knowledge AI in bio-data analytics

Nicos Angelopoulos

The Pirbright Institute Head of Computational Biology & co-lead of bioinformatics

https://stoics.org.uk/~nicos nicos@stoics.org.uk

24.05.29

< □ > < 同 > < Ξ > < Ξ > < Ξ > < Ξ < </p>

overview

- Stream 1. prior knowledge for Bayesian machine learning
- Stream 2. applied knowledge representation for biological big data analytics

< □ > < 同 > < Ξ > < Ξ > < Ξ > < Ξ < </p>

- Stream 3. Bayesian networks for cancer and biological datasets
- Stream 4. single cell RNA for pig immunity

Stream 1. (York) Bayesian machine learning

How

can we incorporate existing (biological) knowledge in the analysis of new experimental data

Bayesian

methods allow for the incorporation of prior knowledge and expectations, although often applications use agnostic priors

Bayesian machine learning theory

Bayes' Theorem

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

Metropolis-Hastings

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

・ロト・(四ト・(川下・(日下・))

A probabilistic programming framework for Bayesian machine learning of structured statistical models (classification trees and Bayesian networks).

Allows the encoding of prior information in the form of a probabilistic logic program.

- Theory (York, 2000-5, KR paper 2017)
- Applications (Edinburgh, 2006-8, IAH 2009, NKI 2013)

・ロト・西ト・ヨト・ヨト・ 日・ のへの

Learning binding molecules

Edinburgh: Pyruvate kinase interactors improve chances of discovering binding molecules based on examples from screened library of chemicals

pyruvate kinase affinity data

582 Active and 582 Inactive, with 1100 property descriptors for each molecule. Compared to Feed Forward NNs and SVMs.

< □ > < 同 > < Ξ > < Ξ > < Ξ > < Ξ < </p>

best likelihood model



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣・のへで

ten-fold validation



▲□▶▲圖▶▲圖▶▲圖▶ 圖 のへで

Stream 2. (Imperial) Knowledge-based data analytics

tkSilac: tyrosine kinase screen

- MCF7 cell line
- 33 SILAC runs
- 65/66 expressed tyrosine kinases
- 4739 proteins quantified in some experiment
- 1000 proteins quantified in 60 or more TK KO
 Molecular and Cellular Proteomics (MCP) 2015

▲ロ ▶ ▲周 ▶ ▲ ヨ ▶ ▲ ヨ ▶ ● の Q @

Tyrosine kinase screen (Imperial)



Fig. 4. Ilterarchical clustering of the 65 TKs expressed in MCCP rells. A, Ilterarchical clustering of the 65 TKs was performed using R Vs helst infration. The complete linkage method which aims to identify similar clusters based on overall cluster messare was used. It 0 distinctive clusters were obtained and the complete dendrogram is shown with the labels colored for three clusters. B, rell list of the TKs included in each cluster. The color coding of the clusters is used throughout to identify the analysis relevant to the corresponding clusters. CH entrange of the protonic quantifications (log2 values of normalized fold changes against control) for the downstream effects (significantly up- or down-regulated proteins, Significant B test p> 0.05), after silencing TKs in cluster 1. D, Number of proteins significantly up or down-regulated in each identified cluster x-axis shows 10 different clusters and y-axis indicates the counts.

Fig. 6. Representatives of defined functional networks in each datafield TK cluster. The functional intervals were generated using (C) analysis: constrained with the TTSD's platform. Proteins in Splate order are use peopleted, whereas heighter order indicates drawn regulation. Access door the interactions helicence commuted proteins, Representative defined functional networks associated with their clusters are shown here. The order order gas for mashes for each during an enderoid is above.

herceptin resistance (BT474HR) — ATG9A / autophagy



proteomics data analytics (Imperial)

tyrosine kinase screen Molecular and Cellular Proteomics (MCP) 2015 KSR1: Breast Cancer Res. and Treat., 2015 ATG9A: Oncotarget 2016 Prolog libraries¹: Real (> 580), proSQLite (> 840), bio_analytics bio_db (currently: 91 tables, 55 M records on human, mouse, chicken)

¹work started at NKI

Stream 3. (Sanger) Bayesian networks in cancer genomics





AUC (1v) for Clinical vs Lasso min model



・ロット (雪) (日) (日) (日)

MPN: myeloproliferative neoplasms



New England Journal of Medicine, October 2018

myeloma structural variations



< □ > < 同 > < Ξ > < Ξ > < Ξ > < Ξ < </p>

Nature Communications, August 2019

BNs in cancer genomics

MPN published in New England J. of Medicine, Oct, 2018

< □ > < 同 > < Ξ > < Ξ > < Ξ > < Ξ < </p>

- multiple myeloma: in Nature Communications (3rd author), Aug, 2019
- colorectal: January 2020 (with Dutch collaborators - J. of Clin. Oncology)
- 1st author methods paper: Communications Biology (April 2022)

Renal carcinoma, Bayesian estimate





Figure 7. Mathematical Modeling of Clear Cell Renal Cell Carcinoma Evolution

(A) Schematic depicting how the age of incidence of renal cell carcinoma may be modeled as the sam of waiting times. Z, representing the time to Z, and the cell cards (representing the time to Z) and the cell cards (representing the time to Z) and the cell cards (representing the time to Z) and the cell cards (representing the time to Z) and the cell cards (representing the time to Z). The time to the time to Z) and the cell cards (representing the time to Z) and the time to Z) and the time to Z) and the time to Z). The time to C and the time to Z) and the time to Z). The time to C and the time to Z) and the time to Z). The time to C and the time to Z) and the time to Z). The time to C and the time to Z) are card to Z and the time to Z) and the time to Z) and the time to Z) and the time to Z). The time to Z and the time to Z) and the time to Z) are card to Z and the time to Z) and the time to Z) and the time to Z). The time to Z and the time to Z) are card to Z and the time to Z) and the time to Z) and the time to Z) and the time to Z). The time to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z). The time to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z) are card to Z and the time to Z are the time to Z are card to Z and the time to Z are the time to Z are card to Z and the time to Z are the time to Z are

Stream 4. single cell RNAseq (scRNA)

Ability to interrogate expression at a single cell level, but ... at major cost to depth

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● ● ● ●

scRNA-seq pig mucosal immunity - all runs



Sample

pig immunity - cell type identification



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへ(?)

Differential expression



Positive Reg of Innate Immune Rensp.

Reg of Defense Response to Virus

э

mucosal immunity - scRNA clusters and cells



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─のへで

generate rich datasets

some loss of depth

some good tools exist

but it is not as straight forward as just applying an R package

truly multi-disciplinary: interactions are necessary

truly multi-disciplinary: appreciation lacking

scientific vision of computational biology

- work in close collaboration with experimental and clinical groups
- get involved early in the formulation of scientific question

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ ―臣 … のへで

- iterative, refining process, (forming a common language)
- properly resource analysis and computational tasks
- data management life cycle
- policies on data and analytics
- harmonious development and use of resources
- simple and robust solutions

Thank you