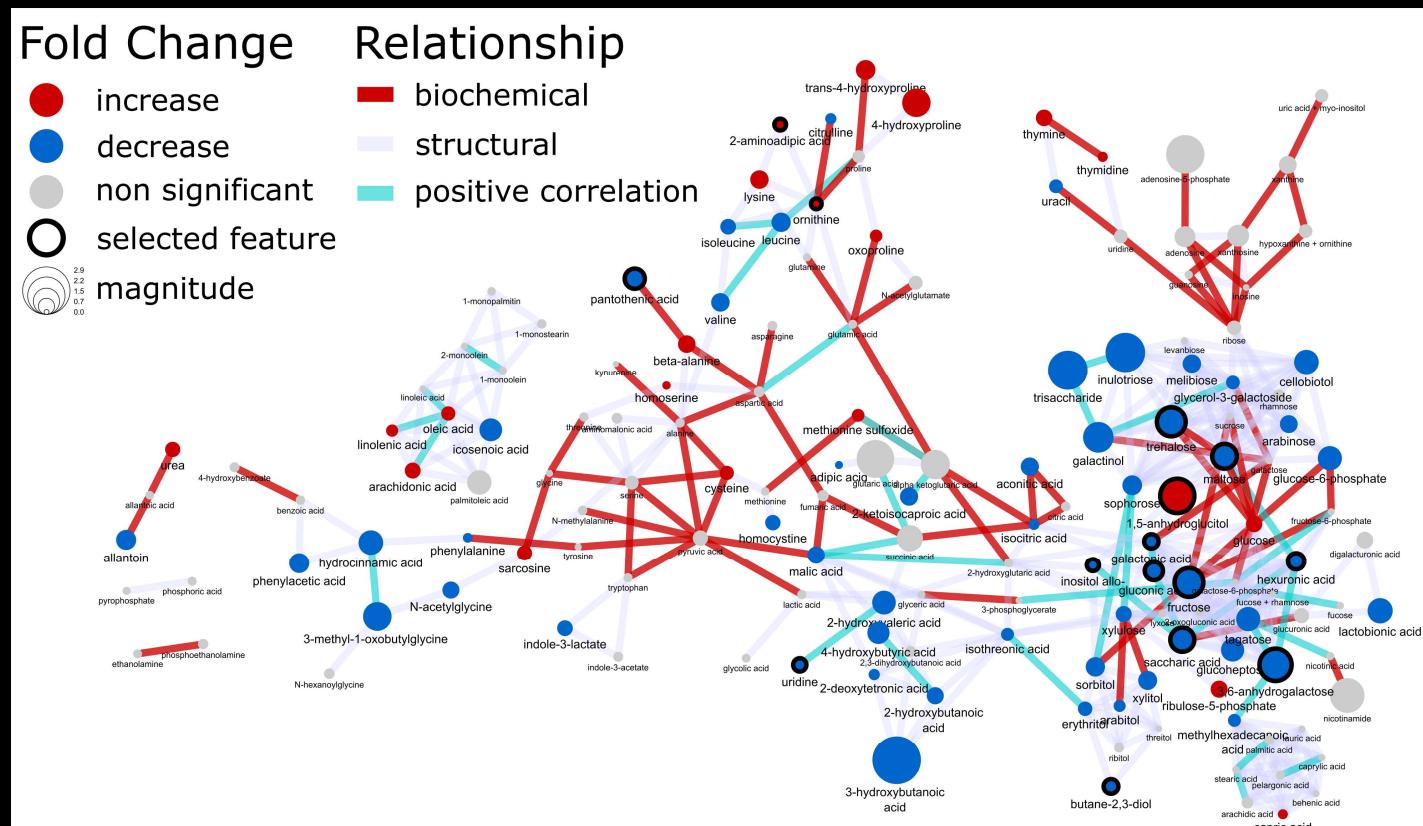


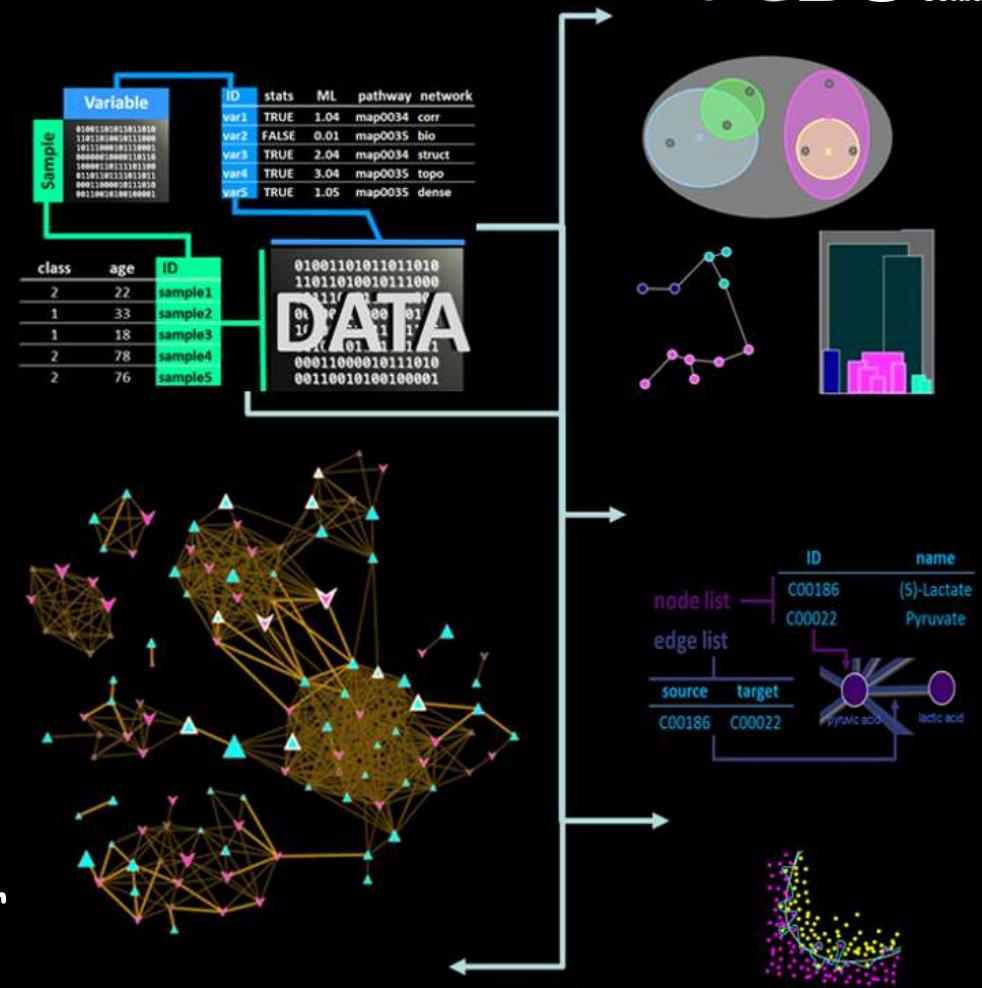
# Welcome to network mapping 101

In the following course you will learn how to integrate statistical, multivariate and machine learning results within a publication quality biochemical network.

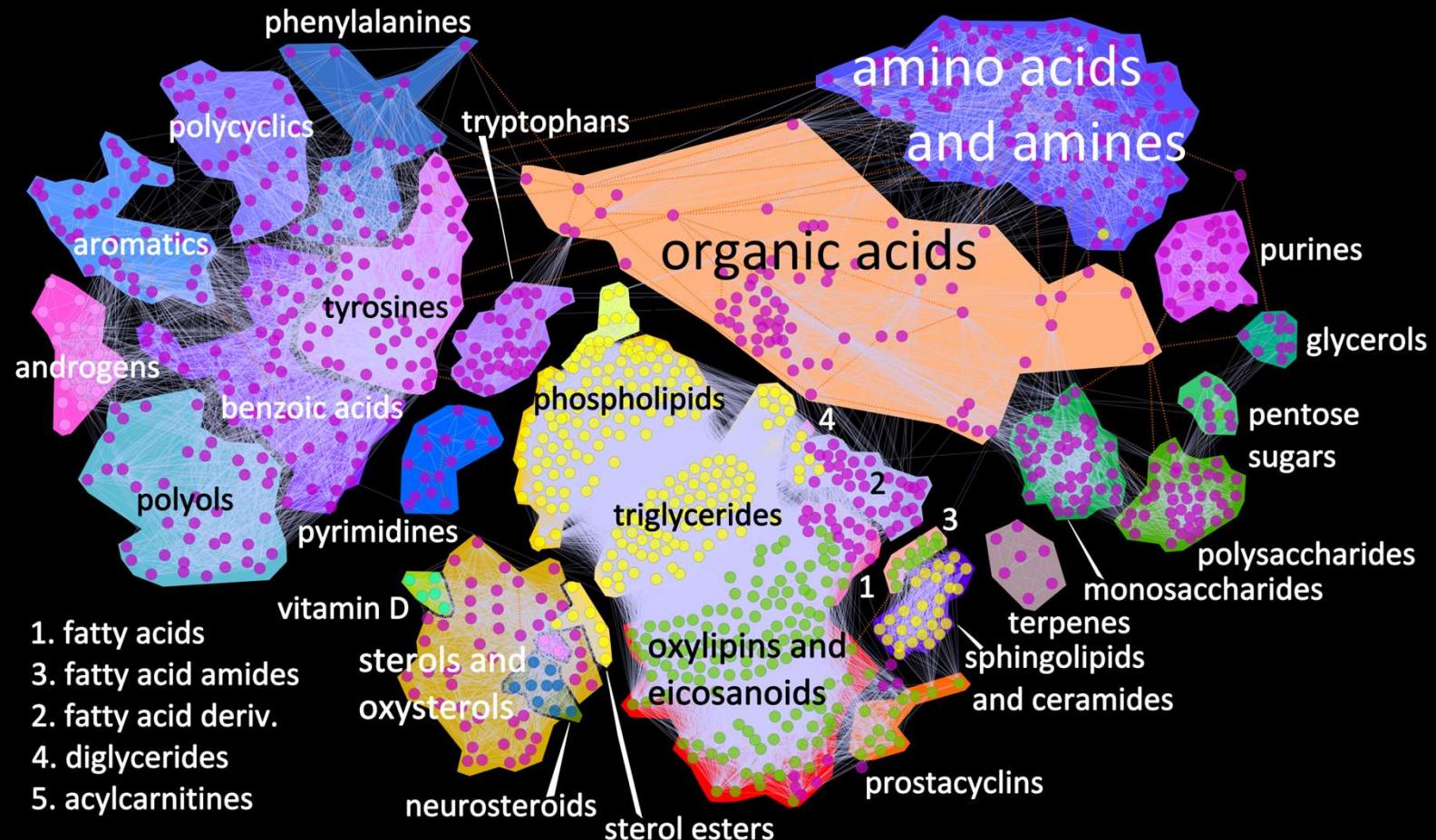


# Tutorials

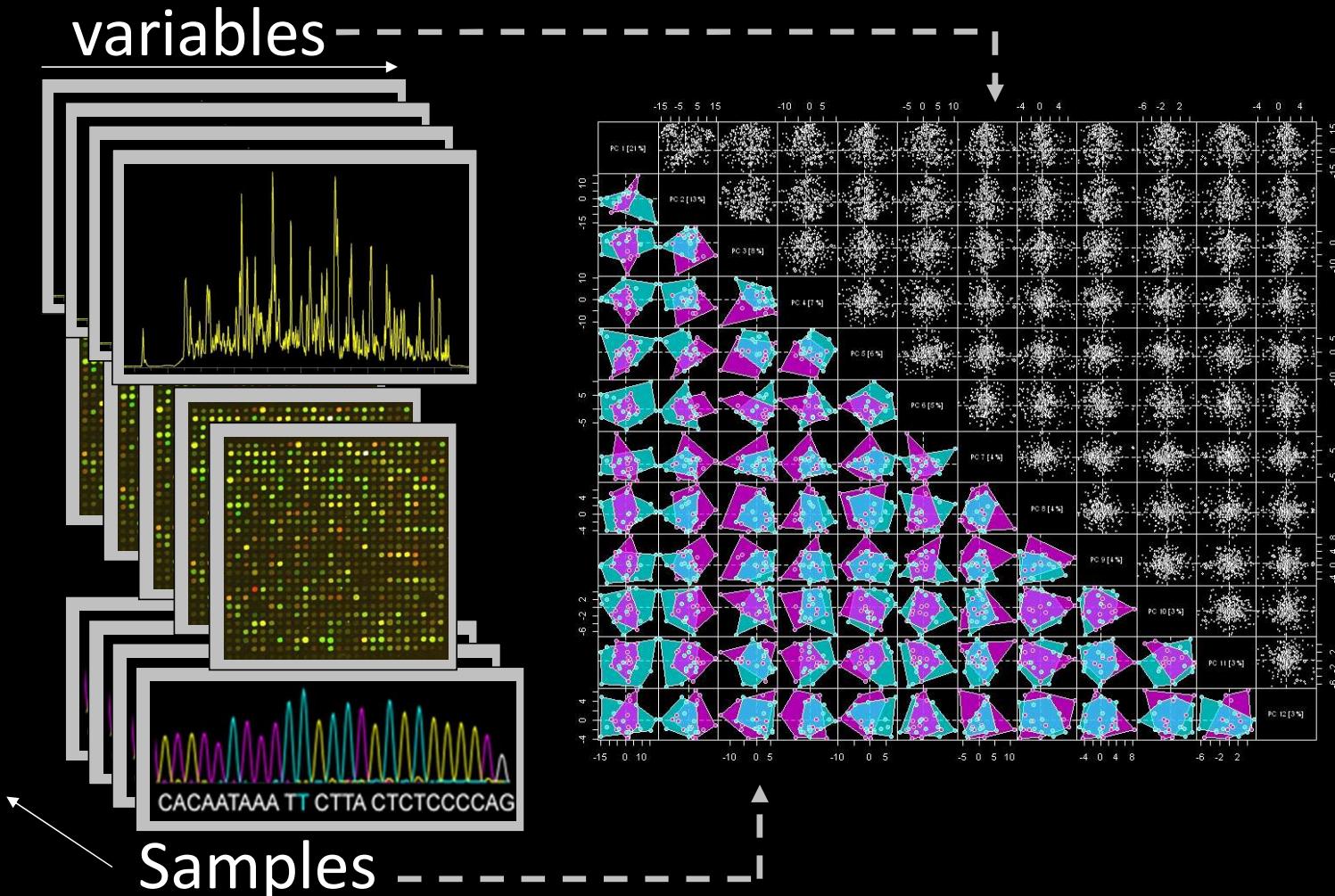
- Preparing raw data for analysis
- Statistical analysis
- Multivariate data exploration
- Supervised clustering
- Machine learning
  - classification
  - model validation
  - feature selection
- Network analysis
  - biochemical
  - structural similarity
  - correlation
- Network mapping - putting it all together



# Analysis at the metabolomic scale



# Integrate high-dimensional data



# Identify what matters



## Univariate

## Multivariate

# Predictive Modeling

Group 1

Group 2

1995-1996  
1996-1997  
1997-1998  
1998-1999  
1999-2000  
2000-2001  
2001-2002  
2002-2003  
2003-2004  
2004-2005  
2005-2006  
2006-2007  
2007-2008  
2008-2009  
2009-2010  
2010-2011  
2011-2012  
2012-2013  
2013-2014  
2014-2015  
2015-2016  
2016-2017  
2017-2018  
2018-2019  
2019-2020  
2020-2021  
2021-2022  
2022-2023  
2023-2024  
2024-2025  
2025-2026  
2026-2027  
2027-2028  
2028-2029  
2029-2030  
2030-2031  
2031-2032  
2032-2033  
2033-2034  
2034-2035  
2035-2036  
2036-2037  
2037-2038  
2038-2039  
2039-2040  
2040-2041  
2041-2042  
2042-2043  
2043-2044  
2044-2045  
2045-2046  
2046-2047  
2047-2048  
2048-2049  
2049-2050  
2050-2051  
2051-2052  
2052-2053  
2053-2054  
2054-2055  
2055-2056  
2056-2057  
2057-2058  
2058-2059  
2059-2060  
2060-2061  
2061-2062  
2062-2063  
2063-2064  
2064-2065  
2065-2066  
2066-2067  
2067-2068  
2068-2069  
2069-2070  
2070-2071  
2071-2072  
2072-2073  
2073-2074  
2074-2075  
2075-2076  
2076-2077  
2077-2078  
2078-2079  
2079-2080  
2080-2081  
2081-2082  
2082-2083  
2083-2084  
2084-2085  
2085-2086  
2086-2087  
2087-2088  
2088-2089  
2089-2090  
2090-2091  
2091-2092  
2092-2093  
2093-2094  
2094-2095  
2095-2096  
2096-2097  
2097-2098  
2098-2099  
2099-20100

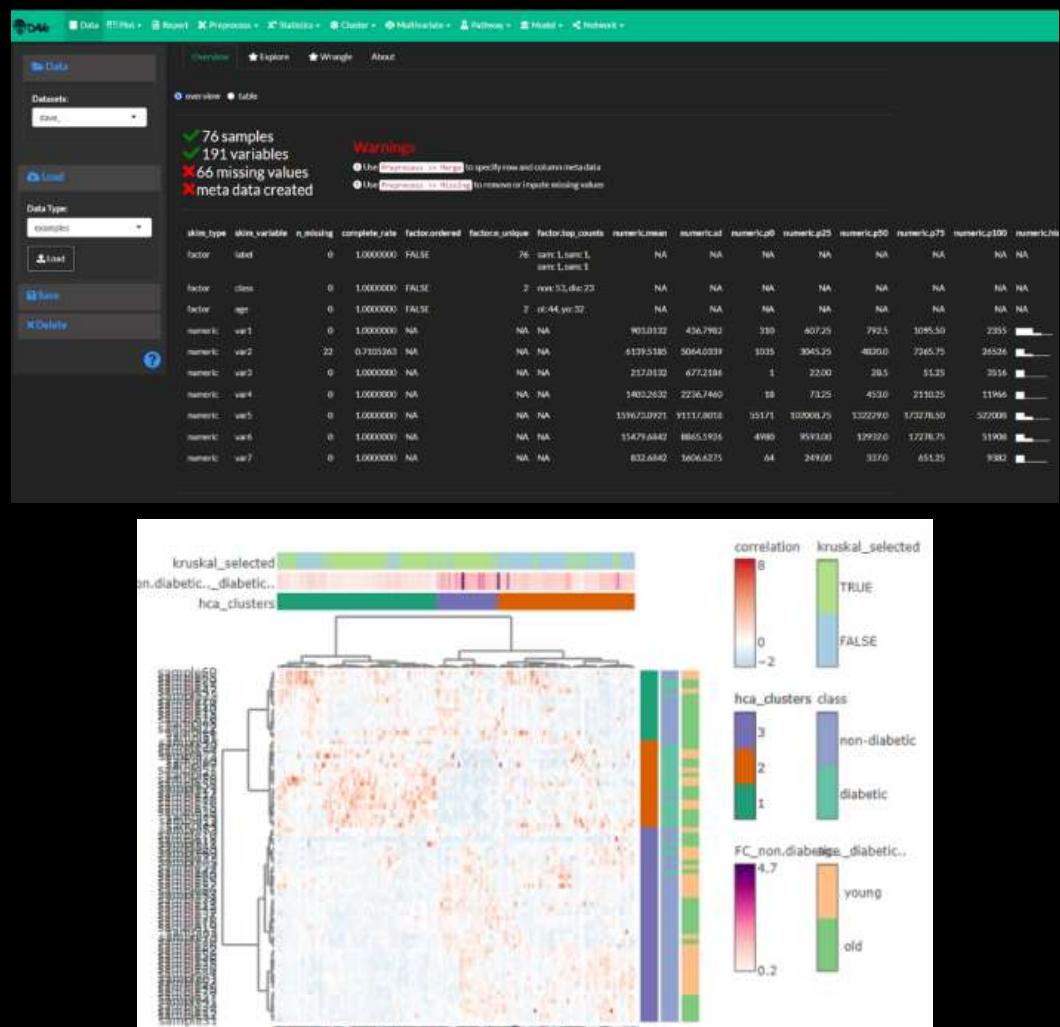
ANOVA

PCA

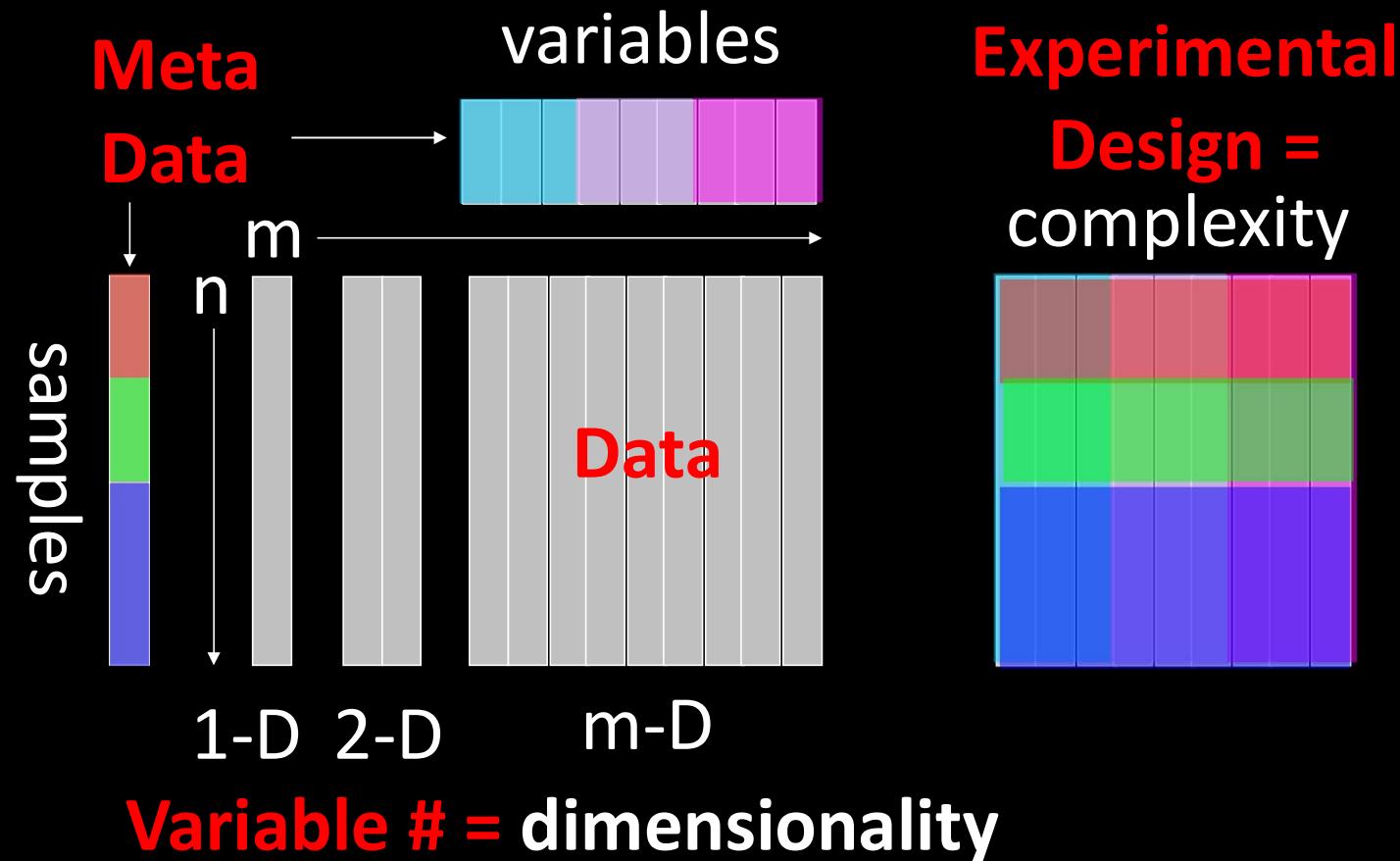
PLS

# Topics

- Data preparation
- Differential expression
- Hierarchical Clustering
- Principal Components Analysis (PCA)
- Statistical analysis
- Machine learning
- Network analysis
- Network mapping



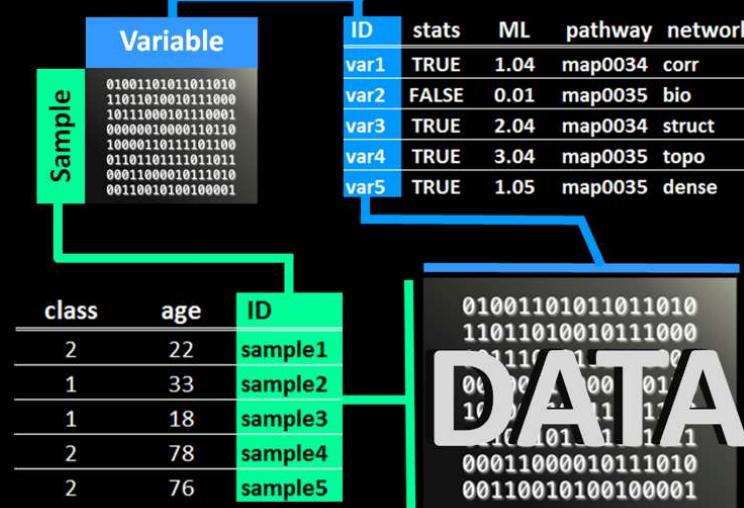
# How to think about data complexity



# Data preprocessing

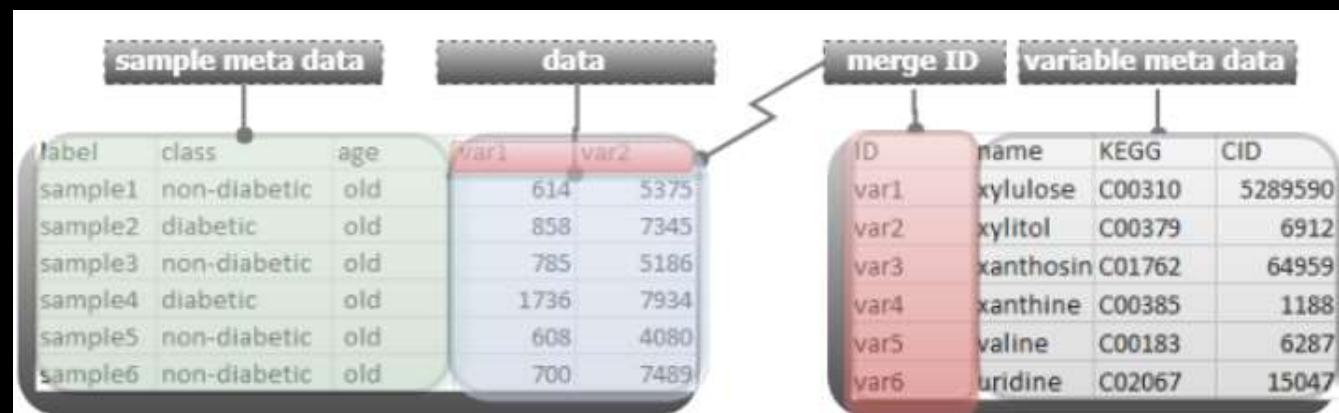
define

- data
- row meta data
- column meta data



remove and/or impute

- missing values



# Your turn

Follow along with the following tutorial:

[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/preprocess/](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/preprocess/)

# Differential expression

compare

- class means

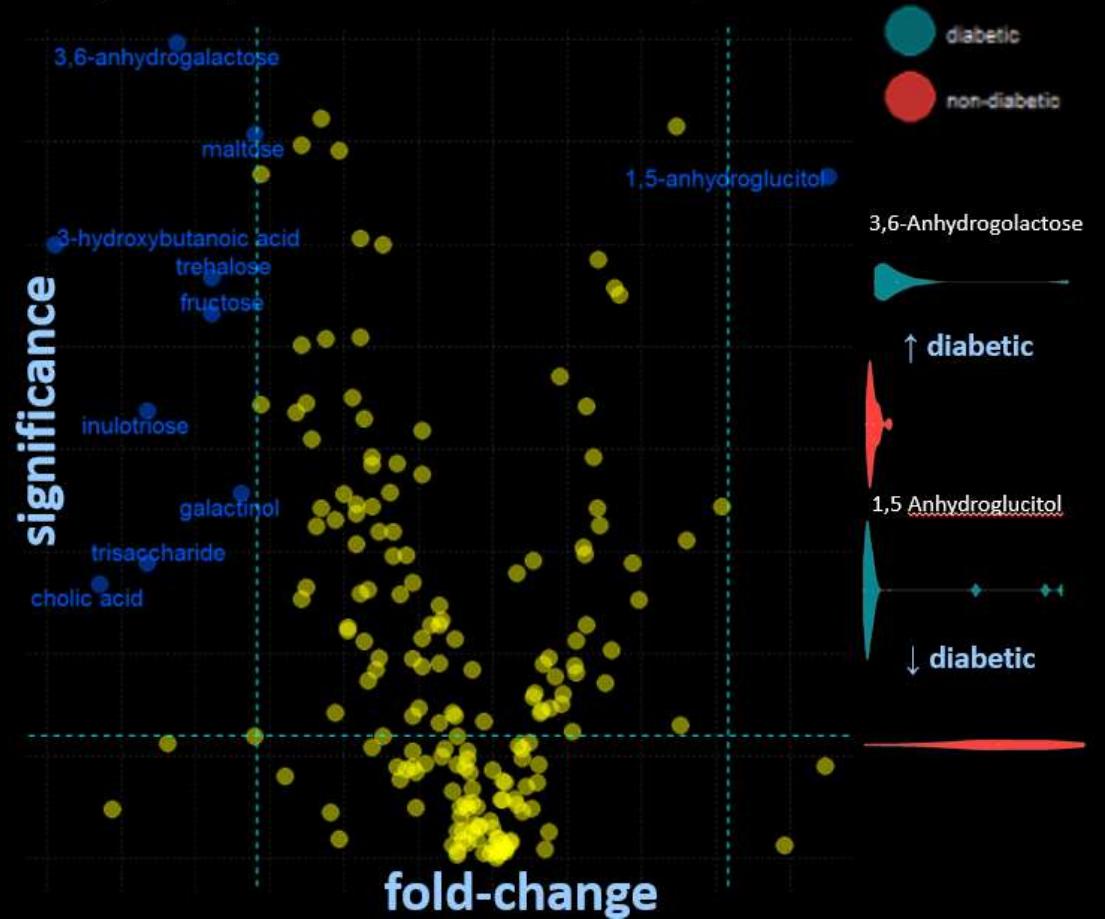
identify

- significant differences

visualize

- volcano plots
- violin plots

Simplest representation: two-class comparison



# Your turn

Follow along with the following tutorial:

[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/statistics/](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/statistics/)

# Hierarchical clustering (HCA)

## group

- samples and/or variables

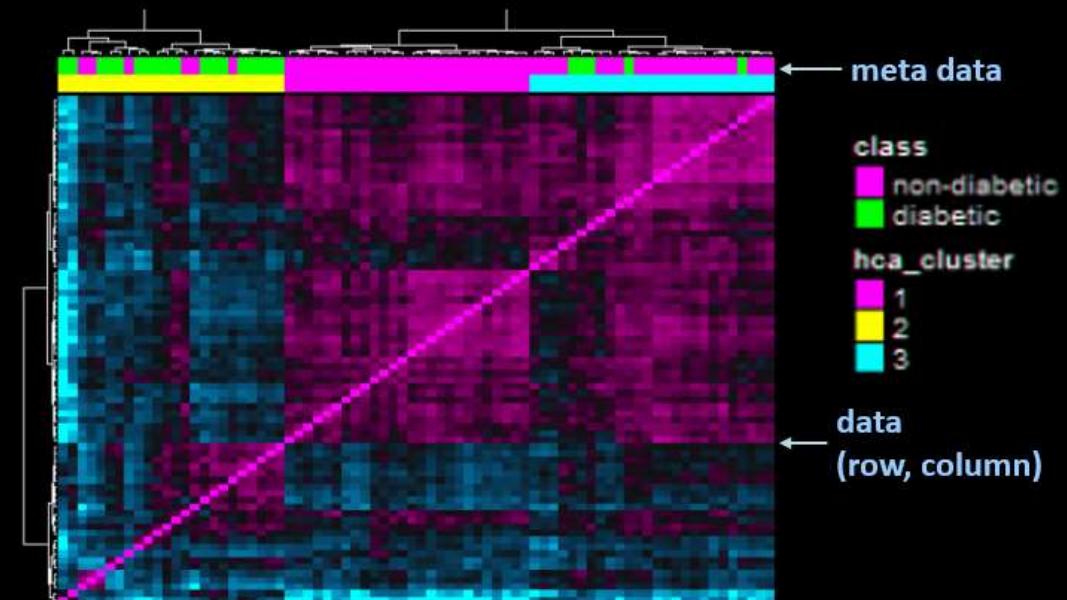
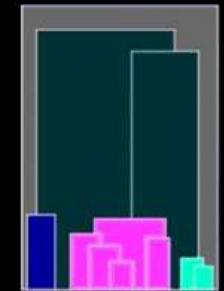
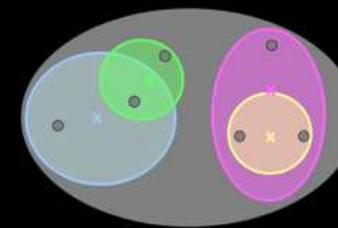
## define similarity

- correlation
- distance
- linkage

## visualize

- heatmaps
- dendrograms

Clustering: distance + linkage



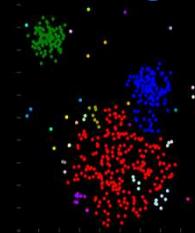
# Clustering basics

Use the concept similarity/dissimilarity to group a collection of samples or variables

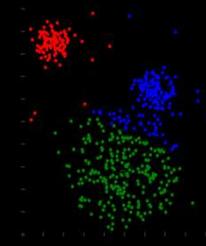
## approaches

- hierarchical (linkage)
- non-hierarchical (k-NN, k-means)
- distribution (mixtures models)
- density (DBSCAN)
- self organizing maps (SOM)

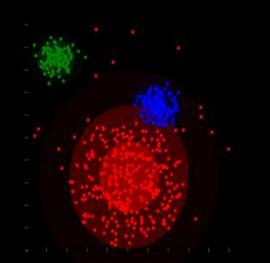
Linkage



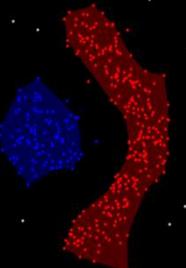
k-means



Distribution



Density



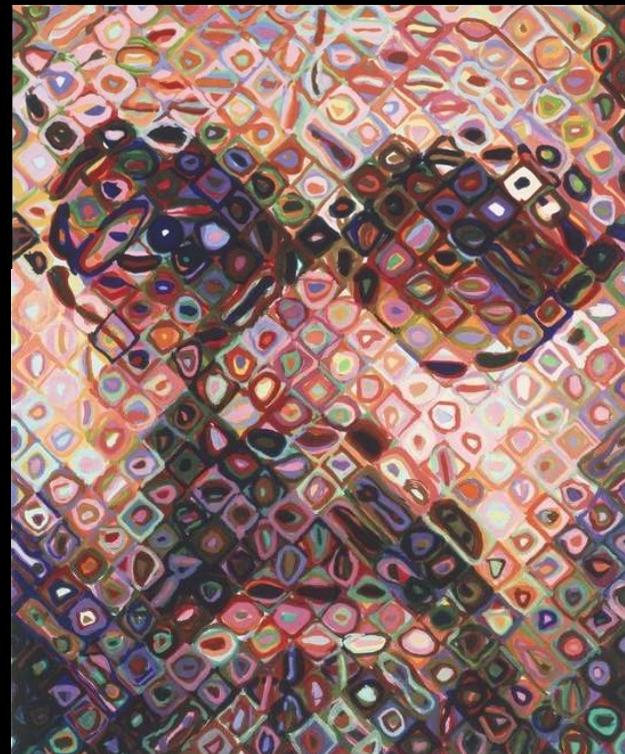
# HCA goals

identify

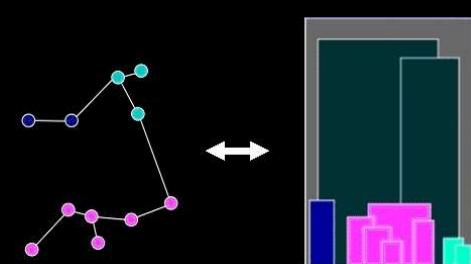
- patterns
- group structure
- relationships

evaluate and refine hypothesis

reduce complexity



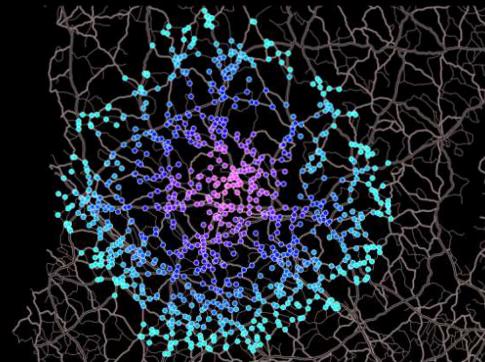
Artist: Chuck Close



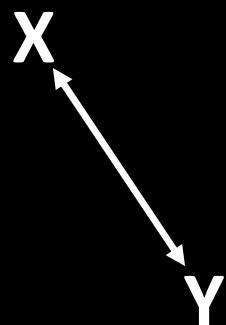
# HCA methods

## distance

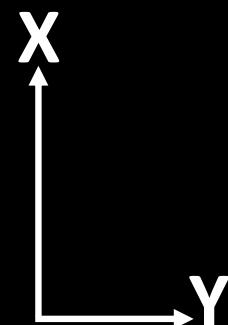
- defines “nearness” or similarity



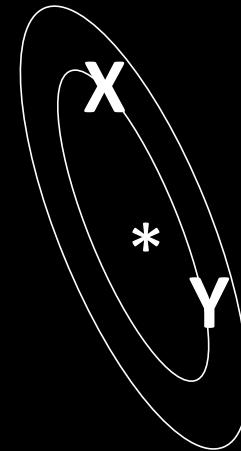
Euclidean



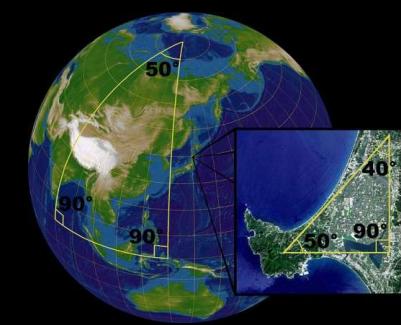
Manhattan



Mahalanobis



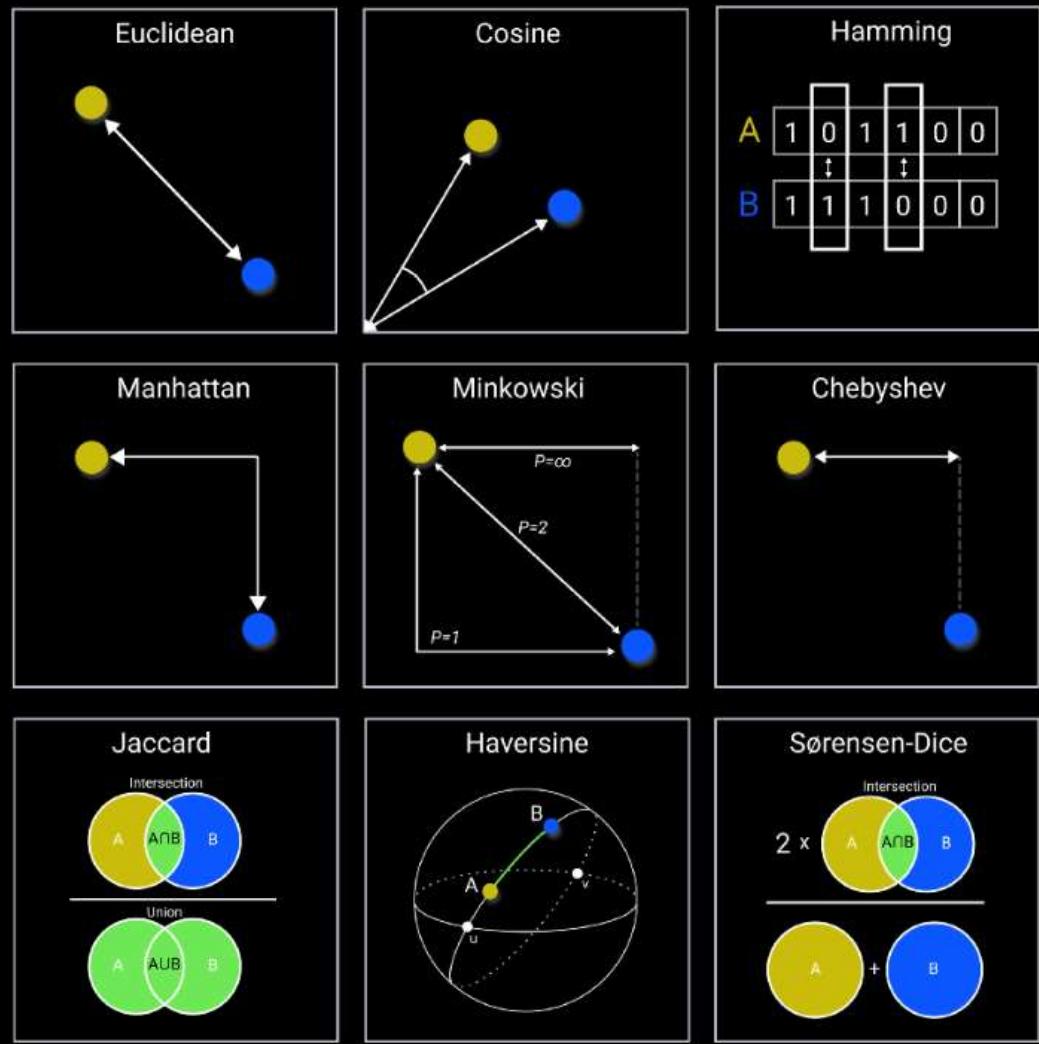
non-Euclidean



# HCA methods

## distance

- defines “nearness” or similarity

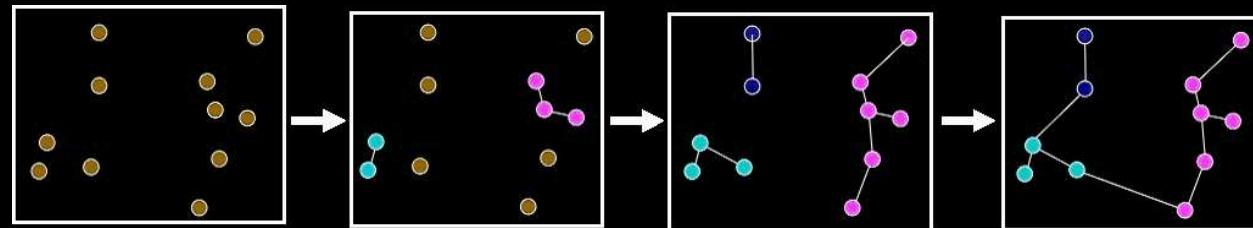


<https://towardsdatascience.com/9-distance-measures-in-data-science-918109d069fa>

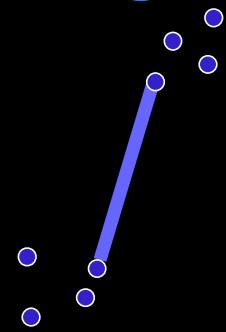
# HCA methods

## linkage or agglomeration

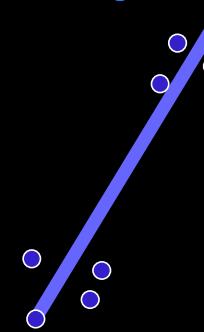
- how samples or variables are connected or grouped



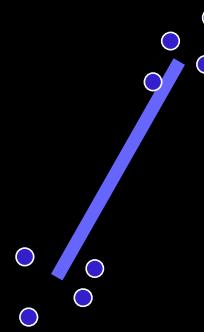
**single**



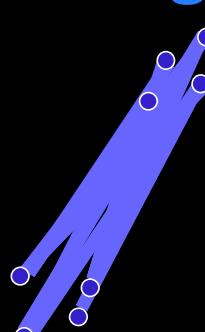
**complete**



**centroid**

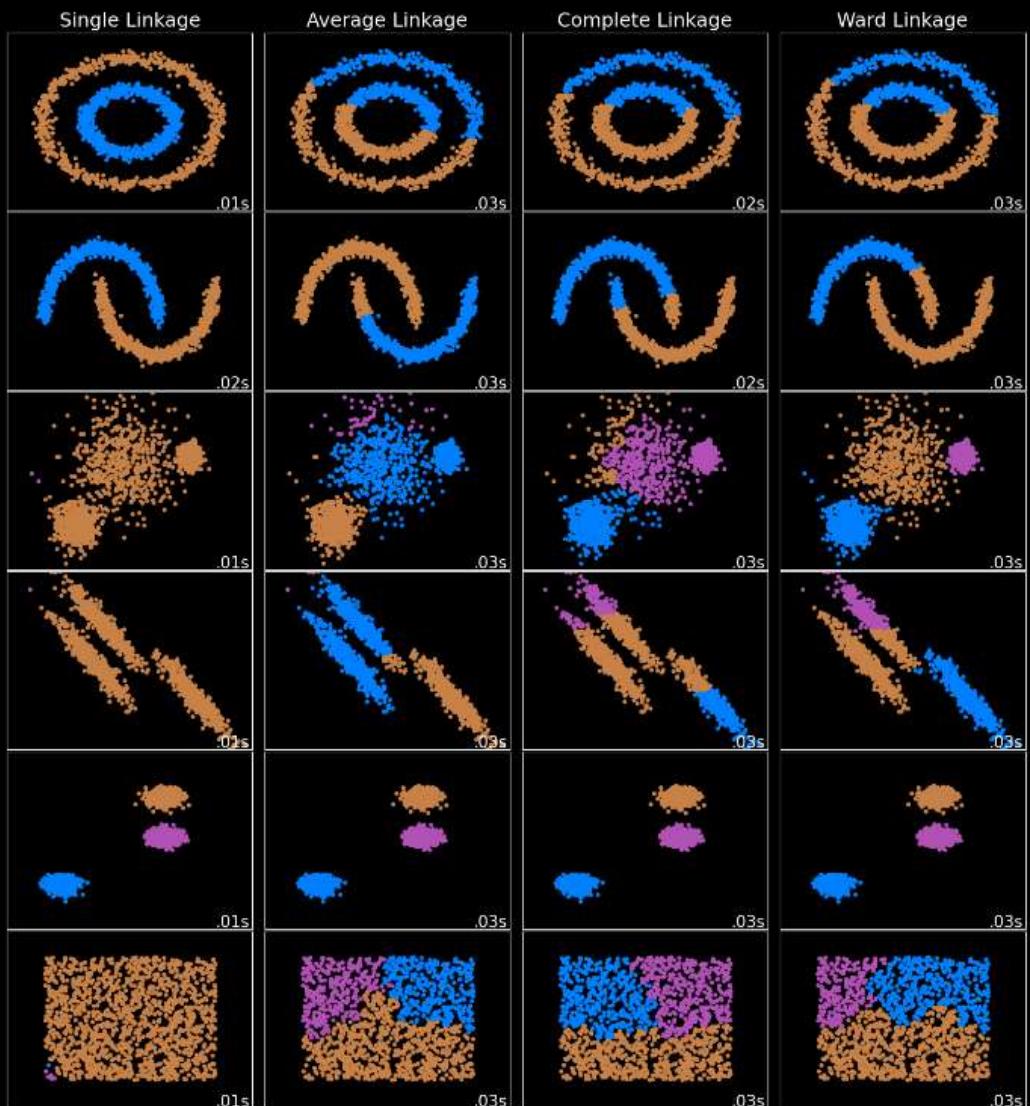


**average**

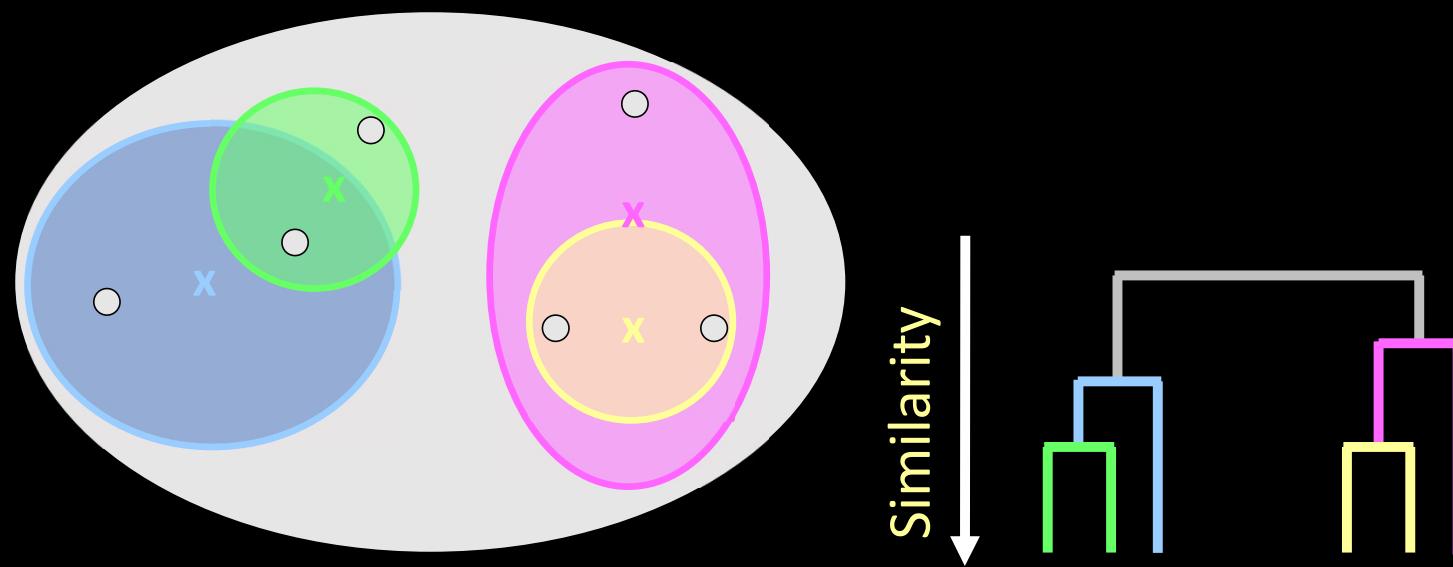


# HCA methods

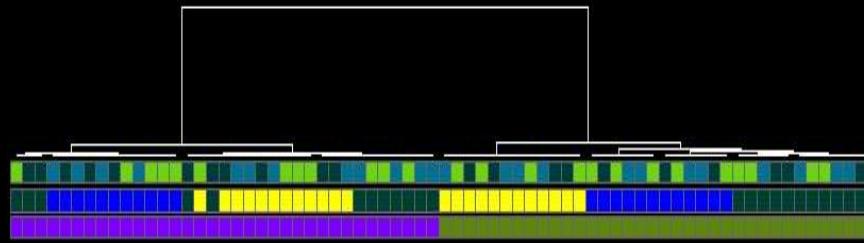
## linkage or agglomeration



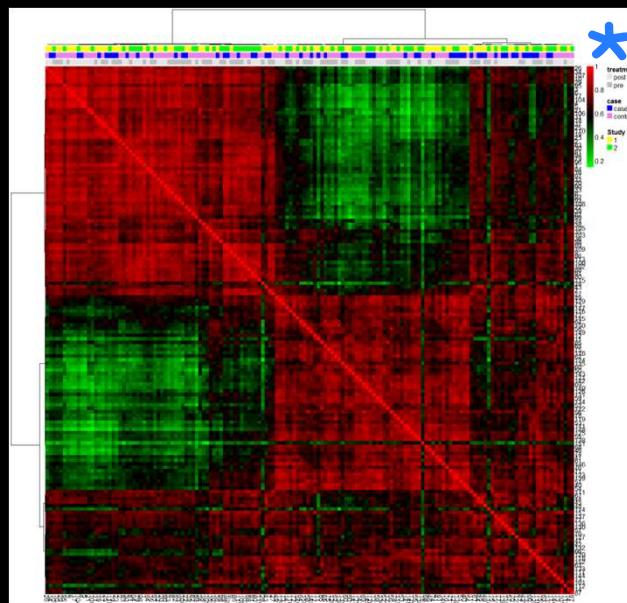
# HCA process



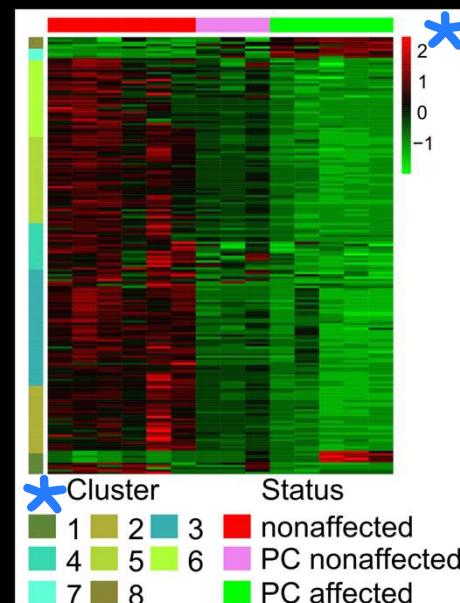
# HCA interpretation



overview



confirmation



How does my metadata \*  
match my data structure?

# Your turn

Follow along with the following tutorial:

[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/clustering/#heirarchical-clustering](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/clustering/#heirarchical-clustering)

# Principal Components Analysis (PCA)

reduce

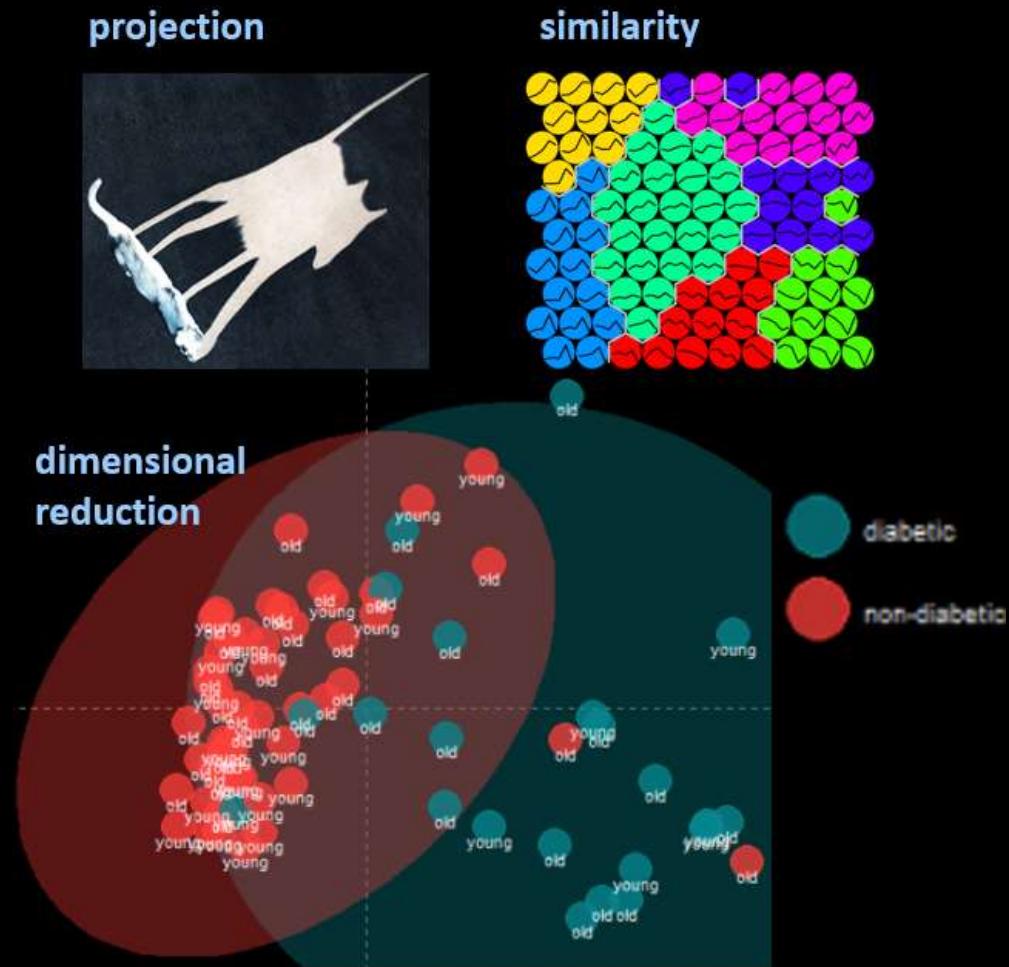
- dimensionality

maximize

- variance explained

visualize

- variance explained
- outliers
- sample scores
- variable loadings



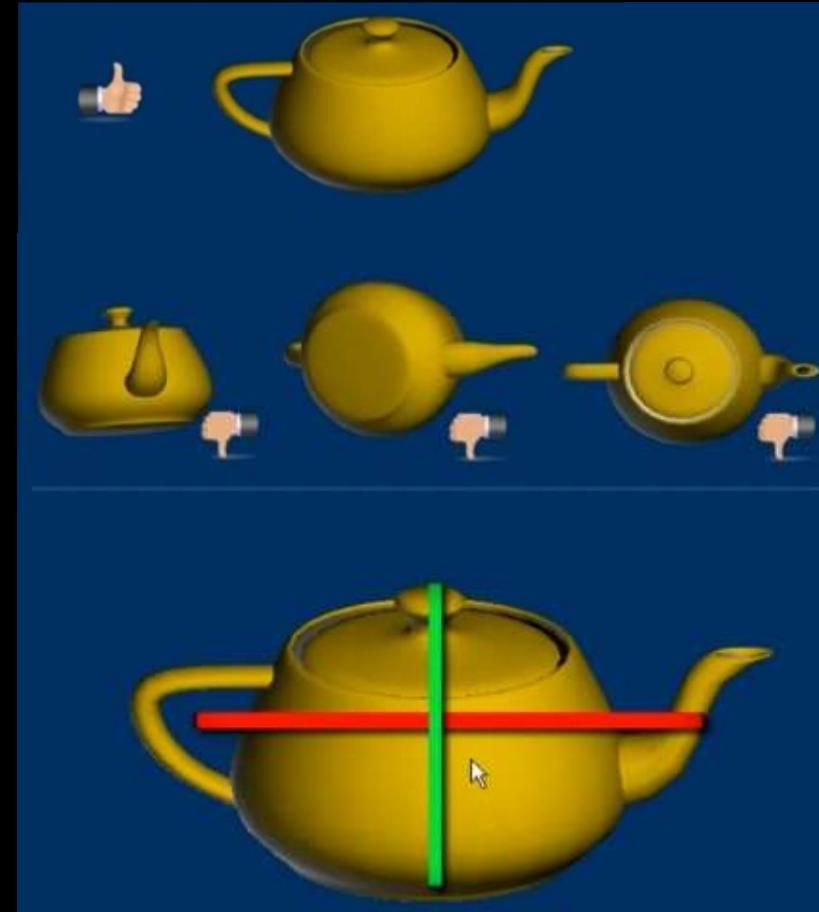
# PCA goals

## Principal Components (PCs)

- non-supervised
- projection of the data which maximize variance explained

## results

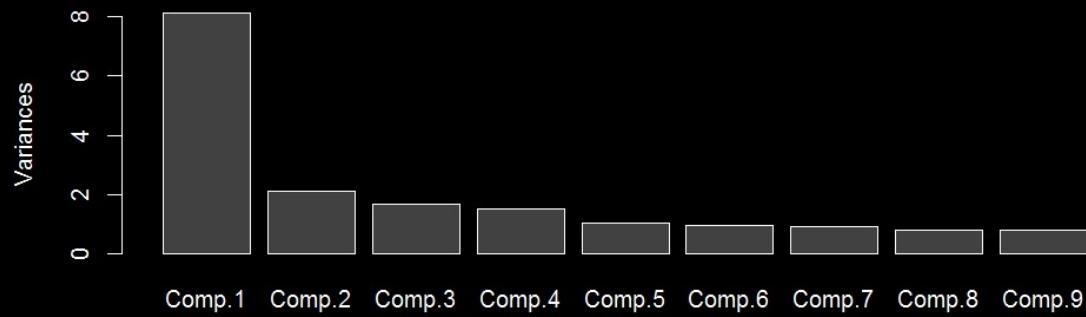
1. eigenvalues = variance explained
2. scores = new coordinates for samples (rows)
3. loadings = linear combination of original variables



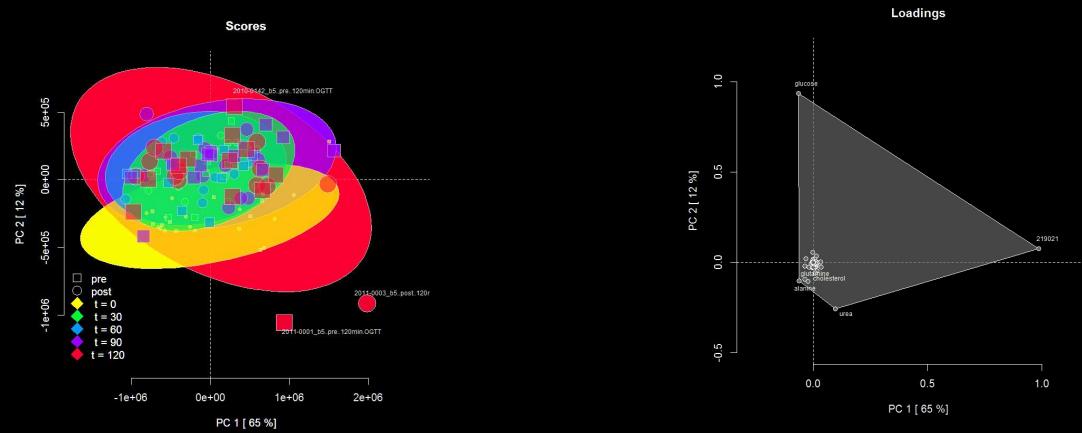
James X. Li, 2009, VisuMap Tech.

# PCA interpretation

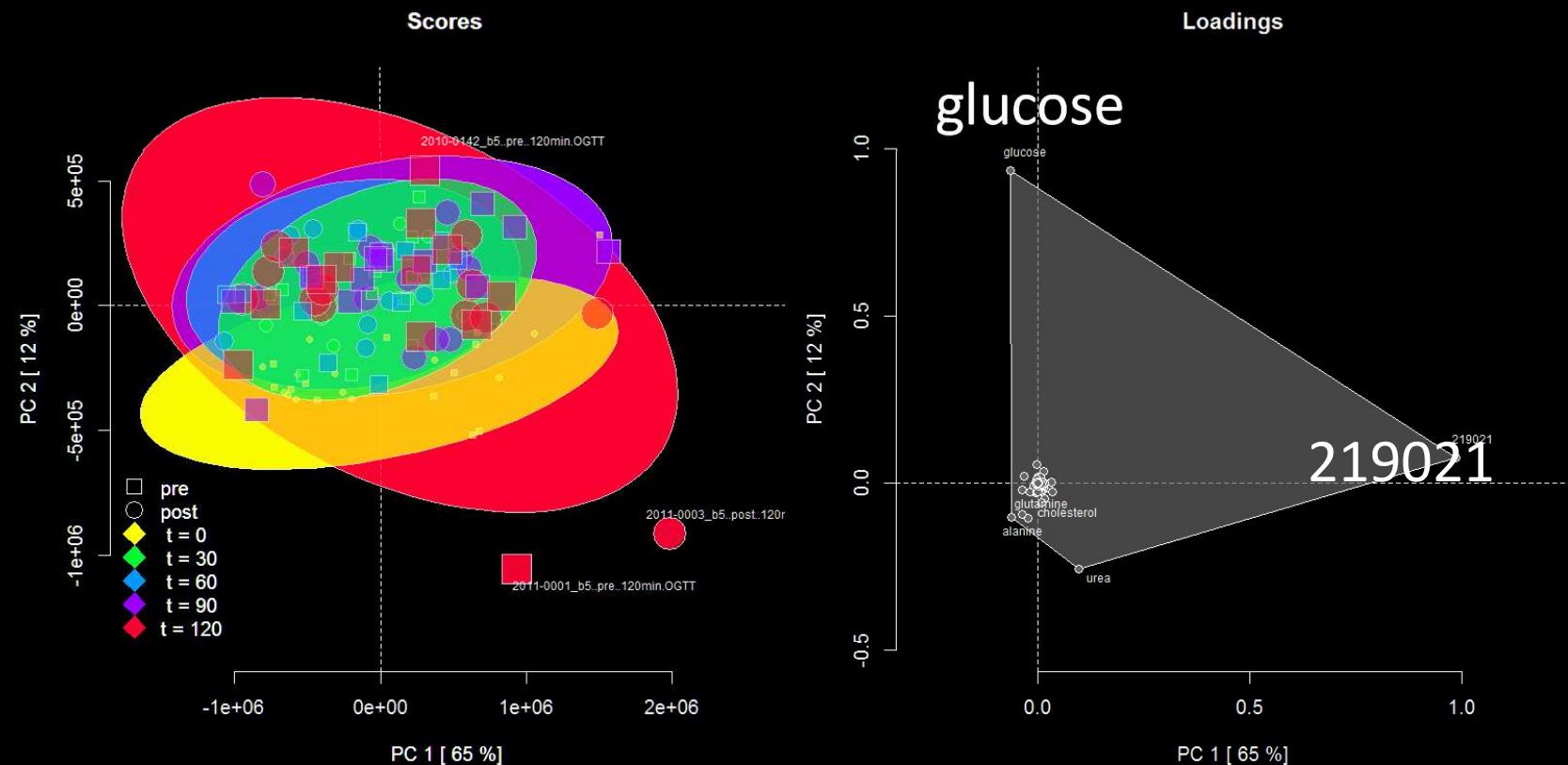
Variance explained (eigenvalues)



Row (sample) scores and column (variable) loadings



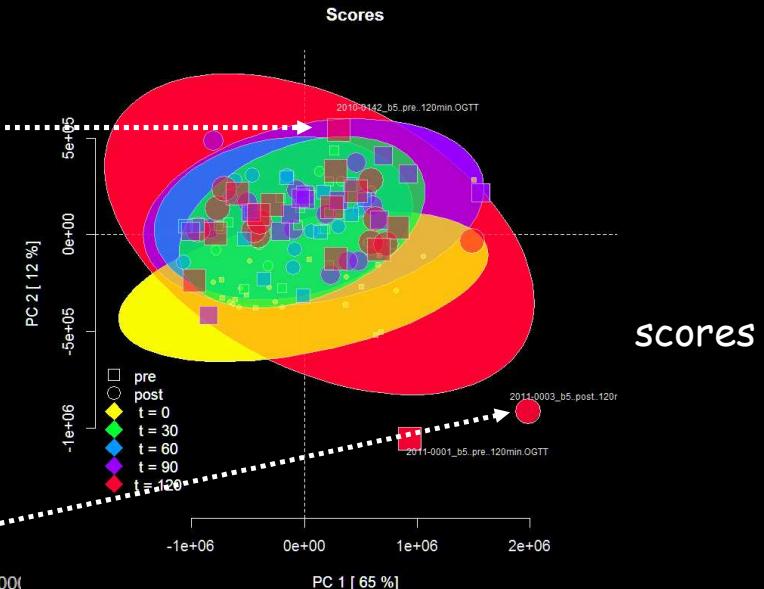
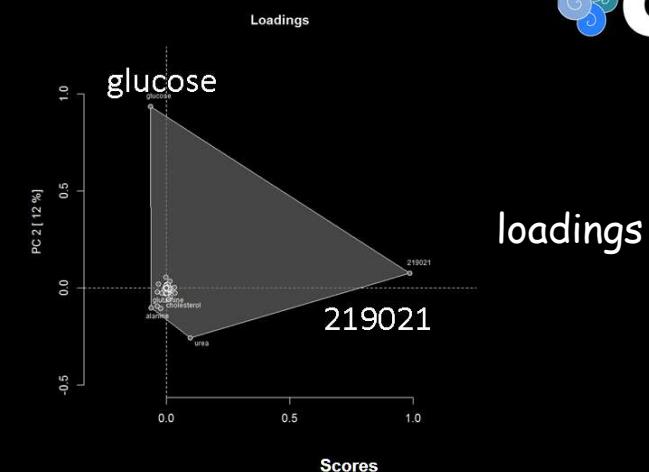
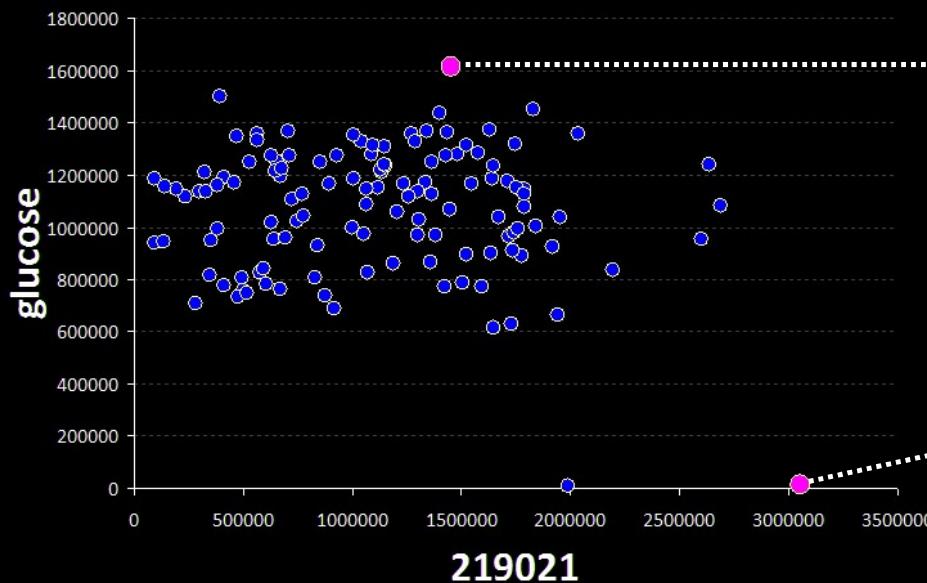
# PCA example



\*no scaling or centering

# Relationship between scores and loadings

top loading variable's scatterplot



# Your turn

Follow along with the following tutorial:

[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/multivariate/#pca](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/multivariate/#pca)

# Machine learning

**predict**

- sample classification

**optimize**

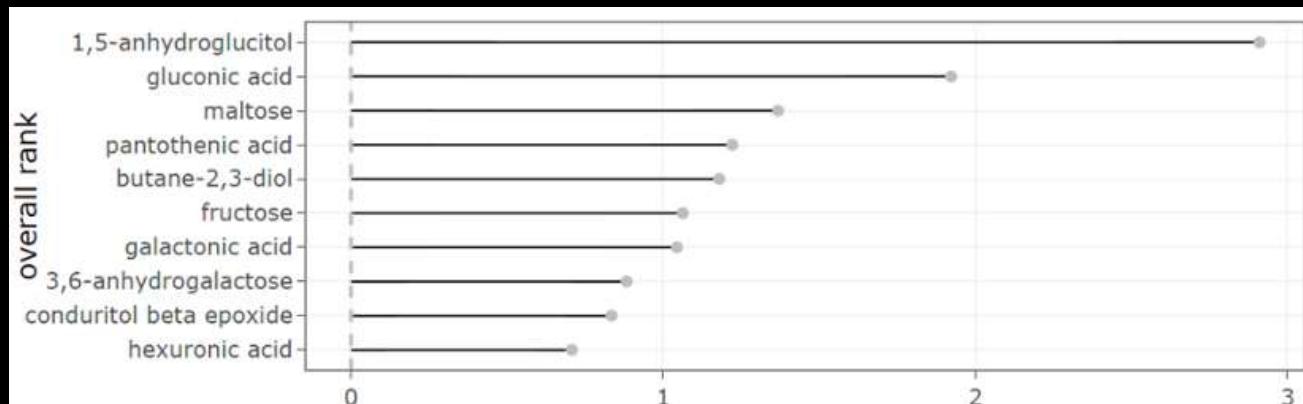
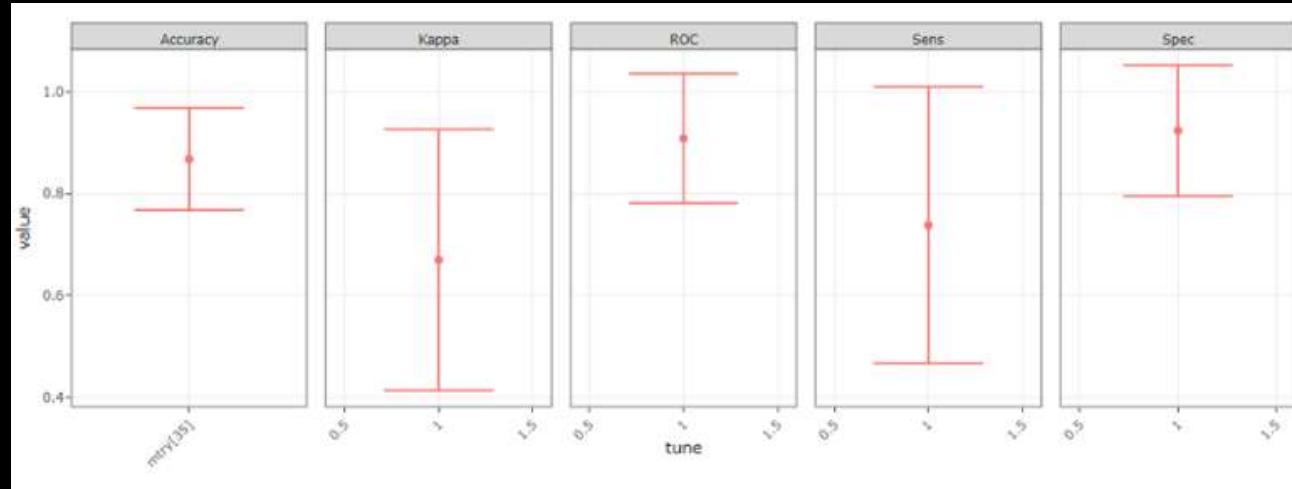
- model performance

**select**

- important features

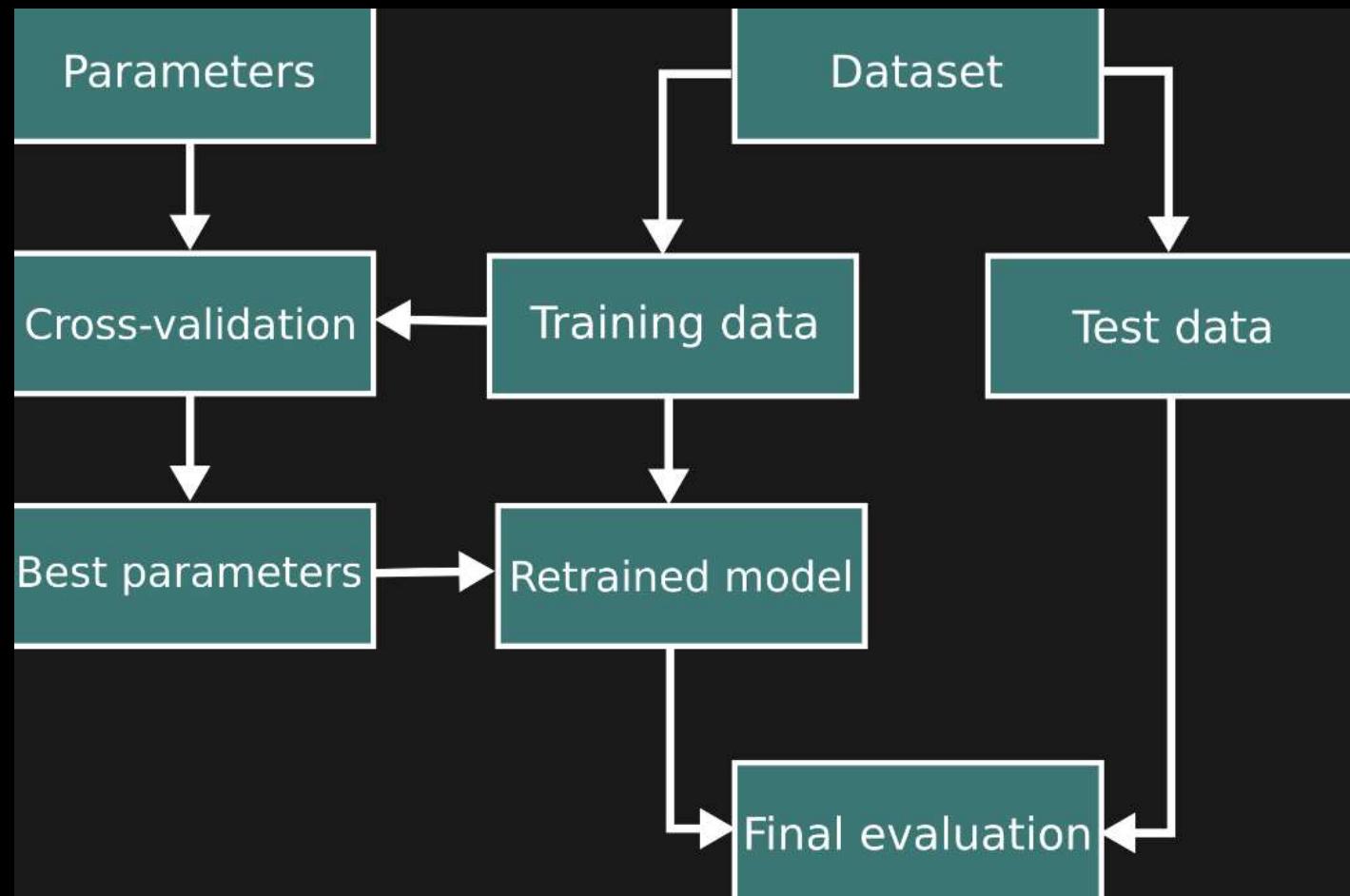
**visualize**

- model performance
- feature importance



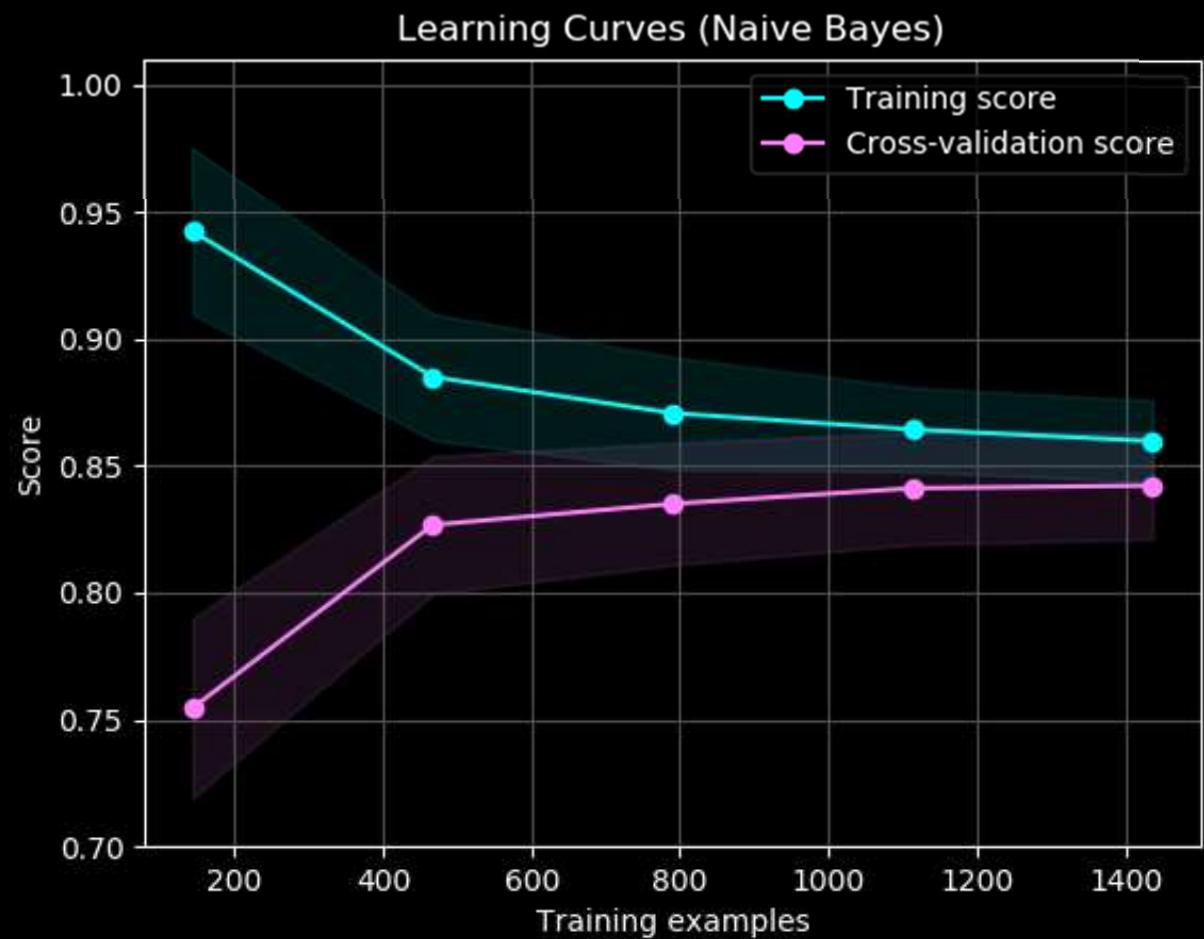
# Machine learning

model validation



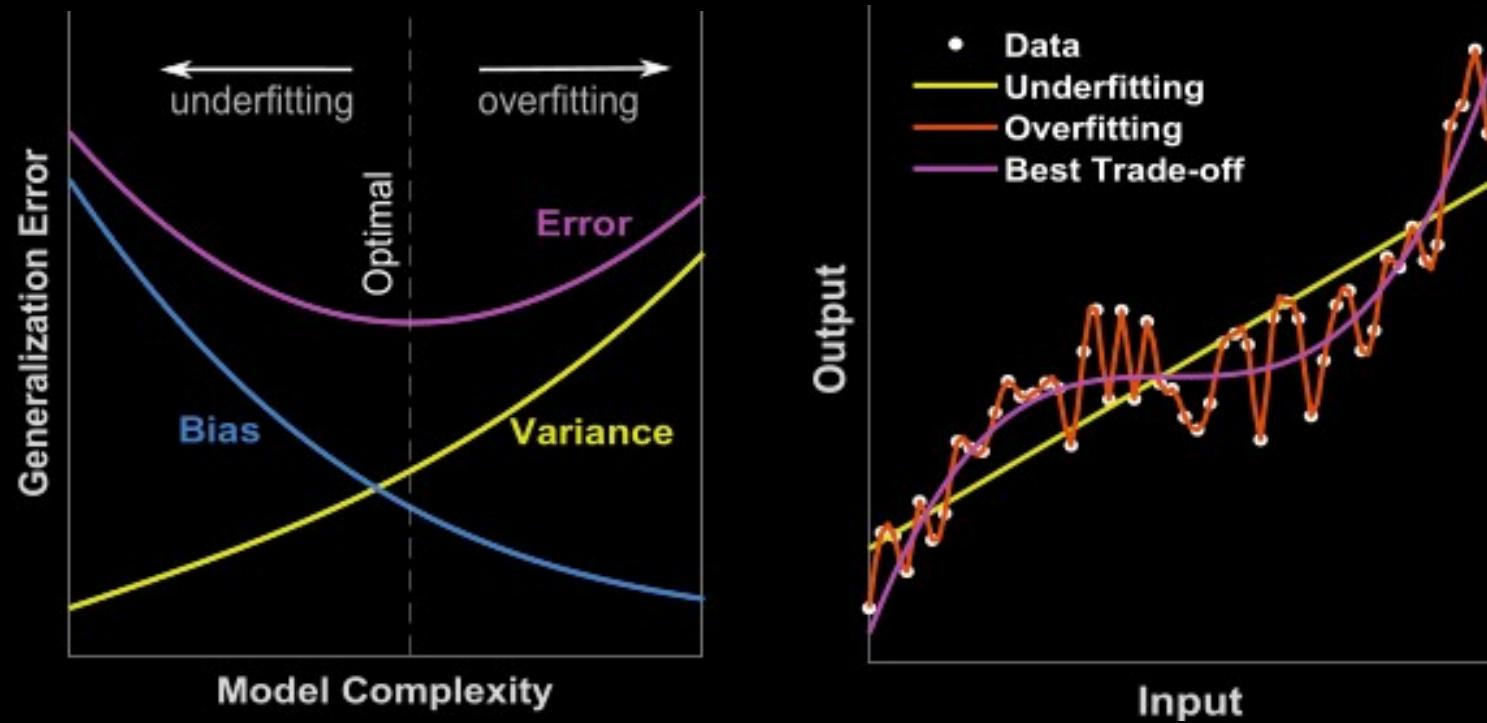
# Machine learning

model training  
and validation



# Machine learning

## bias vs. variance



# Machine learning

## classification performance

		Predicted condition	
		Positive (PP)	Negative (PN)
Total population = P + N			
Actual condition	Positive (P)	True positive (TP), hit	False negative (FN), type II error, miss, underestimation
	Negative (N)	False positive (FP), type I error, false alarm, overestimation	True negative (TN), correct rejection

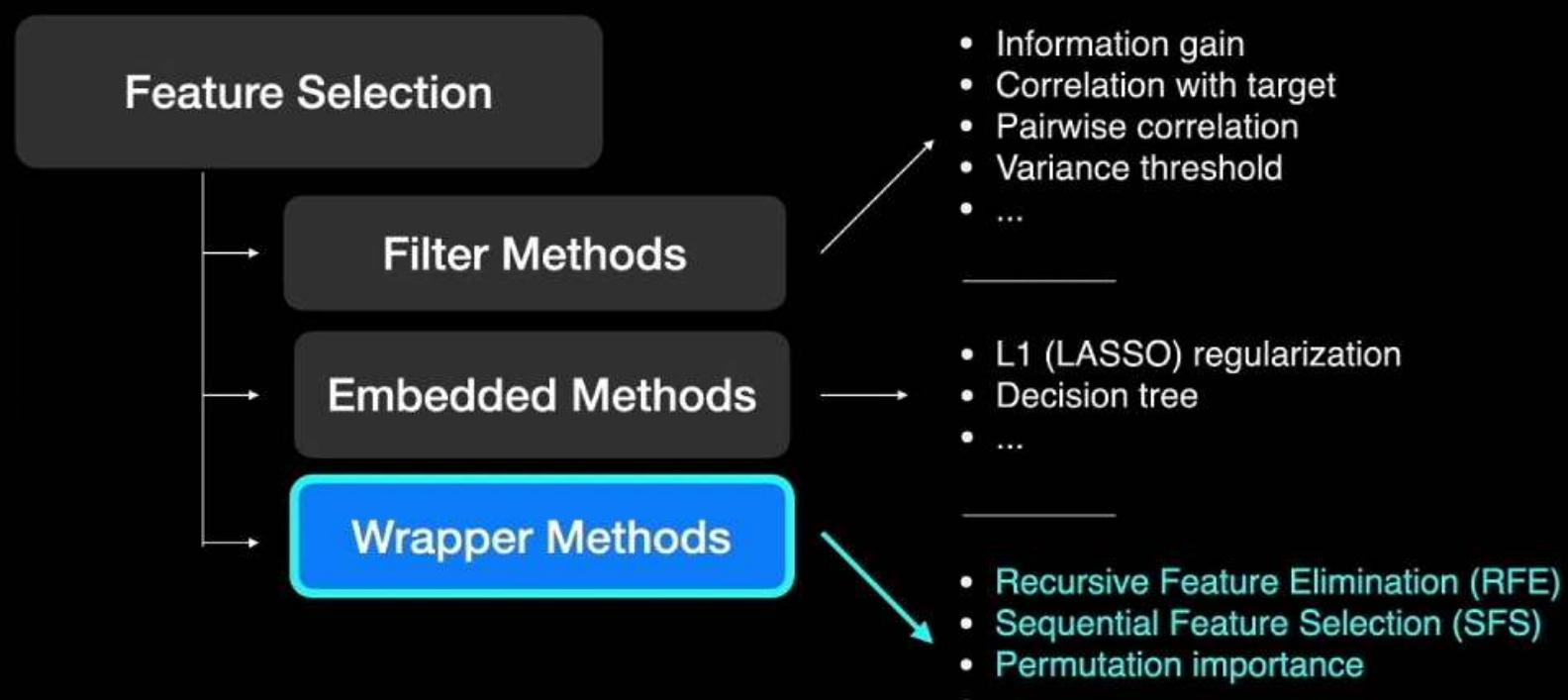
# Machine learning

## cross-validation



# Machine learning

## feature selection

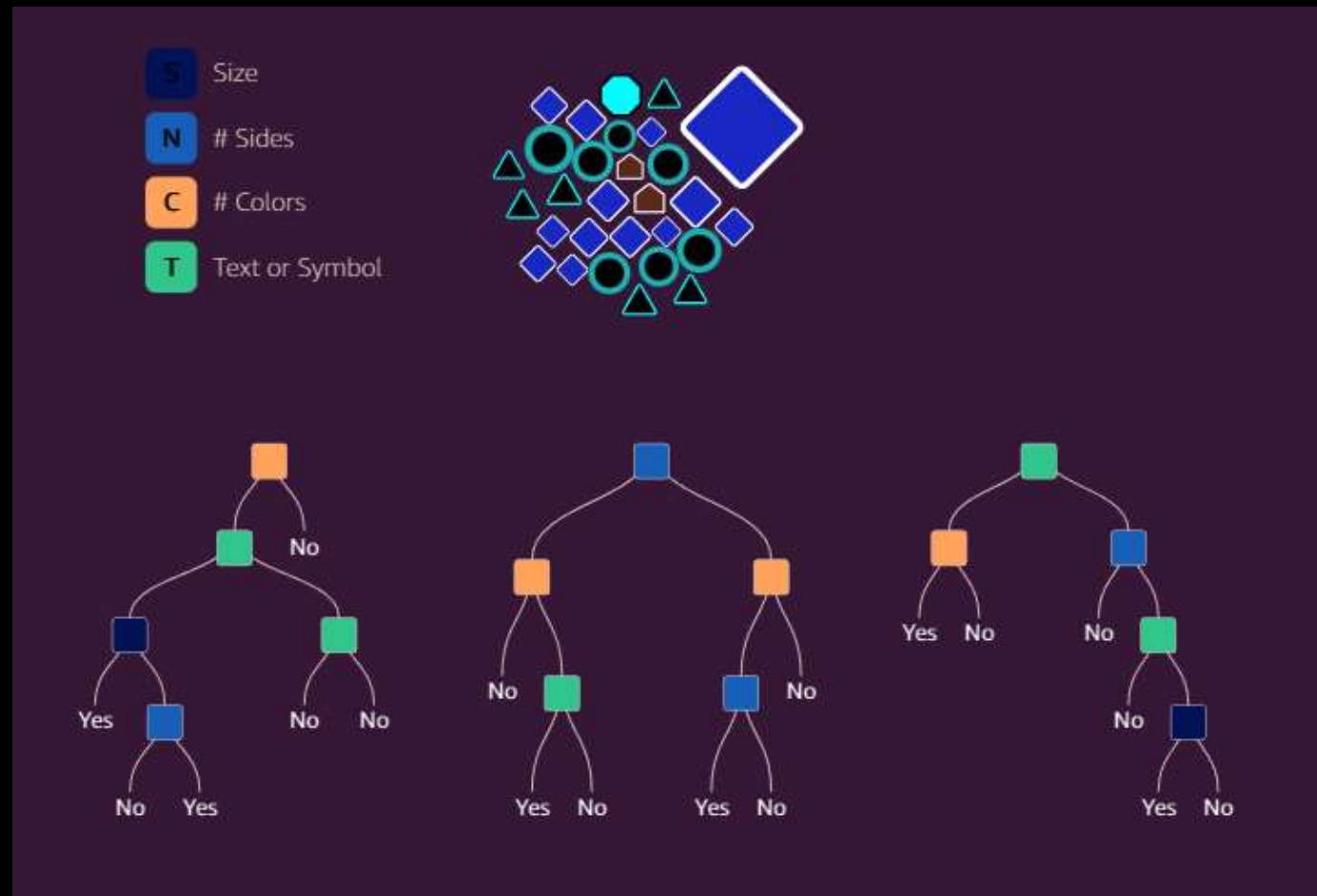


[https://scikit-learn.org/stable/modules/feature\\_selection.html#feature-selection](https://scikit-learn.org/stable/modules/feature_selection.html#feature-selection)

# Machine learning

## Random Forest

- cross-validation
- feature selection

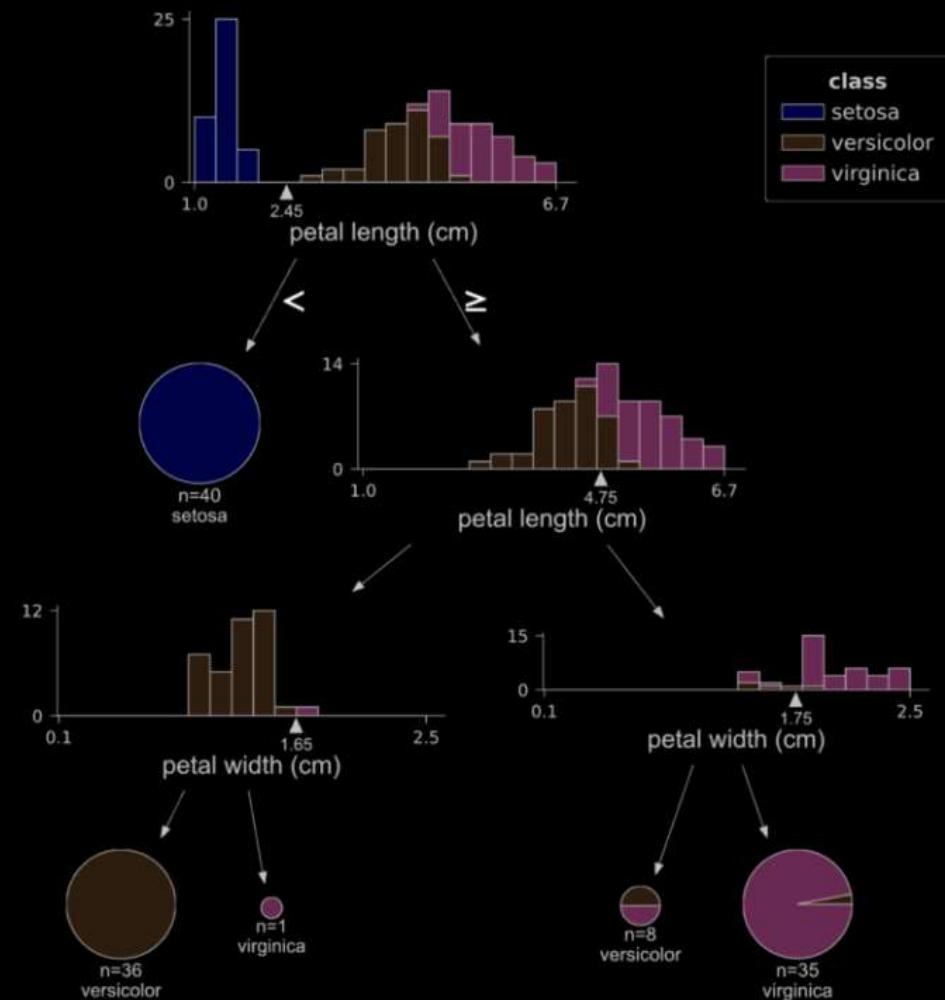


<https://mlu-explain.github.io/random-forest/>

# Machine learning

## Random Forest

- decision path

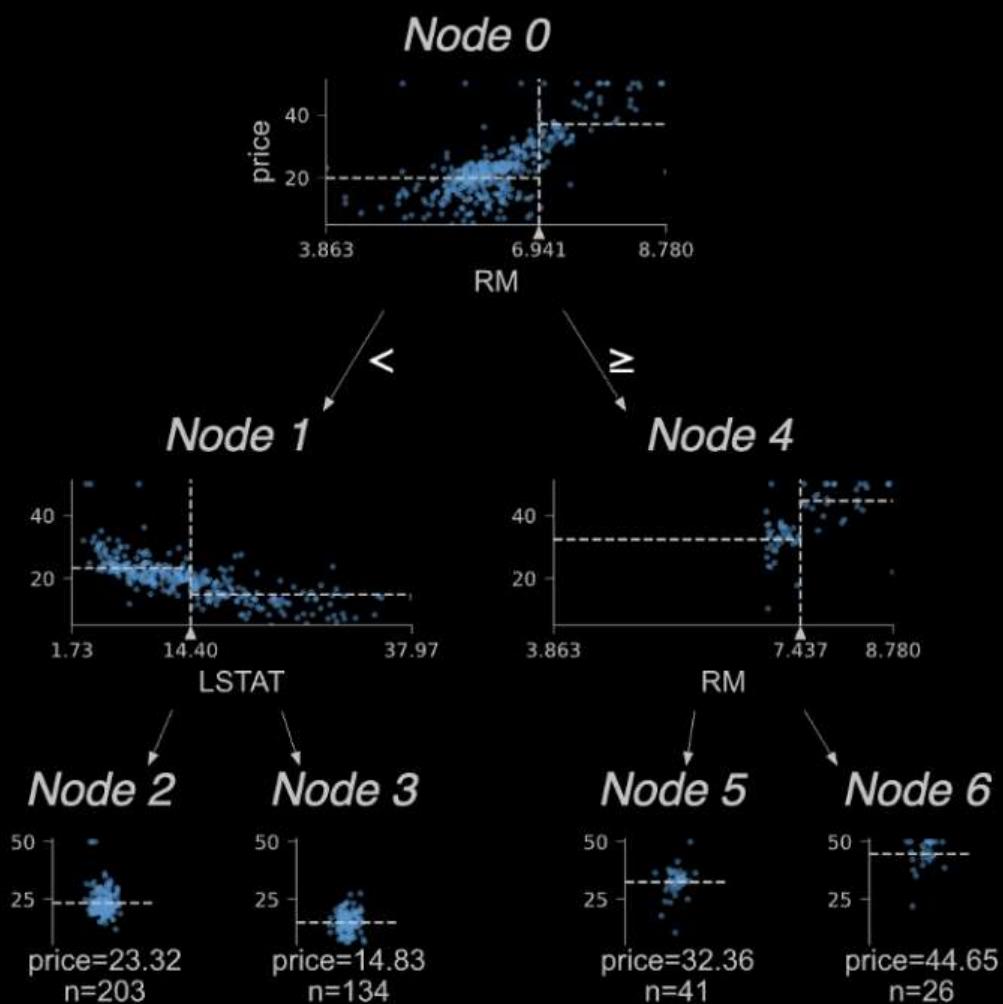


<https://github.com/parrt/dtreeviz>

# Machine learning

## Random Forest

- decision path



# Machine learning

autoML

	Model	Accuracy	AUC	Recall	Prec.	F1	Kappa	MCC	TT (Sec)
<b>catboost</b>	CatBoost Classifier	0.7767	0.8309	0.6056	0.7114	0.6413	0.4823	0.4950	1.3870
<b>lr</b>	Logistic Regression	0.7564	0.8043	0.5056	0.6941	0.5786	0.4145	0.4285	1.1230
<b>gbc</b>	Gradient Boosting Classifier	0.7562	0.8239	0.5667	0.6731	0.6031	0.4314	0.4431	0.0900
<b>ada</b>	Ada Boost Classifier	0.7526	0.8016	0.5889	0.6524	0.6091	0.4310	0.4394	0.0800
<b>lightgbm</b>	Light Gradient Boosting Machine	0.7524	0.8028	0.5778	0.6614	0.6086	0.4299	0.4381	0.1430
<b>rf</b>	Random Forest Classifier	0.7488	0.8035	0.5111	0.6849	0.5740	0.4023	0.4182	0.2350
<b>ridge</b>	Ridge Classifier	0.7452	0.0000	0.4722	0.6844	0.5492	0.3816	0.3997	0.0150
<b>lda</b>	Linear Discriminant Analysis	0.7452	0.7912	0.4833	0.6783	0.5563	0.3859	0.4017	0.0130
<b>xgboost</b>	Extreme Gradient Boosting	0.7449	0.7896	0.5722	0.6442	0.5984	0.4140	0.4207	0.2640
<b>knn</b>	K Neighbors Classifier	0.7153	0.7261	0.5111	0.5962	0.5405	0.3379	0.3467	0.0220
<b>et</b>	Extra Trees Classifier	0.7134	0.7573	0.4333	0.6079	0.4968	0.3072	0.3204	0.1810
<b>dt</b>	Decision Tree Classifier	0.7075	0.6741	0.5722	0.5635	0.5630	0.3445	0.3481	0.0130
<b>nb</b>	Naive Bayes	0.6817	0.7064	0.2389	0.5527	0.3288	0.1657	0.1905	0.0110
<b>svm</b>	SVM - Linear Kernel	0.6015	0.0000	0.3611	0.3419	0.3251	0.0851	0.0924	0.0170
<b>qda</b>	Quadratic Discriminant Analysis	0.5759	0.5889	0.4833	0.4062	0.3705	0.1011	0.1281	0.0180

# Your turn

Follow along with the following tutorial:

[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/model/](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/model/)

# Network analysis

## network mapping

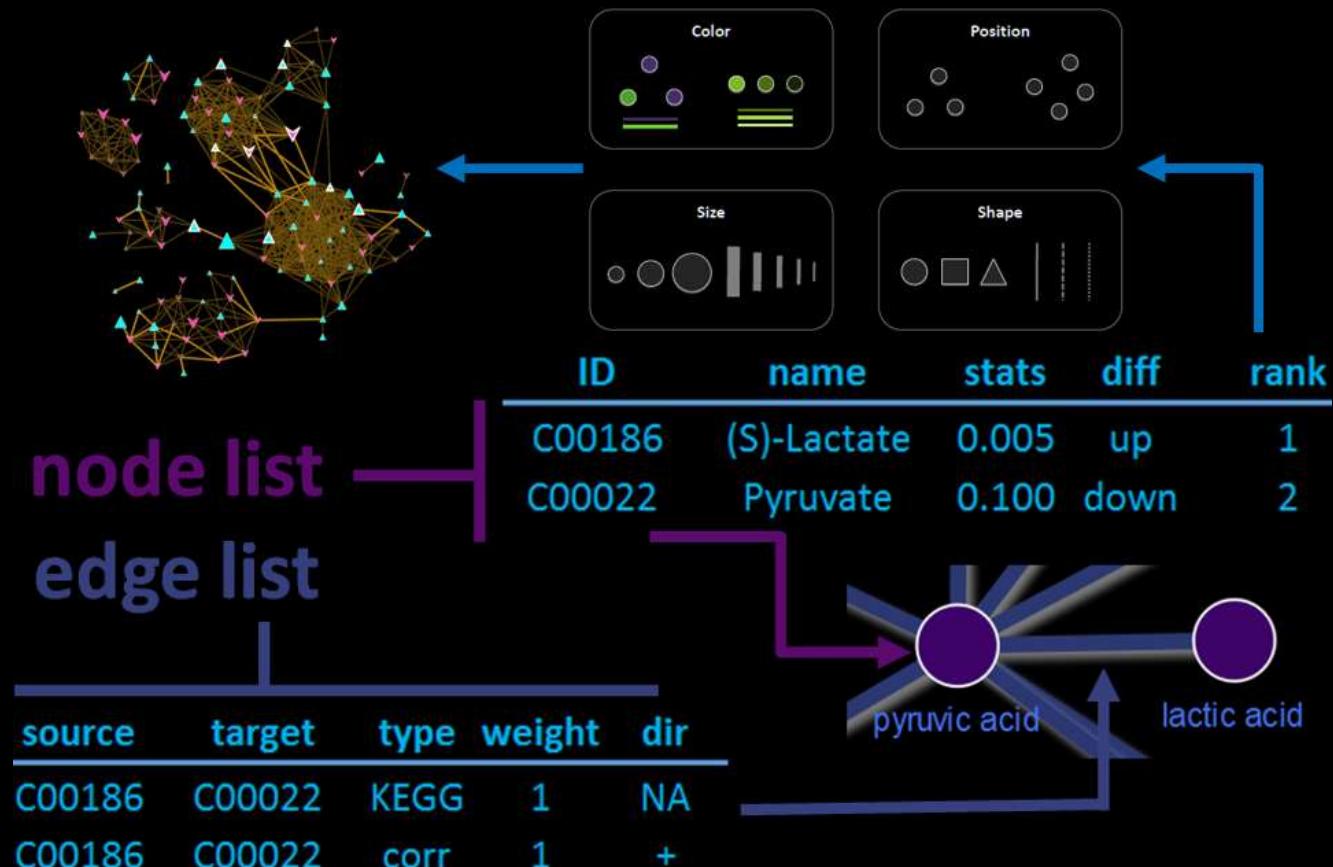
- transform variables

## network calculation

- regularized correlation
- biochemical
- structural similarity
- model performance

## visualize

- interactive networks



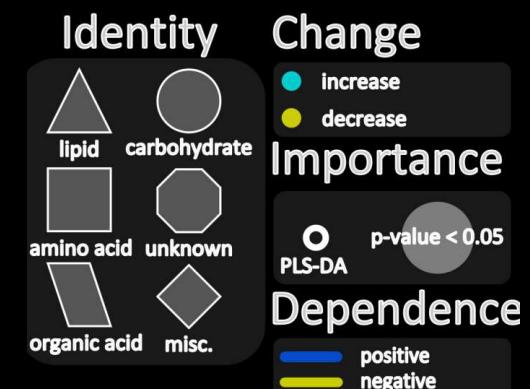
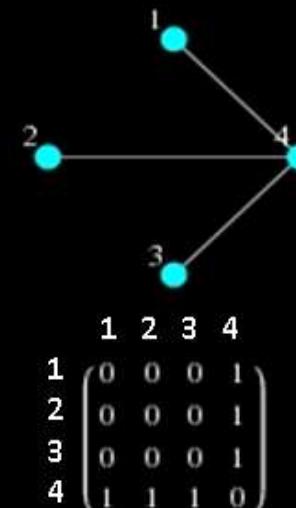
# Components for network mapping

## connections (edges)

- empirical dependency (correlation)
- biochemical (substrate/product)
- chemical similarity
- ...

## nodes (vertices)

- magnitude
- importance
- direction
- relationships
- ...



# Network data structures

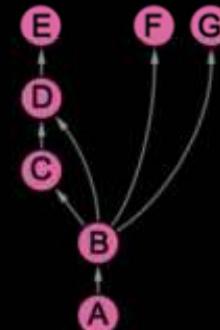
## adjacency matrix

Undirected



	A	B	C	D	E	F	G	Degree
A	0	1	1	1	1	1	0	5
B	1	0	0	0	0	1	0	2
C	1	0	0	0	0	0	0	1
D	1	0	0	0	0	0	0	1
E	1	0	0	0	0	0	0	1
F	1	1	0	0	0	0	1	3
G	0	0	0	0	0	1	0	1

Directed



	A	B	C	D	E	F	G	Out-degree
A	0	1	0	0	0	0	0	1
B	0	0	1	1	0	1	1	4
C	0	0	0	1	0	0	0	1
D	0	0	0	0	1	0	0	1
E	0	0	0	0	0	0	0	0
F	0	0	0	0	0	0	0	0
G	0	0	0	0	0	0	0	0

Weighted

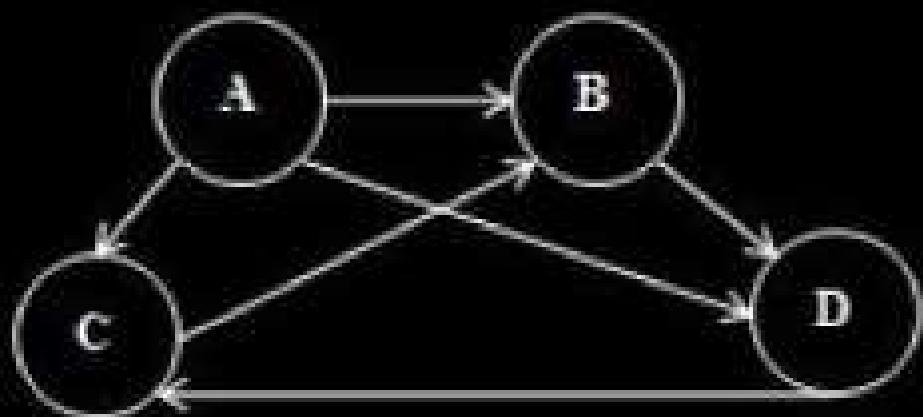


	A	B	C	D	E	F	G	Degree
A	0	8	12	12	12	16	12	72
B	8	0	0	0	0	4	0	12
C	12	0	0	0	0	0	0	12
D	12	0	0	0	0	0	0	12
E	12	0	0	0	0	0	0	12
F	16	4	0	0	0	0	12	32
G	12	0	0	0	0	12	0	24

<https://www.steveclarkapps.com/graphs/>

# Network data structures

## adjacency list



Node	Adjacent Node(s)
A	B C D
B	D
C	B
D	C

# Network data structures

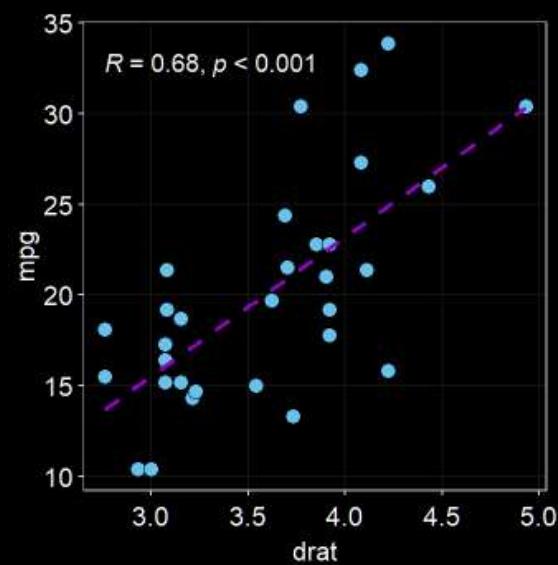
edge list and  
vertices list  
(node attributes)

edgesList		verticesList	index	...
source	target			
1	2	1		
1	4	2		
1	5	3		
2	3	4		
2	5	5		
2	6	6		
3	6			
4	5			
5	6			

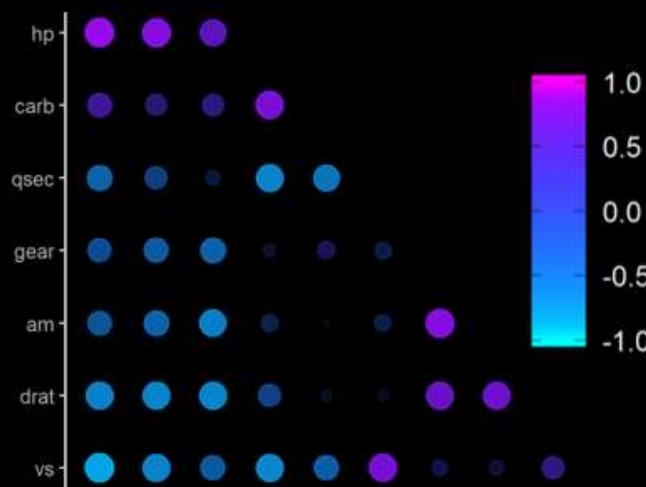
# Correlation networks

Connect molecules based on strength of their correlation or partial-correlation

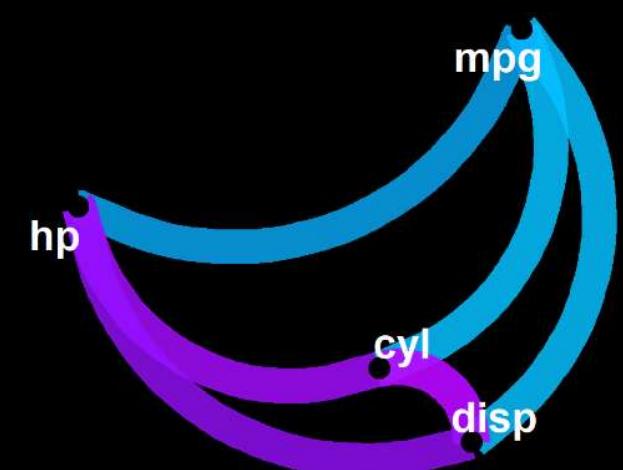
bivariate



multivariate

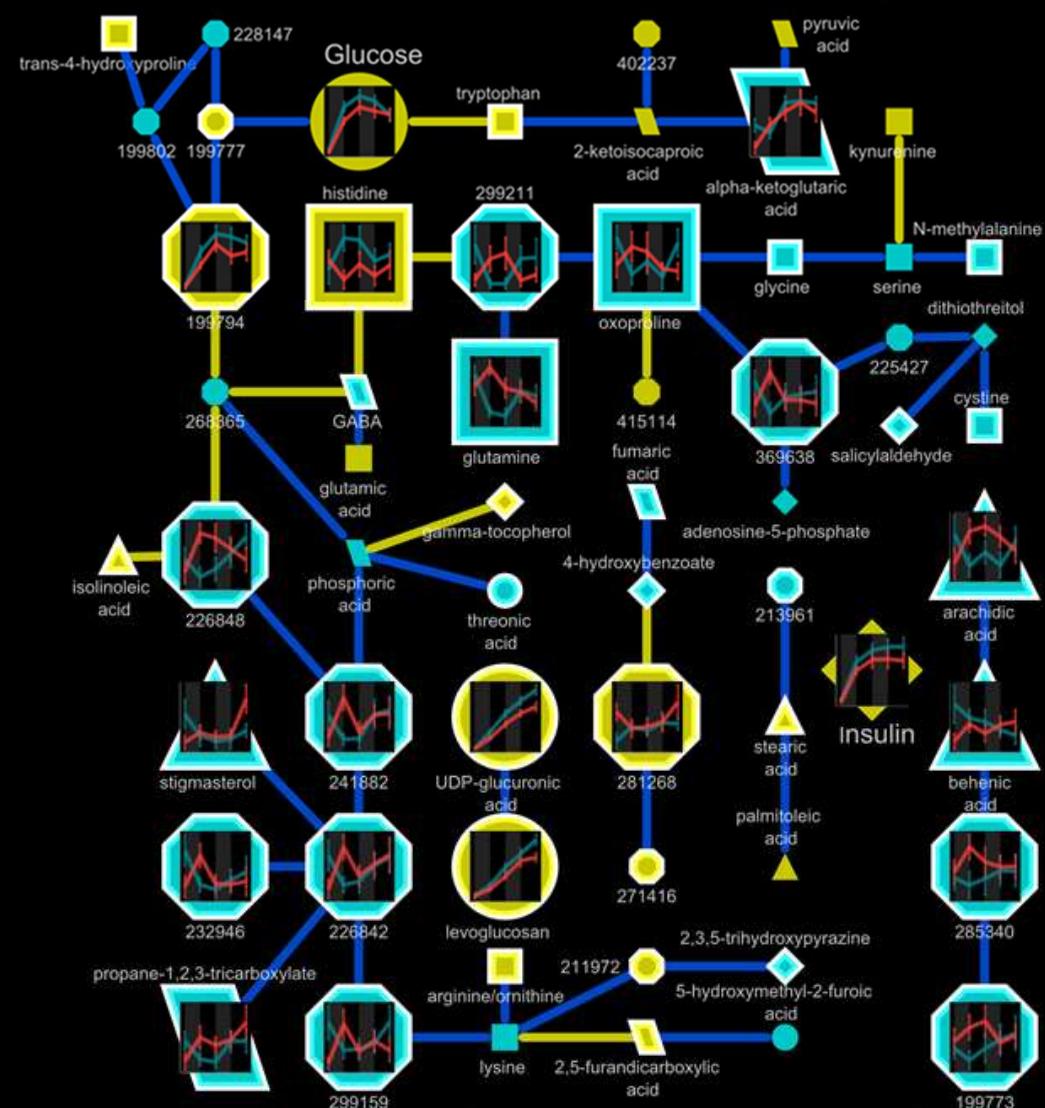


network

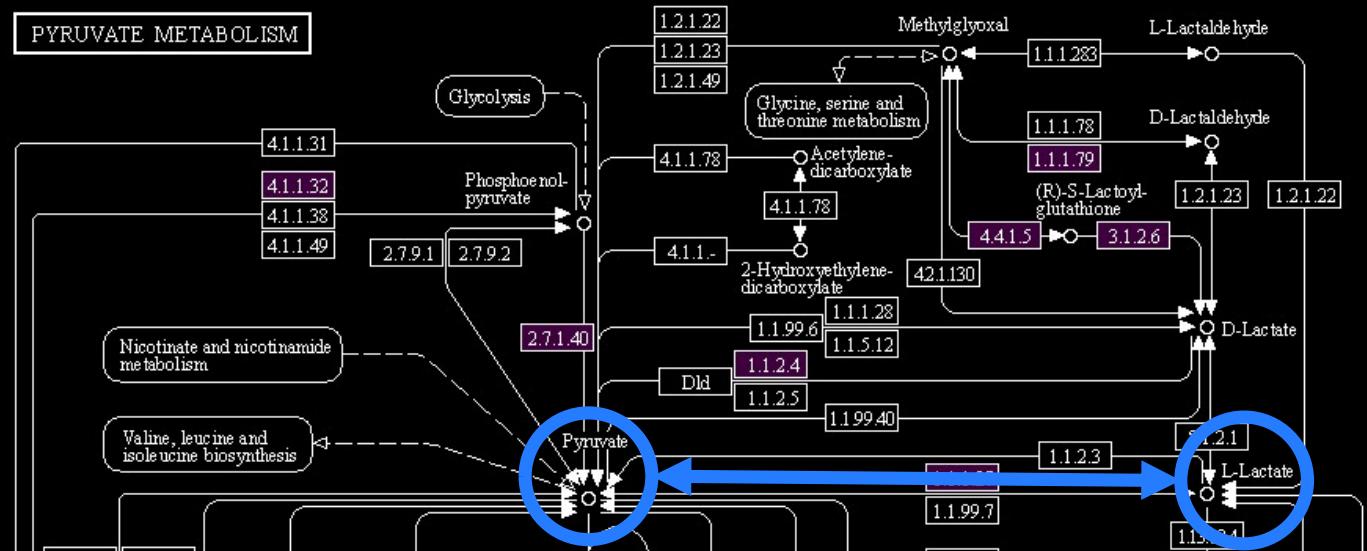
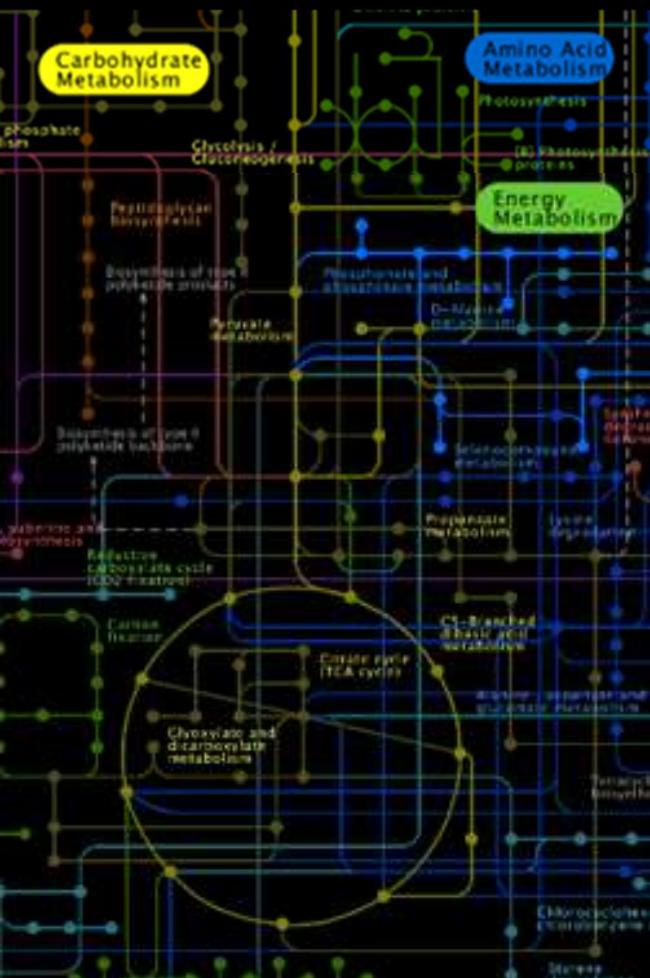


# Correlation example

regularized correlation network showing relationships in metabolic timeseries measurements for two classes of samples



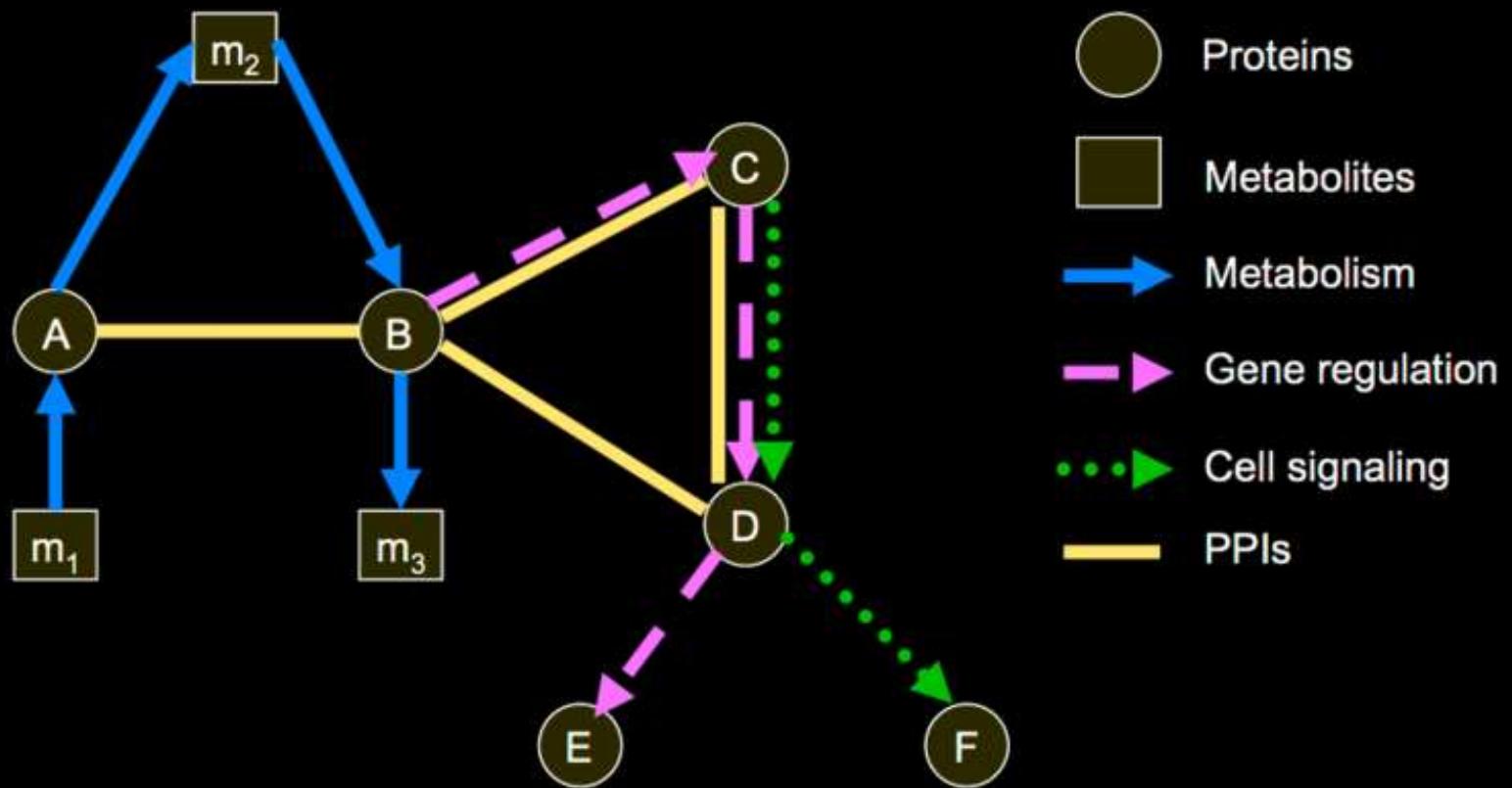
# Biochemical networks



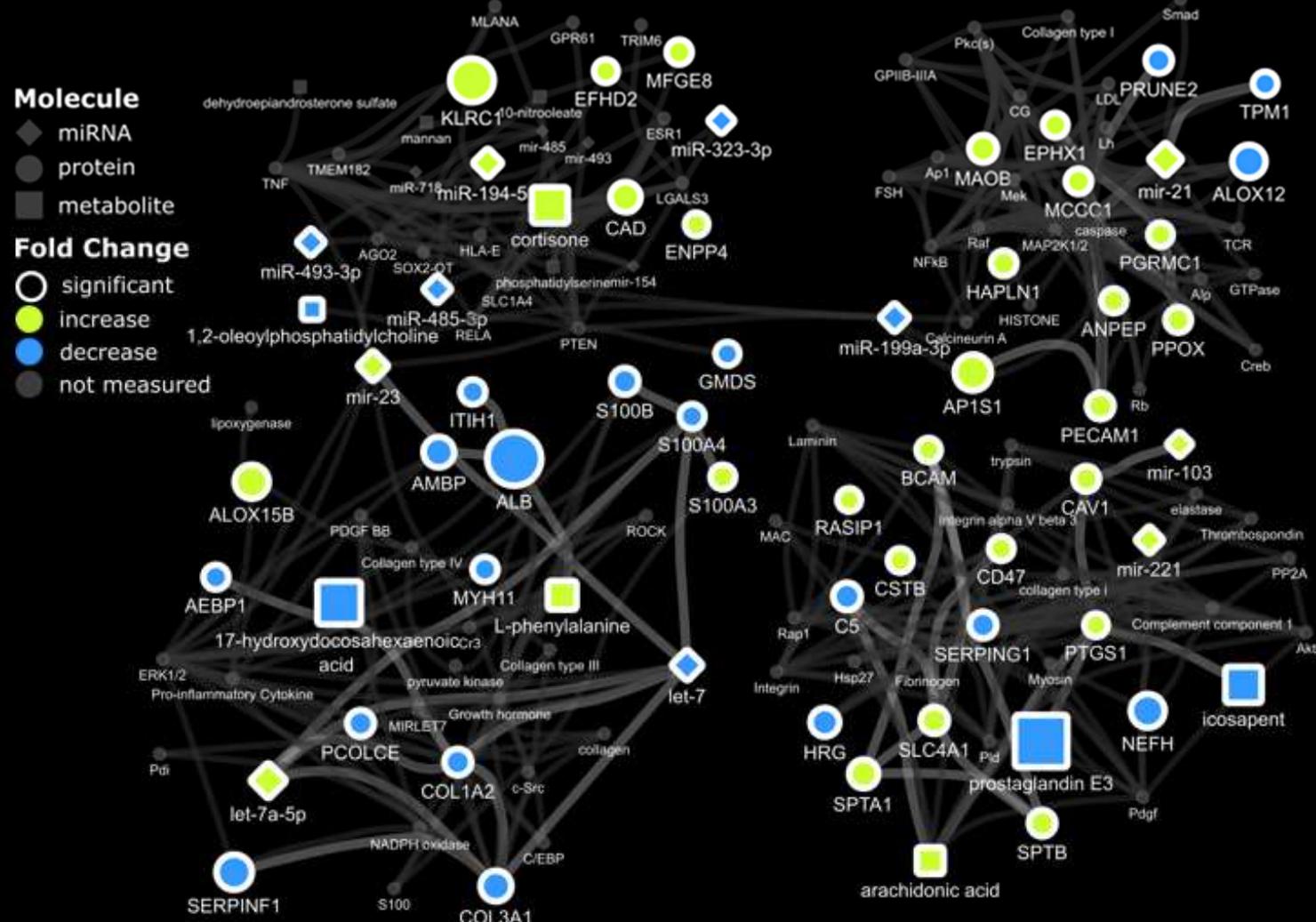
	ID	name	stats	diff	rank
nodes	C00186	(S)-Lactate	0.005	up	1
	C00022	Pyruvate	0.100	down	2
edges	source	target	type	weight	dir
	C00186	C00022	KEGG	1	NA



# Multi-Omic networks



# Multi-Omic networks



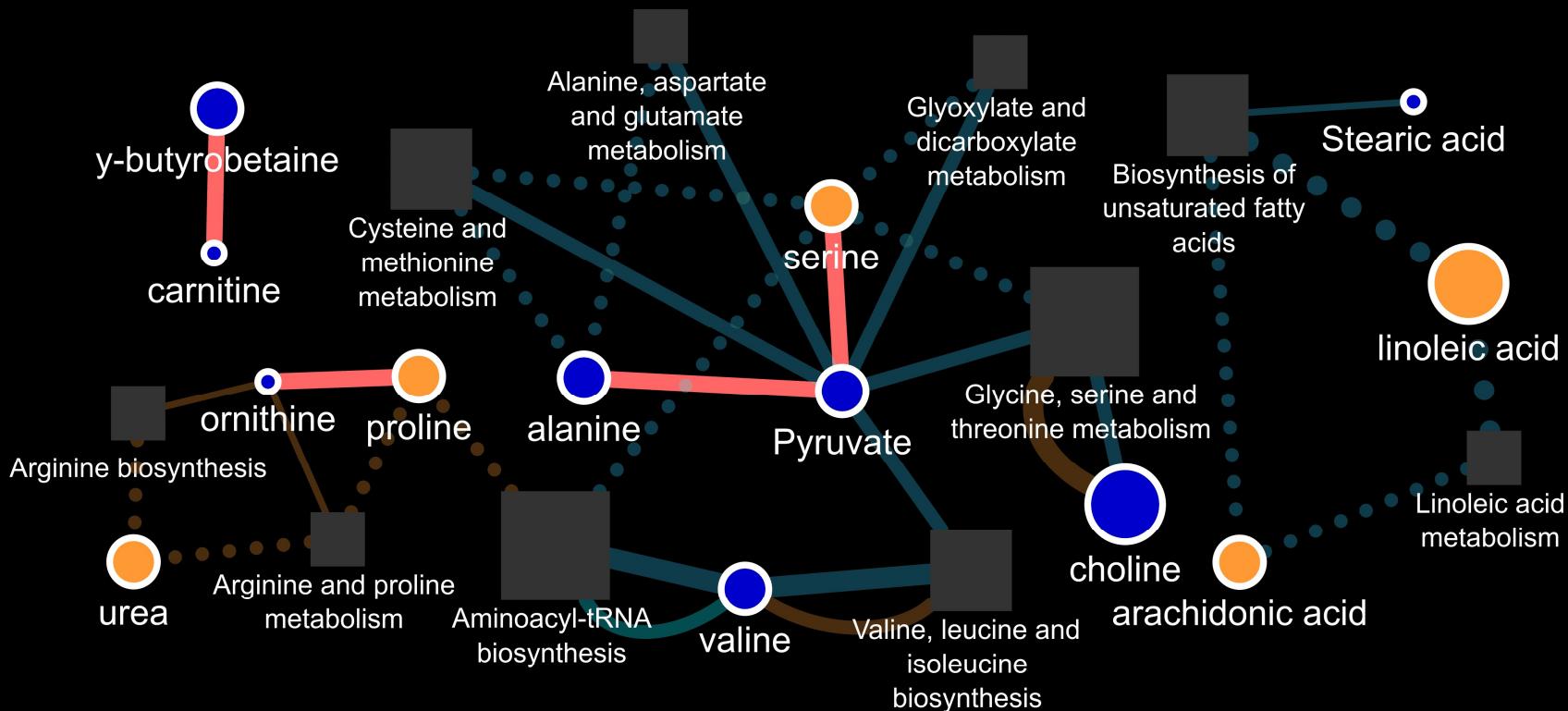
# Multi-Omic networks

## Type

- metabolite
- pathway
- increase
- decrease

## Relationship

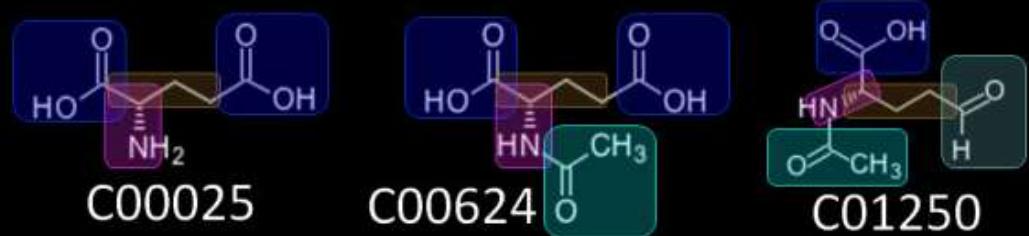
- biochemical
- blood
- saliva
- decrease
- increase



# Structural similarity networks

- Use structure to generate molecular fingerprint
- Calculate similarities between metabolites based on fingerprint

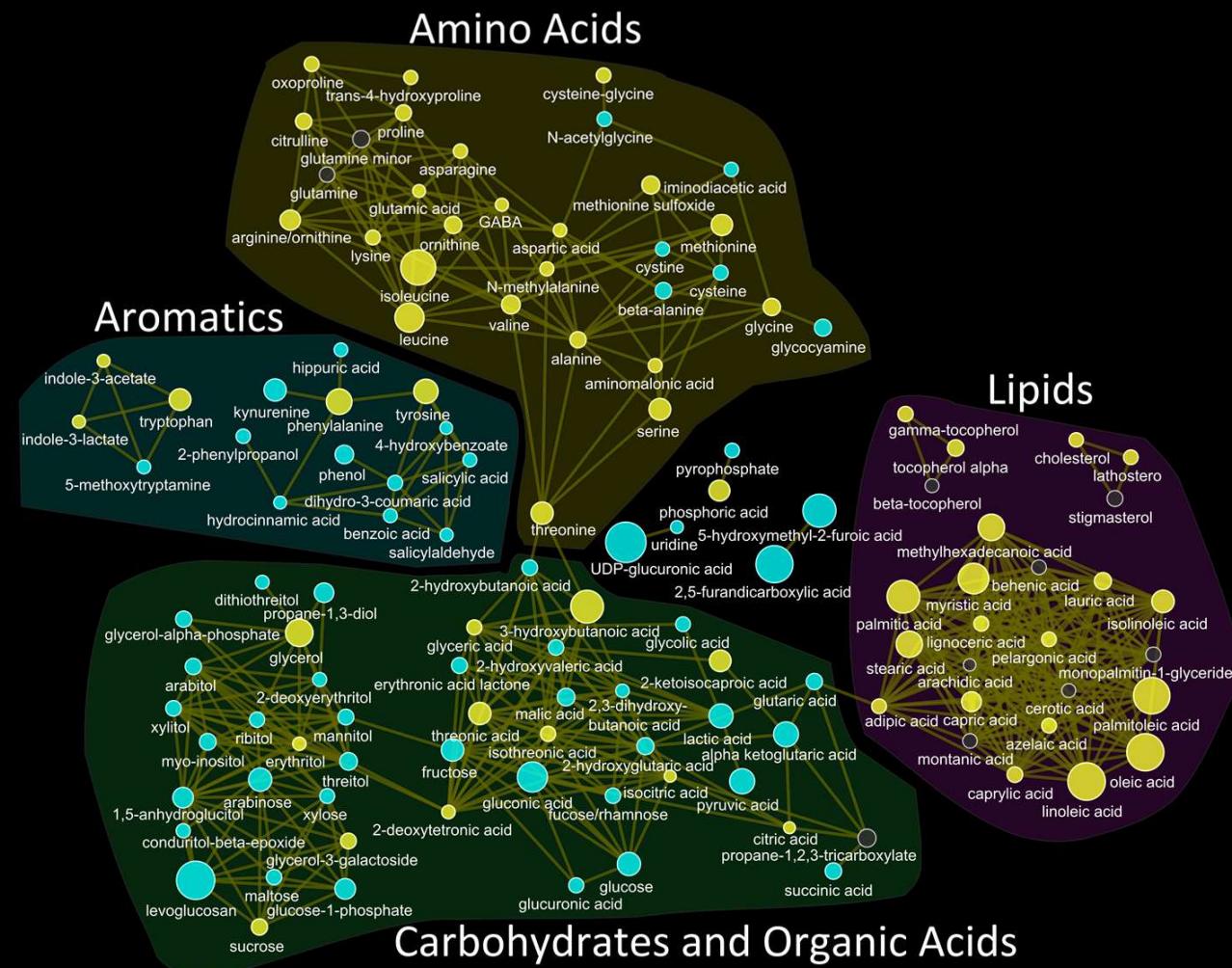
Chemical mapping  
of substructure comparison  
using PubChem



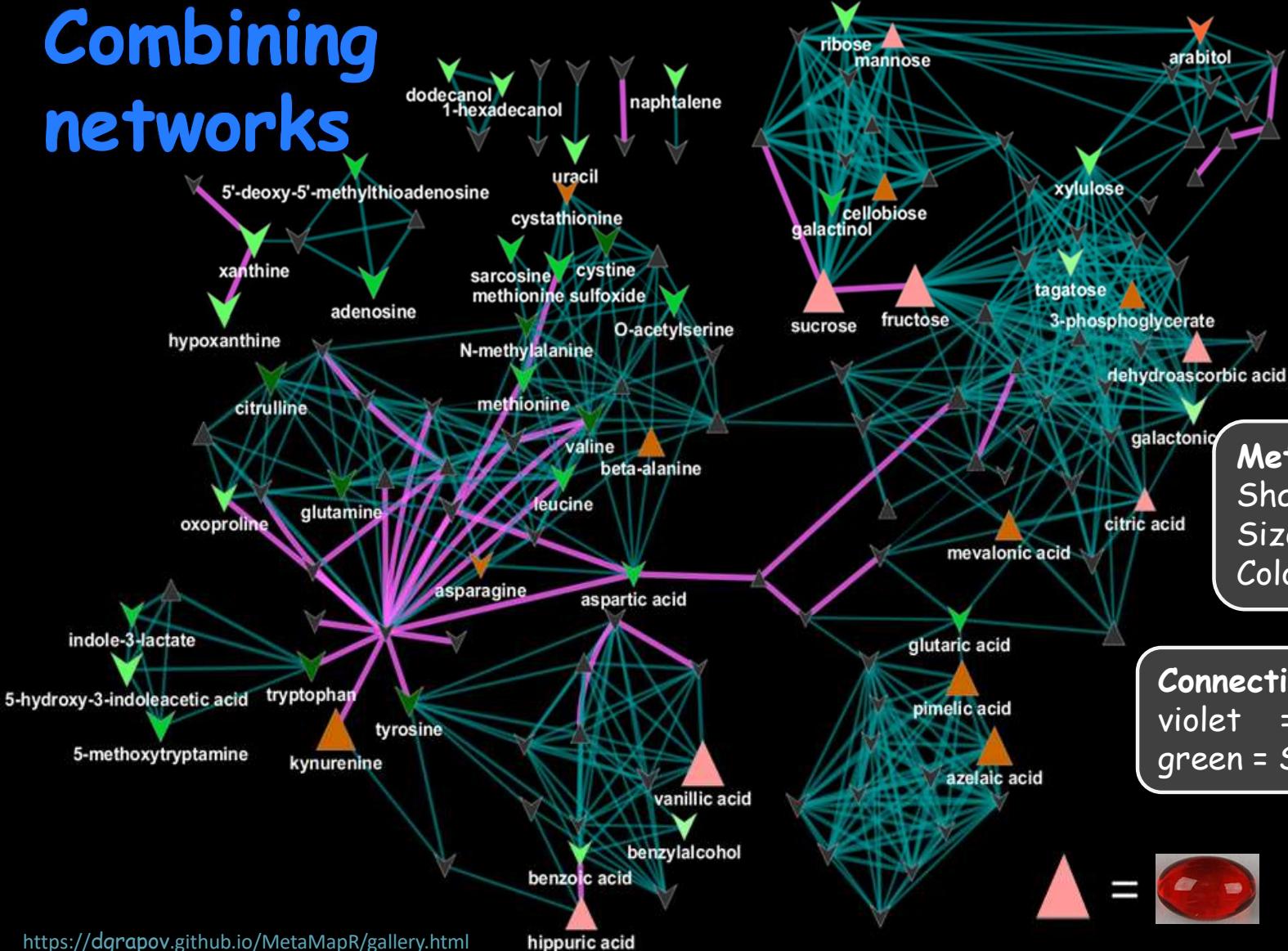
substructure matrix decomposition and  
Tanimoto chemical similarity calculations

BMC Bioinformatics 2012, **13**:99 doi:10.1186/1471-2105-13-99

# Structural similarity example



# Combining networks



## Metabolites

Shape = increase/decrease  
 Size = importance (loading)  
 Color = correlation

## Connections

violet = Biochemical relationships  
 green = Structural similarity



# Your turn

Follow along with the following tutorial:

[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/network/](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/network/)

# Network refinement and visualization

learn

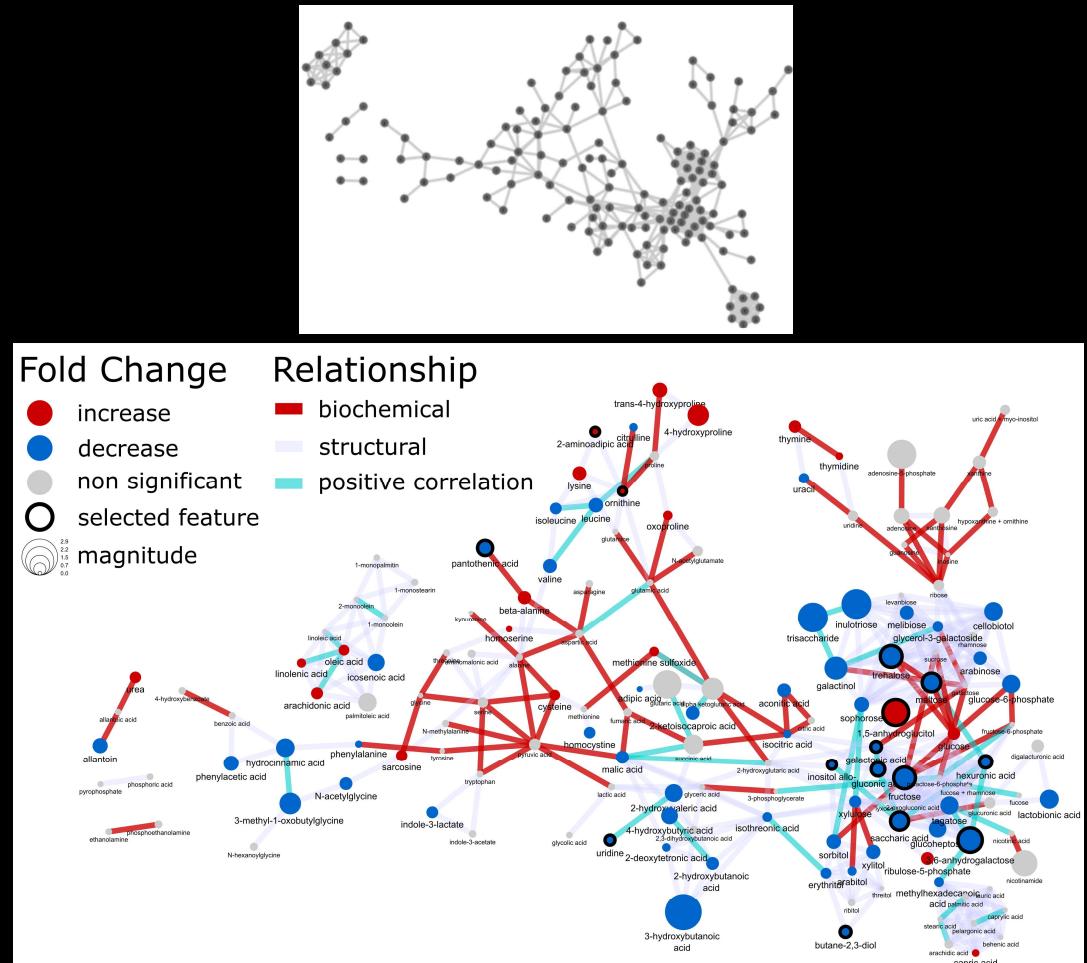
- Cytoscape basics

map variables to

- node attributes
- edge attributes

optimize

- layout
- legend
- publication quality figure



[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/cytoscape/](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/cytoscape/)

# Your turn

Follow along with the following tutorial:

[https://creativedatasolutions.github.io/CDS.courses/courses/network\\_mapping\\_101/docs/partial/cytoscape/](https://creativedatasolutions.github.io/CDS.courses/courses/network_mapping_101/docs/partial/cytoscape/)