

Multi-Omics Data Integration via Machine Learning

Oslo Bioinformatics Workshop Week 2022, Oslo, Norway, 9.12.2022

Nikolay Oskolkov, Lund University, NBIS SciLifeLab, Sweden

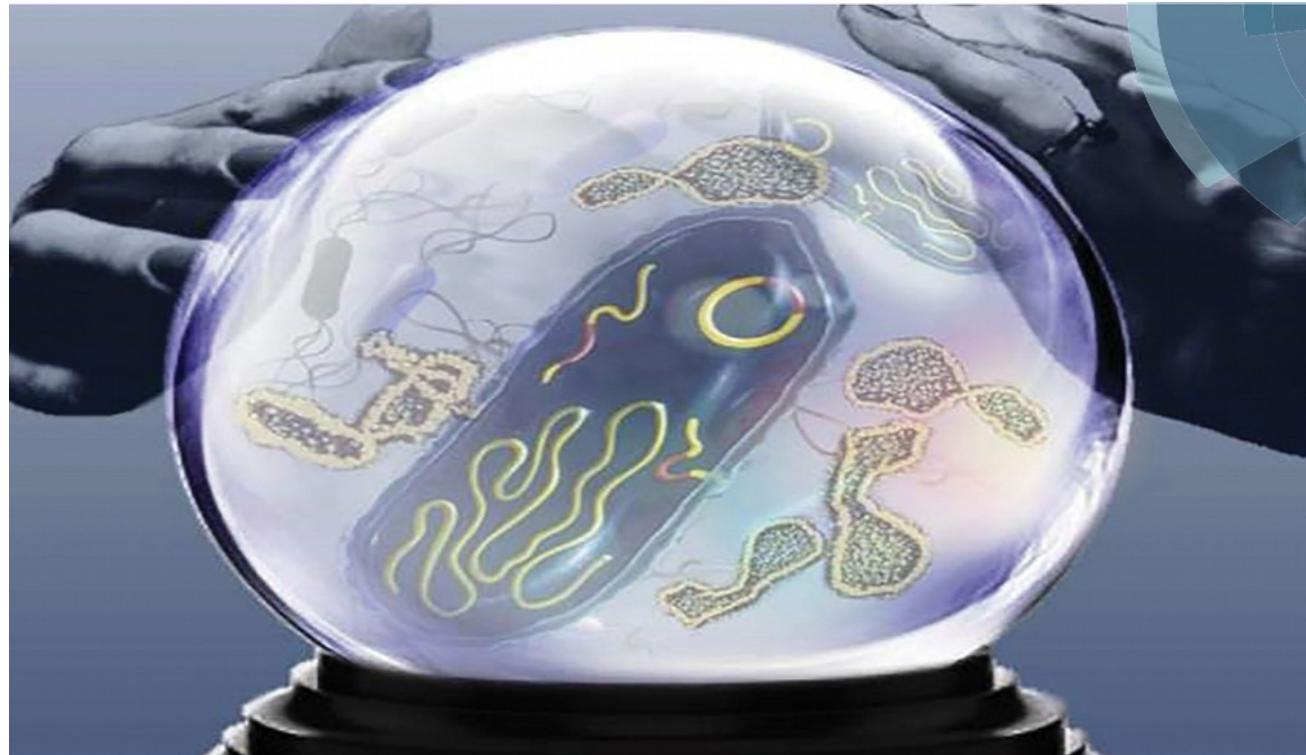
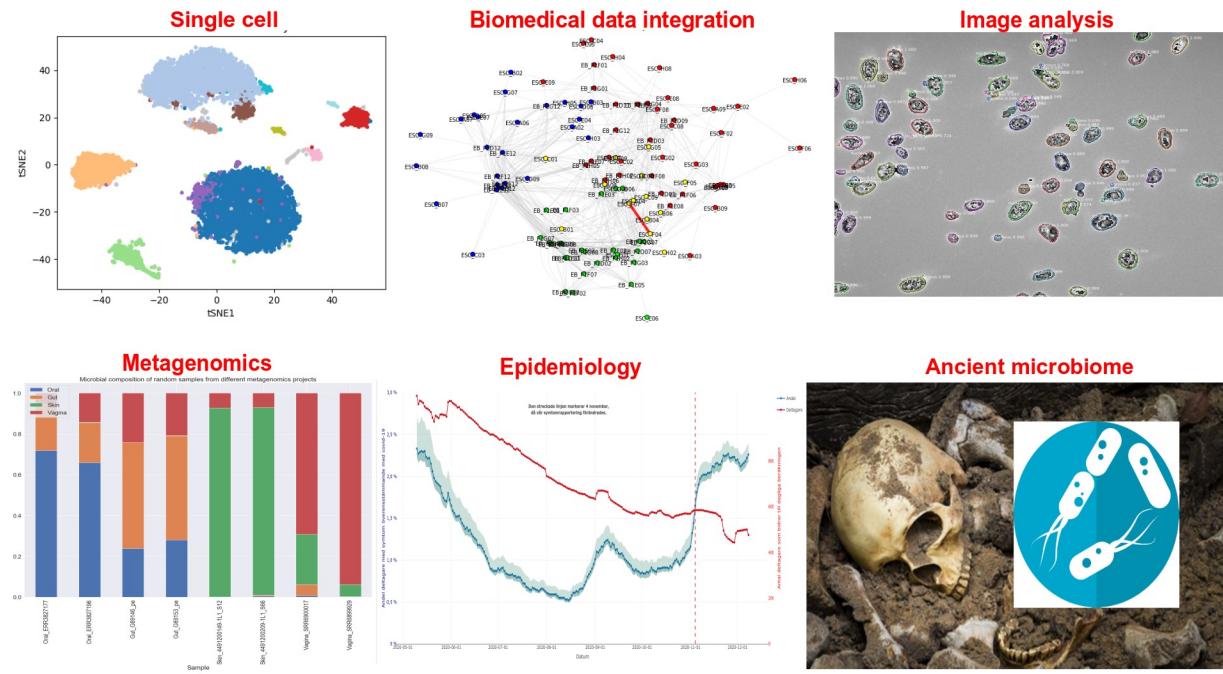
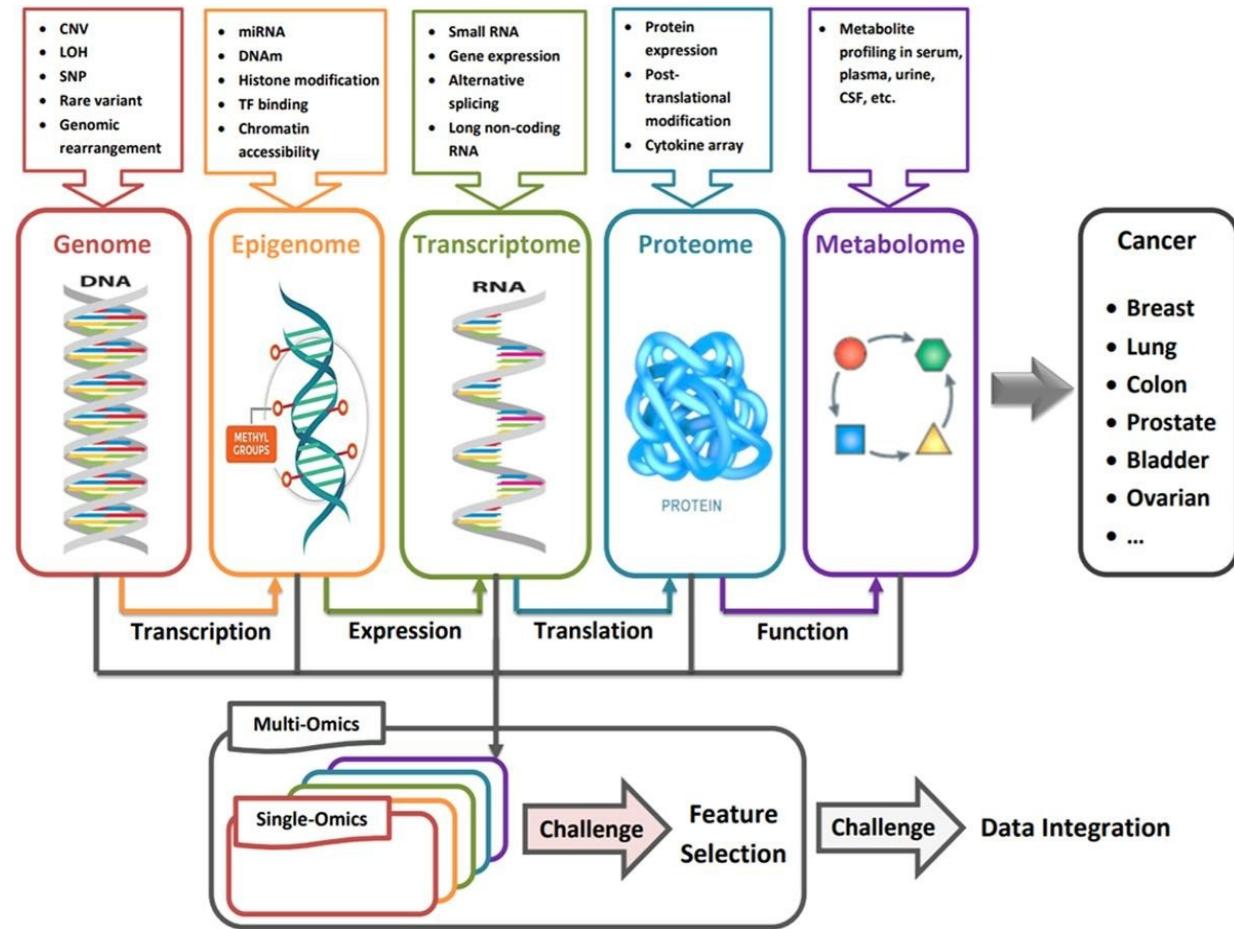


Image adapted from Molecular Omics, Issue 1, 2018

- 2007 PhD in theoretical physics
- 2011 medical genetics at Lund University
- 2016 working at NBIS SciLifeLab, Sweden





Syllabus: ELIXIRSE_OMICSENT_H21 > Syllabus

Workshop Overview: Aimed at providing an integrated view of data-driven hypothesis generation through biological network analysis, constraint-based modelling, and supervised and unsupervised integration methods. General descriptions of different methods for analysing different omics data (e.g. transcriptomics and genomics) will be presented with some of the lectures discussing key methods and pitfalls in their integration. The techniques will be discussed in terms of their rationale and applicability.

Covered topics:

- Data pre-processing and cleaning prior to integration;
- Application of machine learning methods for multi-omics analysis including deep learning;
- Multi-omics integration, clustering and dimensionality reduction;
- Biological network inference, community and topology analysis and visualization;
- Condition-specific and personalized modeling through Genome-scale Metabolic models for integration of transcriptomic, proteomic, metabolomic and fluxomic data;
- Identification of key biological functions and pathways;
- Identification of potential biomarkers and targetable genes through modeling and biological network analysis;
- Application of network approaches in meta-analyses;
- Similarity network fusion and matrix factorization techniques;
- Intricated data visualization techniques;

GitHub Repository: https://github.com/NBISweden/workshop_omics_integration

Contributors: Various contributors listed, including NBISweden, elixirse, and others.

Languages: HTML (67.7%), Python (24.5%), JavaScript (4.5%), and others.

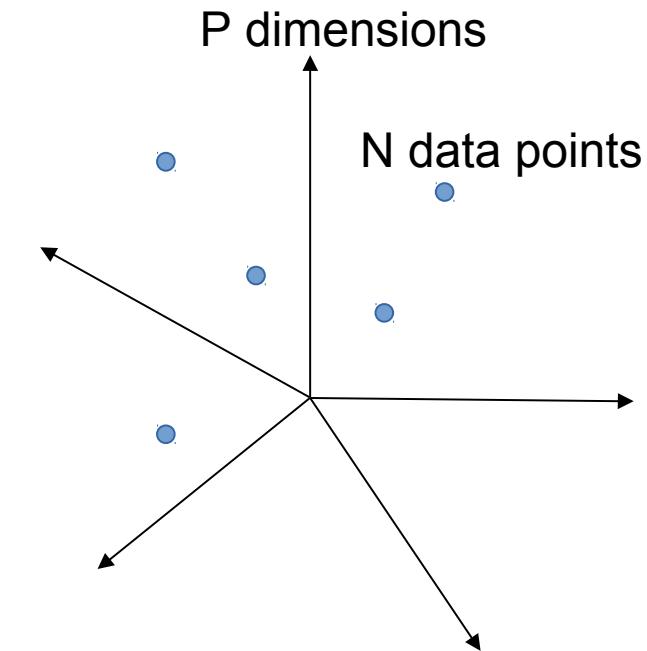
Introduction: High Dimensional Biological Data

Statistical observations:
e.g. samples, cells etc.

Features: genes, proteins,
microbes, metabolites etc.

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



High Dimensional Data:
P >> N

For a robust statistical analysis, one should properly “sample” the P-dimensional space, hence large sample size is required, $N \gg P$

Types of Statistical Analysis

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

Biology / Biomedicine

Bayesianism



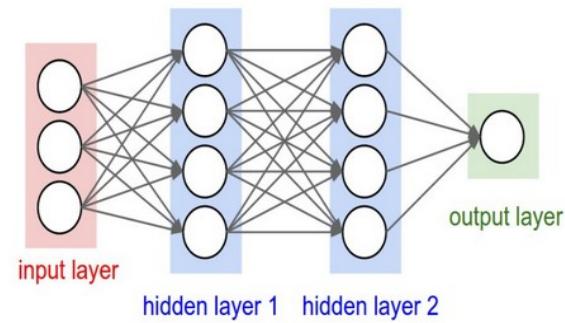
$P \gg N$

Frequentism



$P \sim N$

Deep Learning



$P \ll N$

The Curse of Dimensionality



Amount of Data

Ex.1

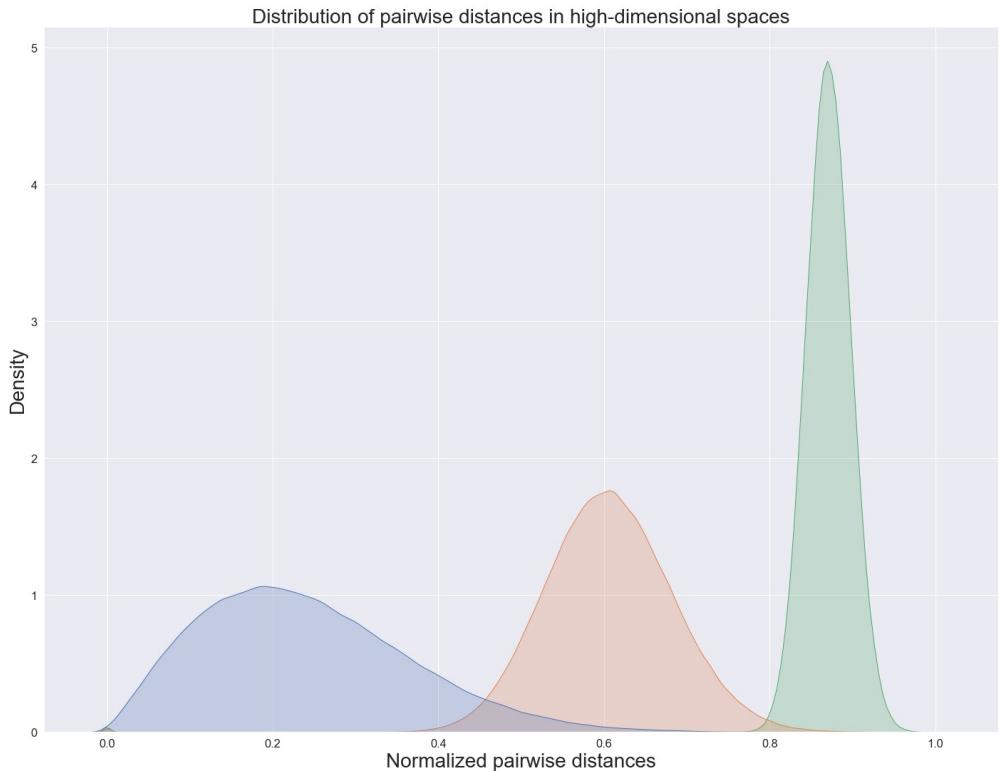
$$Y = \alpha + \beta X$$

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

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

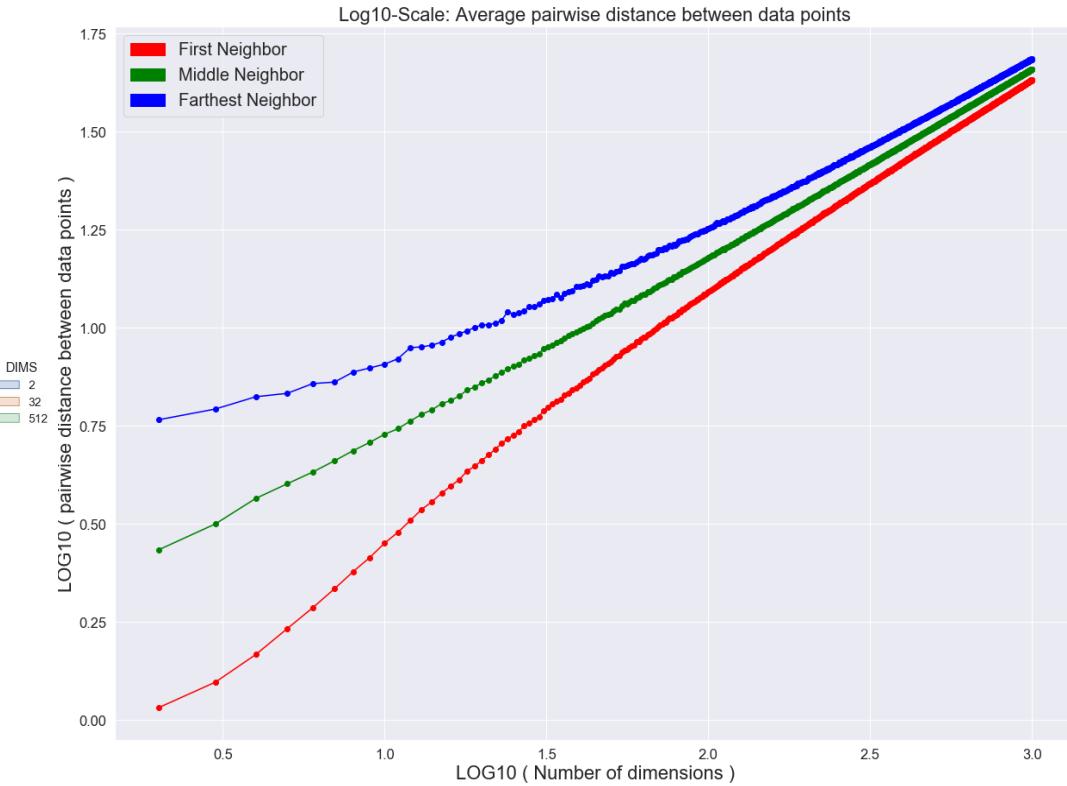
$$\text{Ex.2} \quad E[\hat{\sigma}^2] = \frac{n-p}{n} \sigma^2$$

Biased ML variance estimator in HD-space



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

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



The differences between closest and farthest data point neighbours disappears in high-dimensional spaces:
can't run cluster analysis

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

$N \approx P$

Metabolomics

N ≈ P

$$\begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \end{pmatrix} \quad \text{Proteomics} \quad N \approx P$$

Proteomics N ≈ P

$$\begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \end{pmatrix} \text{Metagenomics} \\ N \approx P$$

Metagenomics

$N \approx P$

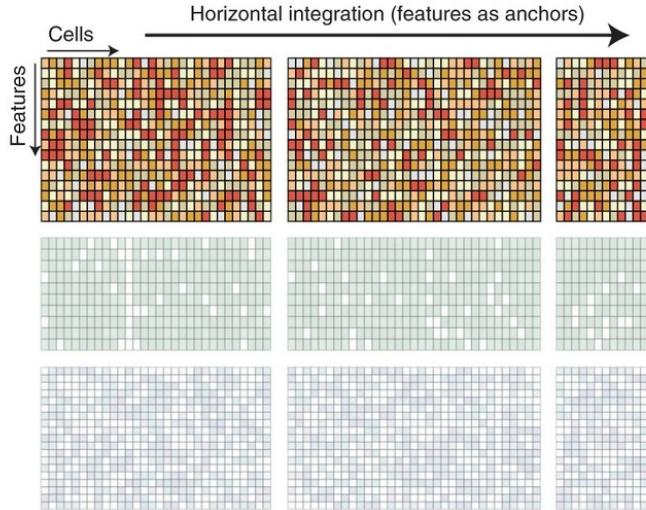
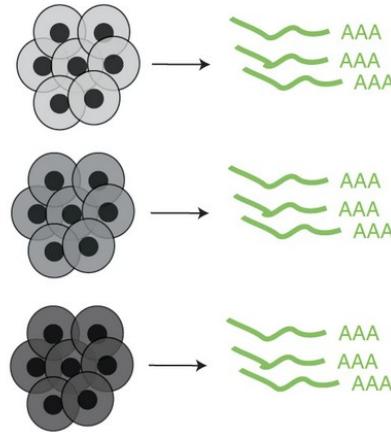
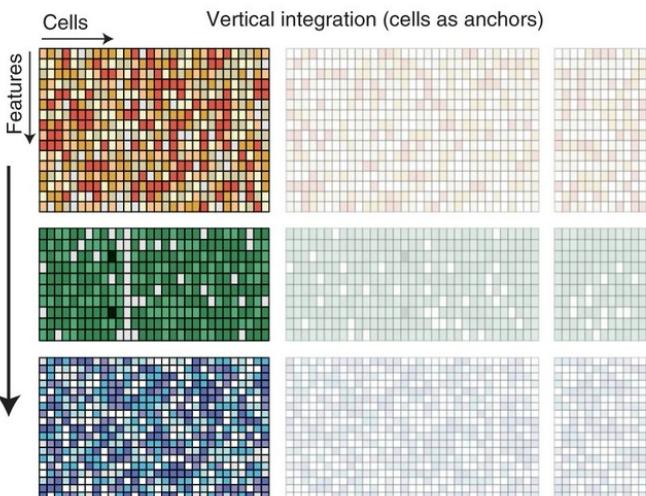
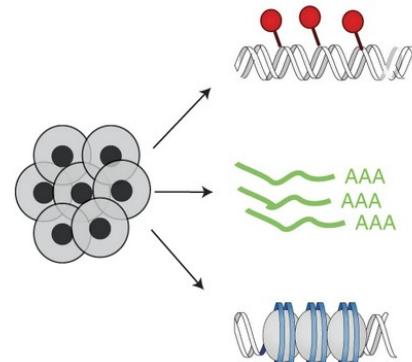
— manageable

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

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

challenging

Multi-Omics Data Integration

a**b**

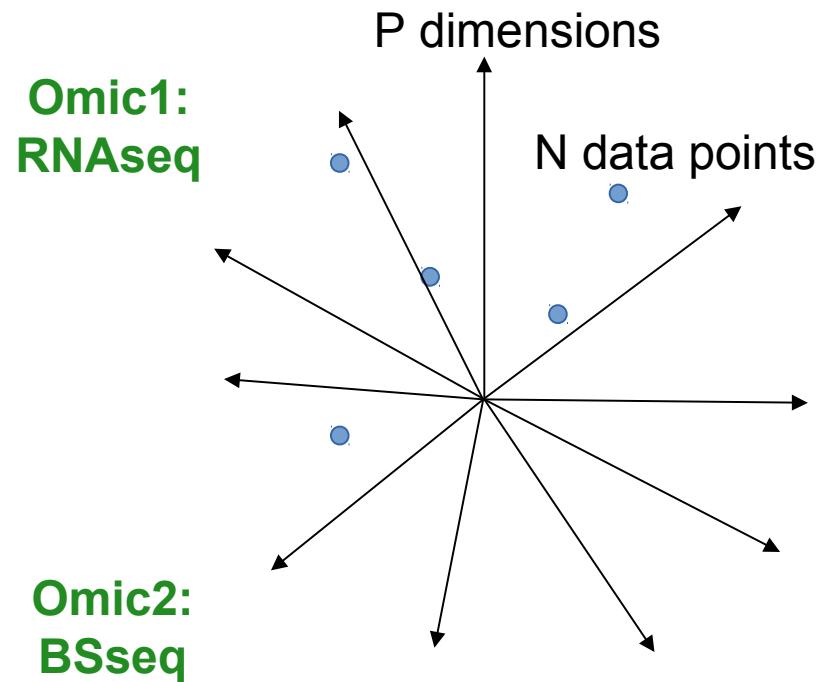
Statistical observations:
e.g. samples, cells etc.

	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	

Features: genes, proteins,
microbes, metabolites etc.

	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	

P₂



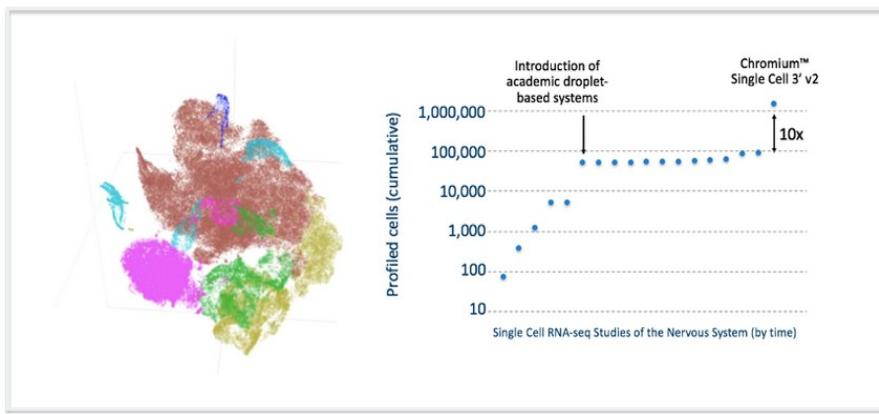
P₁ + P₂ >> N integration across features leads to even more high-dimensional data

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



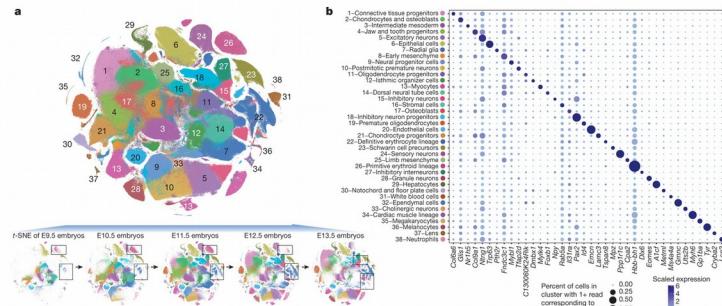
POSTED BY: grace-10x, on Feb 21, 2017 at 2:28 PM

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](#).

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

EXPLORER NOW

A next-generation platform to re-analyze published single-cell sequencing data

Single Cell Analysis

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

by @biomarkers • August 30, 2019

f | t | in

Human Cell Atlas, single-cell data

We are glad to announce that we will upscale 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.

RECENT POSTS

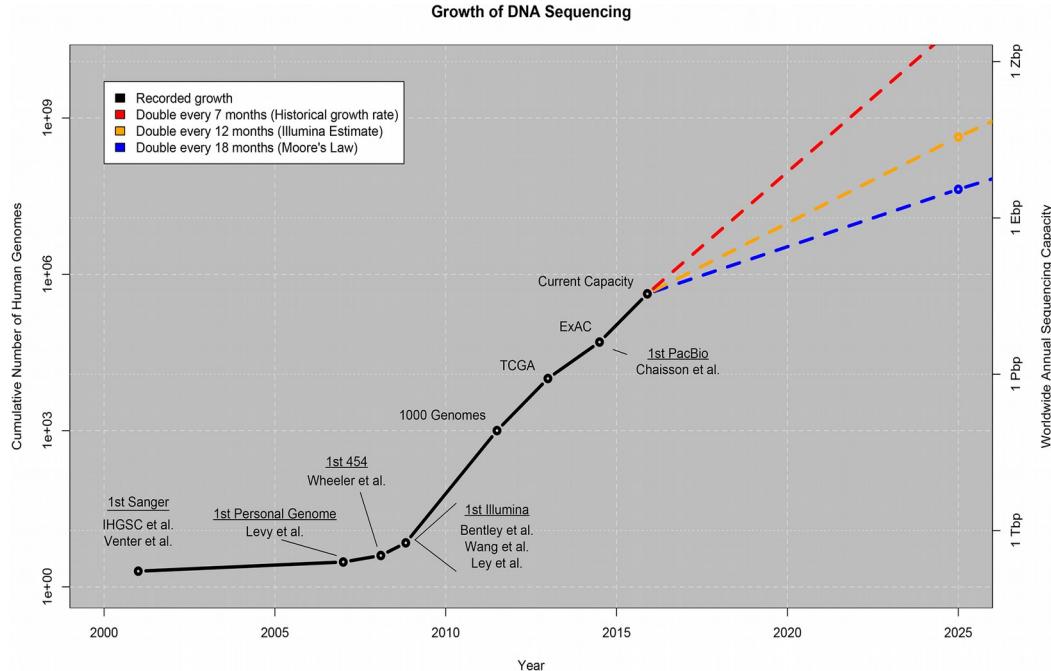
A new tool to interactively visualize single-cell objects (Seurat, Scanpy, SingleCellExperiments,...)

September 26, 2019

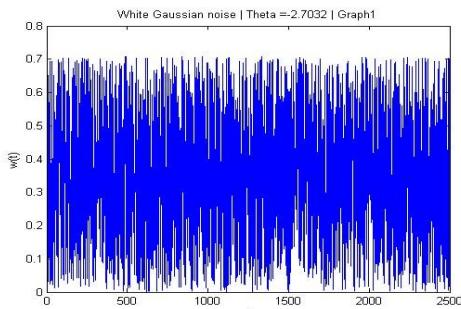
5,500,000 cells will be indexed into BioTuring Single-cell Data Repository this September

August 30, 2019

Big in Size or Sample Size?



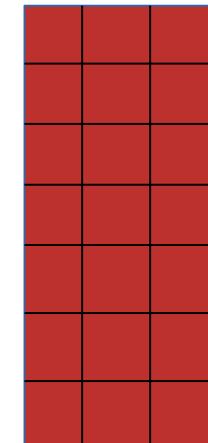
Stephens et al., (2015). Big Data: Astronomical or Genomical? PLoS Biology 13(7)



A file with **White Noise** can also take a lot of disk space

Genomics / WGS: Little Data

$$N_1 \sim 10^3$$

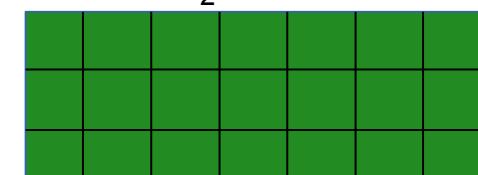


$$P_1 \sim 10^6$$

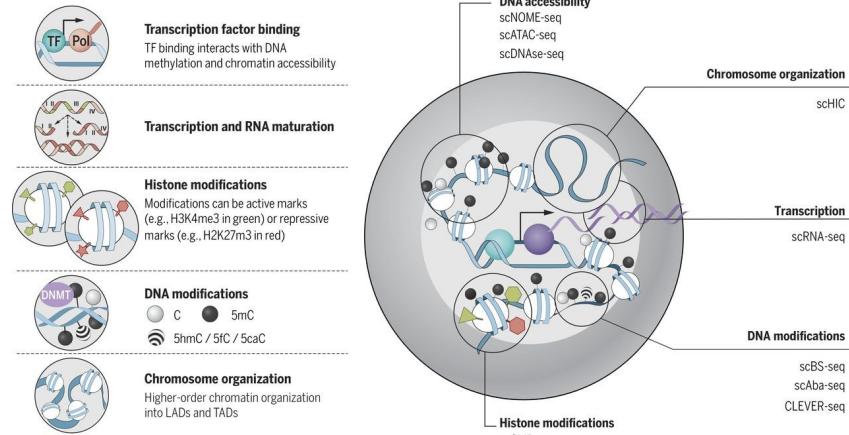
$$N_1 * P_1 = N_2 * P_2 = 10^9$$

scRNAseq: Big Data

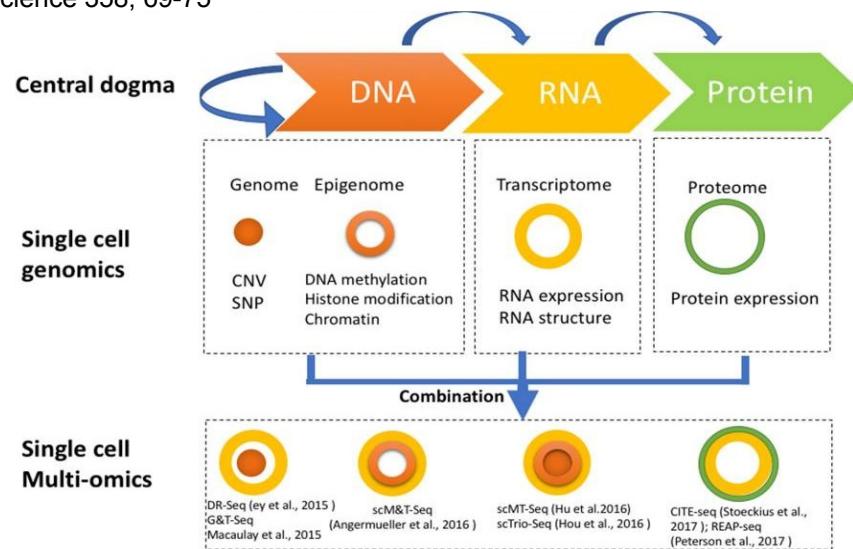
$$N_2 \sim 10^6$$



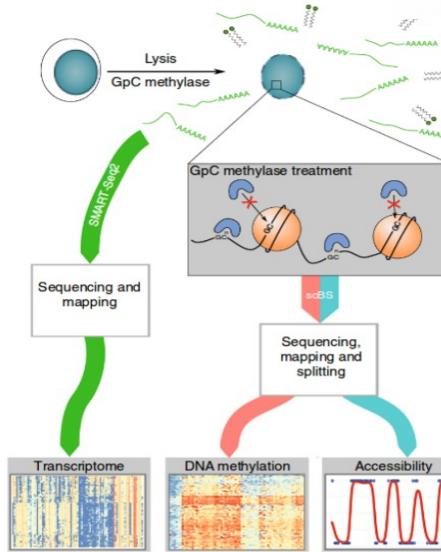
$$P_2 \sim 10^3$$



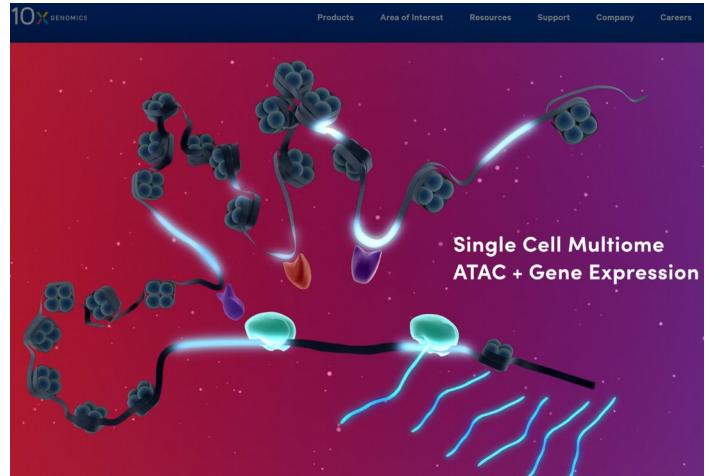
Kelsey et al., 2017, Science 358, 69-75



Hu et al., 2018, Frontier in Cell and Developmental Biology 6, 1-13

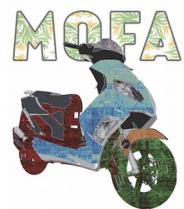
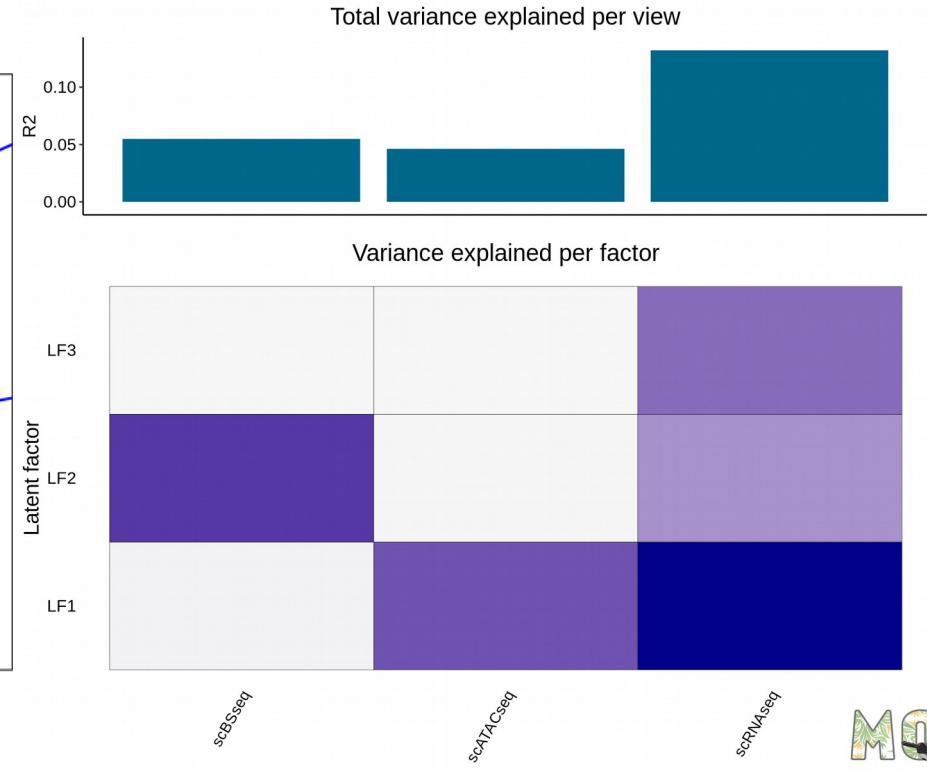
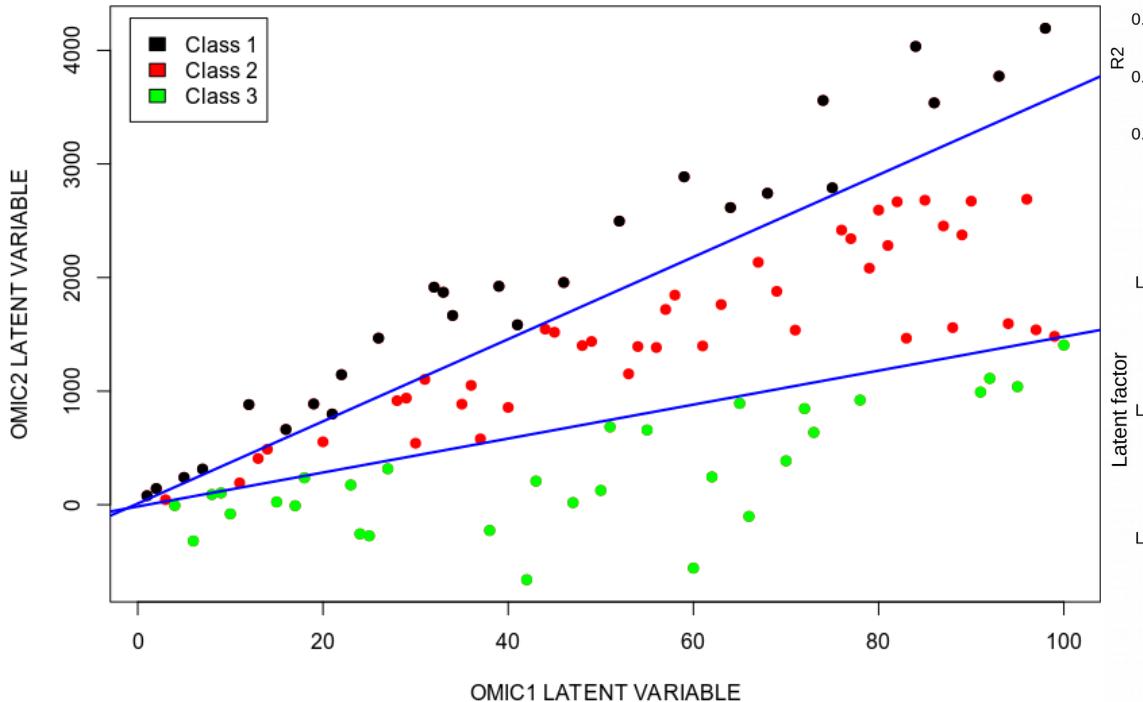


Clark et al., 2018, Nature Communications 9, 781



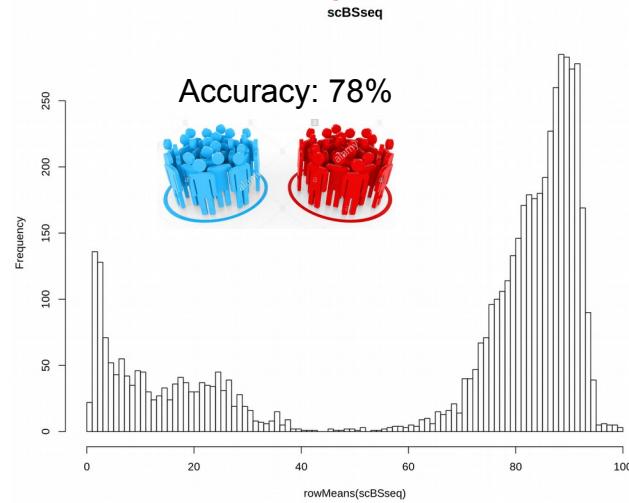
How to define and evaluate multi-Omics data integration?

Idea Behind OMICs Integration:
See Patterns Hidden in Individual OMICS

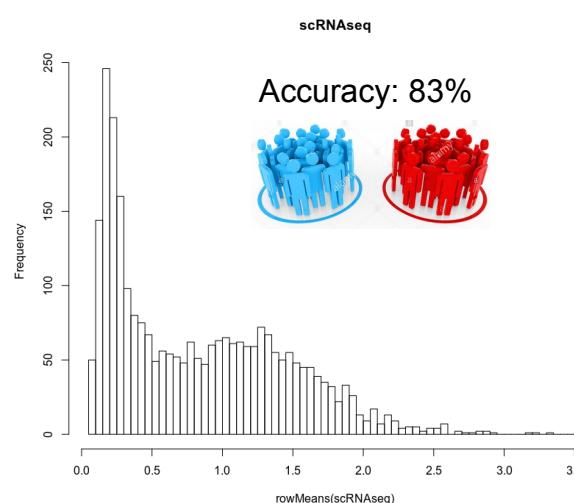


How I Evaluate Omics Integration

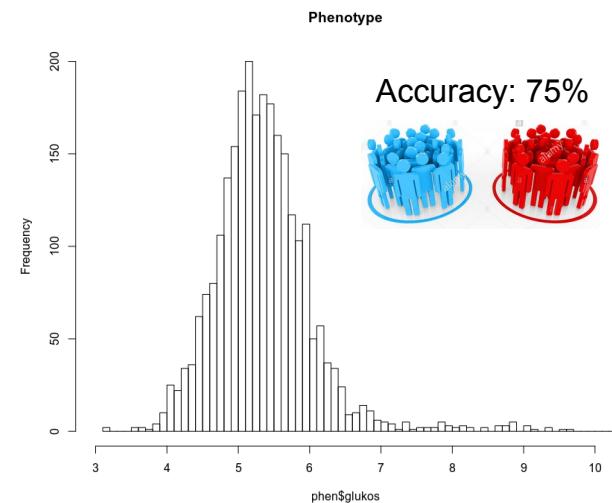
Methylation



Gene Expression



Clinical variable



1) Convert to common space:
Neural Networks, SNF, UMAP

2) Explicitly model distributions:
MOFA, Bayesian Networks

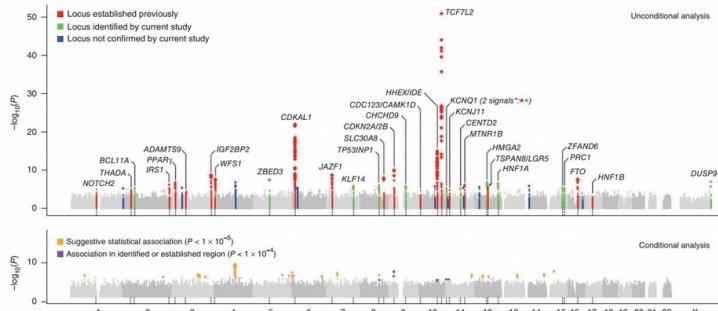
3) Extract common variation:
PLS, CCA, Factor Analysis



Data Integration Accuracy: 96%

Prediction as a Criterion of Success

Statistics searches for candidates



Consequence

NEWS FEATURE PERSONAL GENOMES

NATURE/Vol 456/6 November 2008



The case of the missing heritability

B. Maher, Nature 456, 18-21 (2008)

Machine Learning optimizes prediction



Letter | Published: 31 July 2019

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

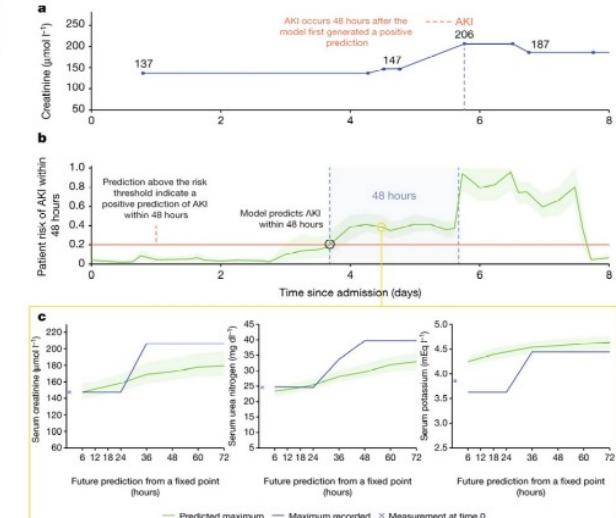
Nenad Tomalec, Kavir Glorot, [...] Shahril Mohamed

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

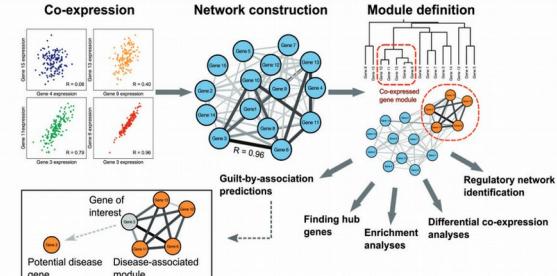
Abstract

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

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



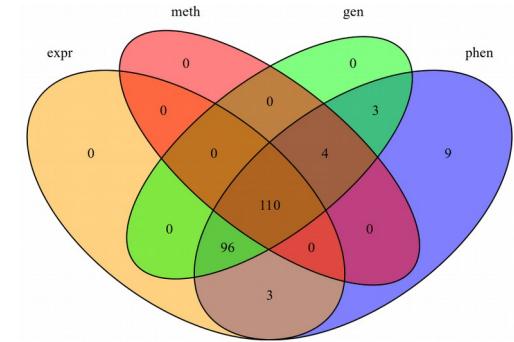
Consequence



	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

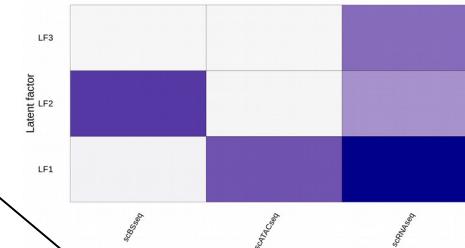
For Example:

- 1) With ~100 samples it is a good idea to do **linear** Omics integration
- 2) T2D is a phenotype of interest, therefore **supervised** integration



Data Set (4 Omics)
110 overlapping individuals

Check covariance



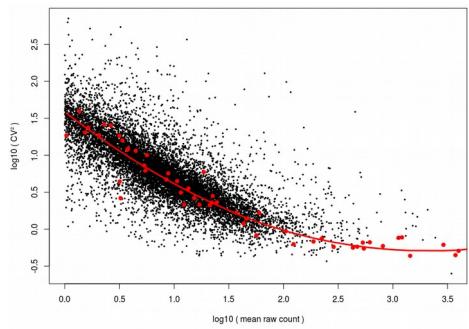
Train Set (n = 80)

Supervised:
LASSO

Test Set (n = 30)

Evaluation

Unsupervised:
remove low-variance

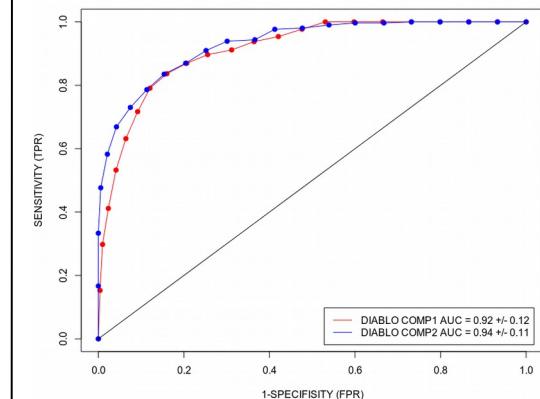


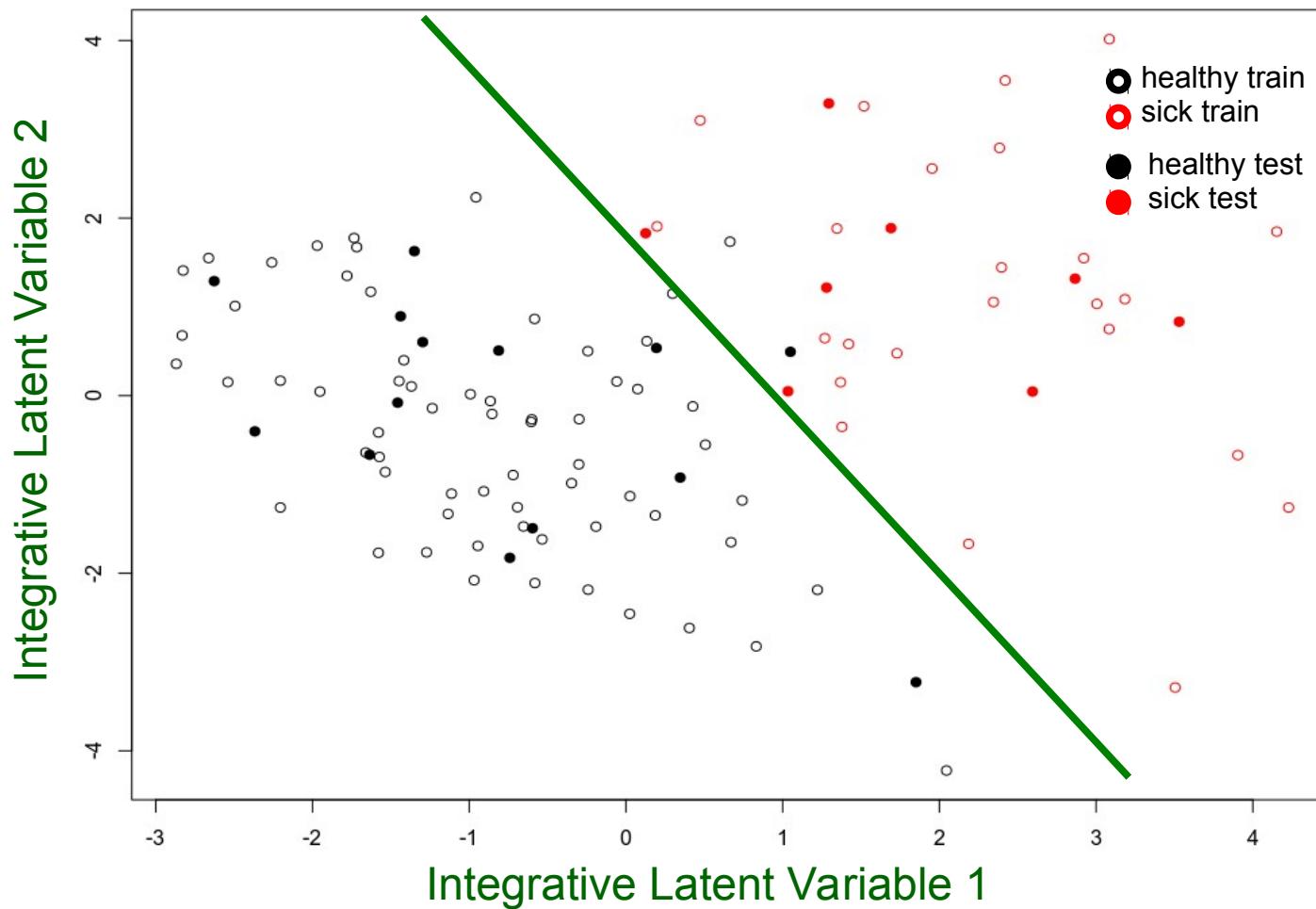
Feature Selection

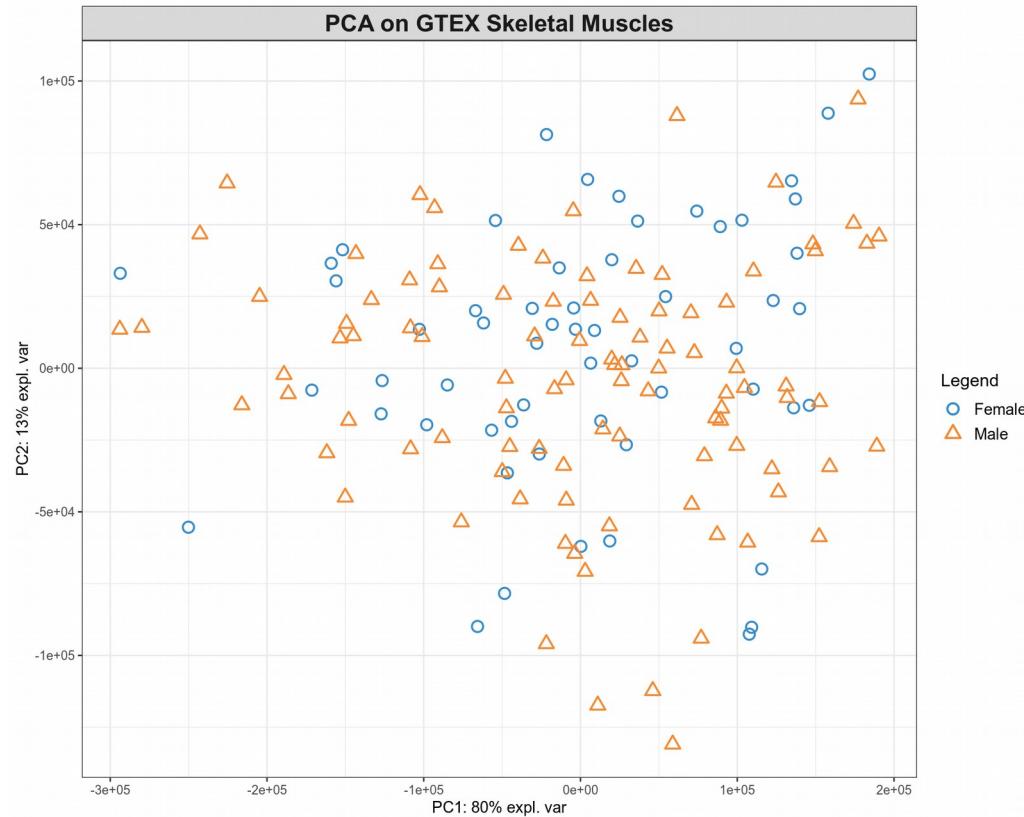


Omics Integration

Trained Model







Test feature-by-feature for correlation with phenotype of interest: aka DGE

```

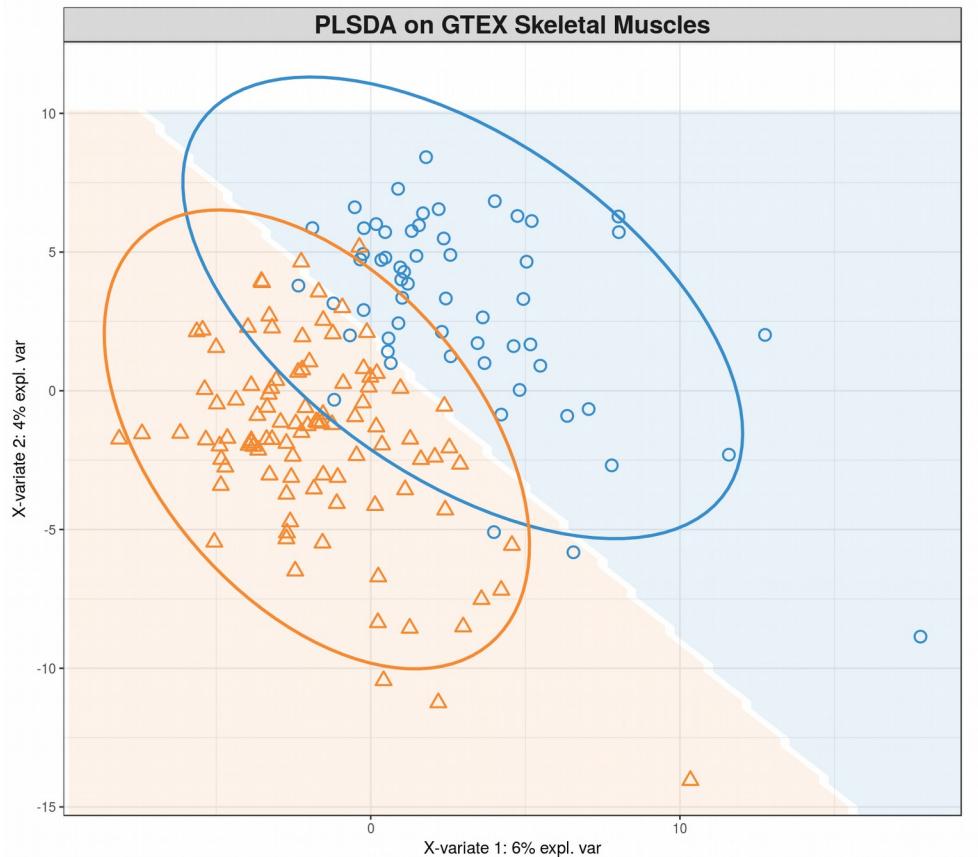
1 rho <- vector()
2 p <- vector()
3 a <- seq(from=0, to=dim(x)[2], by=100)
4 for(i in 1:dim(x)[2])
5 {
6   corr_output <- cor.test(X[,i], as.numeric(Y), method="spearman")
7   rho <- append(rho,as.numeric(corr_output$estimate))
8   p <- append(p,as.numeric(corr_output$p.value))
9   if(iTRUE(i %in% a)==TRUE){print(paste("FINISHED ",i," FEATURES",sep=""))}
10 }
11 output <- data.frame(GENE=colnames(X), SPEARMAN_RHO=rho, PVALUE=p)
12 output$FDR <- p.adjust(output$PVALUE, method="fdr")
13 output <- output[order(output$FDR, output$PVALUE, -output$SPEARMAN_RHO), ]
14 head(output,10)

```

UnivarFeatureSelect.R hosted with ❤ by GitHub

[view raw](#)

##	GENE	SPEARMAN_RHO	PVALUE	FDR
## 256	ENSG00000184368.11_MAP7D2	-0.5730196	4.425151e-15	2.416132e-12
## 324	ENSG00000110013.8_SIAE	0.3403994	1.288217e-05	3.516833e-03
## 297	ENSG00000128487.12_SPECCI	-0.3003621	1.323259e-04	2.408332e-02
## 218	ENSG00000162512.11_SDC3	0.2945390	1.807649e-04	2.467441e-02
## 38	ENSG00000129007.10_CALML4	0.2879754	2.549127e-04	2.783647e-02
## 107	ENSG00000233429.5_HOTAIRM1	-0.2768054	4.489930e-04	4.085836e-02
## 278	ENSG00000185442.8_FAM174B	-0.2376098	2.731100e-03	2.130258e-01
## 421	ENSG00000234585.2_CCT6P3	-0.2322268	3.426233e-03	2.338404e-01
## 371	ENSG00000113312.6_TTC1	0.2284351	4.007655e-03	2.431310e-01
## 269	ENSG00000226329.2_AC005682.6	-0.2226587	5.064766e-03	2.523944e-01

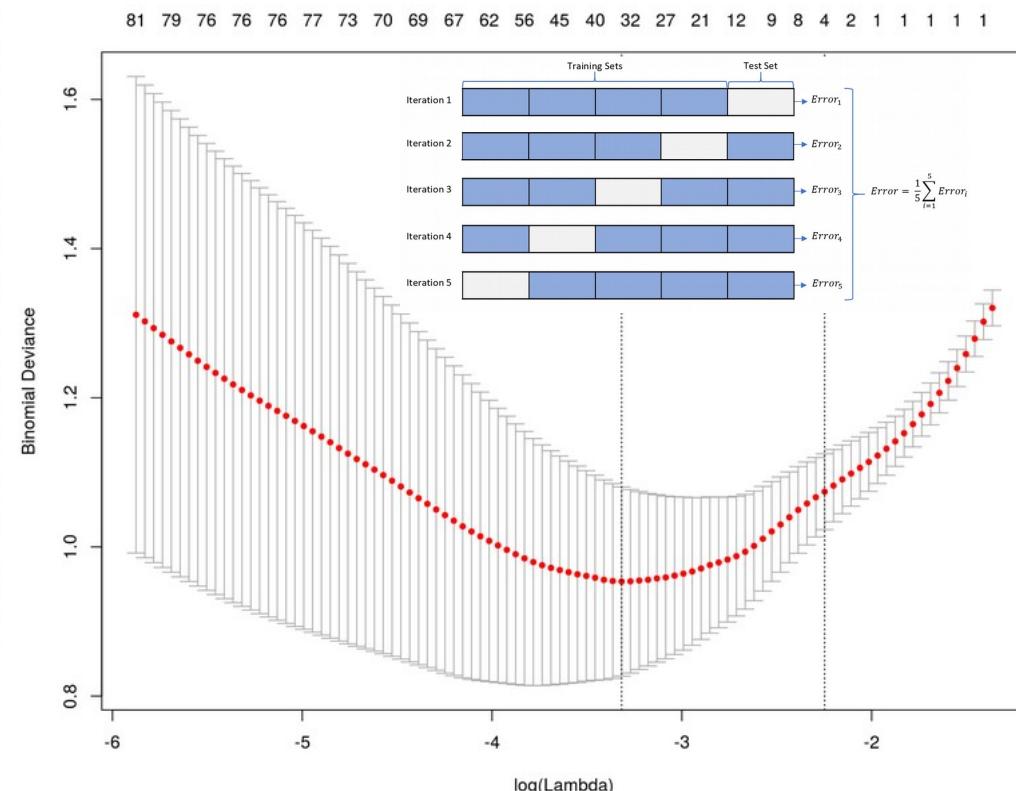


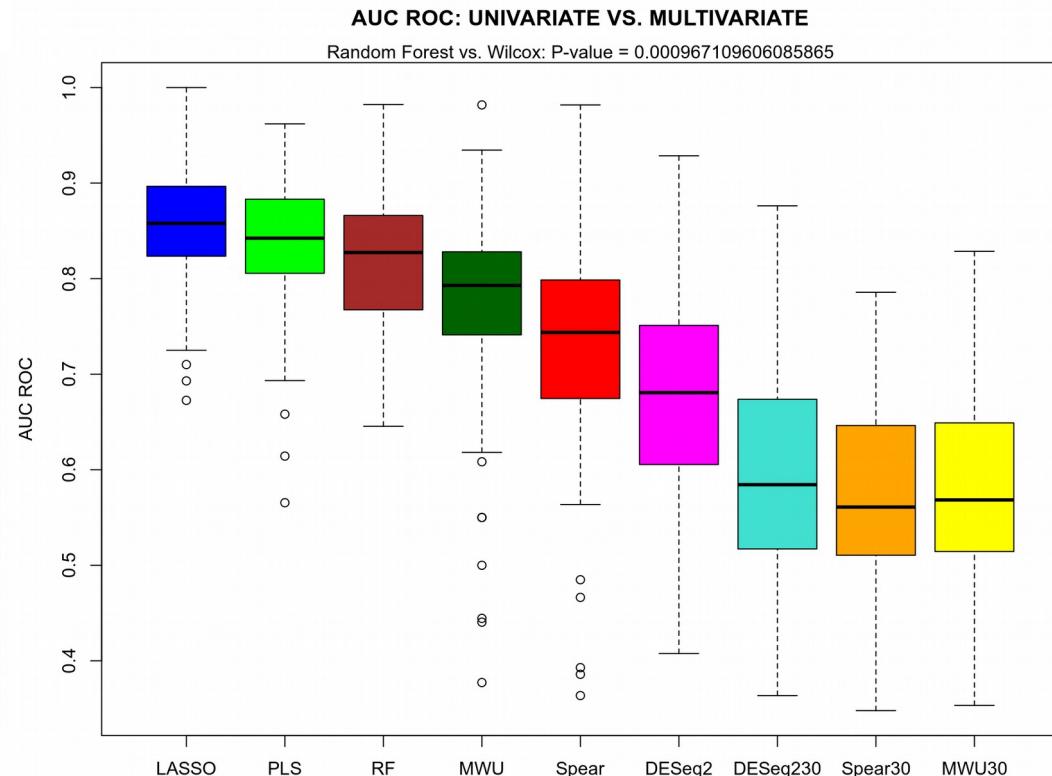
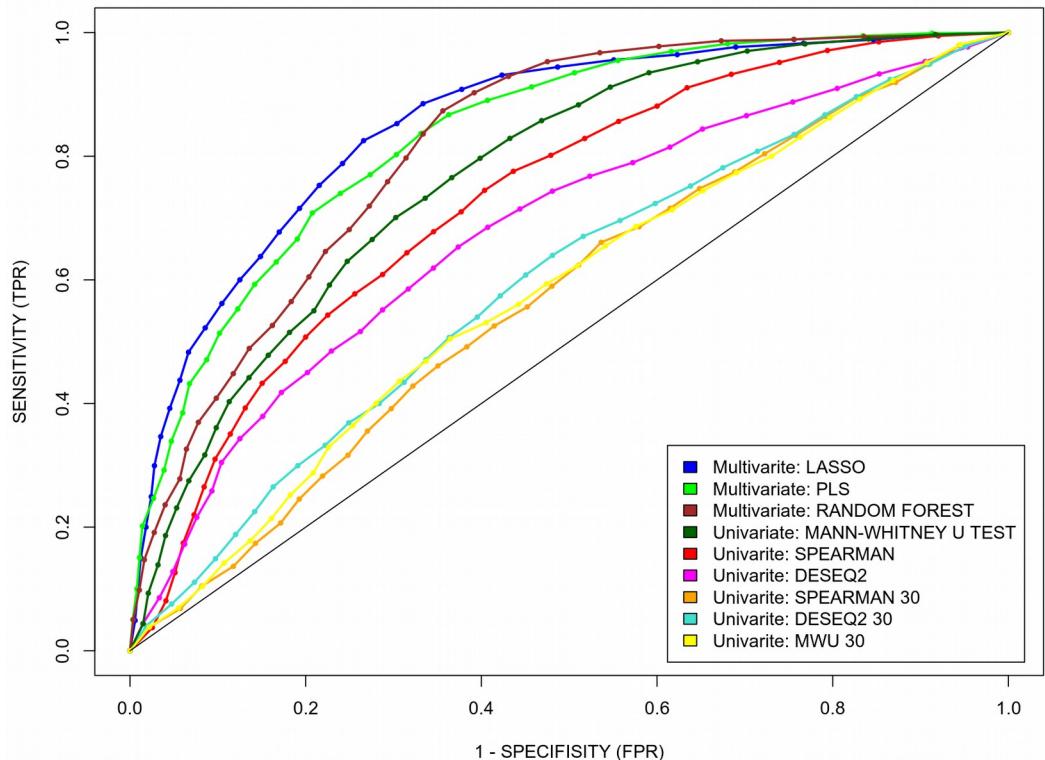
Select groups of features correlated with phenotype of interest: PLS, LASSO, LDA

$$Y = \beta_1 X_1 + \beta_2 X_2 + \epsilon$$

$$\text{OLS} = (y - \beta_1 X_1 - \beta_2 X_2)^2$$

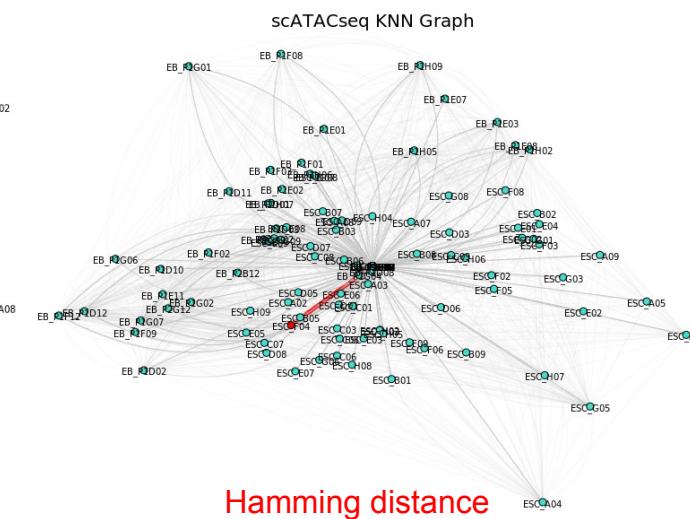
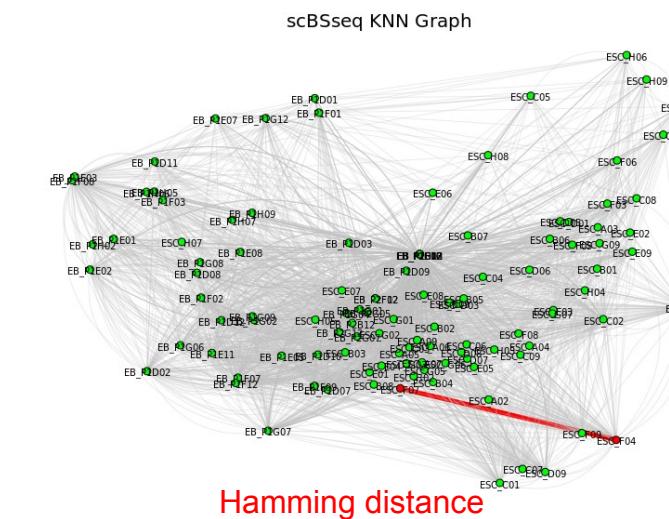
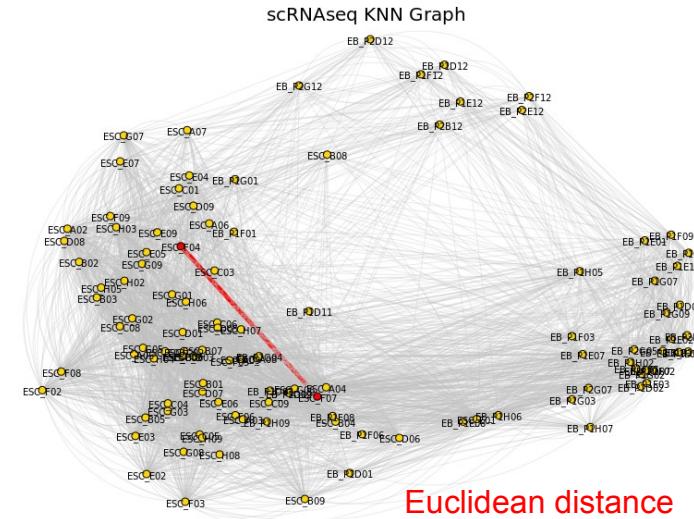
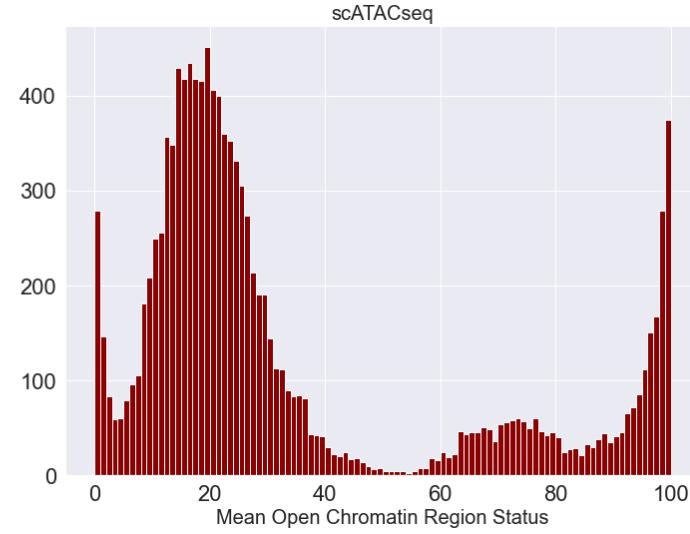
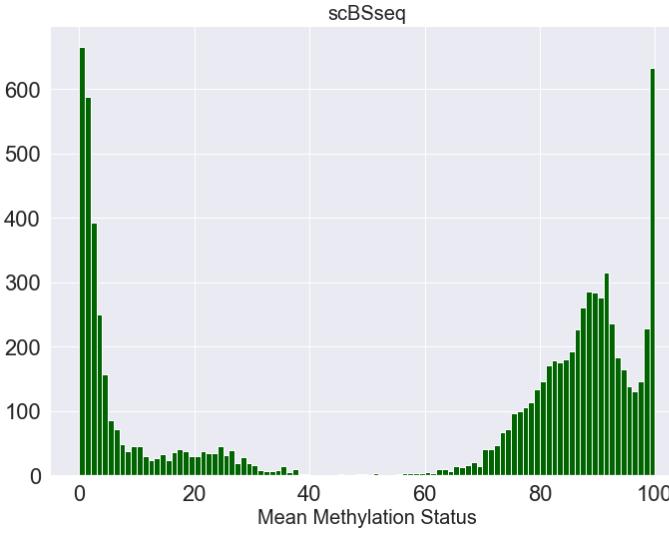
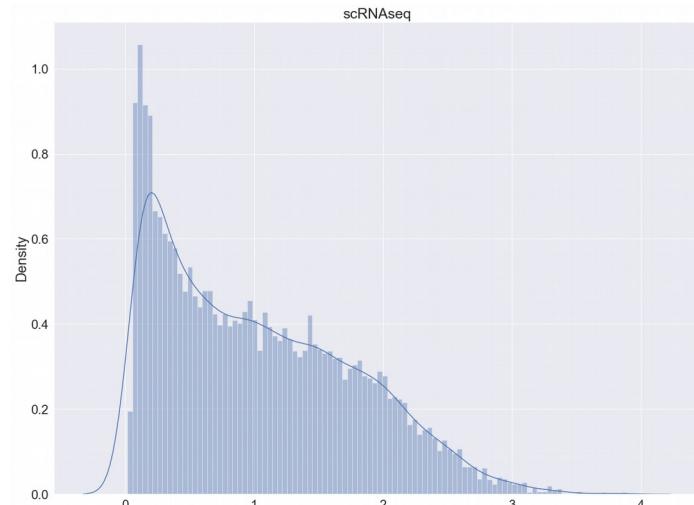
$$\text{Penalized OLS} = (y - \beta_1 X_1 - \beta_2 X_2)^2 + \lambda(|\beta_1| + |\beta_2|)$$



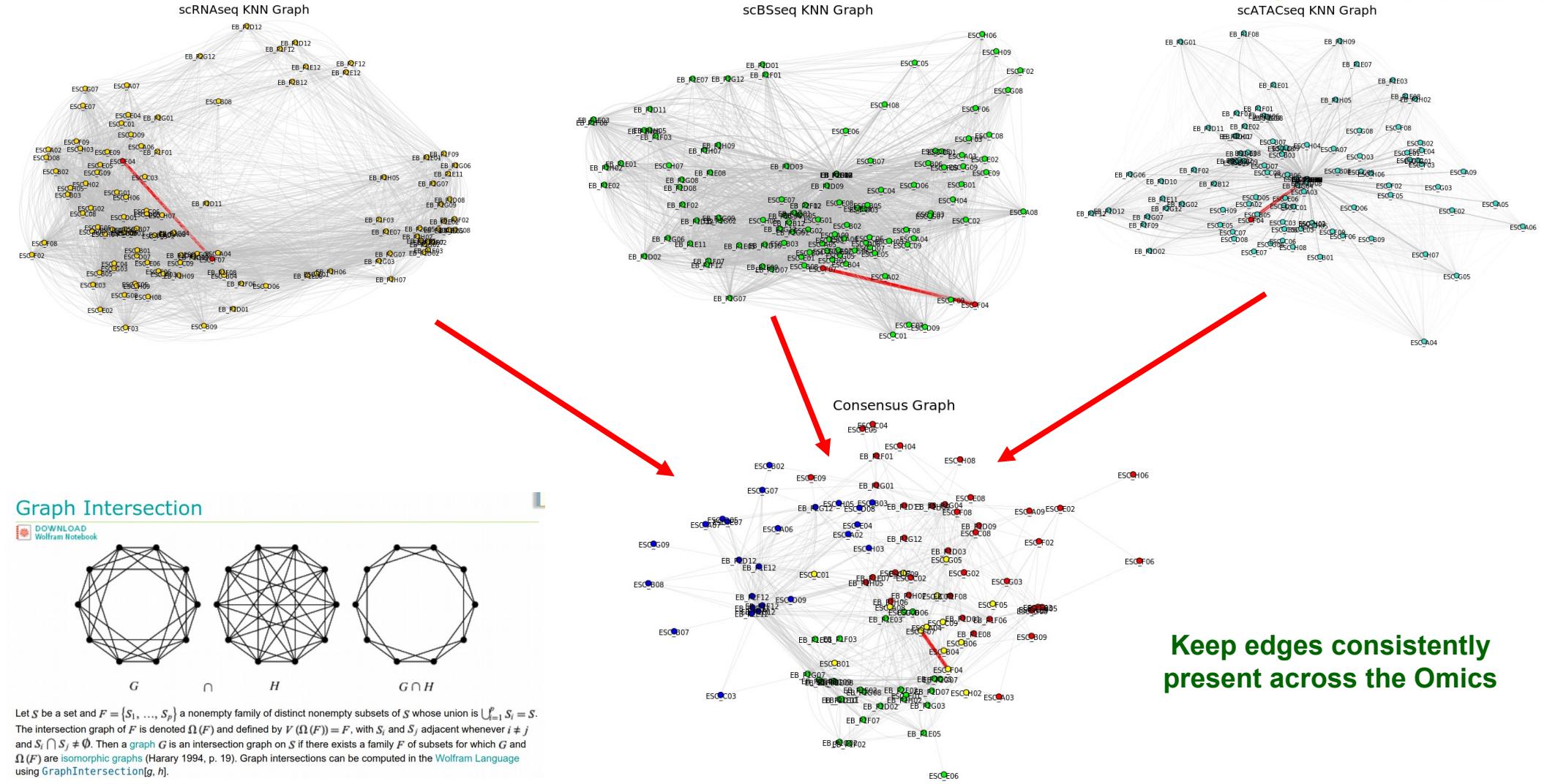


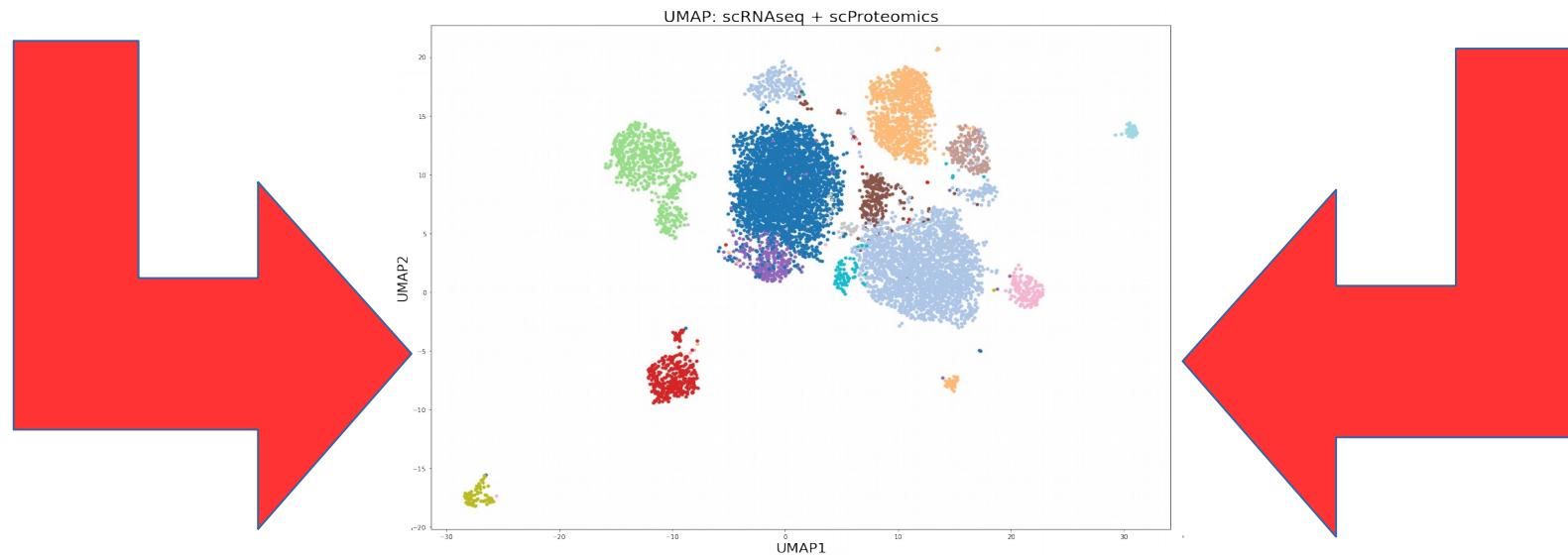
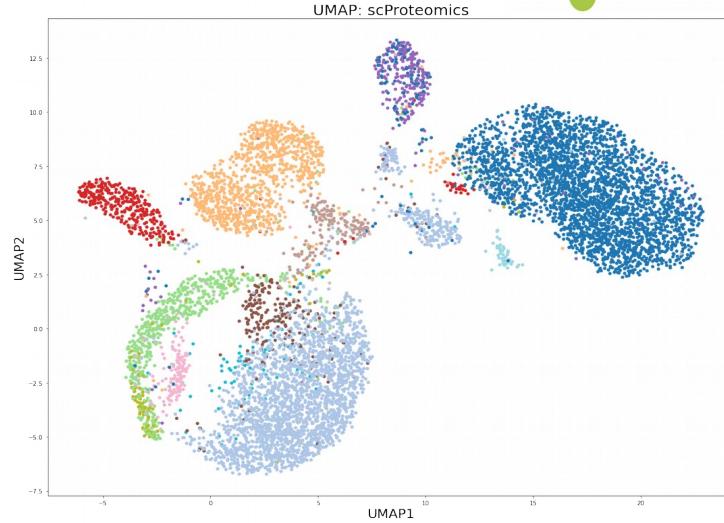
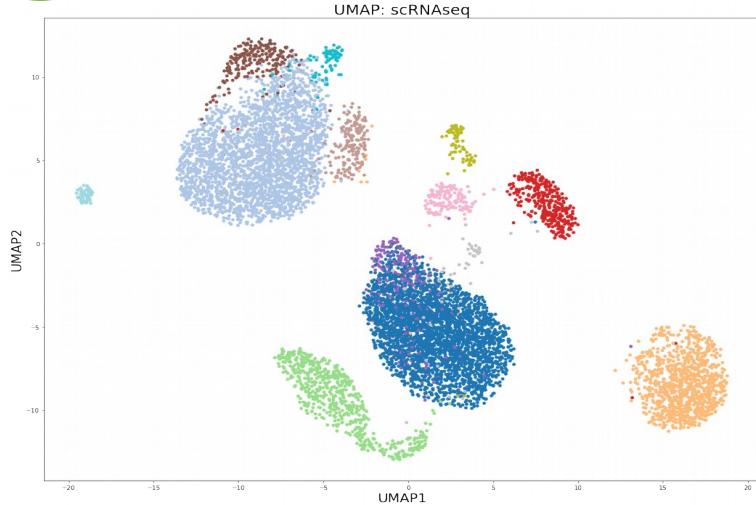
Multivariate feature selection generally has higher predictive capacity compared to univariate feature selection

How Does This Work in Practice? (an example from unsupervised integration)

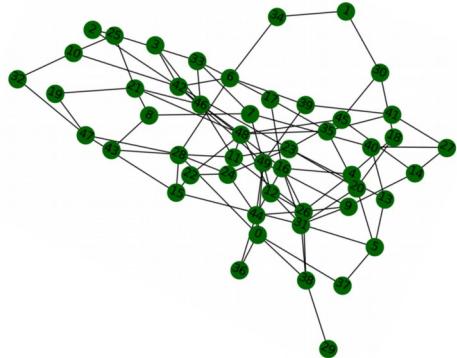


Graph Intersection Method

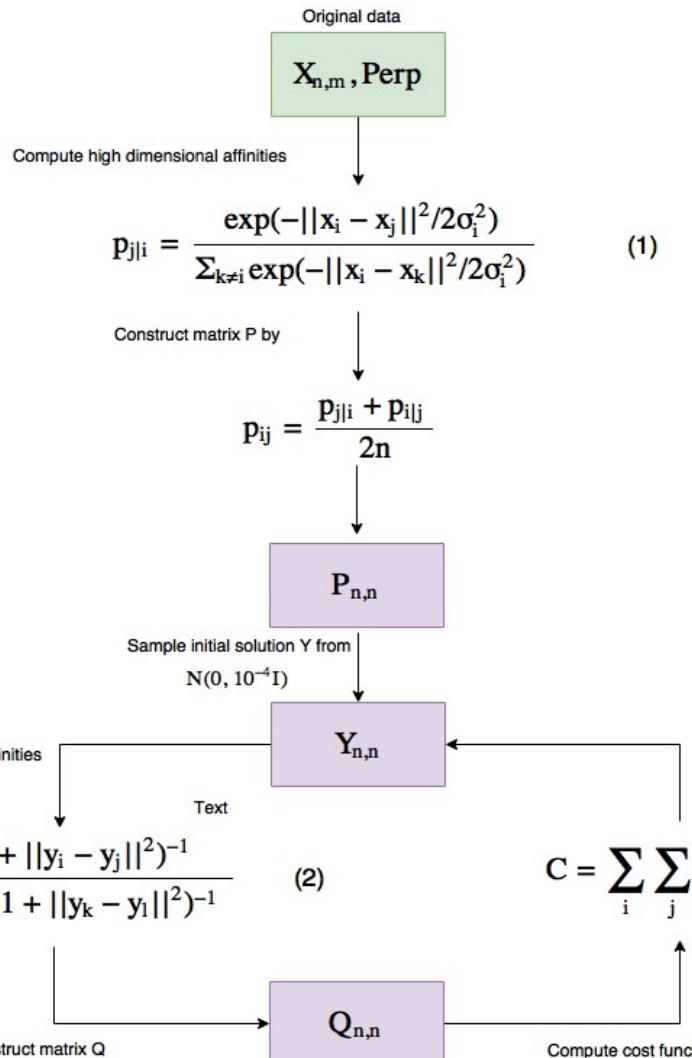
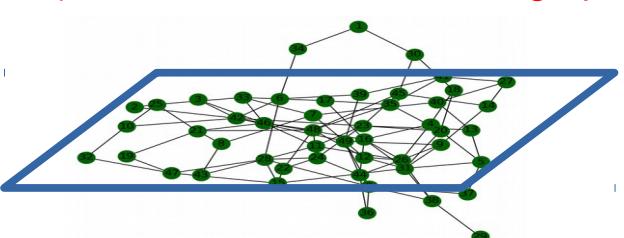




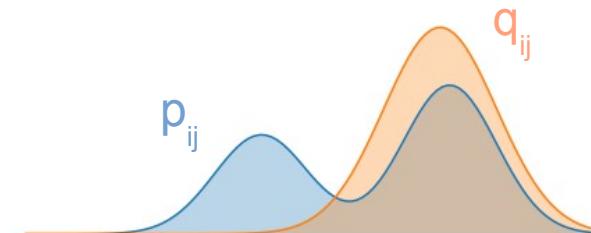
1) Construct high-dimensional graph



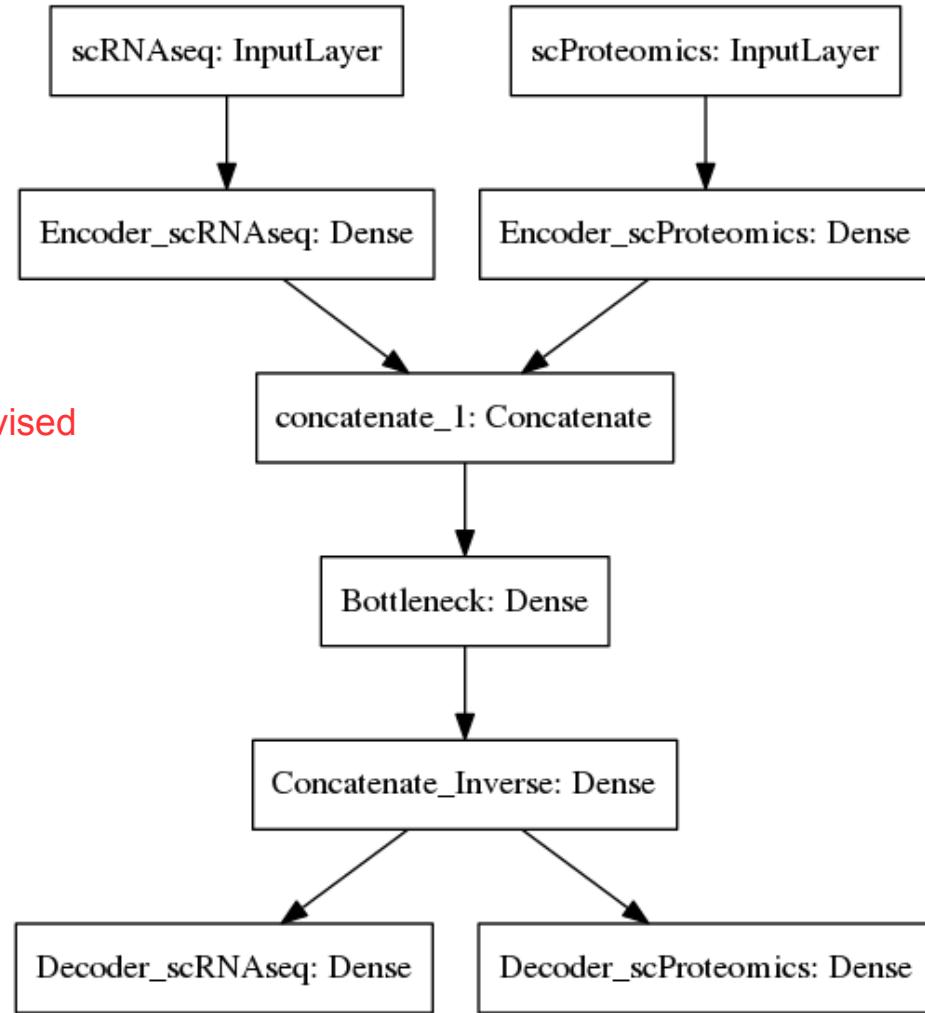
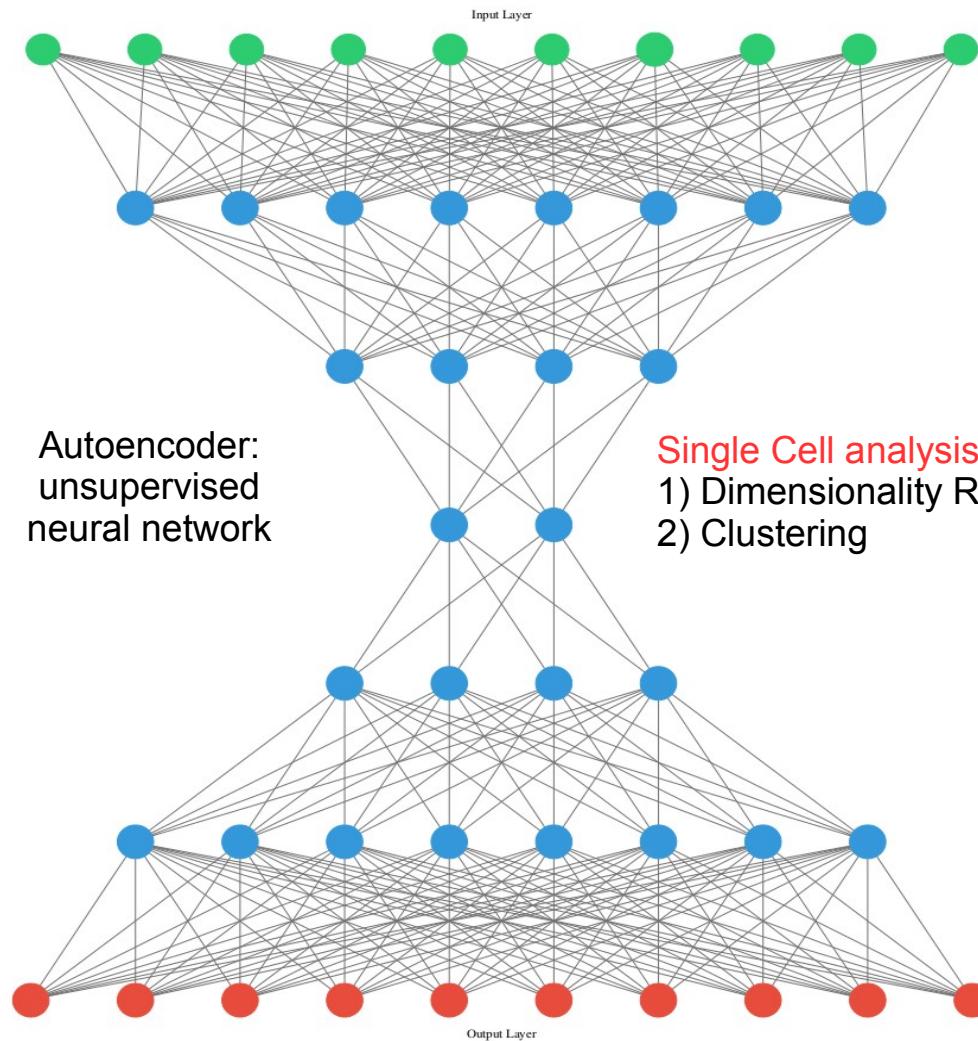
2) Construct low-dimensional graph

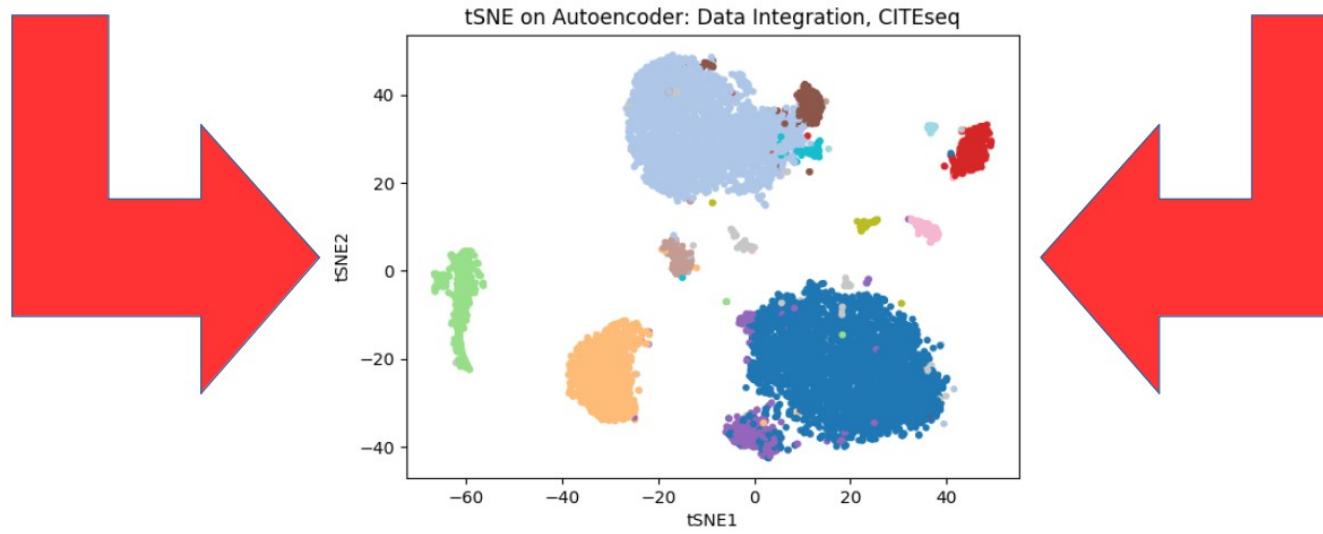
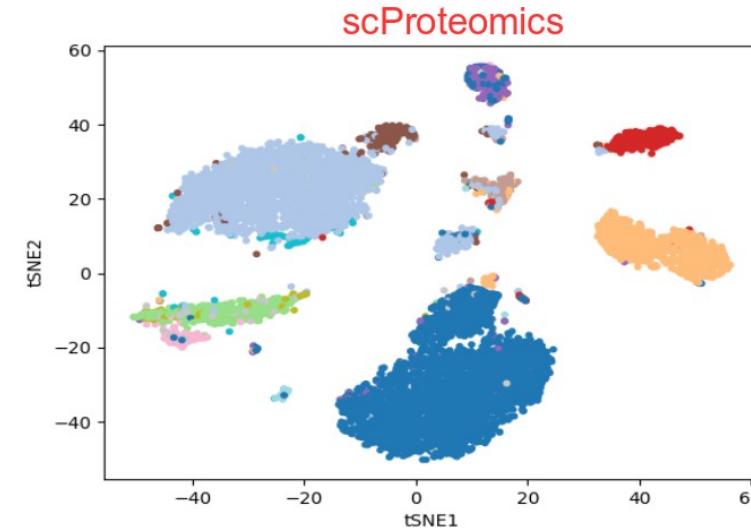
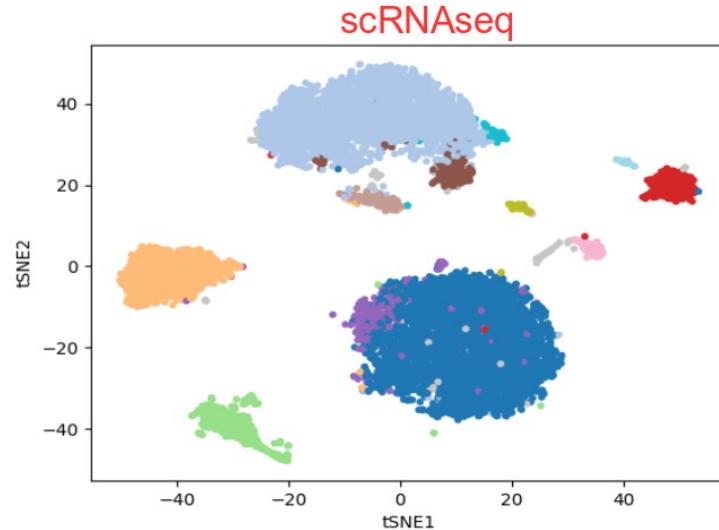


3) Collapse the graphs together



$$C = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}} \quad (3)$$







*Knut och Alice
Wallenbergs
Stiftelse*



**LUNDS
UNIVERSITET**

<https://github.com/NikolayOskolkov/OsloBioinfoWeek2022>

The screenshot shows a GitHub repository page. At the top, there's a navigation bar with links for Pull requests, Issues, Codespaces, Marketplace, and Explore. Below the navigation bar, the repository name 'NikolayOskolkov / OsloBioinfoWeek2022' is displayed, along with a 'Public' badge. To the right of the repository name are buttons for Pin, Unwatch (with a count of 1), Fork (with a count of 0), and Star (with a count of 0). The main content area shows a list of files uploaded by 'NikolayOskolkov'. The files listed are:

- CITESeq.zip
- DeepLearningDataIntegration.html
- DeepLearningDataIntegration.ipynb
- ENSEMBLE_TO_GENE_SYMBOL_...
- GTEX_SkeletalMuscles_157Sample...
- GTEX_SkeletalMuscles_157Sample...
- GTEX_Data_2014-01-17_Annotatio...
- OmicsIntegration_FeatureSelection....
- OmicsIntegration_FeatureSelection....
- TMM_NormalizedCounts_157_Sam...
- UMAP_DataIntegration.html
- UMAP_DataIntegration.ipynb
- UnsupervisedOMICsIntegration_MO...
- UnsupervisedOMICsIntegration_MO...
- scNMT.zip

Each file entry includes a small preview icon, the file name, the action 'Add files via upload', and the time it was added (26 seconds ago). On the right side of the page, there are sections for 'About', 'Releases', and 'Packages'. The 'About' section notes 'No description, website, or topics provided.' The 'Releases' section says 'No releases published' and 'Create a new release'. The 'Packages' section says 'No packages published' and 'Publish your first package'. At the bottom of the page, there's a call to action 'Help people interested in this repository understand your project by adding a README.' followed by a green 'Add a README' button.