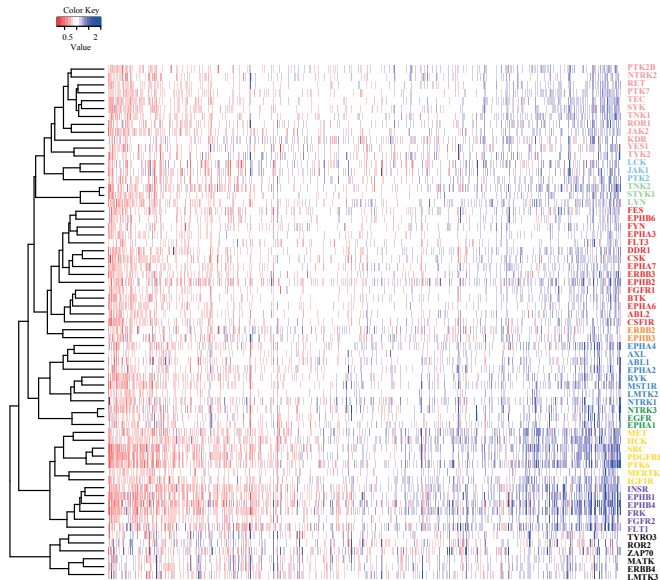- ▶ 1000 quantified in 60 or more TK KO

Figure 2



**Fig. 2. Heatmap of quantified proteins after TK silencing.** The overall pattern of regulation is shown in the heat-map of quantified values. After normalized to siControl, values of fold changes are all above 0, with value 1 show-ing that the expression levels of the specific protein are not altered after silencing TKs. For each knockdown (rows)
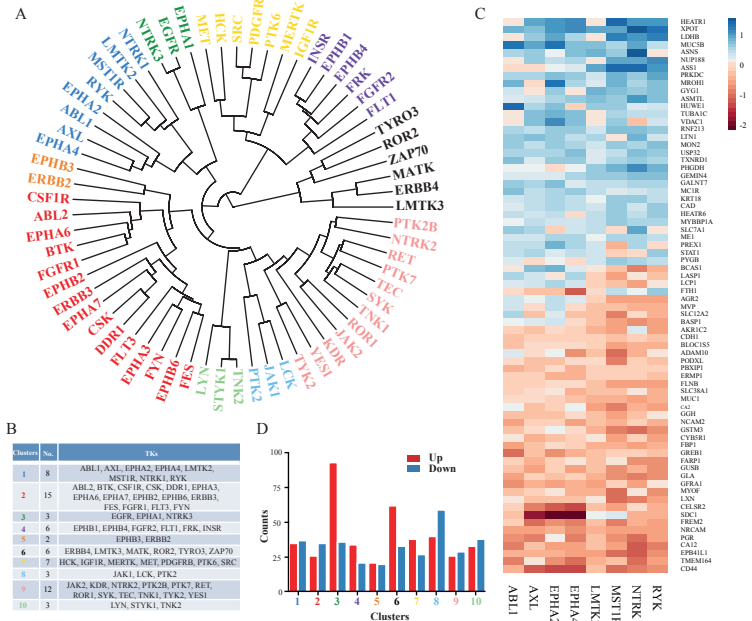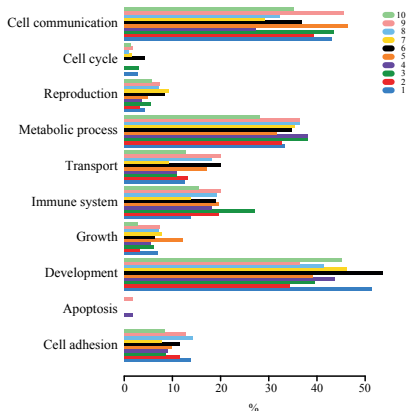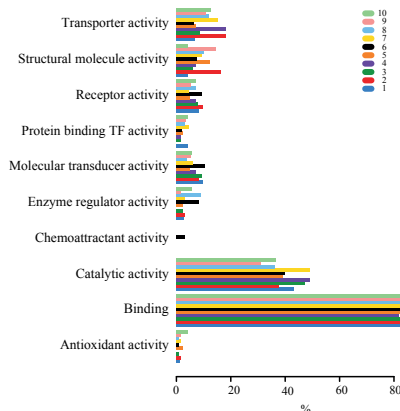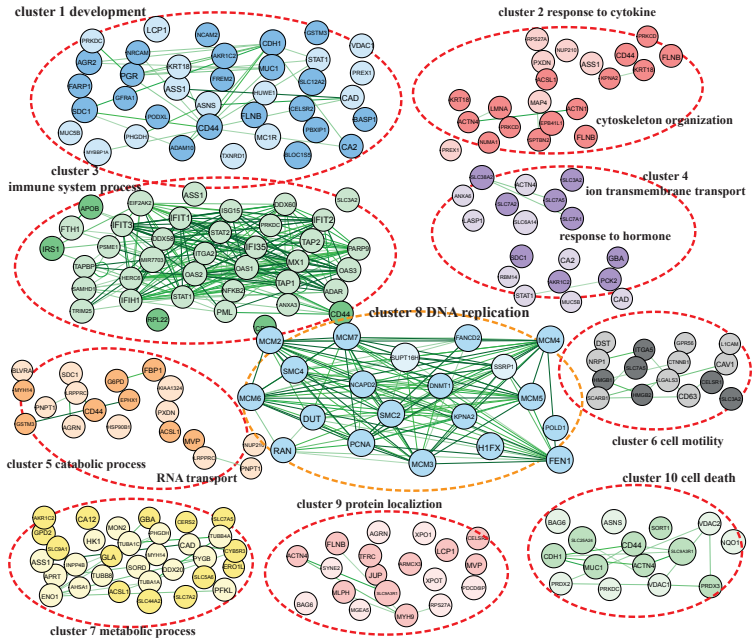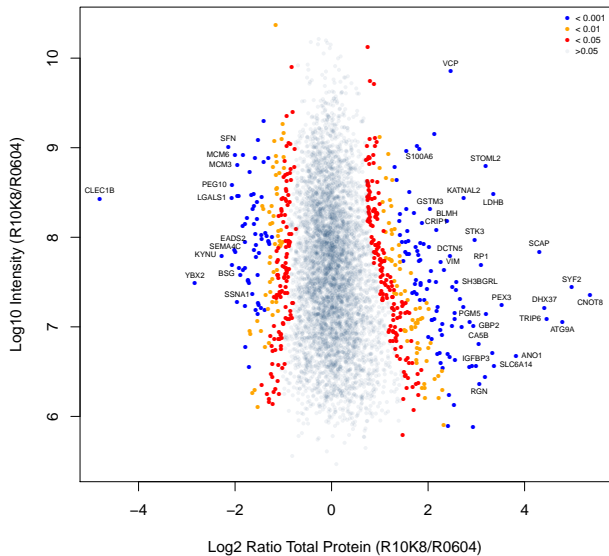
Figure 4

Figure 5



**Fig. 5. Characterization of a functional portrait for each cluster.** A, A functional profile of top GO biologic processes that the up- and downregulated proteins belong to is presented. x-axis shows the percentage of hits in each cluster that belong to a GO biologic process term. The color coding and the number for each cluster are indicated as above. B, A functional profile of top GO molecular functions that the up- and downregulated proteins belong to is presented. x-axis shows the percentage of hits in each cluster that belong to a GO molecular function term.
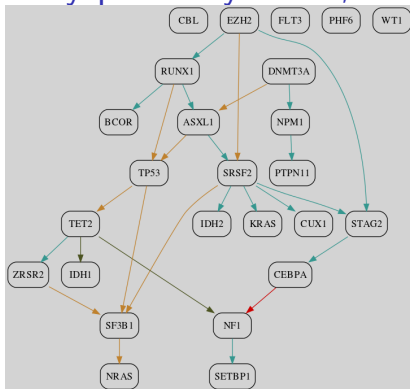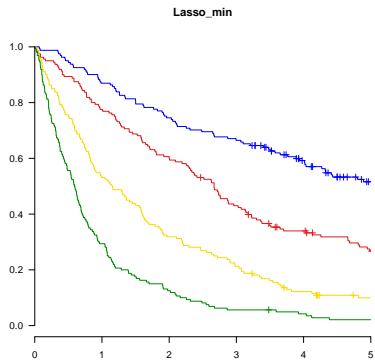
Figure 6

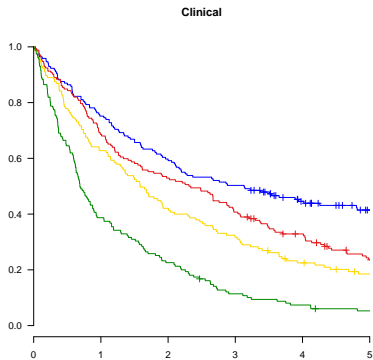# volcano plot (BT474HR H/M)

# Autophagy

# Myelodysplastic syndrom, NGS somatic mutations profiling



AUC (1y) for Clinical vs Lasso_min model

# 5 year: Clinical vs Lasso (Optimal)

# myeloma structural variations



Fisher test odds
Mut.excl (shown odds=0.25)
Co-occur (shown odds=4)
# events (shown med=225,max=643)

# logic programming for (biological) data analytics

Positives

- ▶ interpreted
- ▶ memory management
- ▶ clean and high level
- ▶ probabilistic ML & reasoning (Prism,Bims,Pepl)
- ▶ intuitive database integration (db_facts,bio_db)
- ▶ multi-threaded and web-capable
- ▶ talking to other systems (R:Real,ODBC,proSQLite)
- ▶ (largely) OS independence

Negatives

- ▶ graphics
- ▶ SWI-Prolog, at core a one-person project
- ▶ code sharing in toddler stage (but showing promise)
- ▶ in-browser interaction with other technologies

# KR bottom line

(probabilistic) logic programming and Bayesian networks are powerful tools for

explainable, accountable, open and **shareable** AI & ML

# KR bottom line

(probabilistic) logic programming and Bayesian networks are powerful tools for

explainable, accountable, open and **shareable** AI & ML

symbolic AI education, can be a central player in

contributing tangibly to the current AI resurgence, while
managing expectations of modern AI
see media coverage of Facebook/Cambridge-Analytica
& Uber/Tesla driveless accidents

# KR bottom line

(probabilistic) logic programming and Bayesian networks are powerful tools for

 explainable, accountable, open and **shareable** AI & ML

symbolic AI education, can be a central player in

 contributing tangibly to the current AI resurgence, while
  managing expectations of modern AI
      see media coverage of Facebook/Cambridge-Analytica
                    & Uber/Tesla driveless accidents

biology presents a unique application area, where

  unprecedented volumes of data are generated
    knowledge is a crucial concept, currently being shaped
    transferable to other big data areas