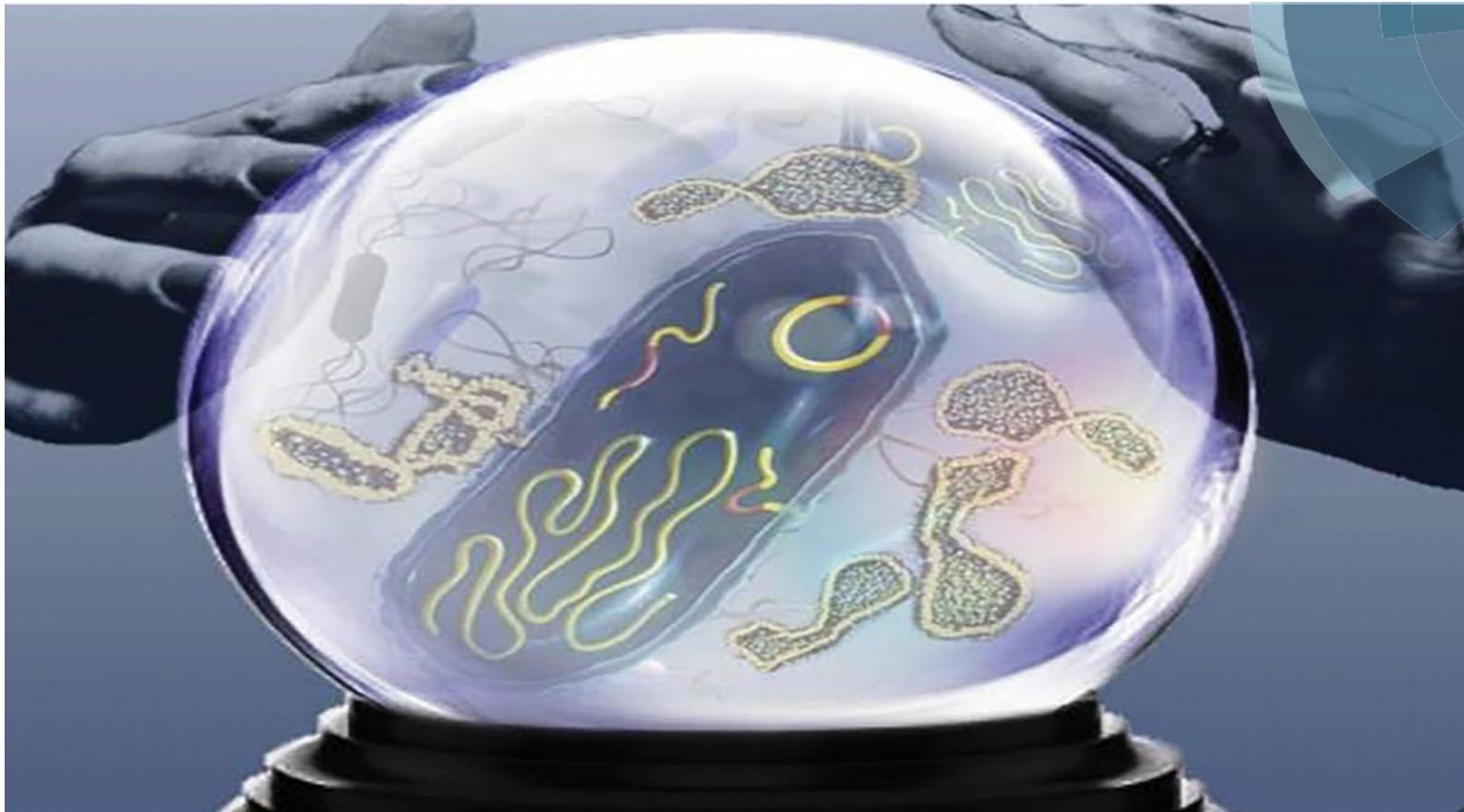


Machine Learning View of OMICs Integration

OMICs Integration and Systems Biology course

Nikolay Oskolkov, NBIS SciLifeLab

Lund, 5.10.2020

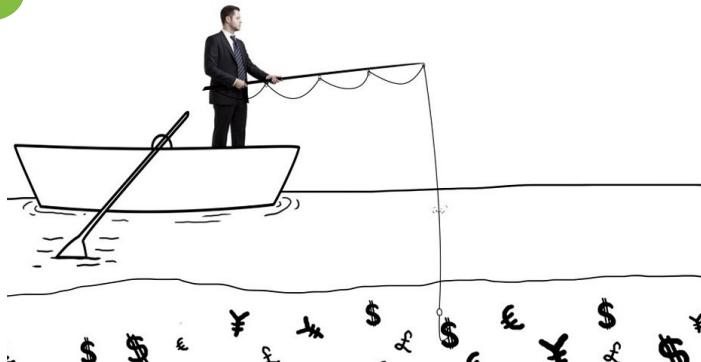




Rätt är rätt och
snett är snett

Peculiarities of OMICs Integration

Fishing expedition



BMC Part of Springer Nature

Search

Genome Biology

Home About Articles Submission Guidelines

Editorial | Open Access | Published: 03 September 2020

A hypothesis is a liability

Itai Yanai & Martin Lercher

Genome Biology 21, Article number: 231 (2020) | [cite this article](#)

12k Accesses | 619 Altmetric | [Metrics](#)

"'When someone seeks,' said Siddhartha, 'then it easily happens that his eyes see only the thing that he seeks, and he is able to find nothing, to take in nothing. [...] Seeking means: having a goal. But finding means: being free, being open, having no goal.' " Hermann Hesse

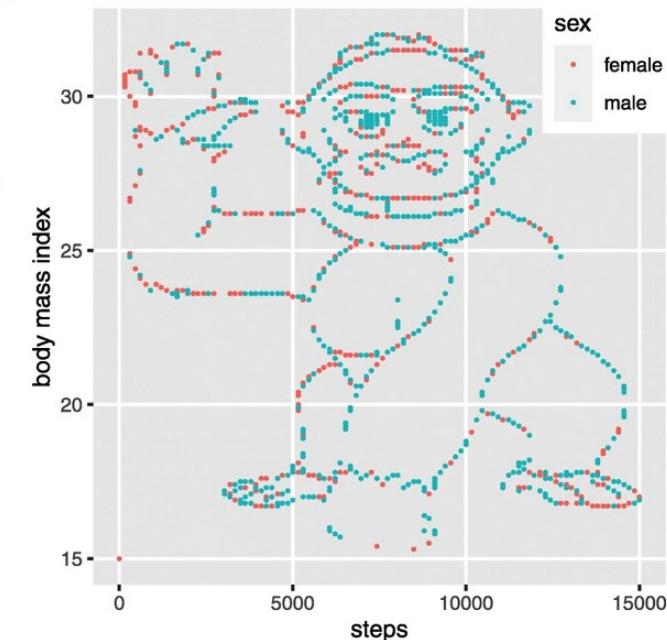
There is a hidden cost to having a hypothesis. It arises from the relationship between night science and day science, the two very distinct modes of activity in which scientific ideas are generated and tested, respectively [1, 2]. With a hypothesis in hand, the impressive strengths of day science are unleashed, guiding us in designing tests, estimating parameters, and throwing out the hypothesis if it fails the tests. But when we analyze the results of an experiment, our mental focus on a specific hypothesis can prevent us from exploring other aspects of the data, effectively blinding us to new ideas. A hypothesis then becomes a liability for any night science explorations. The corresponding limitations on our creativity, self-imposed in hypothesis-driven research, are of particular concern in the context of modern biological datasets, which are often vast and likely to contain hints at multiple distinct and potentially exciting discoveries. Night science has its own liability though, generating many spurious relationships and false hypotheses. Fortunately, these are exposed by the light of day science, emphasizing the complementarity of the two modes, where each overcomes the

- I do not understand your biological hypothesis
- I do not have any
- Then I reject your manuscript

a

ID	steps	bmi
3	15000	17.8
4	14861	17.2
5		
14		
15	15000	16.9
16	15000	16.9
21	6	14861
23	7	14861
26	8	14699
28	10	14560
31	11	14560
33	13	14560
34	17	14560
35	18	14560
36	19	14560
38	20	14560
39	22	14560
41	24	14560
44	25	14560
45	27	14560
29	14560	19.6
30	14560	19.8
32	14398	20.9
37	14398	17.5
40	14398	17.1
42	14259	21.1
43	14259	21.1
44	14259	19.0

b



c

	Gorilla not discovered	Gorilla discovered
Hypothesis-focused	14	5
Hypothesis-free	5	9

a An artificial dataset given to students with and without explicit hypotheses on the relationship between BMI and the steps taken on a particular day, for men and women. b A plot of the dataset. c The contingency table for students in the two groups ("hypothesis-focused," "hypothesis-free") that discovered the gorilla or not [6]

$$N \left(\begin{array}{cccccccccc} 0 & 3 & 1 & 0 & 2 & 3 & 8 & 1 & 1 & 3 \\ 1 & 1 & 0 & 0 & 7 & 1 & 2 & 2 & 3 & 3 \\ 1 & 2 & 2 & 0 & 0 & 6 & 7 & 1 & 2 & 2 \\ 1 & 2 & 3 & 10 & 0 & 4 & 6 & 1 & 0 & 5 \\ 3 & 2 & 2 & 1 & 4 & 3 & 2 & 1 & 6 & 0 \\ 7 & 4 & 4 & 5 & 3 & 9 & 6 & 1 & 6 & 1 \\ 7 & 1 & 1 & 5 & 2 & 8 & 9 & 1 & 3 & 6 \\ 5 & 0 & 1 & 6 & 2 & 0 & 0 & 0 & 1 & 5 \\ 1 & 6 & 3 & 3 & 4 & 6 & 2 & 0 & 1 & 1 \\ 1 & 2 & 2 & 4 & 1 & 1 & 3 & 0 & 8 & 2 \end{array} \right) P_1$$

OMIC1

$$N \left(\begin{array}{cccccccccc} 0 & 3 & 1 & 0 & 2 & 3 & 8 & 1 & 1 & 3 \\ 1 & 1 & 0 & 0 & 7 & 1 & 2 & 2 & 3 & 3 \\ 1 & 2 & 2 & 0 & 0 & 6 & 7 & 1 & 2 & 2 \\ 1 & 2 & 3 & 10 & 0 & 4 & 6 & 1 & 0 & 5 \\ 3 & 2 & 2 & 1 & 4 & 3 & 2 & 1 & 6 & 0 \\ 7 & 4 & 4 & 5 & 3 & 9 & 6 & 1 & 6 & 1 \\ 7 & 1 & 1 & 5 & 2 & 8 & 9 & 1 & 3 & 6 \\ 5 & 0 & 1 & 6 & 2 & 0 & 0 & 0 & 1 & 5 \\ 1 & 6 & 3 & 3 & 4 & 6 & 2 & 0 & 1 & 1 \\ 1 & 2 & 2 & 4 & 1 & 1 & 3 & 0 & 8 & 2 \end{array} \right)$$

OMIC2

N										
0	3	1	0	2	3	8	1	1	3	
1	1	0	0	7	1	2	2	3	3	
1	2	2	0	0	6	7	1	2	2	
1	2	3	10	0	4	6	1	0	5	
3	2	2	1	4	3	2	1	6	0	
7	4	4	5	3	9	6	1	6	1	
7	1	1	5	2	8	9	1	3	6	
5	0	1	6	2	0	0	0	1	5	
1	6	3	3	4	6	2	0	1	1	
1	2	2	4	1	1	3	0	8	2	

OMIC3

$$\begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \end{pmatrix}$$

Metabolomics

N ≈ P

$$\begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \end{pmatrix}$$

Proteomics N ≈ P

– manageable

$$\begin{pmatrix} 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 & 1 & 0 & 1 & 1 \end{pmatrix}$$

Transcriptomics N << P (Single cell: N \leq P)

Genomics

N <<< P

Methylomics N < B

The Curse of Dimensionality complicates OMICs Integration

P is the number of features (genes, proteins, genetic variants etc.)
 N is the number of observations (samples, cells, nucleotides etc.)

Biomedicine

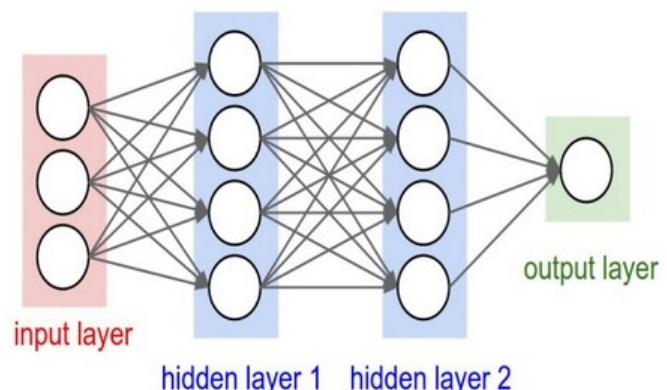
Bayesianism

 $P \gg N$

Frequentism

 $P \sim N$

Deep Learning

 $P \ll N$ 

Amount of Data

Ex.1

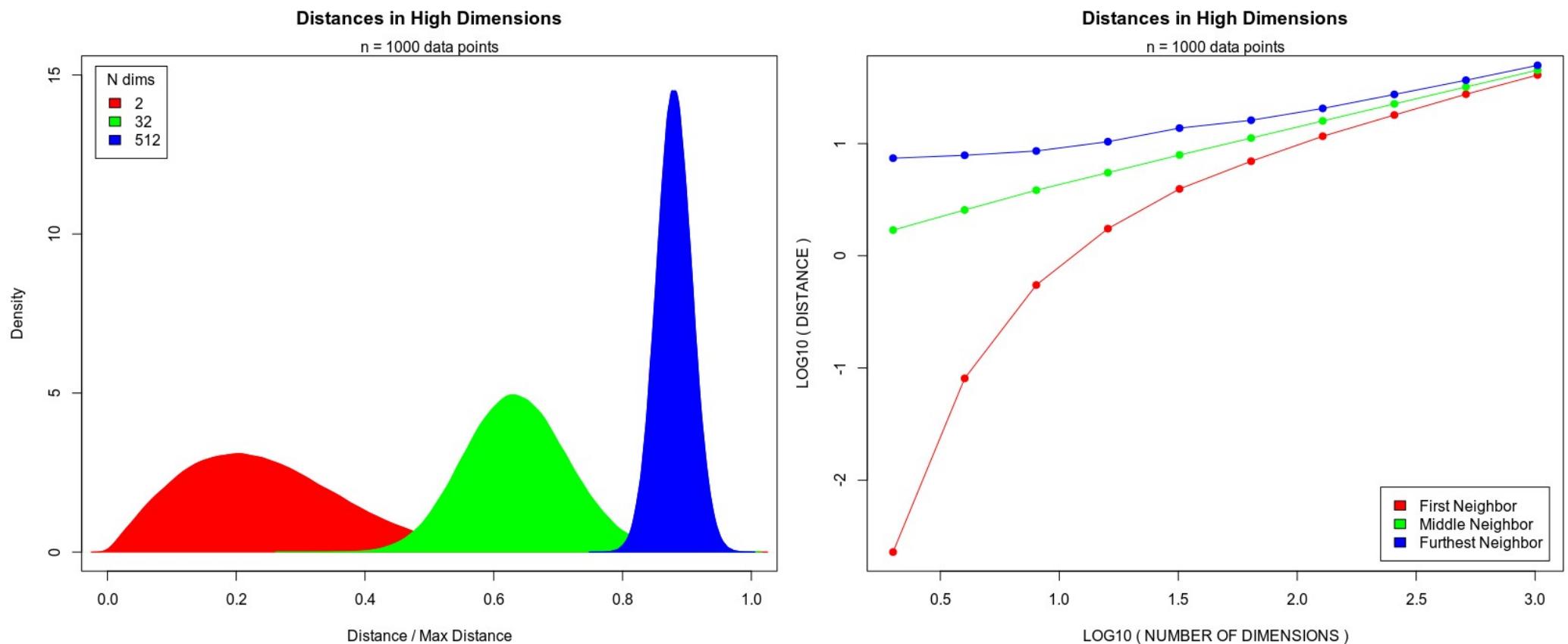
$$Y = \alpha + \beta X$$

$$\beta = (X^T X)^{-1} X^T Y$$

$$(X^T X)^{-1} \sim \frac{1}{\det(X^T X)} \dots \rightarrow \infty, \quad n \ll p$$

$$\text{Ex.2} \quad E[\hat{\sigma}^2] = \frac{n-p}{n} \sigma^2$$

Biased ML variance estimator in HD-space



Data points become far from each other and equidistant from each other in high dimensions

The differences between closest and furthest data point neighbours disappears in high-dimensional spaces – can't cluster

In high-dimensional space we can not separate cases and controls any more

How to define and evaluate OMICs Integration?



Exploration and
Integration of
Omics datasets

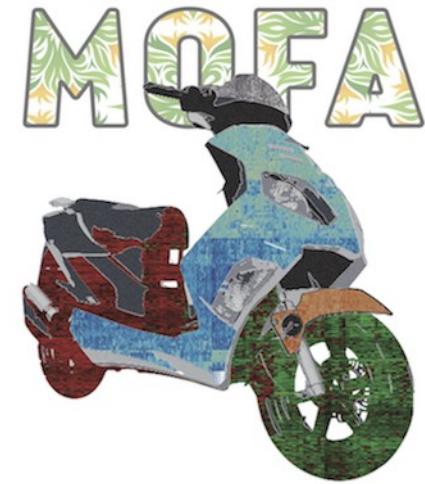
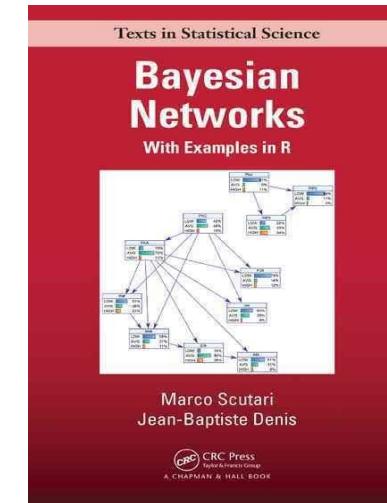
Clustering of Clusters



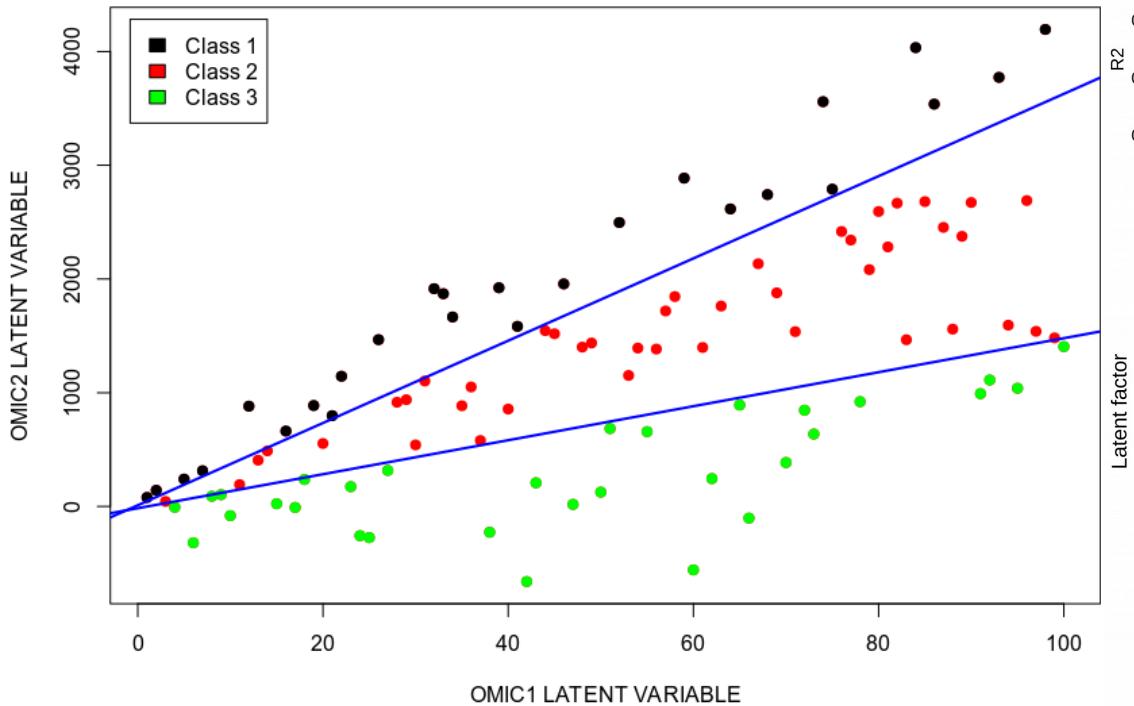
OnPLS

JIVE

DISCO

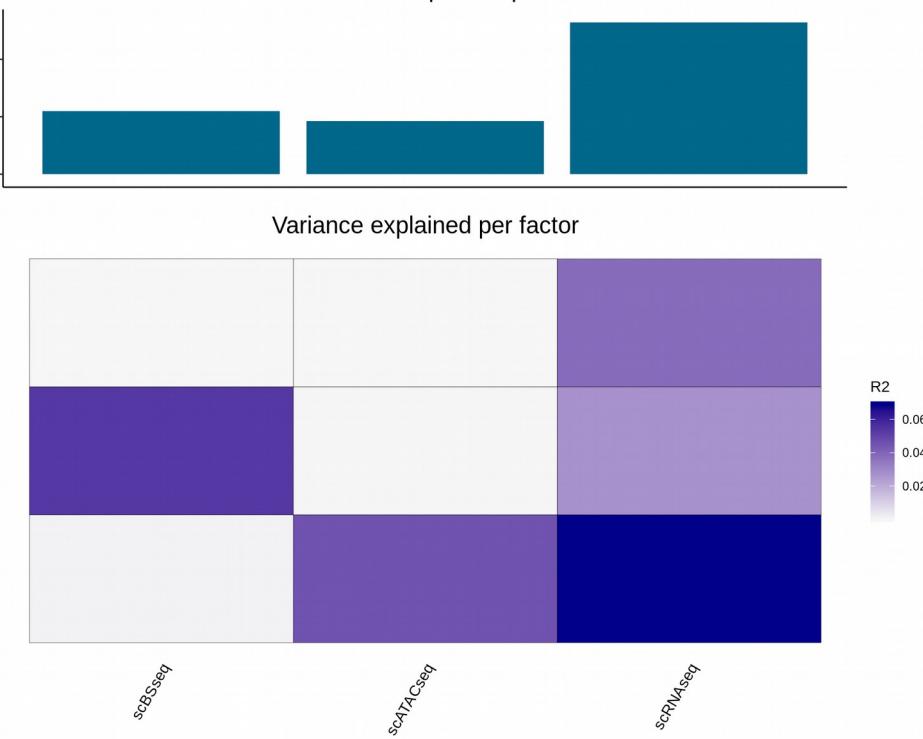


Idea Behind OMICs Integration: See Patterns Hidden in Individual OMICS



Total variance explained per view

Variance explained per factor



How I Evaluate OMICs Integration, Data Science: Boost in Prediction

TEXT (78%)

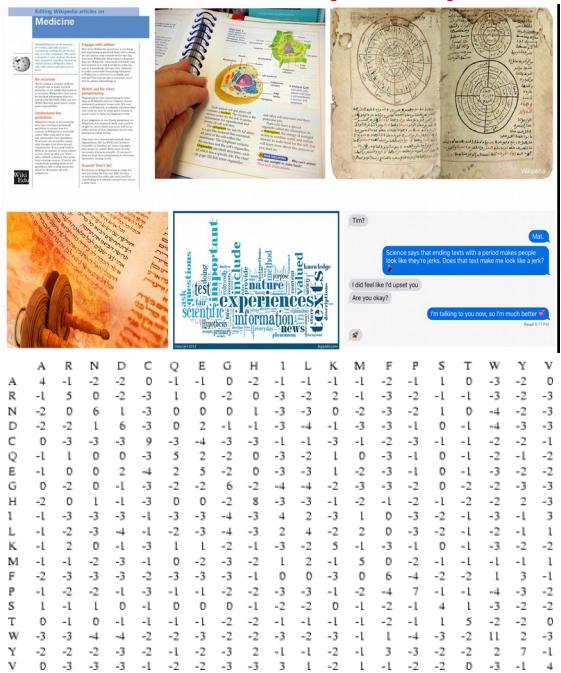
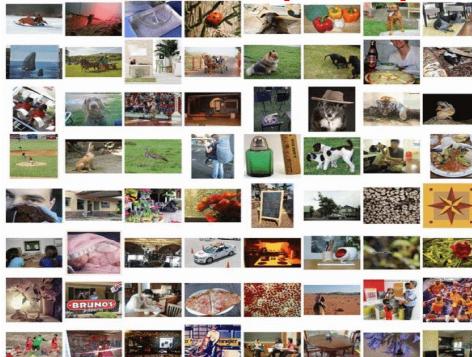
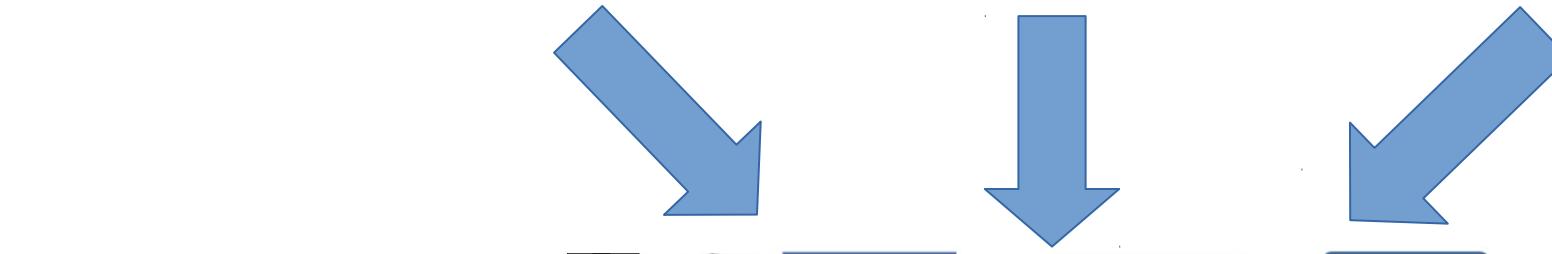
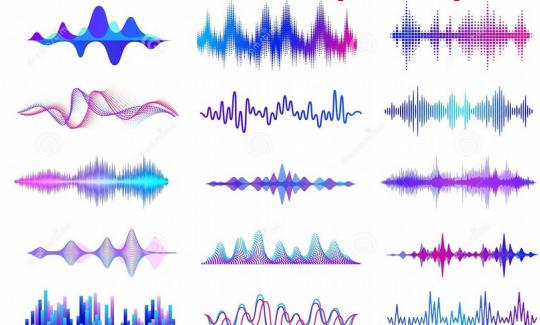


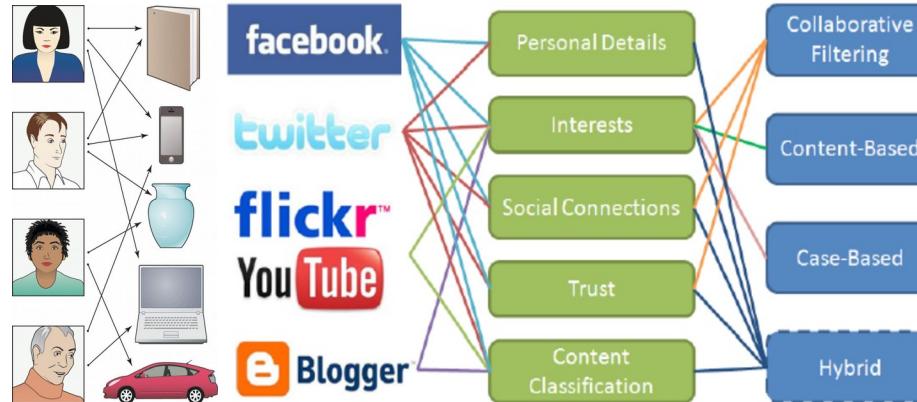
IMAGE (83%)



SOUND (75%)



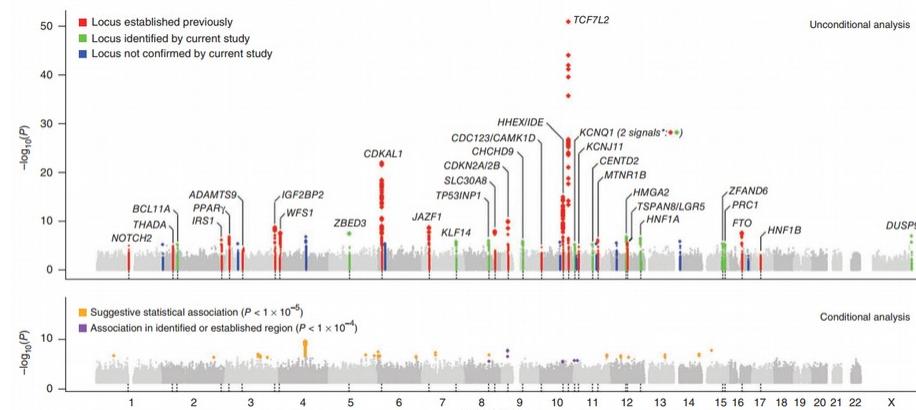
Predict Facebook user interests



Data Integration Accuracy: 96%

Prediction is an Ultimate Criterion of Successful OMICS Integration

Statistics searches for candidates



Consequence



NEWS FEATURE PERSONAL GENOMES NATURE/Vol 456/6 November 2008



The case of the missing heritability

Machine Learning optimizes prediction



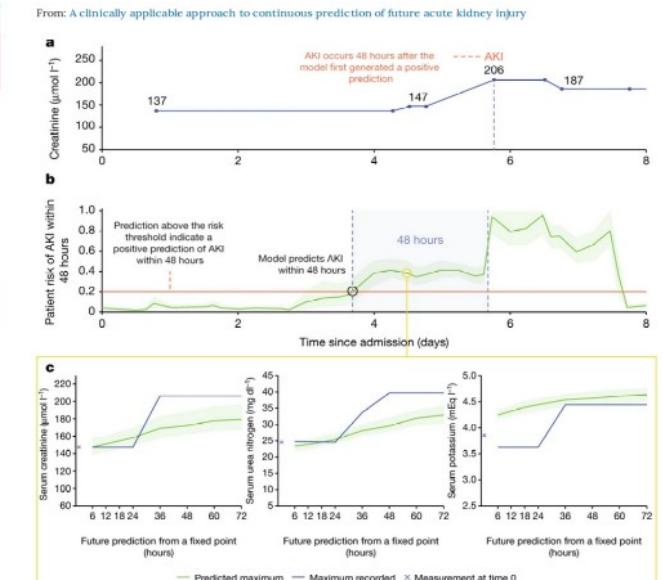
A clinically applicable approach to continuous prediction of future acute kidney injury

Nenad Tomasev, Xavier Glorot, [...] Shakir Mohamed

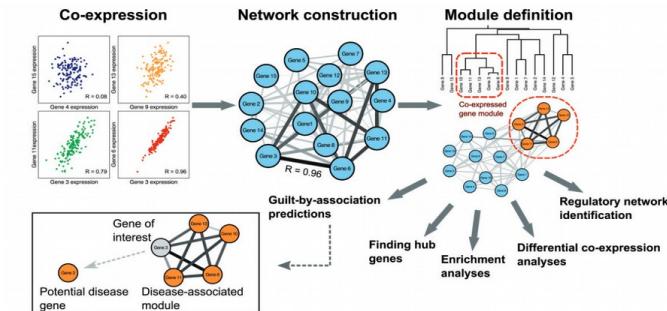
Nature 572, 116–119 (2019) Download Citation ↗

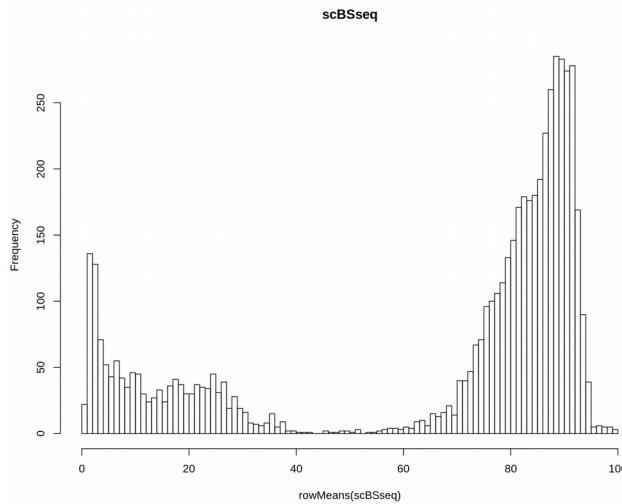
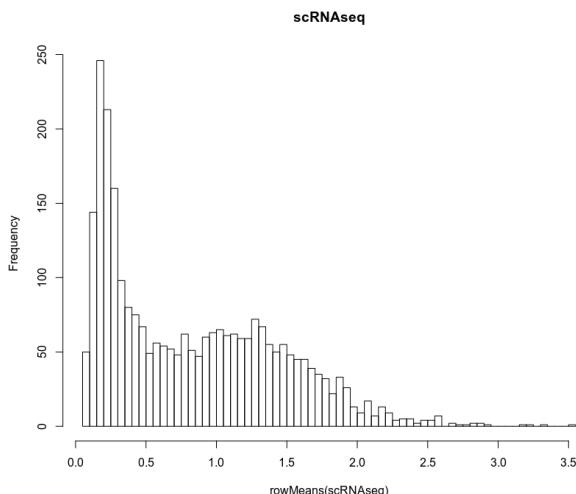
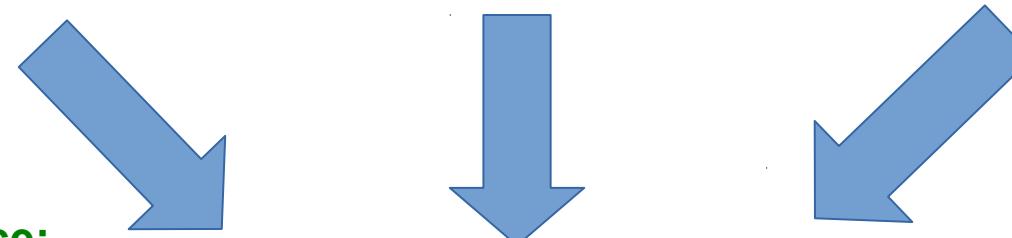
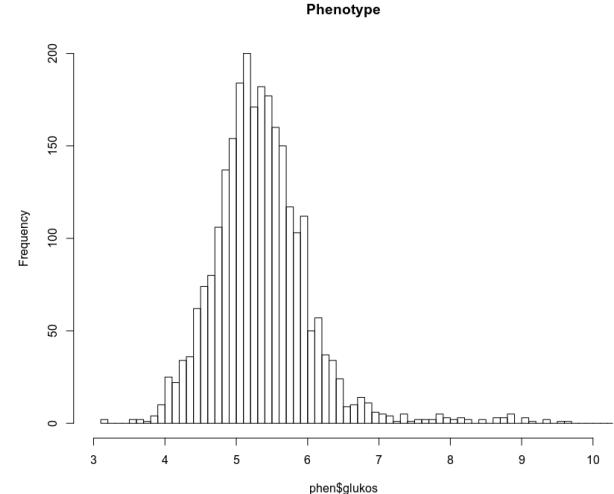
Abstract

The early prediction of deterioration could have an important role in supporting healthcare professionals, as an estimated 11% of deaths in hospital follow a failure to promptly recognize and treat deteriorating patients¹. To achieve this goal requires predictions of patient risk that are continuously updated and accurate, and delivered at an individual level with sufficient context and enough time to act. Here we develop a deep learning approach for the continuous risk prediction of future deterioration in patients, building on recent work that models adverse events from electronic health records^{2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17} and using acute kidney injury—a common and potentially life-threatening condition¹⁸—as an exemplar. Our model was developed on a large, longitudinal dataset of electronic health records that cover diverse



Consequence



Methylation (78%)**Gene Expression (83%)****Phenotype (75%)****1) Convert to common space:**

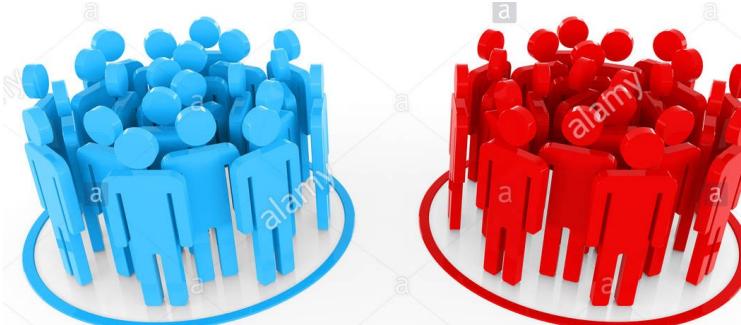
Neural Networks, SNF, UMAP

2) Explicitly model distributions:

MOFA, Bayesian Networks

3) Extract common variation:

PLS, CCA, Factor Analysis

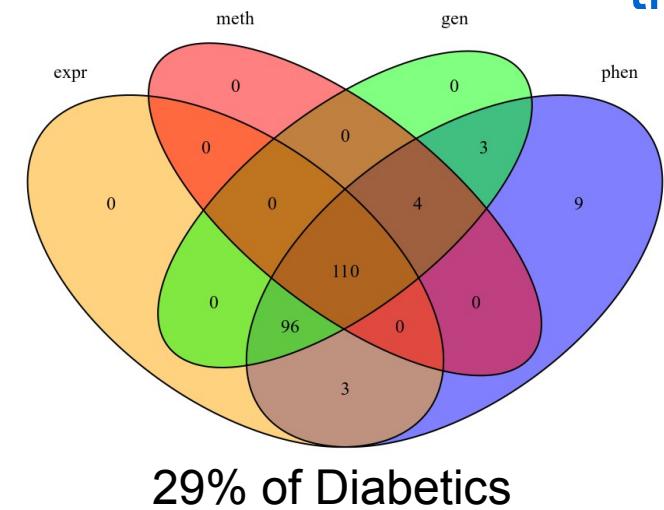
**HEALTHY****SICK****Data Integration
Accuracy: 96%**

Protocol of OMICs Integration

	Linear	Non-Linear
Supervised	PLS / OPLS / mixOmics, LASSO / Ridge / Elastic Net	Neural Networks, Random Forest, Bayesian Networks
Unsupervised	Factor Analysis / MOFA	Autoencoder, SNF, UMAP, Clustering of Clusters

For Example:

- 1) With ~110 samples it is a good idea to do **linear** OMICs integration
- 2) T2D is a phenotype of interest, therefore **supervised** integration



29% of Diabetics

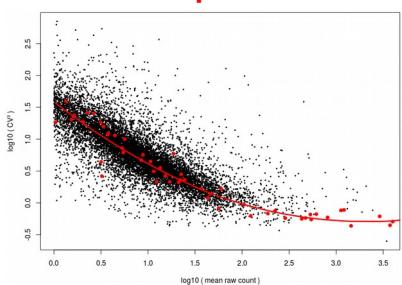
Data Set (4 OMICs)

~~WGS (~30 mln dims)
BSseq (~30 mln dims)~~

Train Set (n = 88)

Test Set (n = 22)

unsupervised



supervised



Feature Pre-Selection

Evaluation

OMICs Integration

Trained Model



- 1) Check that there is a relation between the OMICs (MOFA)
- 2) Choose integrative model based on amount of data and goal (linear, supervised)
- 3) Do feature pre-selection (supervised or unsupervised) on train data set
- 4) Integrate the OMICs using your favorite model chosen in 2) on train data set
- 5) Compare prediction of integrative model with predictions from individual OMICs

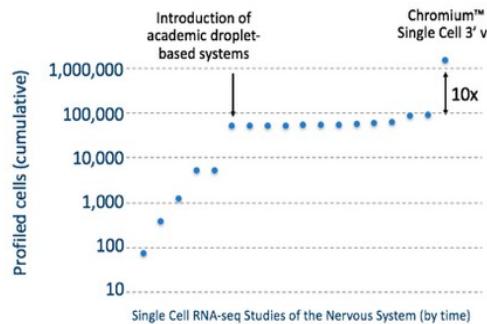
OMICs Integration in Single Cell

CAREERS BLOG 10X UNIVERSITY

10X GENOMICS SOLUTIONS & PRODUCTS RESEARCH & APPLICATIONS EDUCATION & RESOURCES

< Back to Blog

< Newer Article Older Article >



Our 1.3 million single cell dataset is ready 0 KUDOS



POSTED BY: grace-10x, on Feb 21, 2017 at 2:28 PM

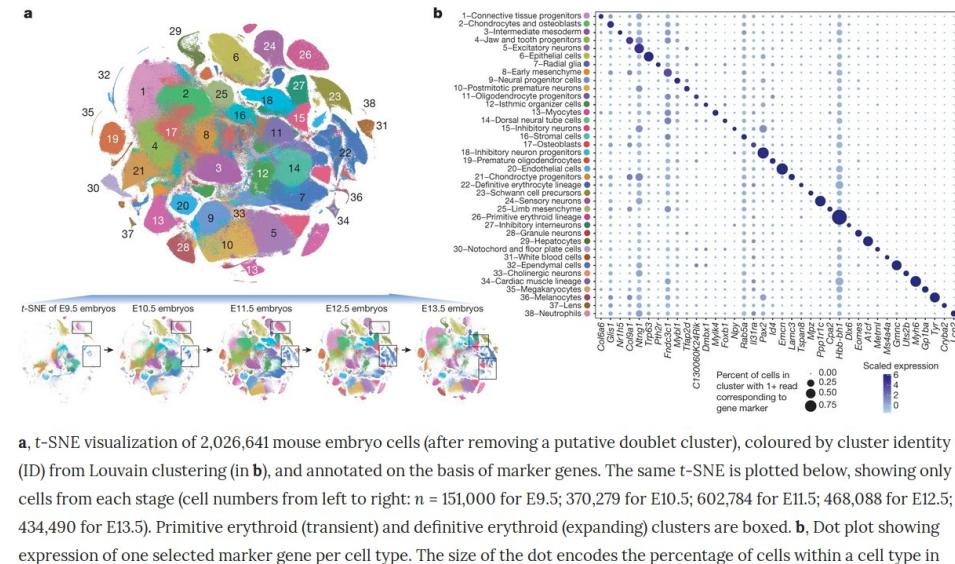
At ASHG last year, we announced our 1.3 Million Brain Cell Dataset, which is, to date, the largest dataset published in the single cell RNA-sequencing (scRNA-seq) field. Using the Chromium™ Single Cell 3' Solution (v2 Chemistry), we were able to sequence and profile 1,308,421 individual cells from embryonic mice brains. Read more in our application note [Transcriptional Profiling of 1.3 Million Brain Cells with the Chromium™ Single Cell 3' Solution](#).

**Watch out Underfitting!
Paradise for Deep Learning!**

MENU nature

Fig. 2: Identifying the major cell types of mouse organogenesis.

From: [The single-cell transcriptional landscape of mammalian organogenesis](#)



BioTuring™ Solutions Resources

Explore 4,000,000 CELLS at ease with BIOTURING BROWSER A next-generation platform to re-analyze published single-cell sequencing data

EXPLORER NOW

Single Cell Analysis

5,500,000 cells will be indexed into BioTuring Single-cell Data Repository this September

by biomembers • August 30, 2019

Human Cell Atlas, single-cell data

We are glad to announce that we will upsize the current single-cell database in BioTuring Single-cell Browser to 5,500,000 cells this September. With this release, we will double the current number of publications indexed in BioTuring Single-cell Browser, and cross the number of cells hosted on available public single-cell data repositories like [Human Cell Atlas \(HCA\)](#) and [Broad Institute's Single-cell Portal](#).

Search

RECENT POSTS

A new tool to interactively visualize single-cell objects (Seurat, Scanpy, SingleCellExperiments, ...)
September 26, 2019

5,500,000 cells will be indexed into BioTuring Single-cell Data Repository this September
August 30, 2019



National Bioinformatics Infrastructure Sweden (NBIS)



*Knut och Alice
Wallenbergs
Stiftelse*



LUNDS
UNIVERSITET