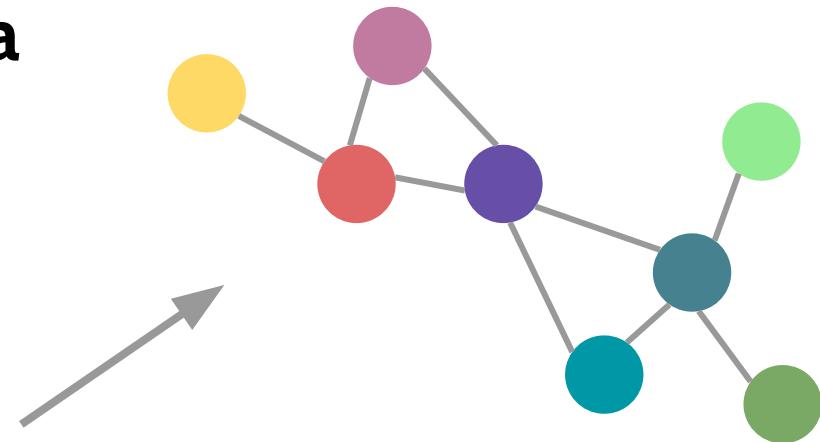
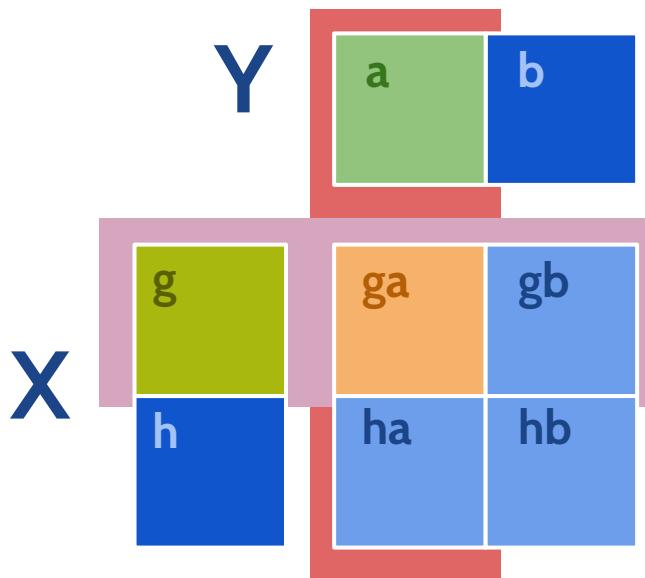
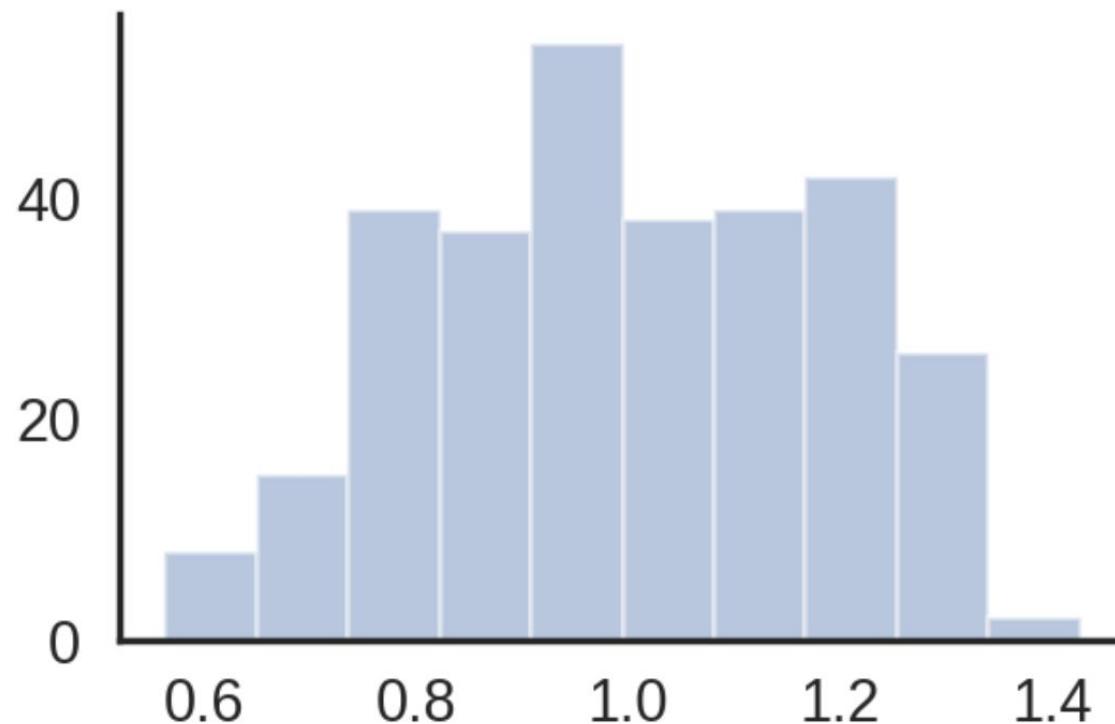


Exploratory Data Visualization for High-dimensional Data

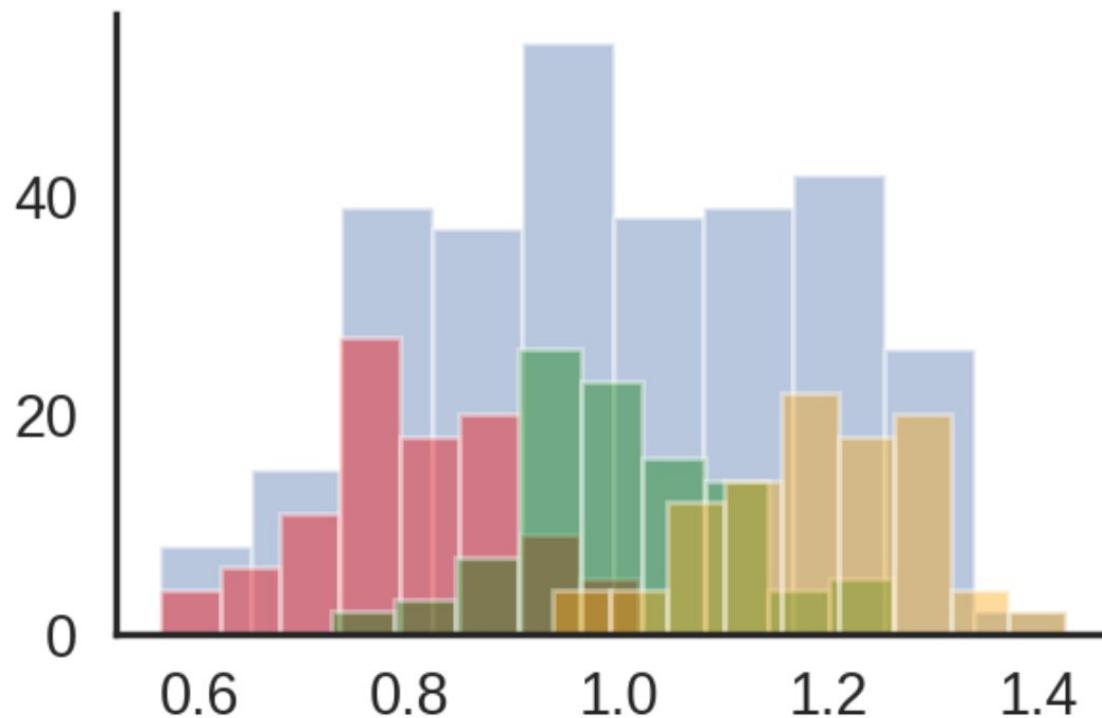


Jane Adams
Data Visualization Artist
Vermont Complex Systems Center
Tuesday March 30th 2021

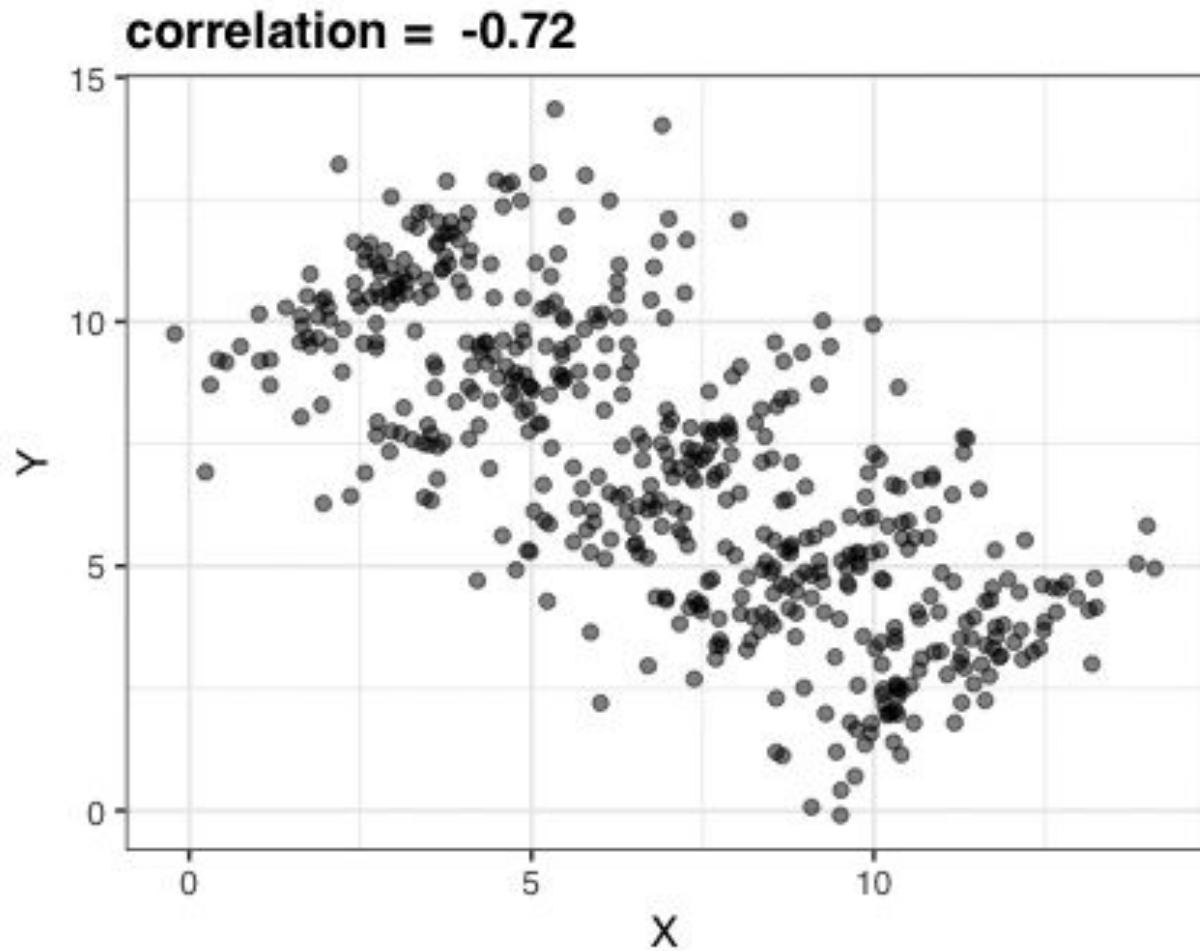
Hidden stories...



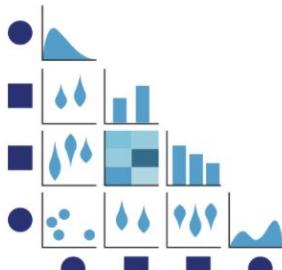
Hidden stories...



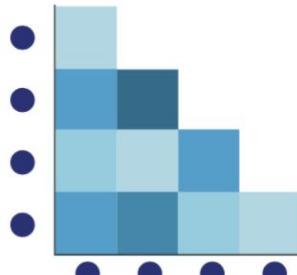
Simpson's Paradox



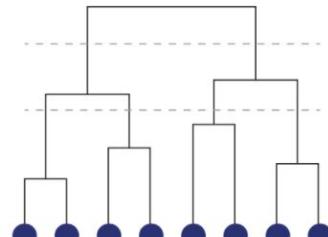
Sirius: Exploratory Analysis Python Package



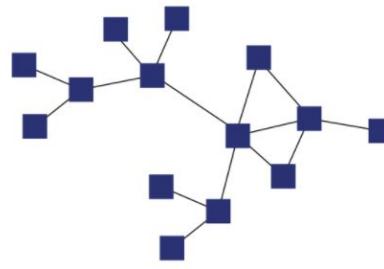
(A) Small Multiples



(B) Matrix

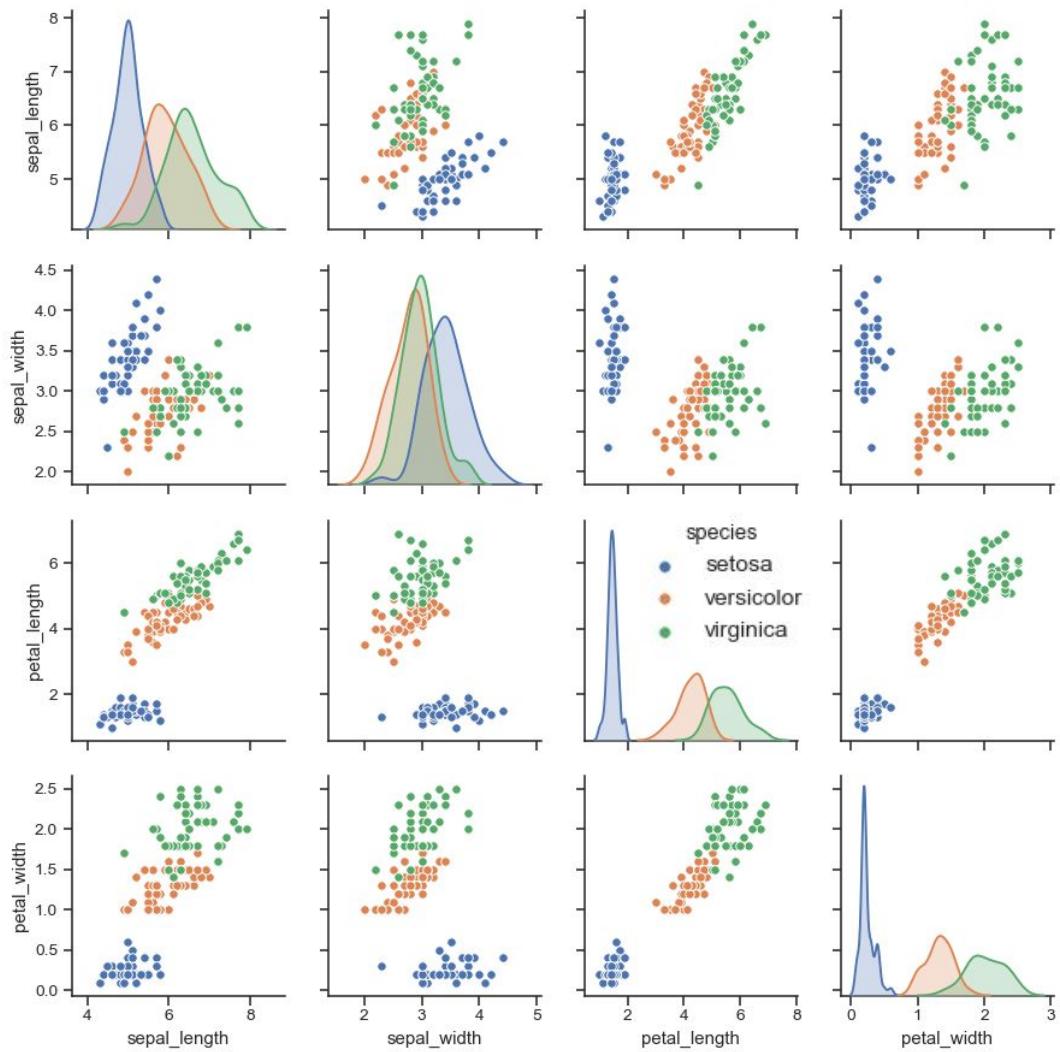


(C) Dendrogram



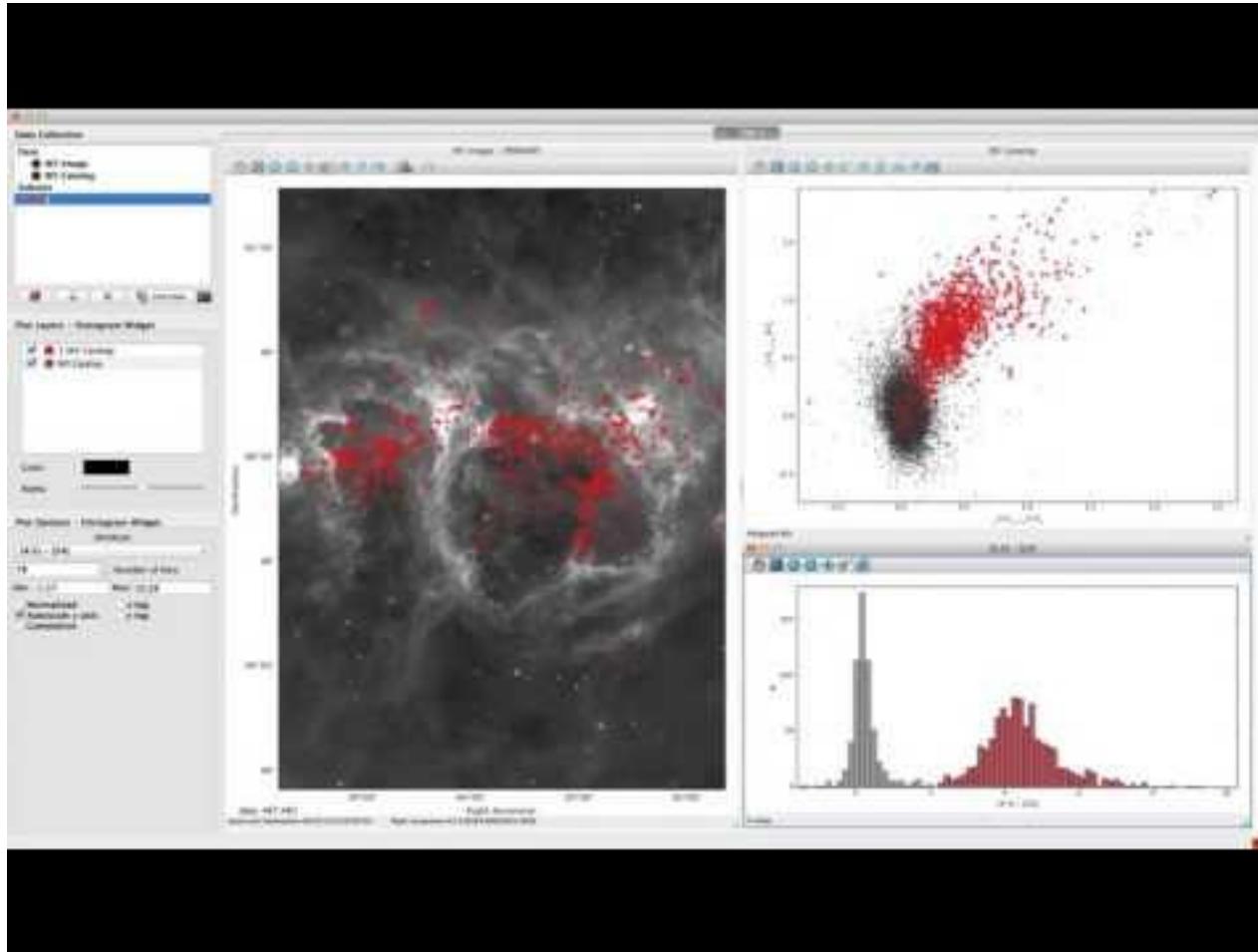
(D) Network

Pair Plots

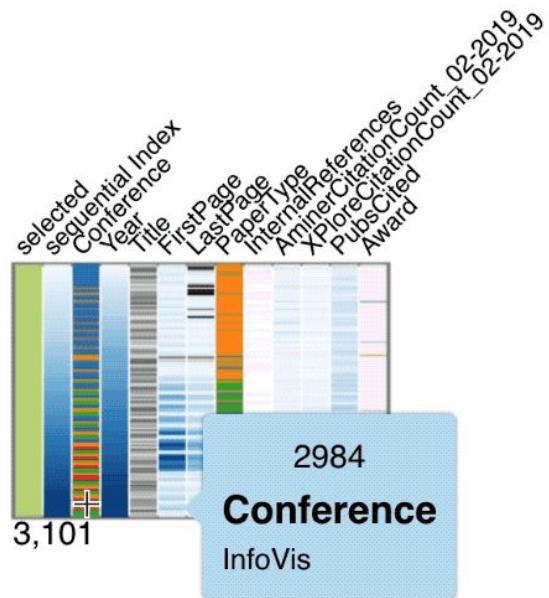


src: seaborn.pydata.org

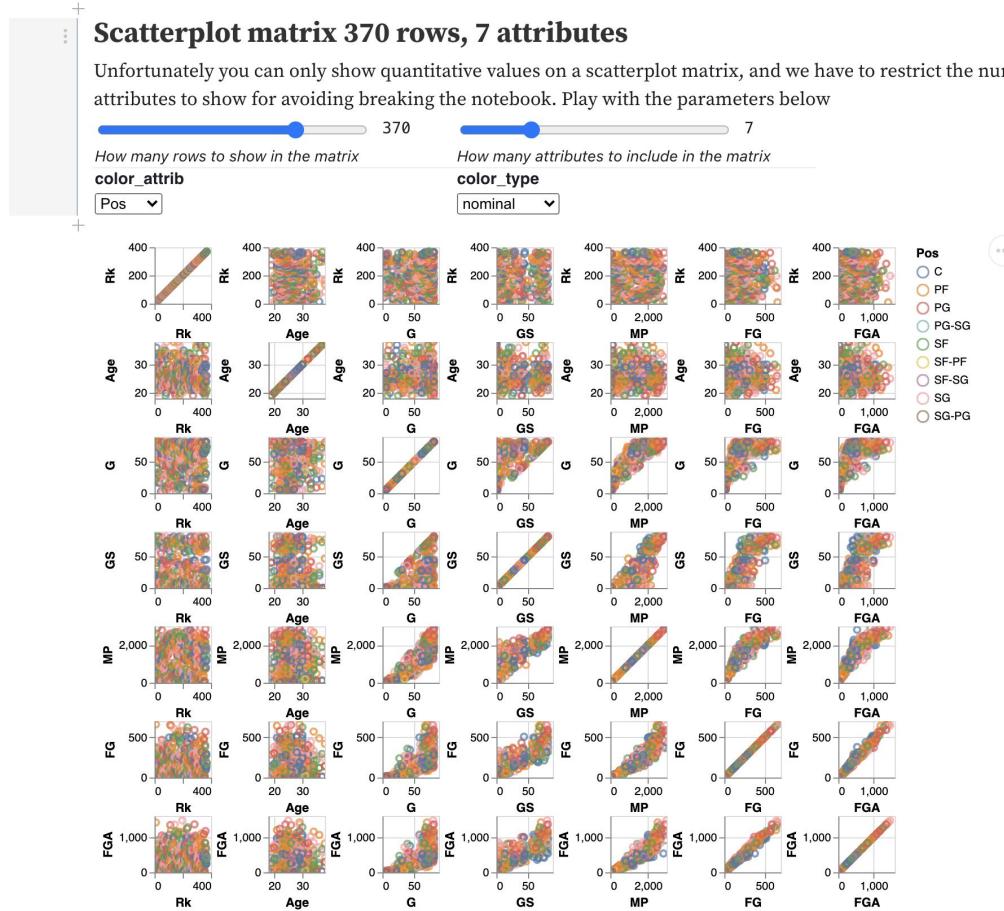
Brushing and Linking: GlueViz



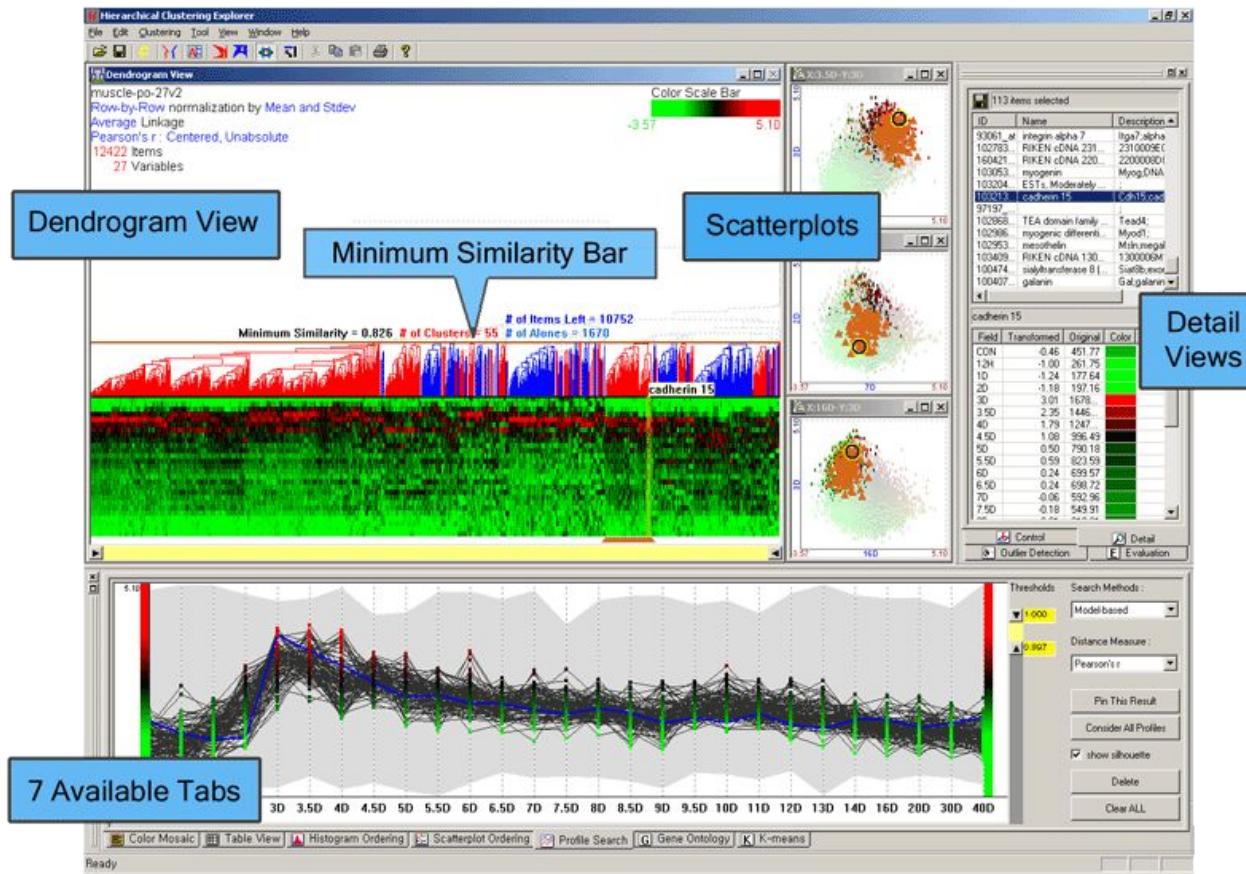
Navio: <https://navio.dev/>



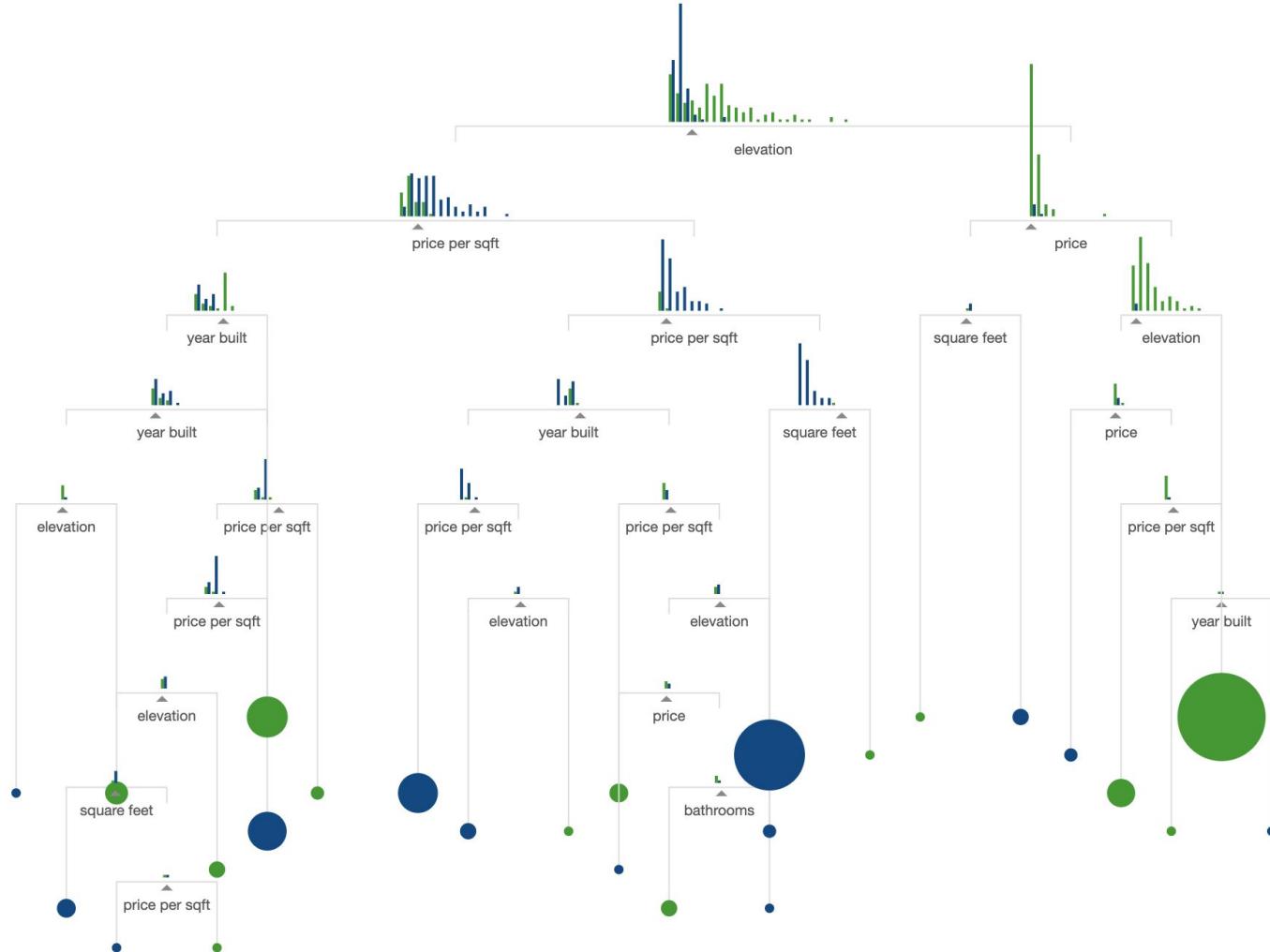
Comparing Navio to a Scatterplot Matrix



Hierarchical Clustering Explorer (HCE)



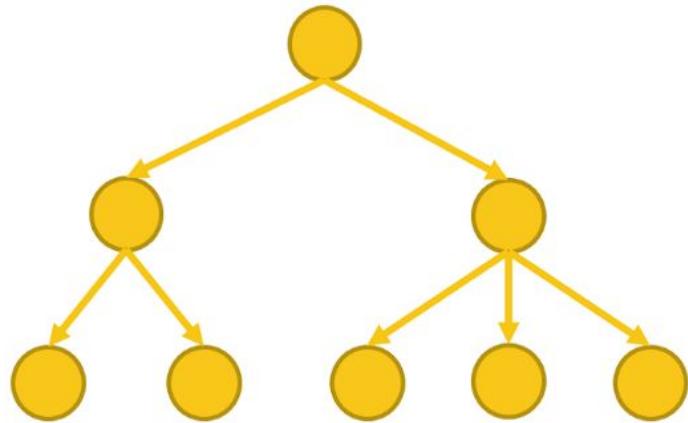
Decision Trees



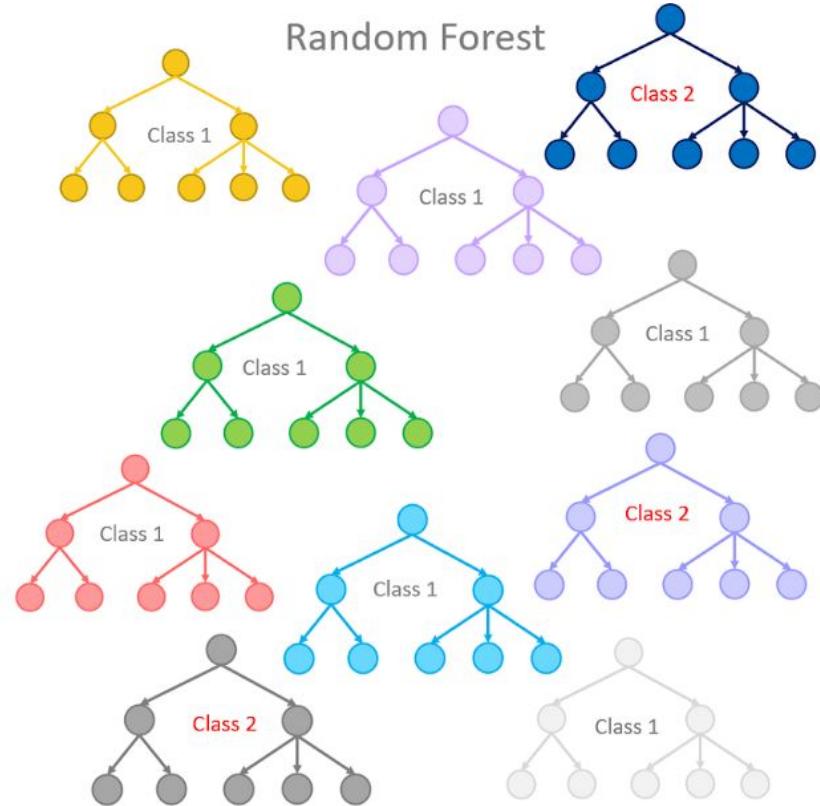
src: r2d3.us

Random Forests

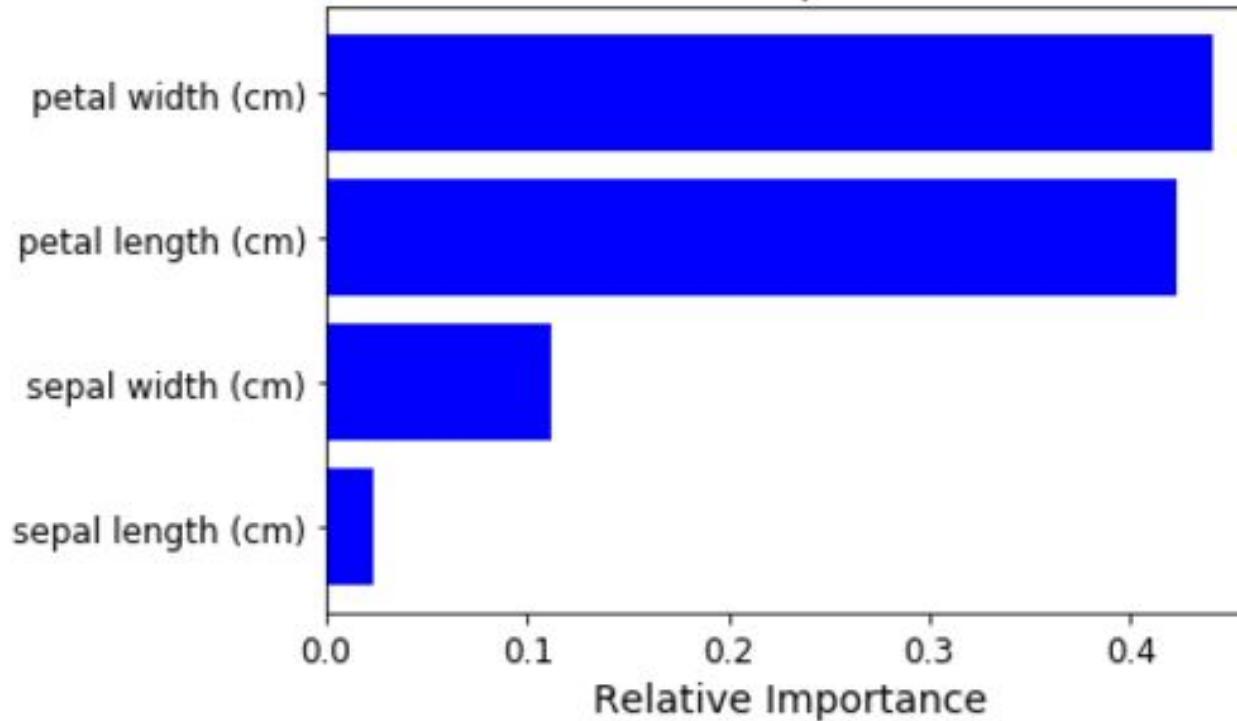
Single Decision Tree



Random Forest

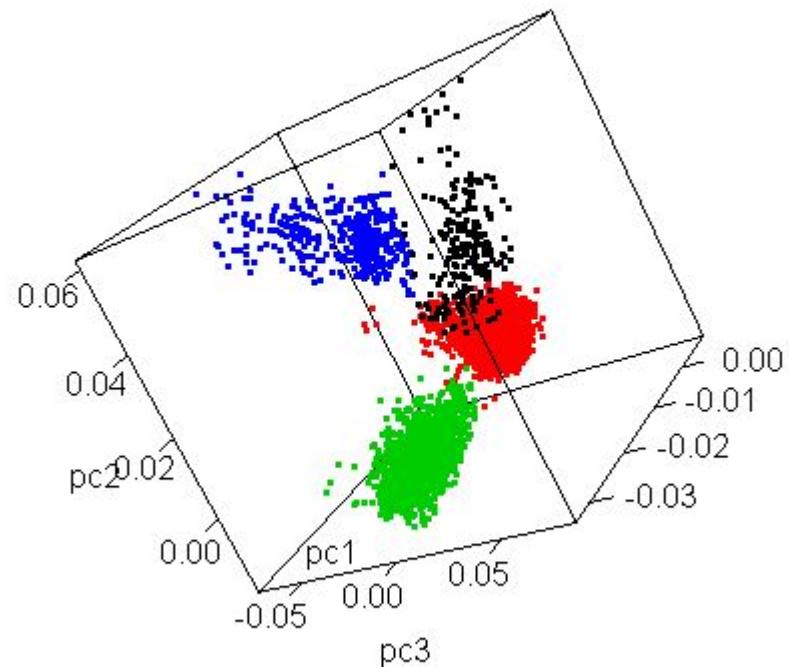
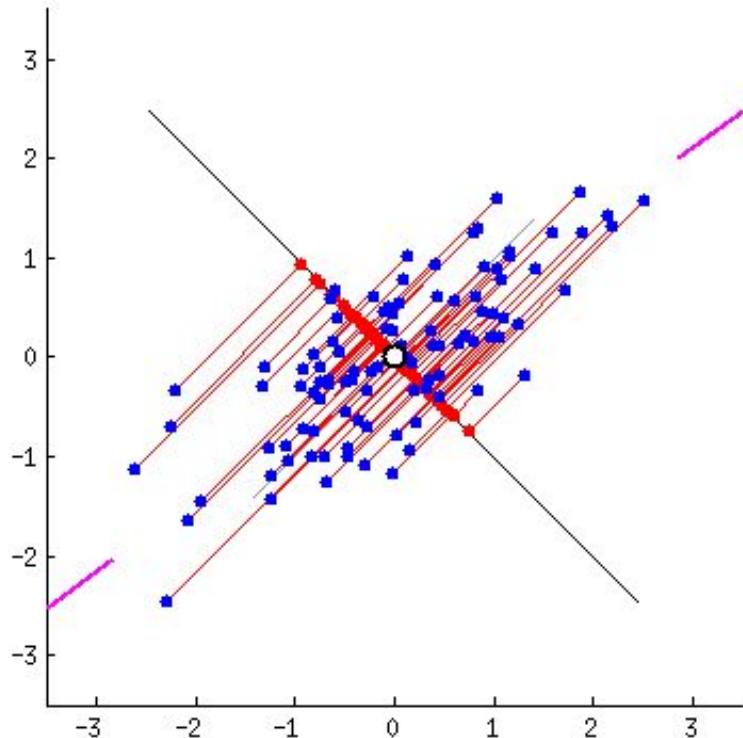


Feature Importances

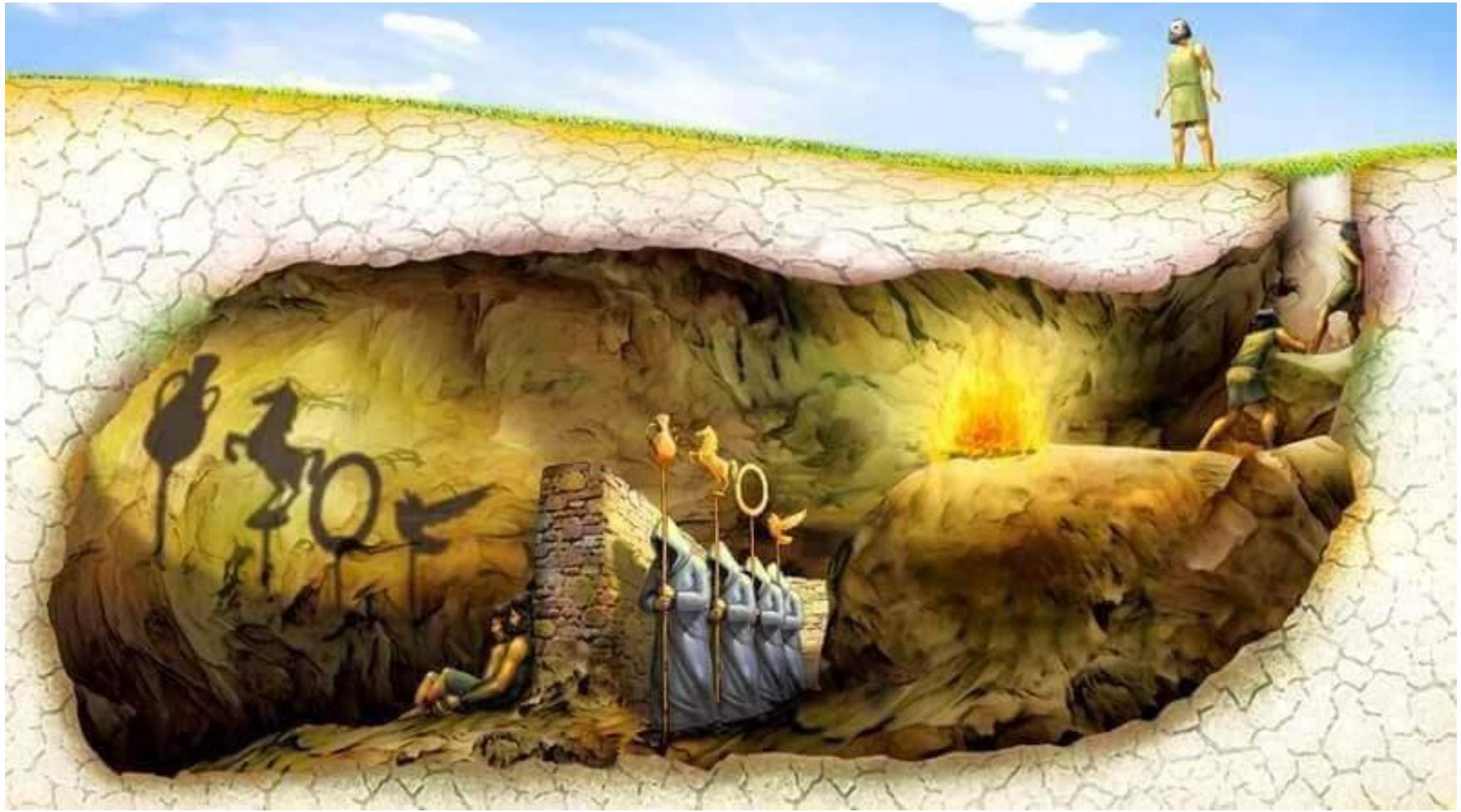


src: [Stack Overflow](#)

Principal Component Analysis



Plato's Cave



Plato's Cave: Modern Version



Ally [REDACTED]
Where's the W?!?! #trippy

Algorithmic Bias

98.7%

68.6%

100%

92.9%



Rekognition
Performance on
Gender Classification



**DARKER
MALES**



**DARKER
FEMALES**

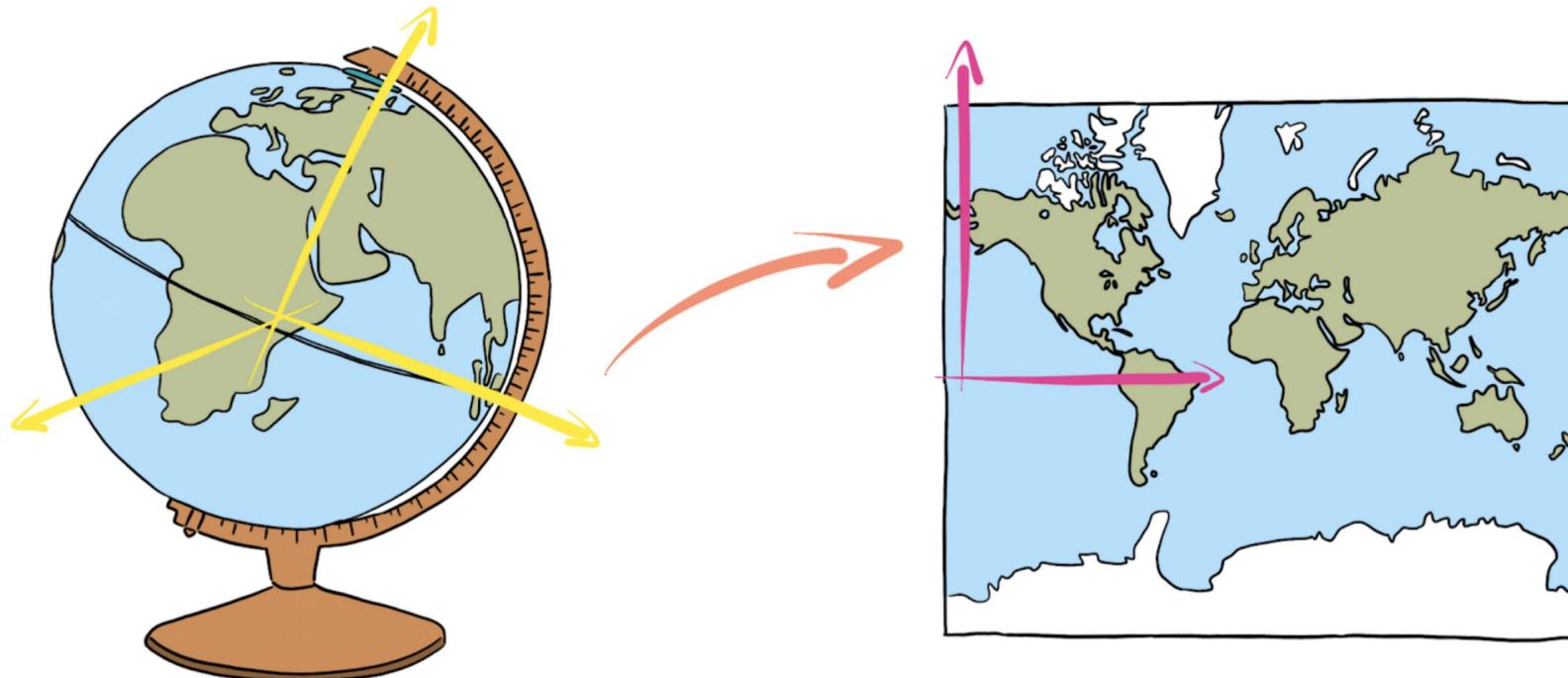


**LIGHTER
MALES**



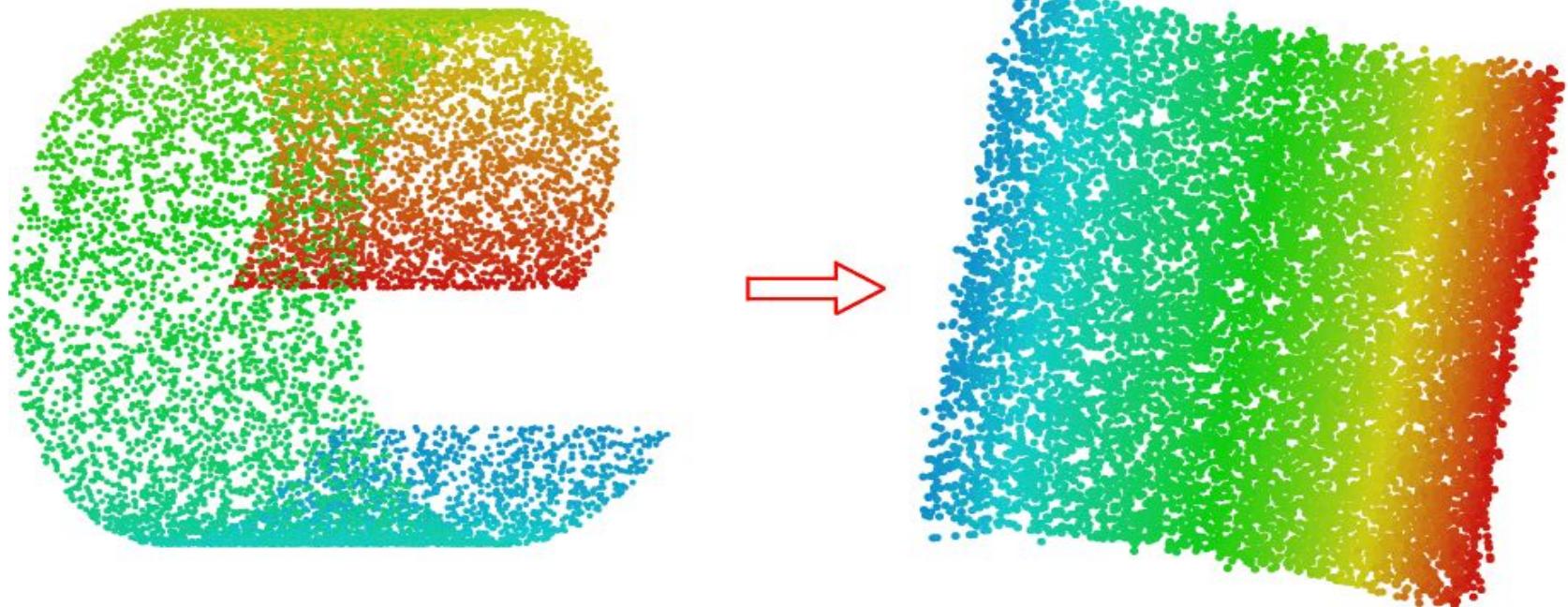
**LIGHTER
FEMALES**

The curse of dimensionality...



src: Fanny Kasapian, "How to Build a Non-Geographical Map"

Locally Linear Mapping



Source: "Nonlinear Dimensionality Reduction via Path-Based Isometric Mapping", Amir Najafi et. al.
<https://www.computer.org/csdl/journal/tp/2016/07/07293680/13rRUwfZC1K>

**How do we narrow the
exploration space to
combinations of variables
that are ‘most interesting’?**

Probability

		a .5	b .5
	Y		
X	g .5	ga .25	gb .25
	h .5	ha .25	hb .25

Probability

Marginal: $P(Y_a)$

a	.5
b	.5

g	.5
h	.5
ga	.25
ha	.25
gb	.25
hb	.25

Probability

Marginal: $P(Y_a)$

a	.5
b	.5

Joint: $P(X_g , Y_a)$

		Y	
X			
	g	.5	
	h	.5	
	ga	.25	gb
	ha	.25	hb

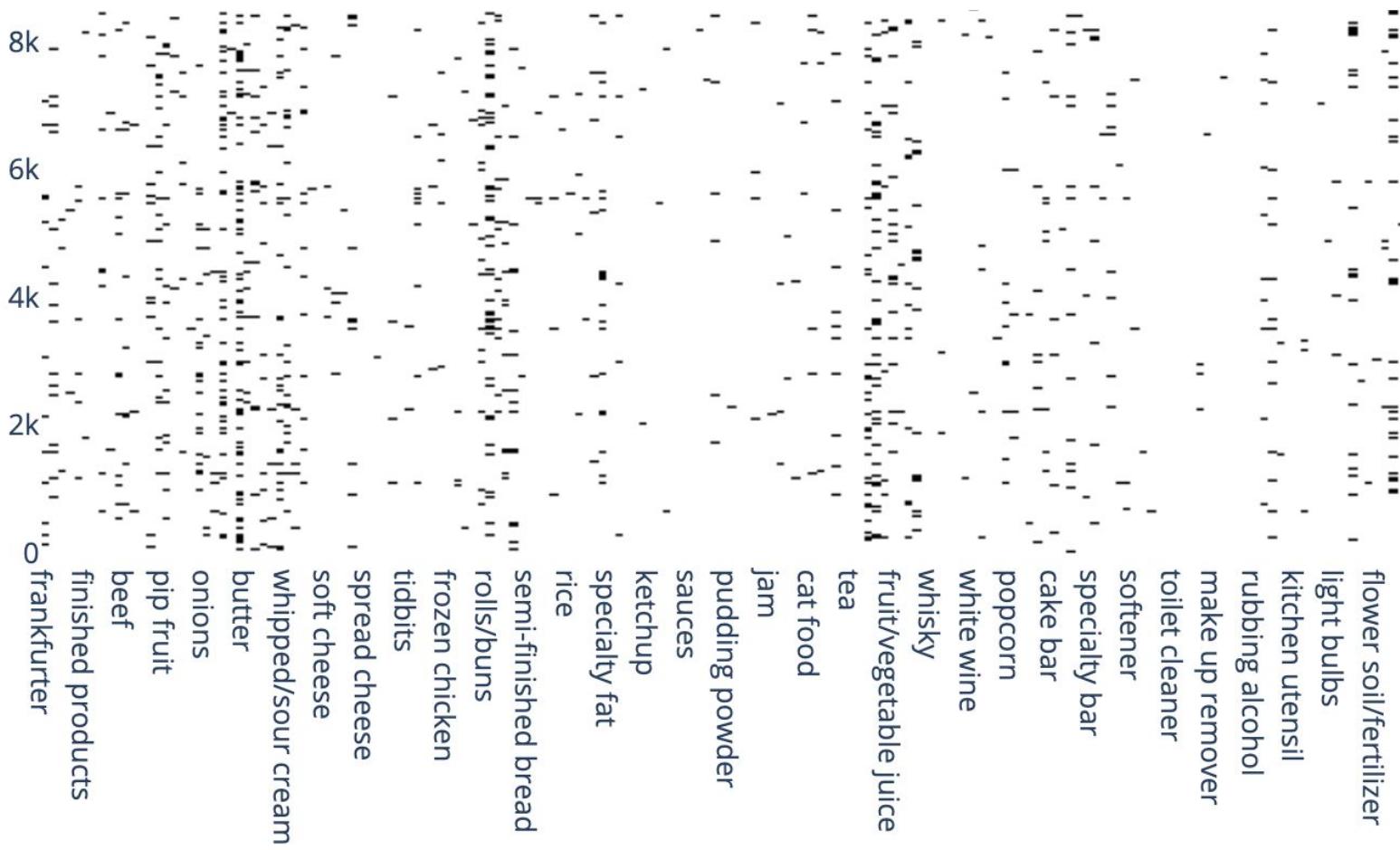
Groceries: Association Networks of Binary Features

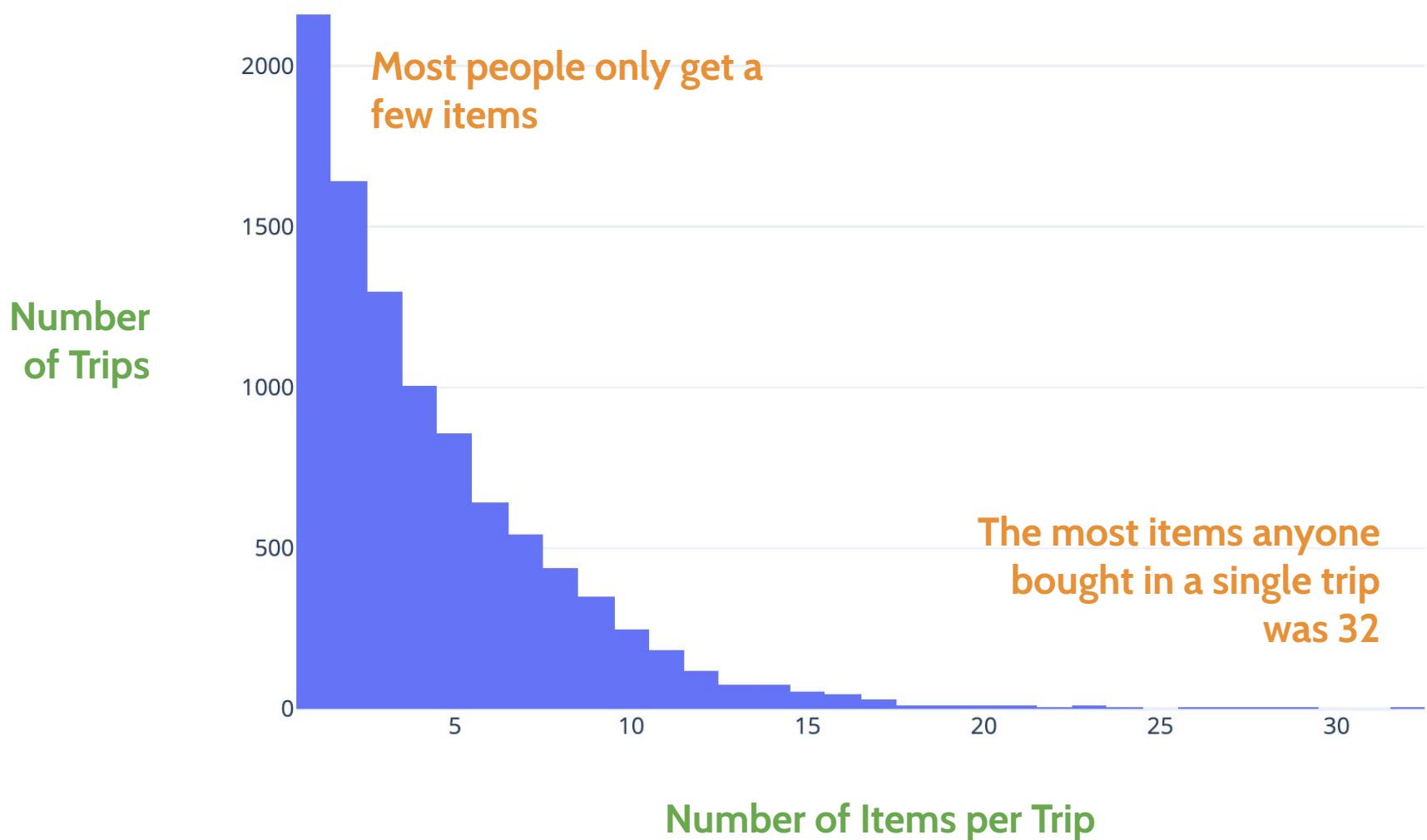
170 products (items, features)

10k customers (trips)

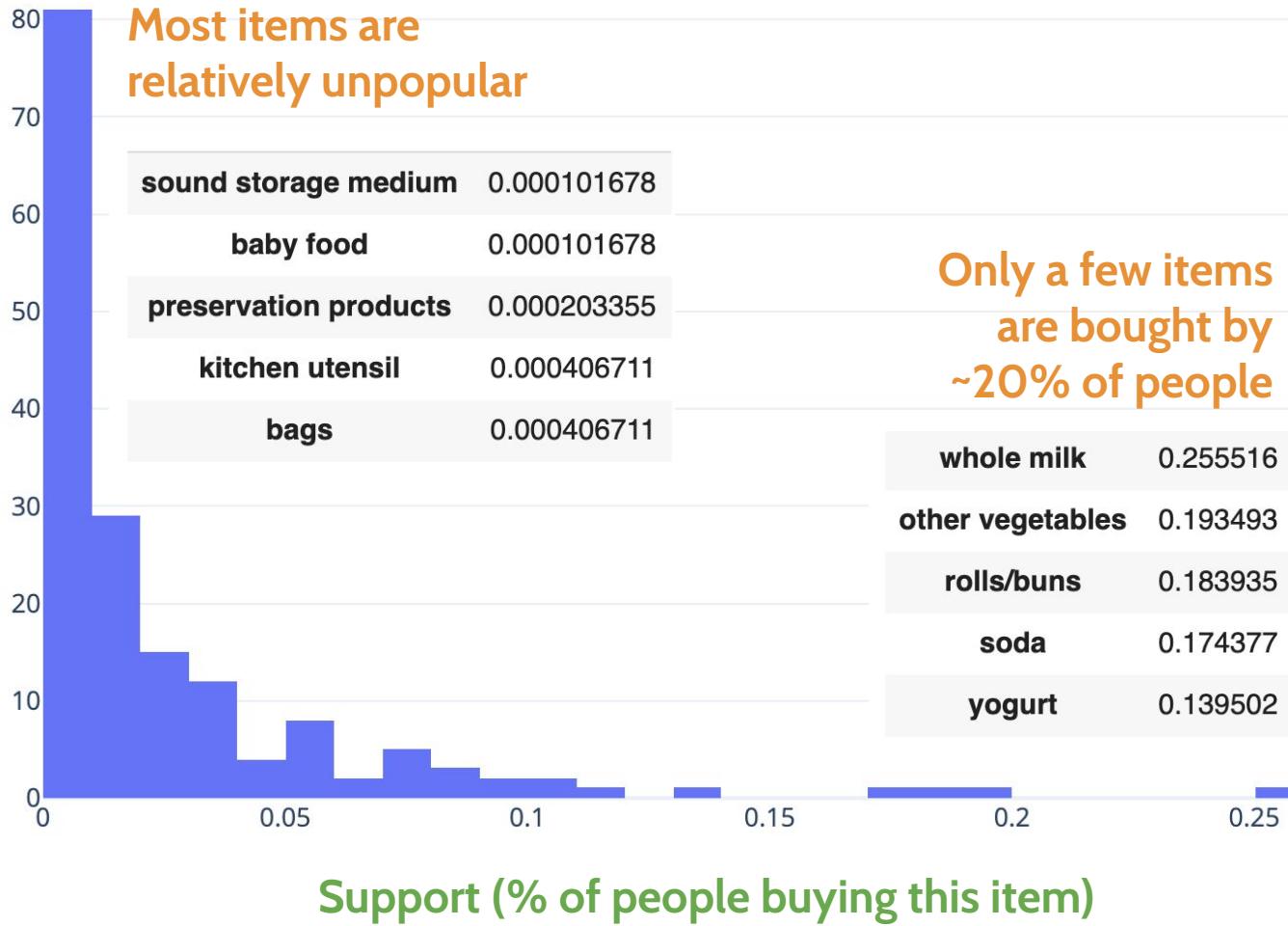
Trips

Items





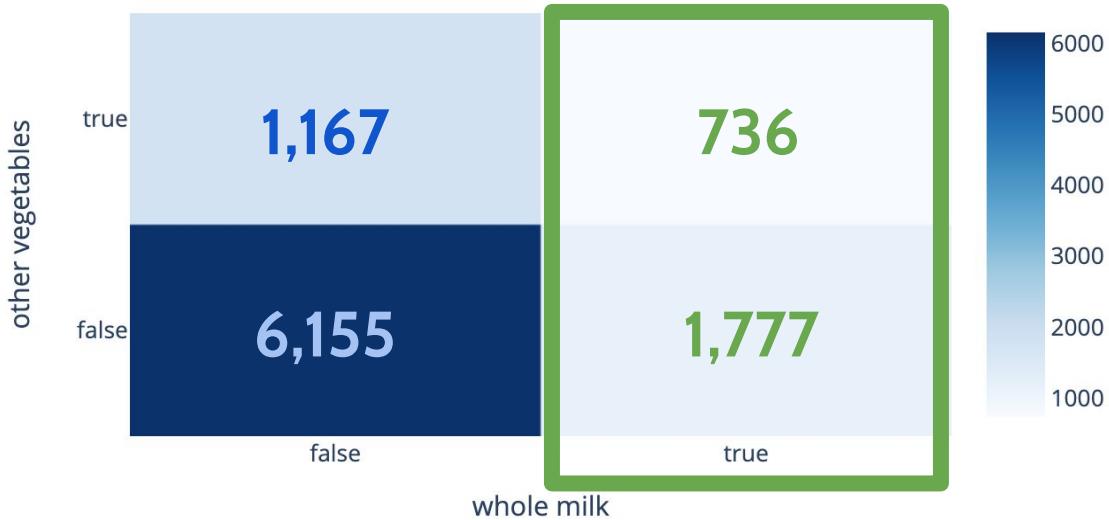
Number
of items



Groceries



Calculating Marginal Probability

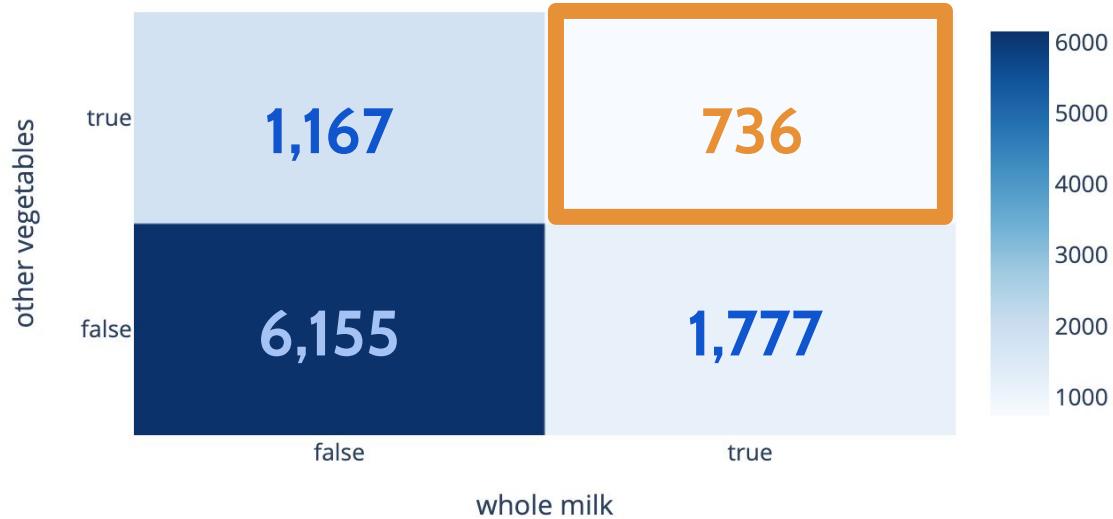


736 + 1,777

$\approx 25.5\%$

$736 + 1,777 + 1,167 + 6,155$

Calculating Joint Probability



736

$$736 + 1,167 + 6,155 + 1,777$$

$\approx 7.5\%$

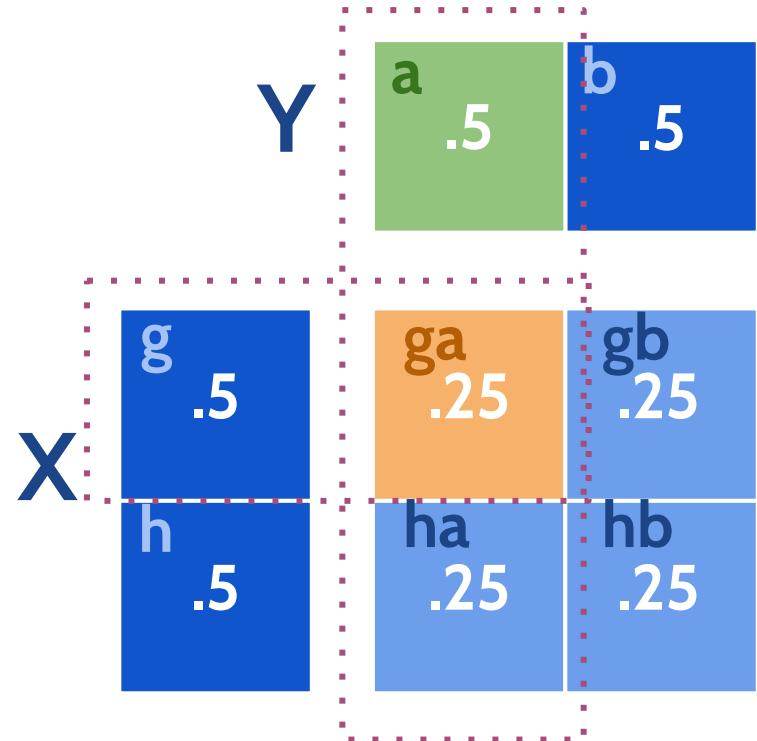
Calculating Conditional Probability

Marginal: $P(Y_a)$

Joint: $P(X_g, Y_a)$

Conditional:
 $P(X_g | Y_a)$

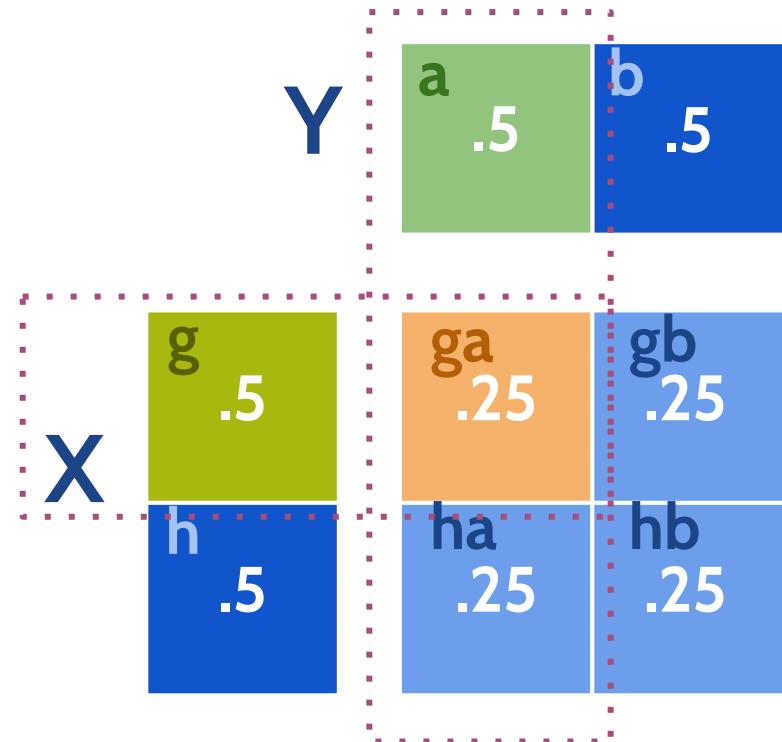
$$\begin{matrix} ga \\ .25 \end{matrix} / \begin{matrix} a \\ .5 \end{matrix} = \begin{matrix} g | a \\ .5 \end{matrix}$$



Calculating “Lift”

$$\begin{matrix} \text{ga} \\ .25 \end{matrix} / \begin{matrix} \text{a} \\ .5 \end{matrix} = \begin{matrix} \text{g|a} \\ .5 \end{matrix}$$

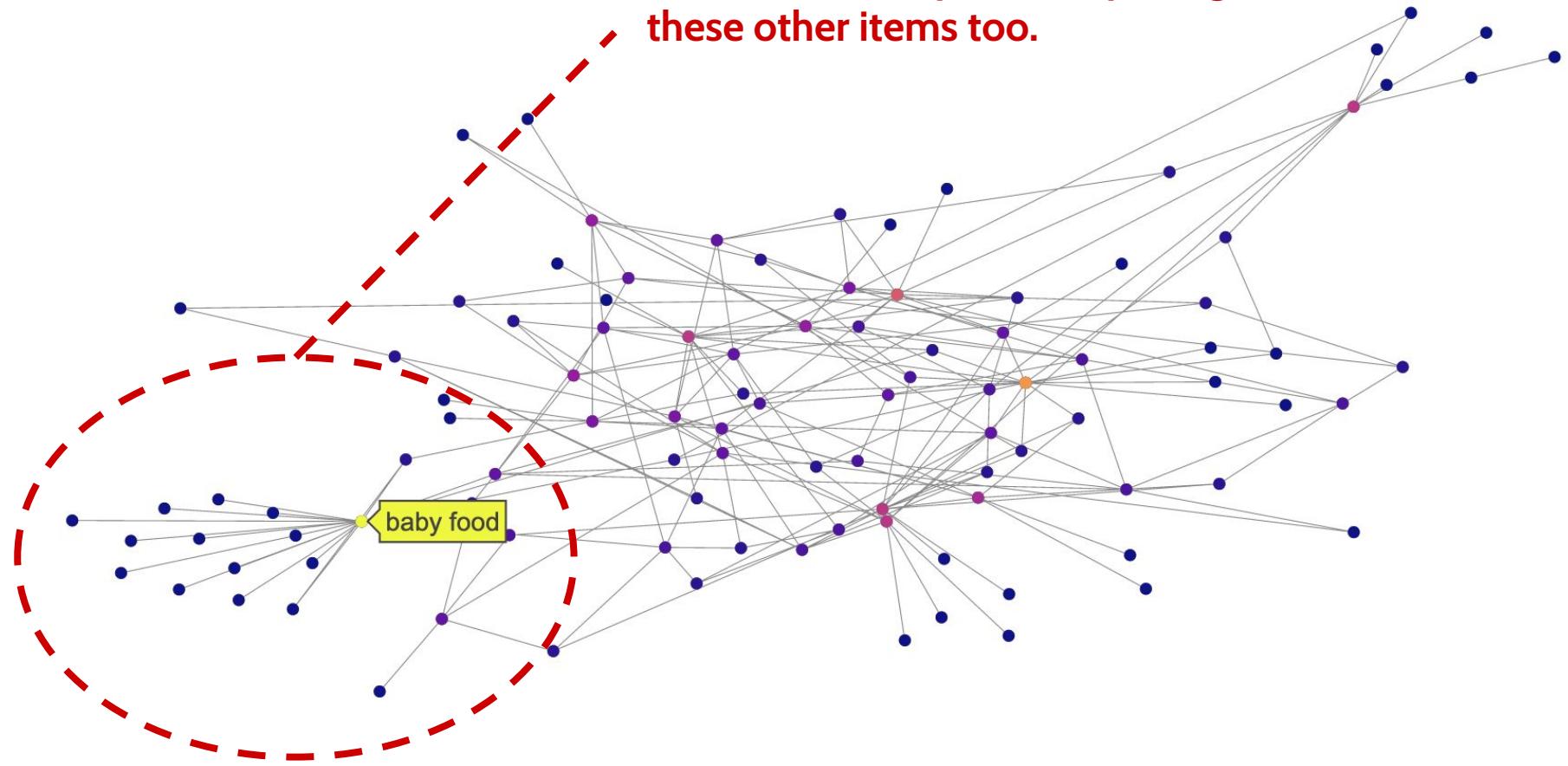
$$\begin{matrix} \text{g|a} \\ .5 \end{matrix} / \begin{matrix} \text{g} \\ .5 \end{matrix} = 1$$



“It is **1** times more likely that **X = g** if I know that **Y = a**.”

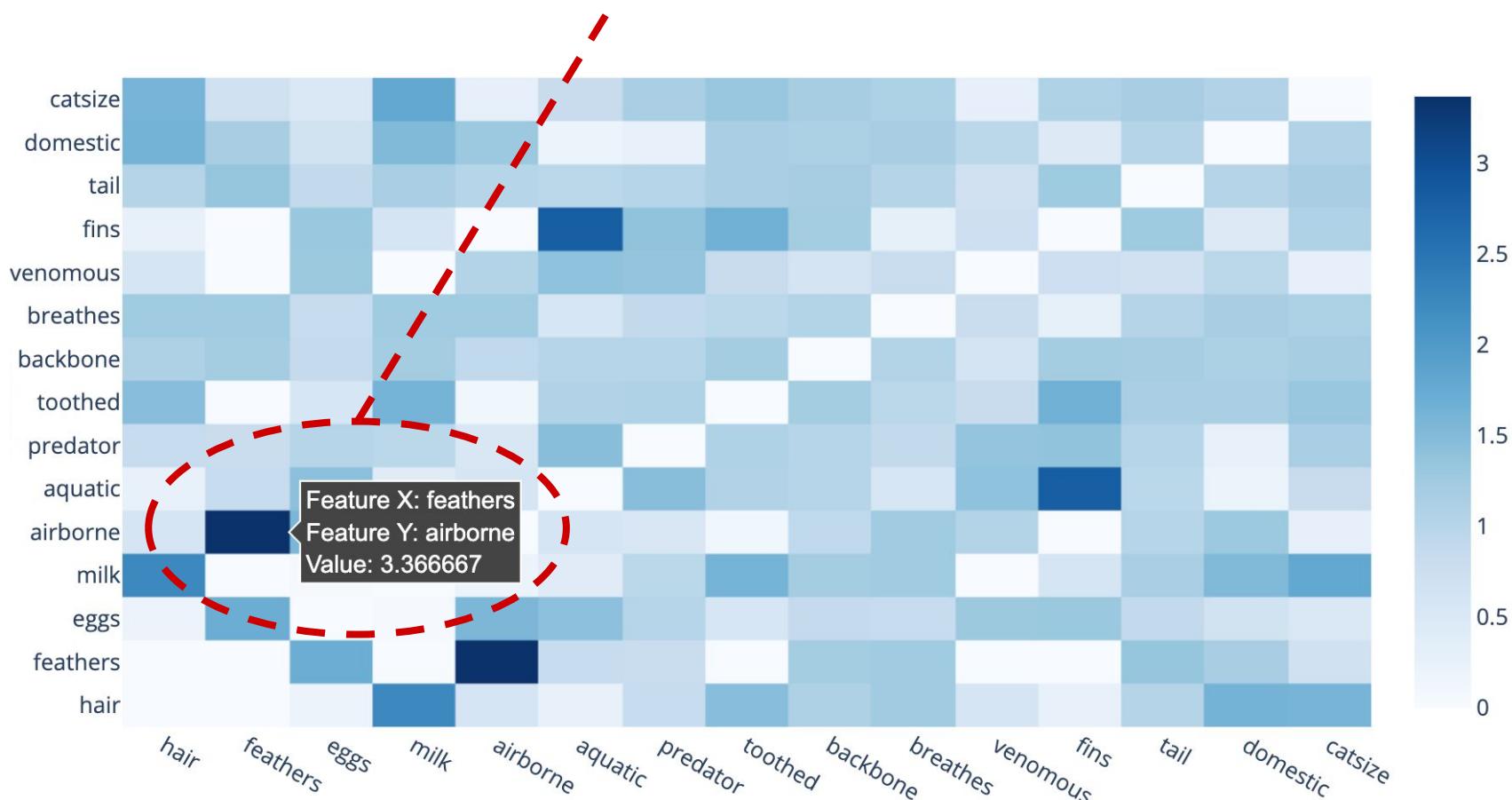
Association Network

If someone bought baby food, it's much more likely that they bought these other items too.

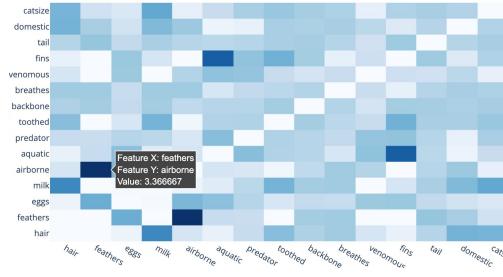
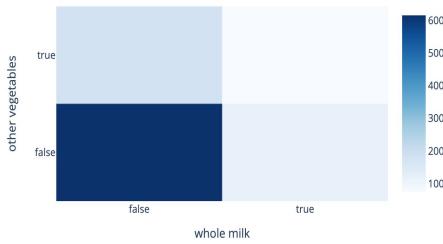


Association Network

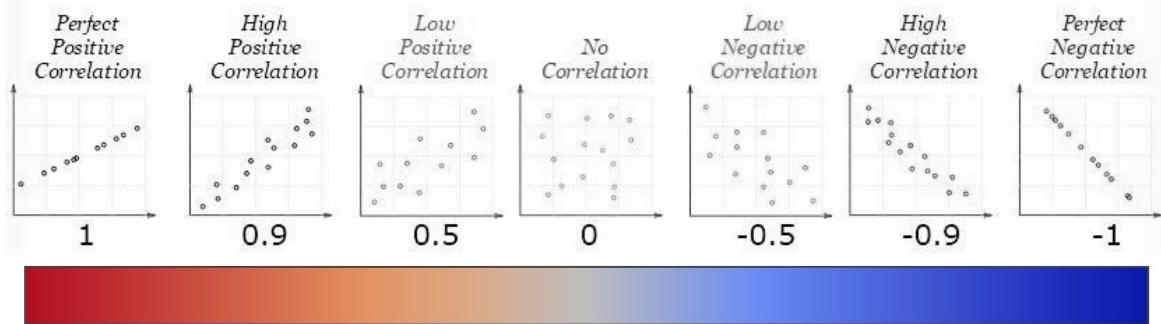
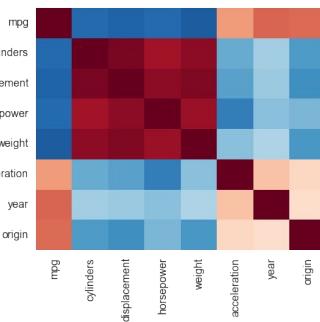
If I know that an animal has feathers, I'm 3x more likely to know that it can be airborne.



Association Network For Binary Features



Correlation Network For Continuous Features



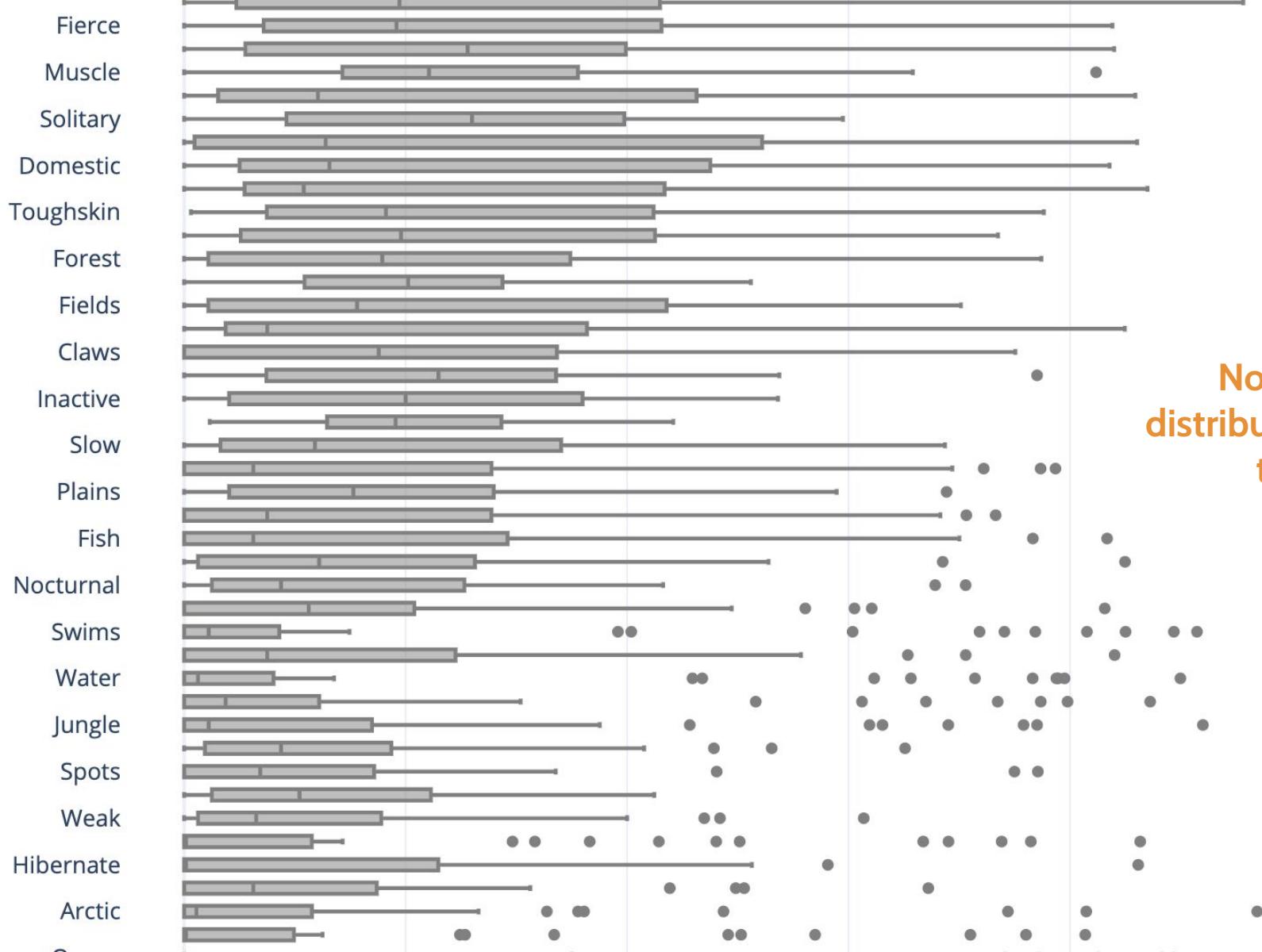
Animals: Correlation Networks of Quantitative Features

85 characteristics (features)

49 animals (observations)

animal	Black	White	Blue	Brown	Gray	Orange	Red	Yellow	Patches	Spots	Stripes	Furry
Grizzly Bear	39.25	1.39	0.00	74.14	3.75	0.00	0.0	0.00	1.25	0.00	0.00	82.37
Killer Whale	83.40	64.79	0.00	0.00	1.25	0.00	0.0	0.00	68.49	32.69	0.00	1.25
Beaver	19.38	0.00	0.00	87.81	7.50	0.00	0.0	0.00	0.00	7.50	0.00	46.25
Dalmatian	69.58	73.33	0.00	6.39	0.00	0.00	0.0	0.00	37.08	100.00	0.00	27.15
Persian Cat	19.38	50.09	29.44	8.98	38.19	0.00	0.0	0.00	17.93	6.25	6.25	90.19
Horse	44.90	42.91	4.44	69.41	35.94	0.00	0.0	0.00	22.29	15.80	0.00	40.58
German Shepherd	43.54	15.88	5.00	54.16	26.82	3.12	2.5	0.38	48.78	11.59	1.56	66.05

Not all of the distributions look the same...



Is it nocturnal?



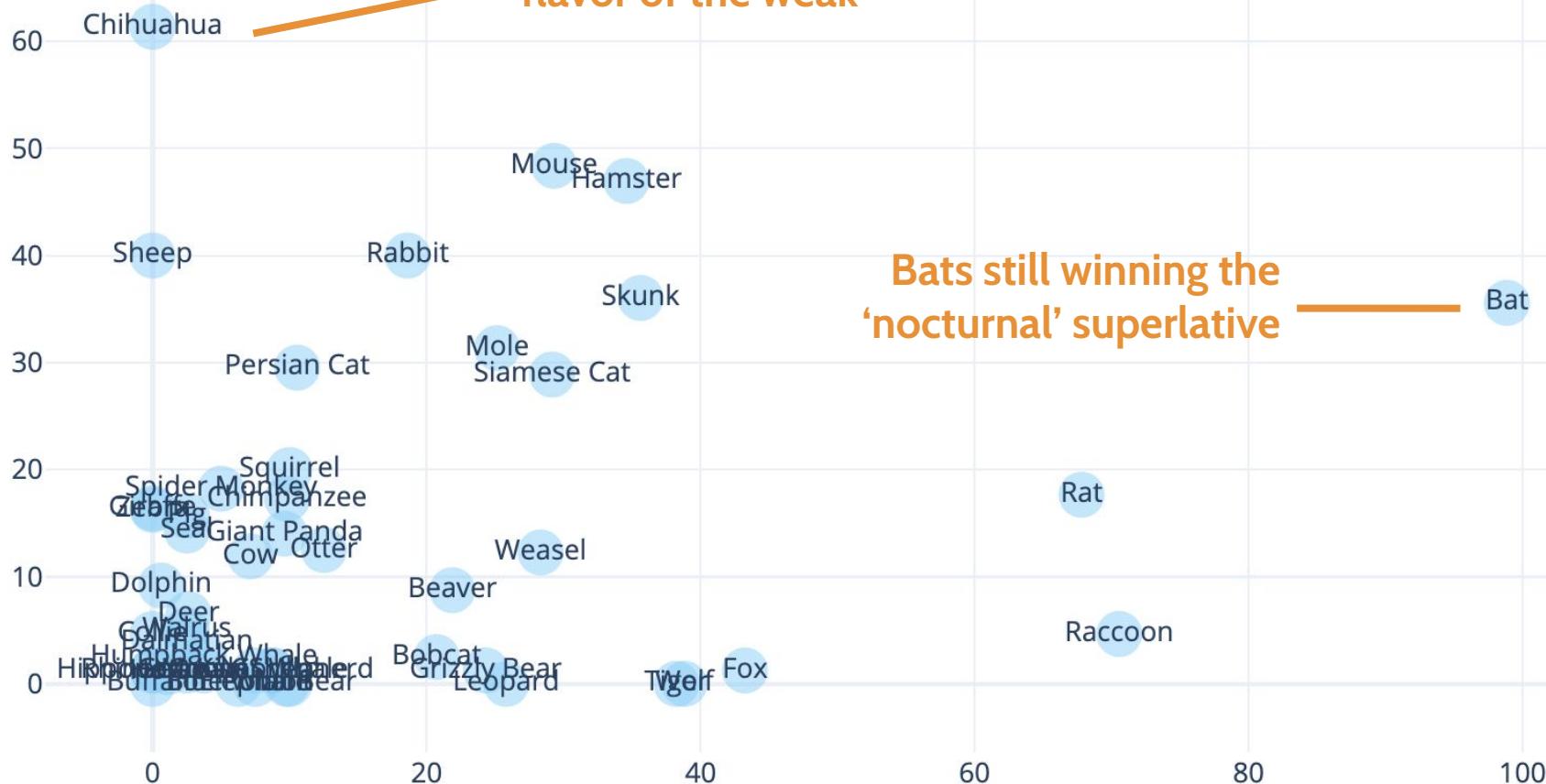
Does it fly?

There are sometimes
clear winners in the
superlative pair matchups

Is it a weak animal?

Chihuahuas winning
'flavor of the weak'

Bats still winning the
'nocturnal' superlative



Is it nocturnal?

Is it a slow animal?



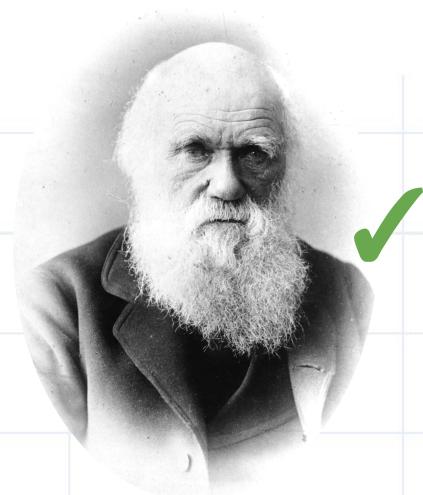
Is it a hunter?

Hunters generally
aren't slow...

Is it a slow animal?

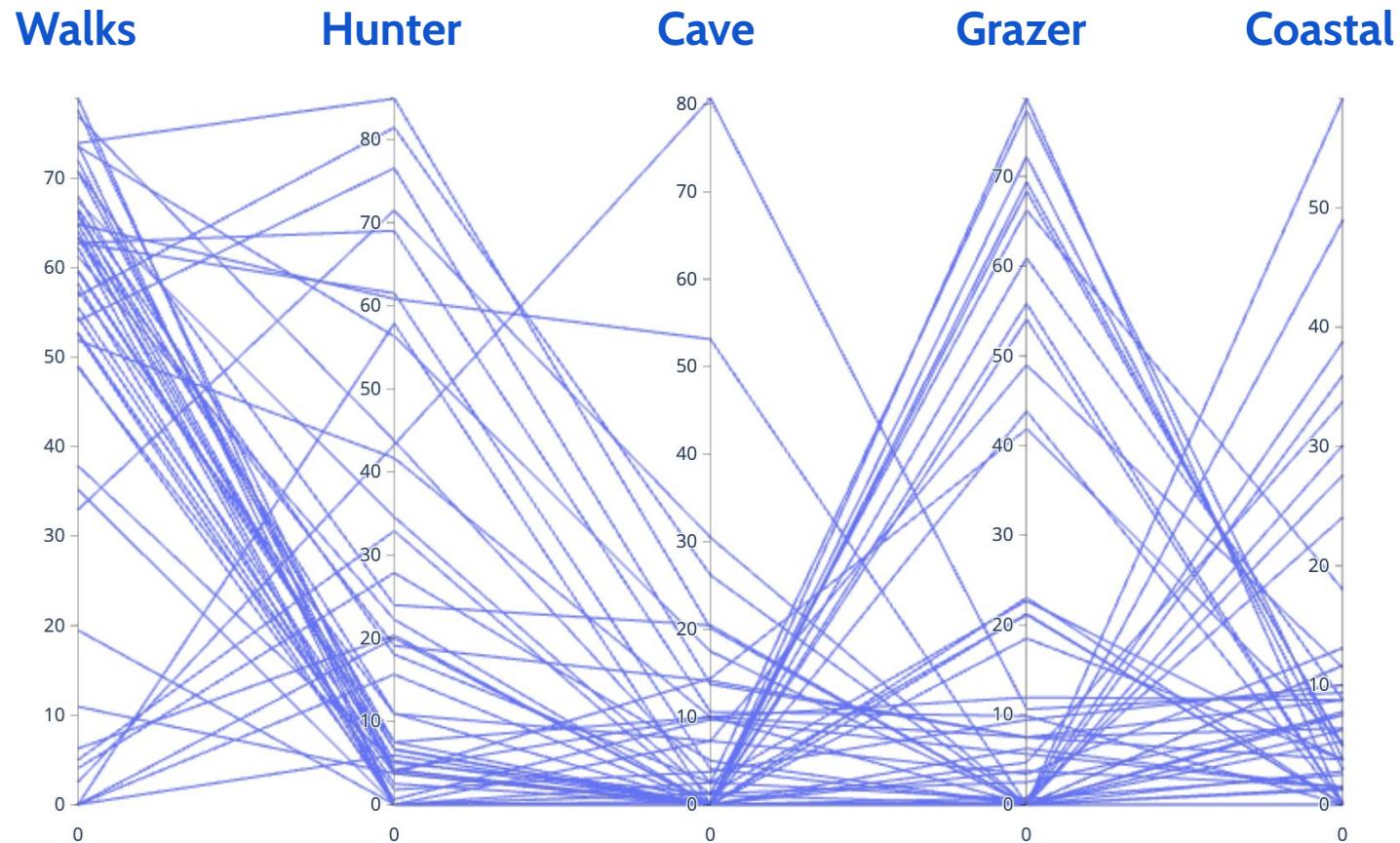


Hunters generally
aren't slow...

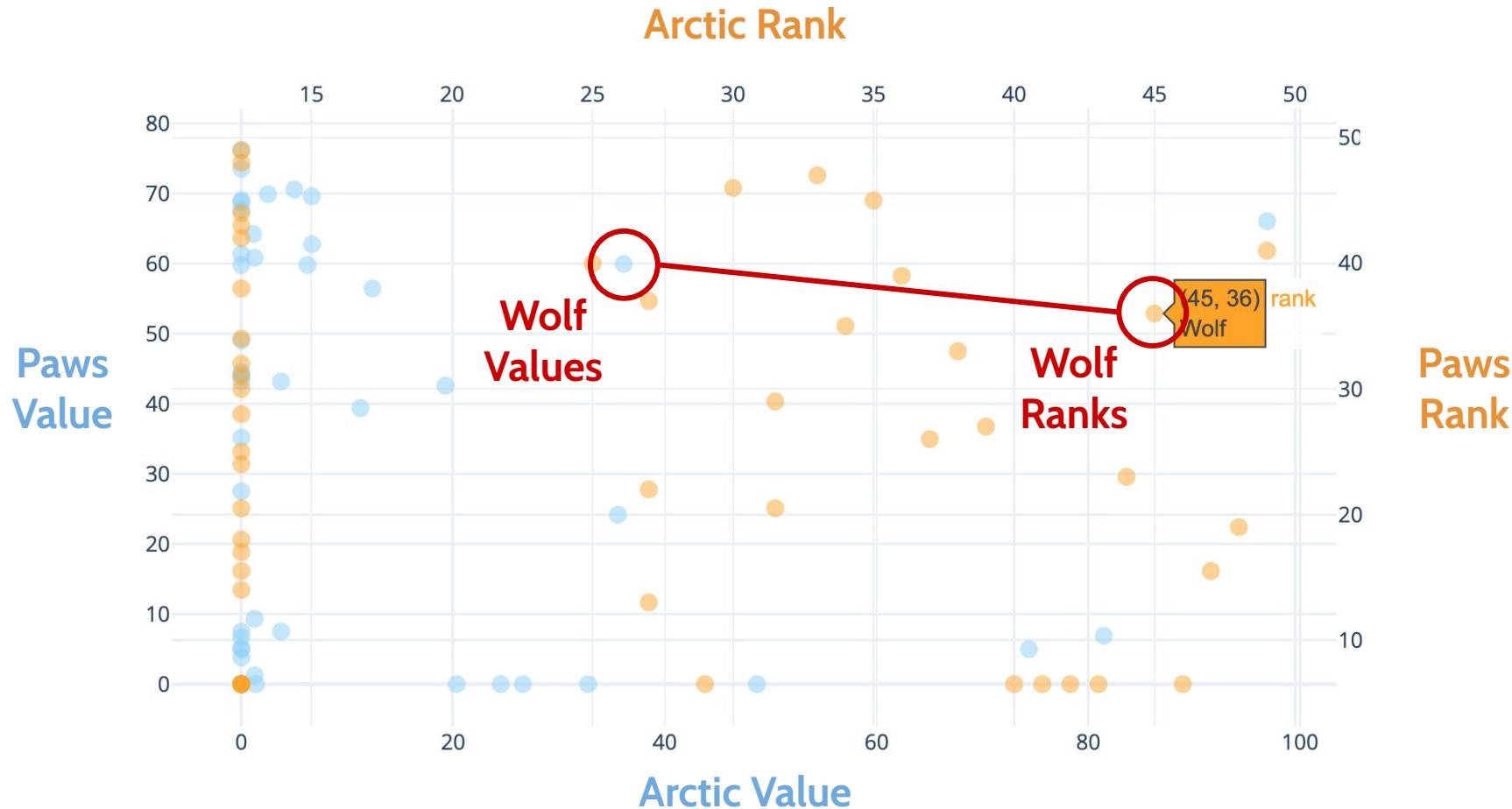


Is it a hunter?

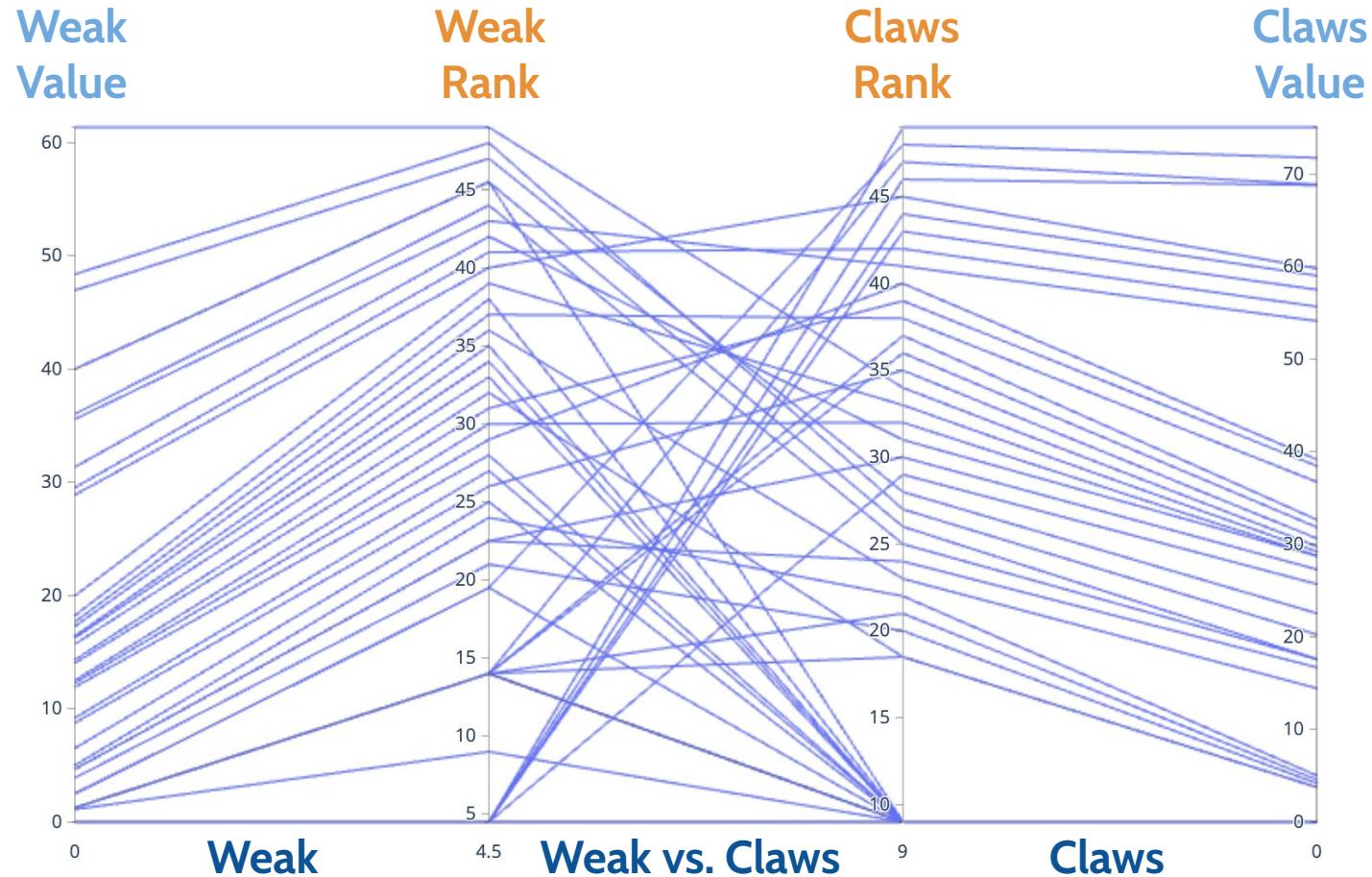
Parallel Coordinate Plot



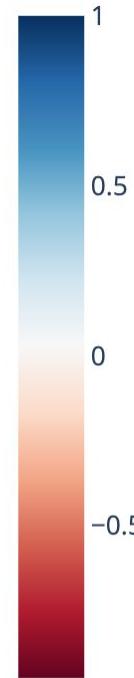
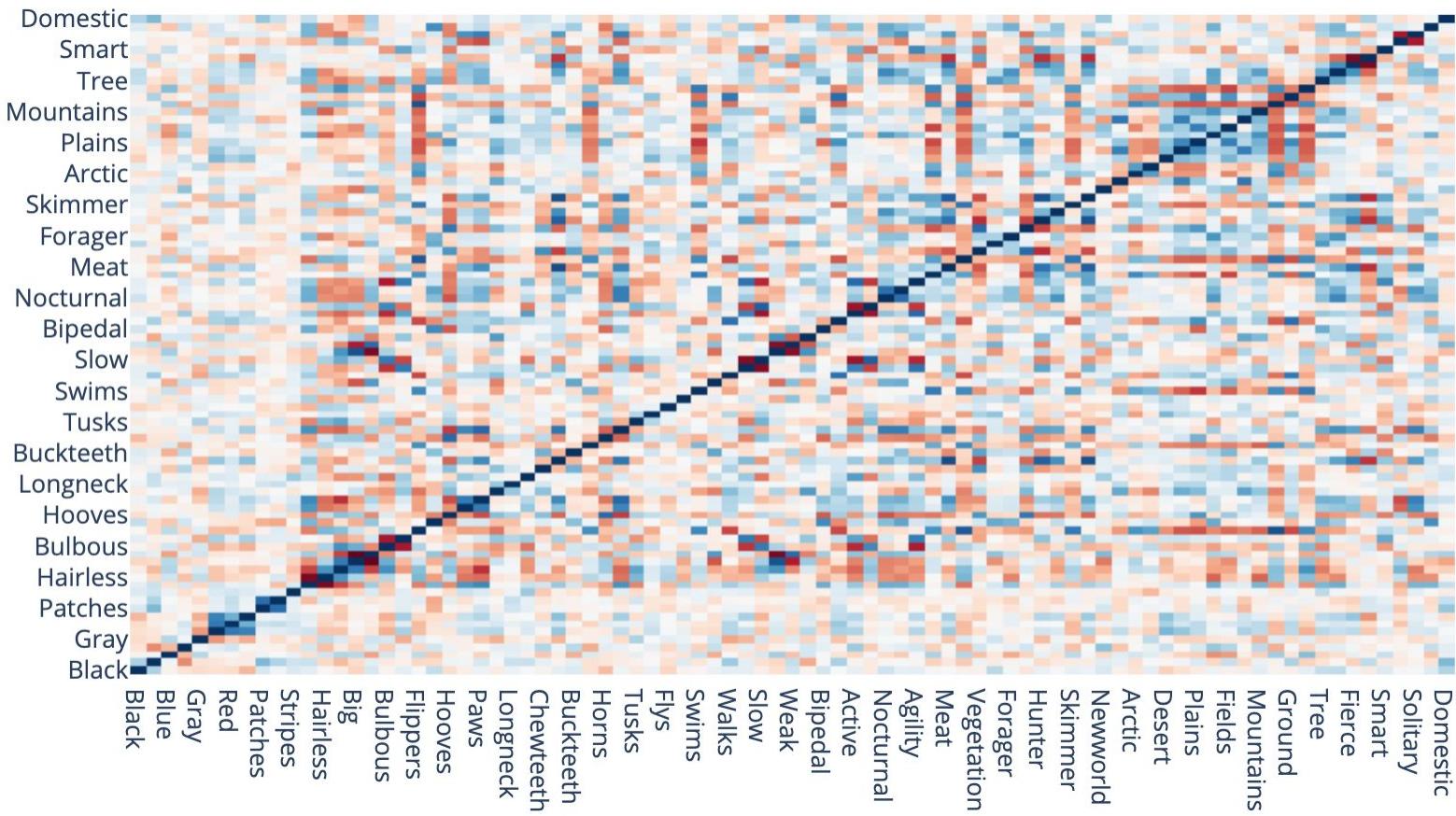
Rank



Parallel Coordinate Plot

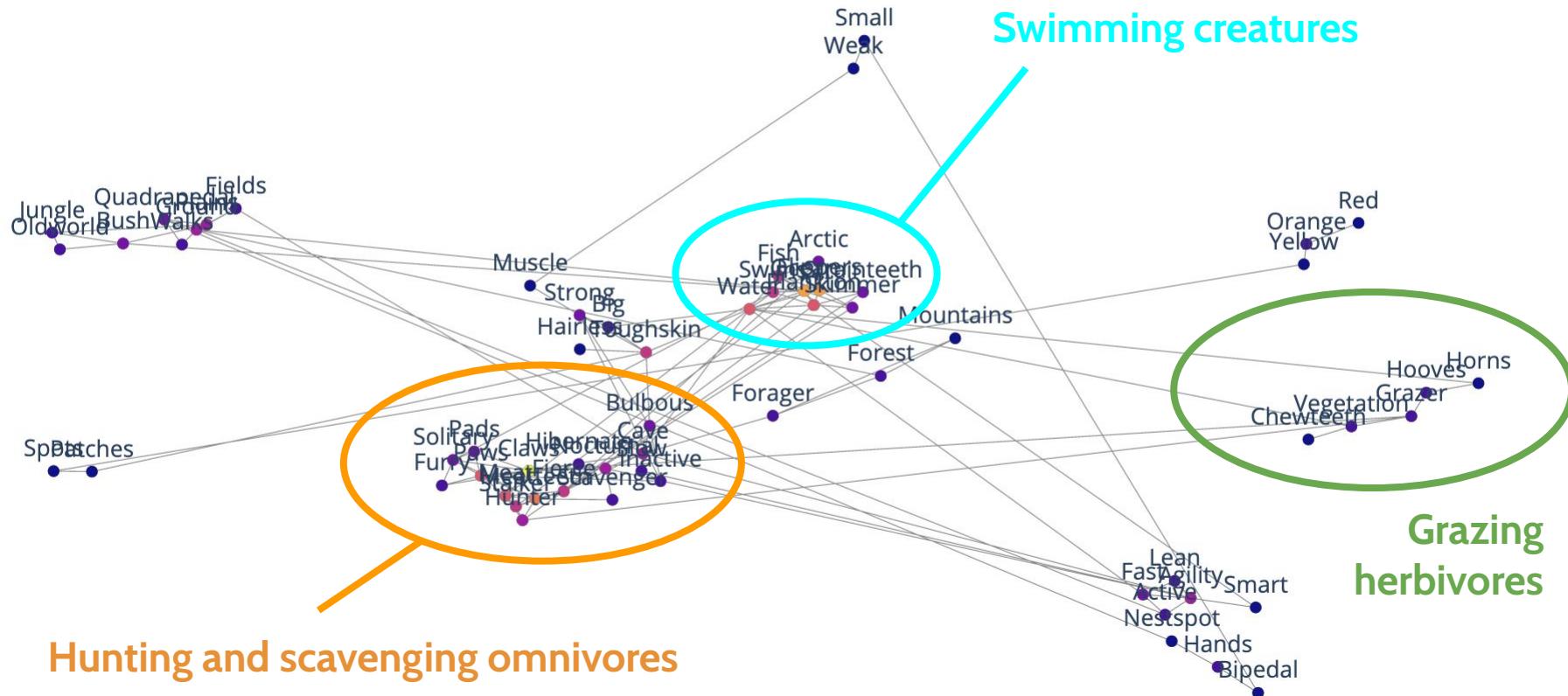


Correlation Matrix (Heatmap)



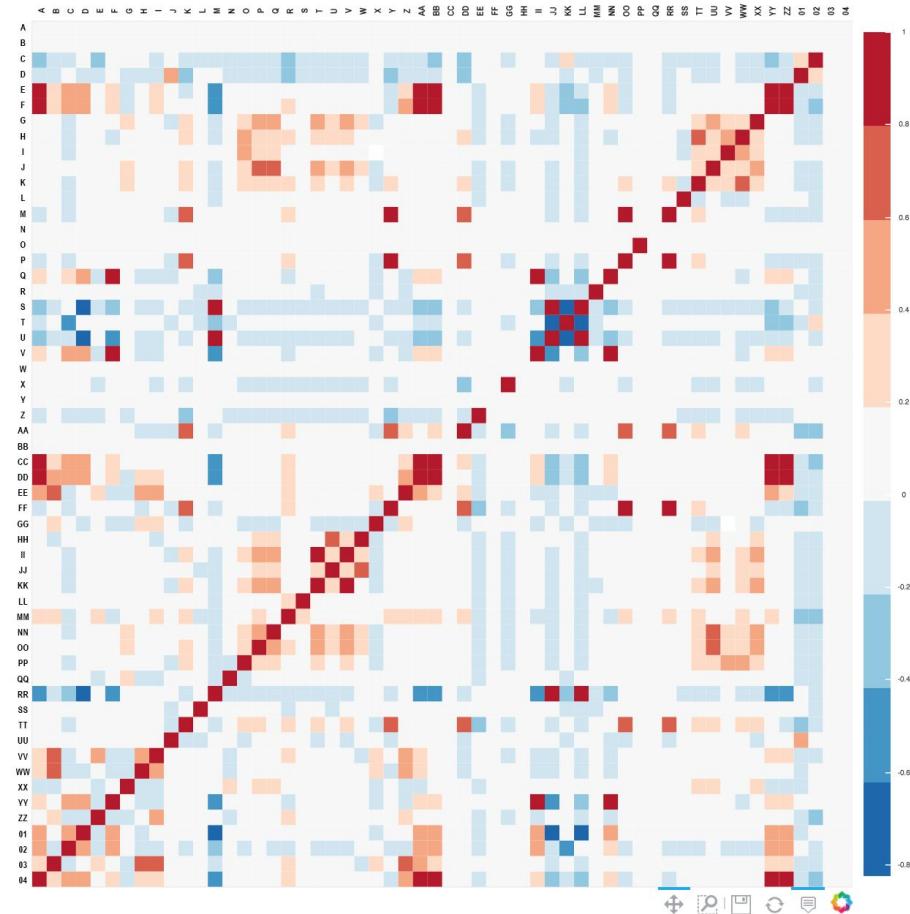
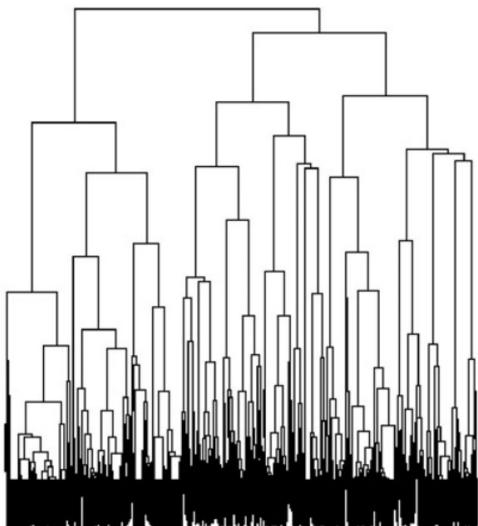
Correlation Network

Showing edges where correlation exceeds 0.6



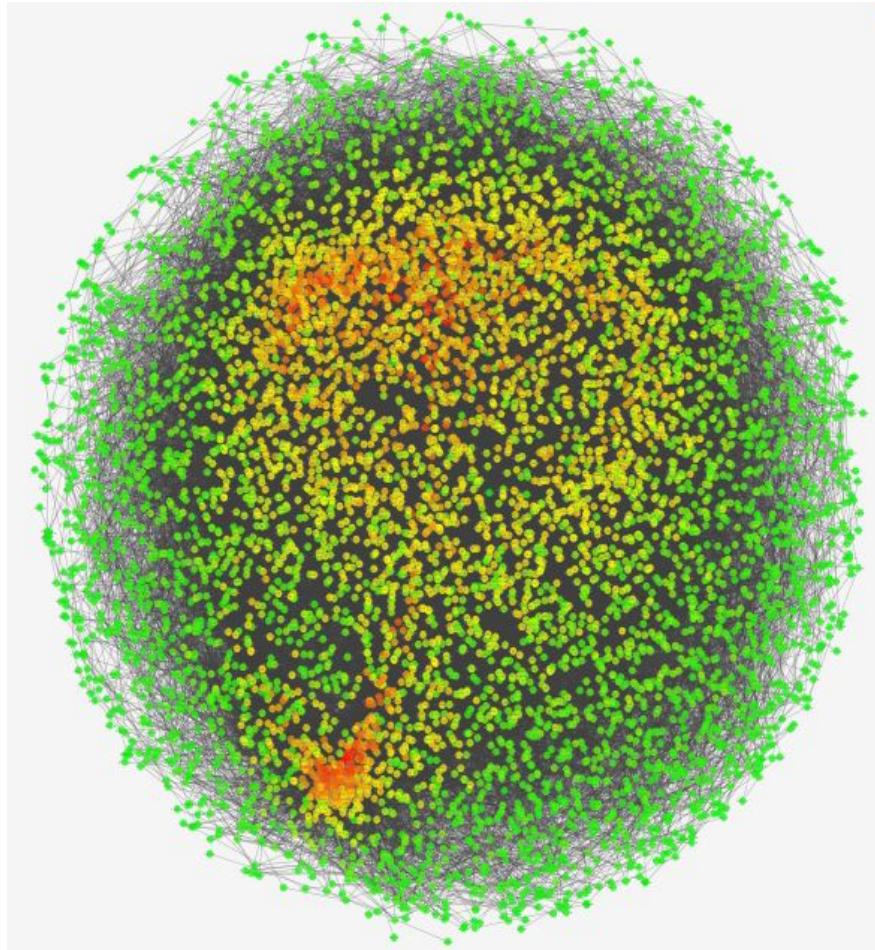
Correlation Network

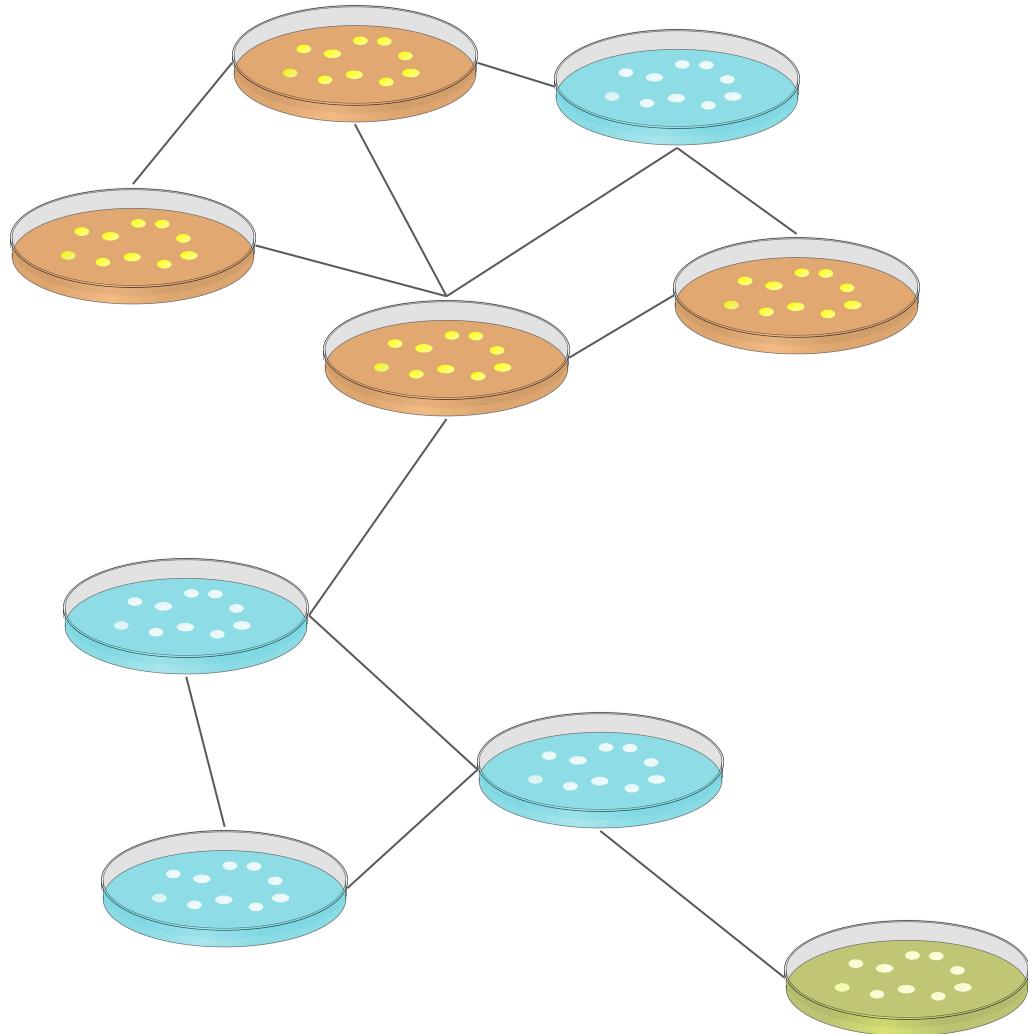
sorted using SciPy's
hierarchical clustering
dendrogram:



Gene Correlation Network

There is a lot of
existing research about
gene correlation /
co-expression
networks...

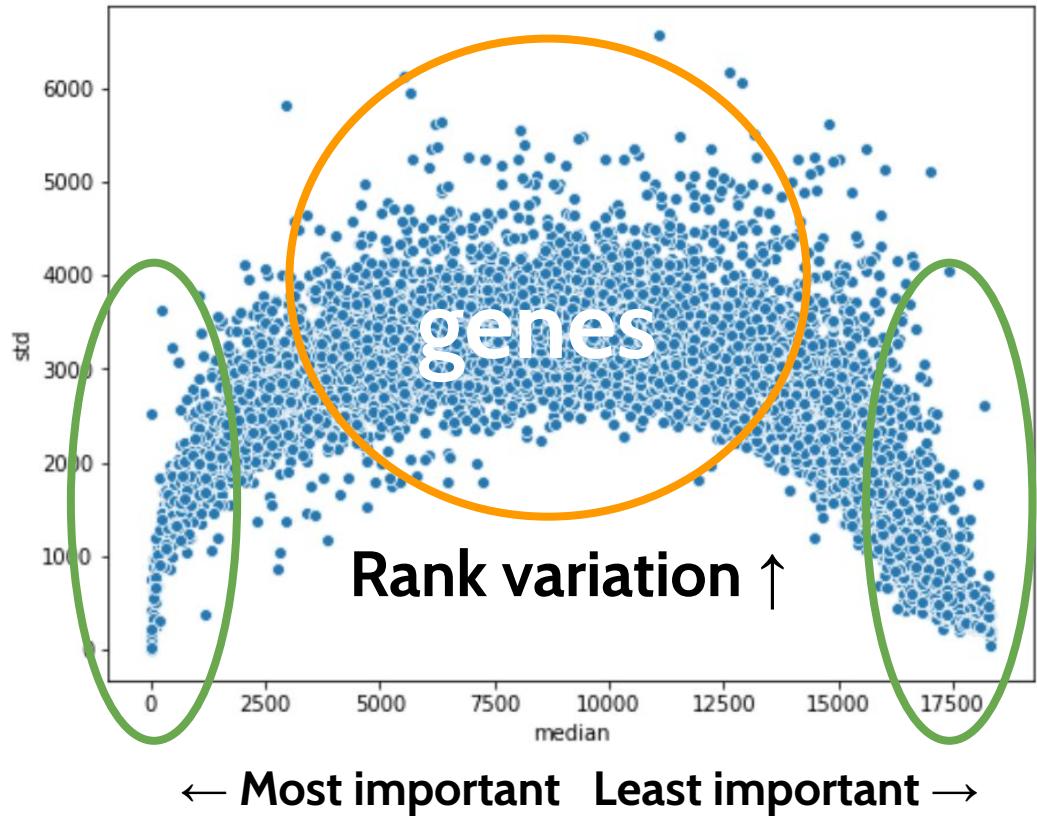




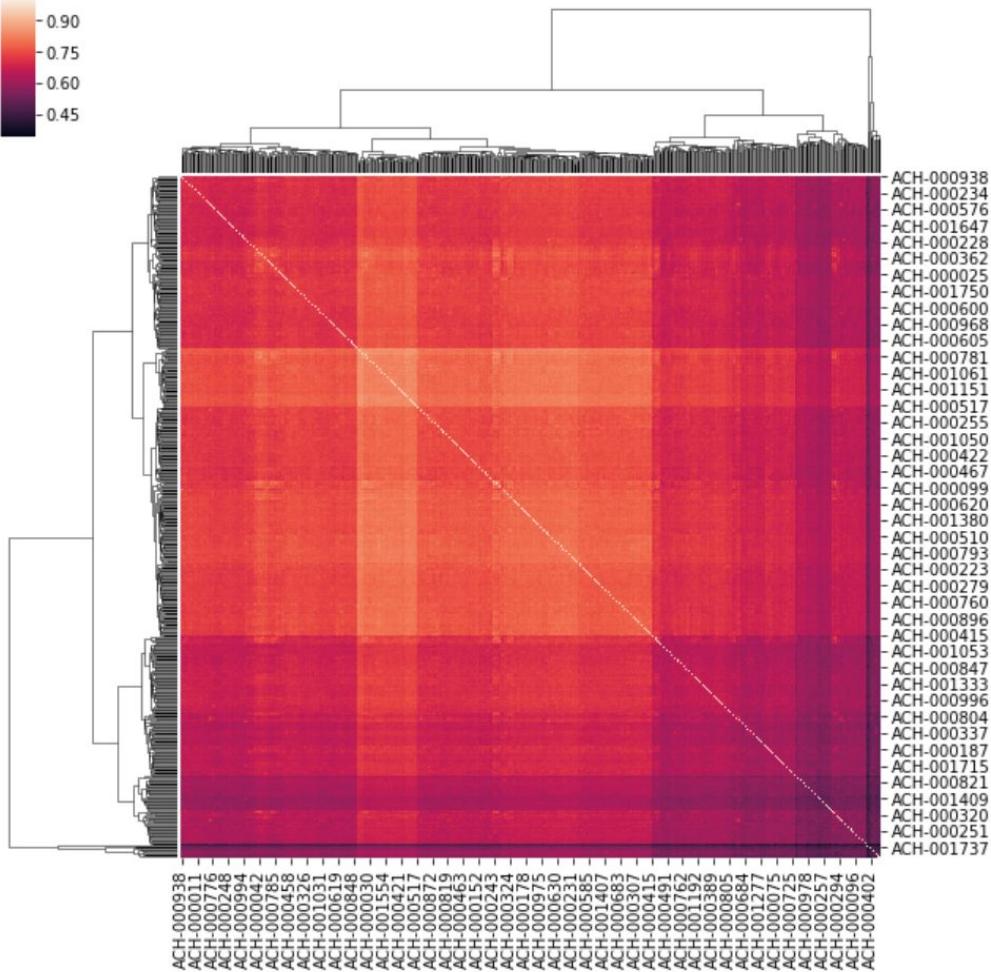
...but how similar are
the cell lines of
different types of
cancer?

There is much
agreement on which
genes are most or least
important...

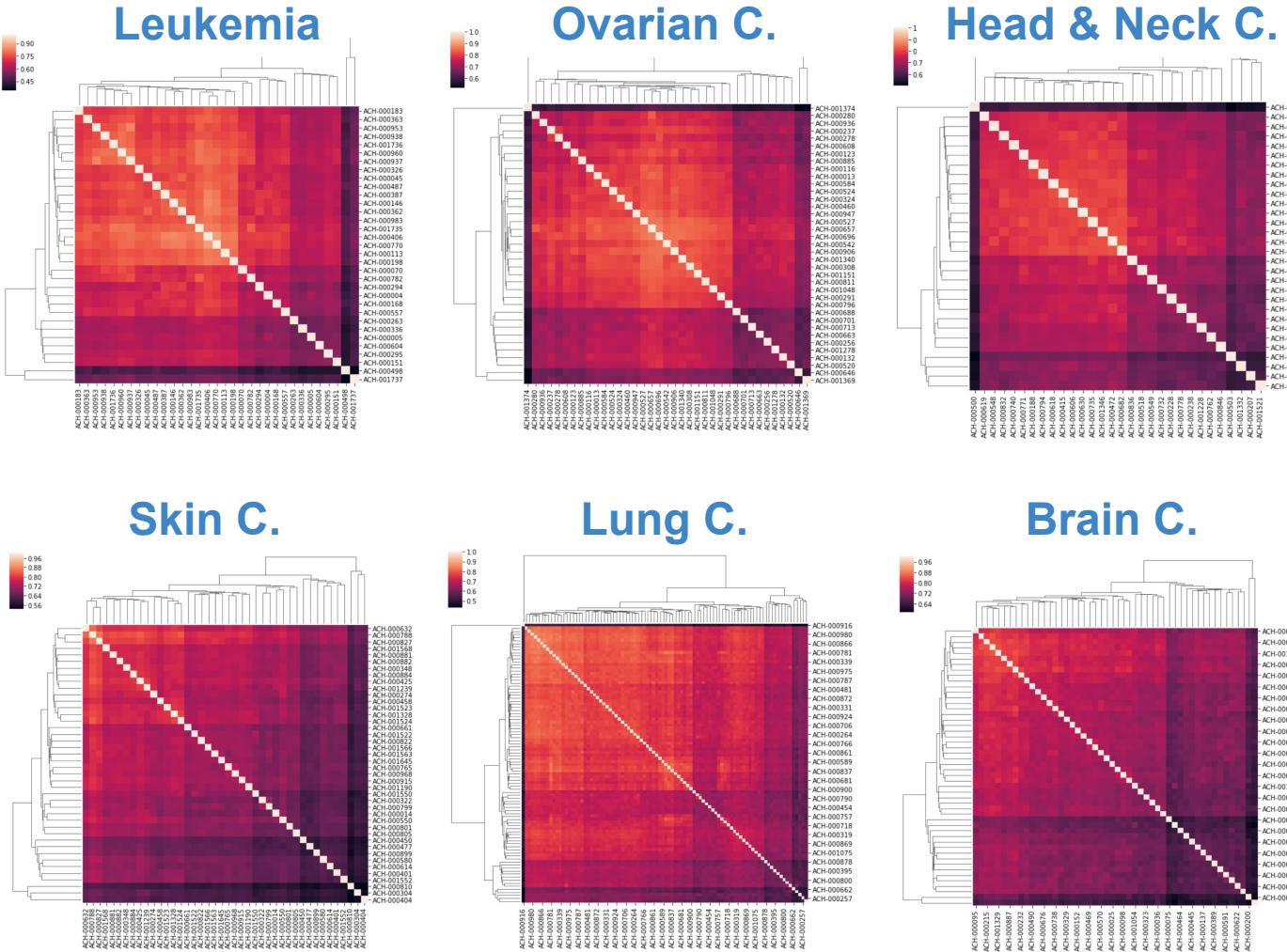
...but there is a lot of
rank variance in the
middle gene dependence
rankings



Gene dependence rank correlations among all cell lines

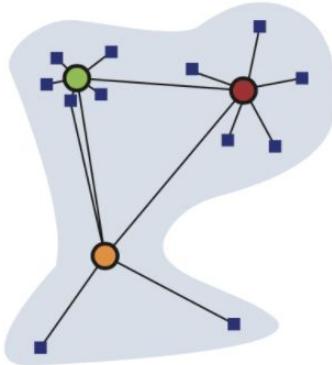


Gene rank correlations between cell lines by primary disease

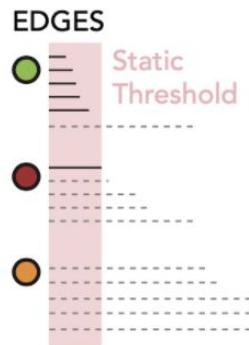
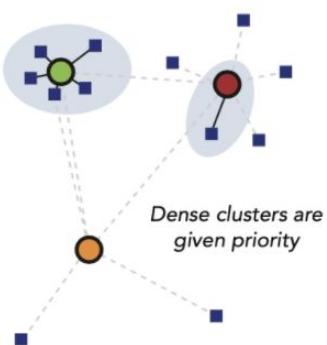


Edge Thresholding

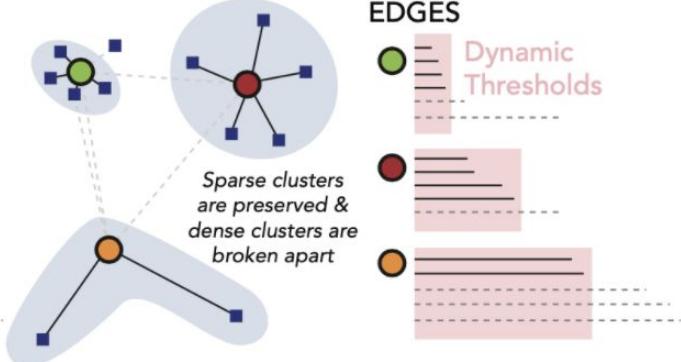
Consider this network:



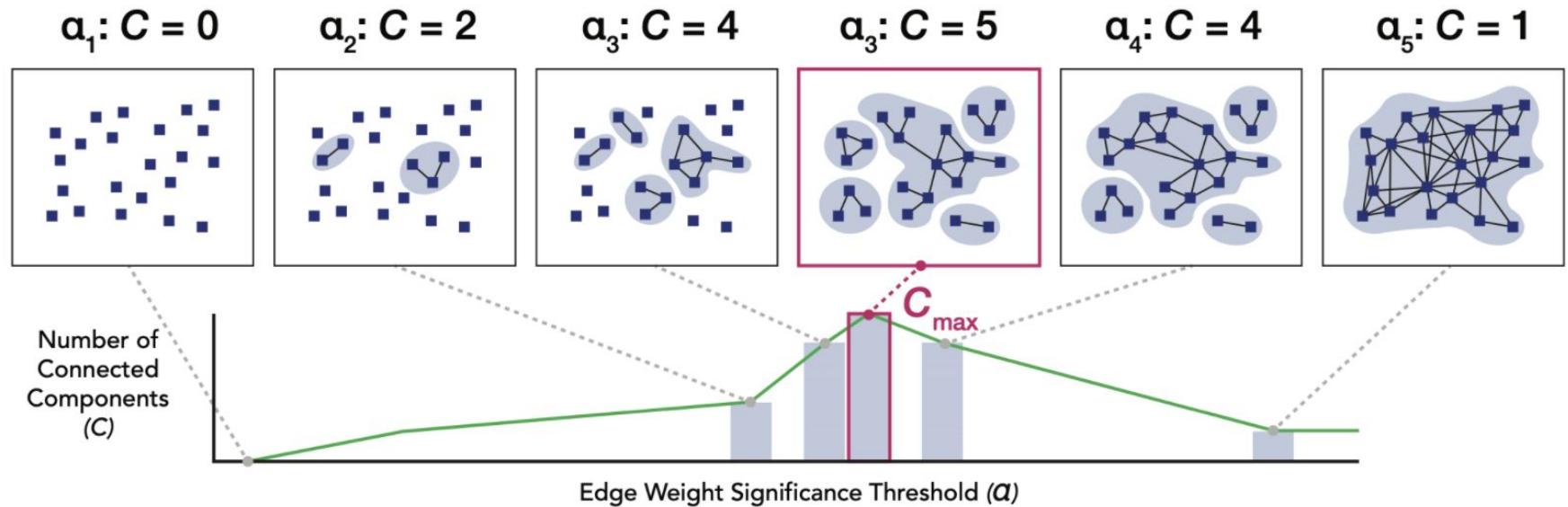
Static thresholding:
Edges below a certain weight are cut



Backbone sparsification:
Edge importance determined by local significance

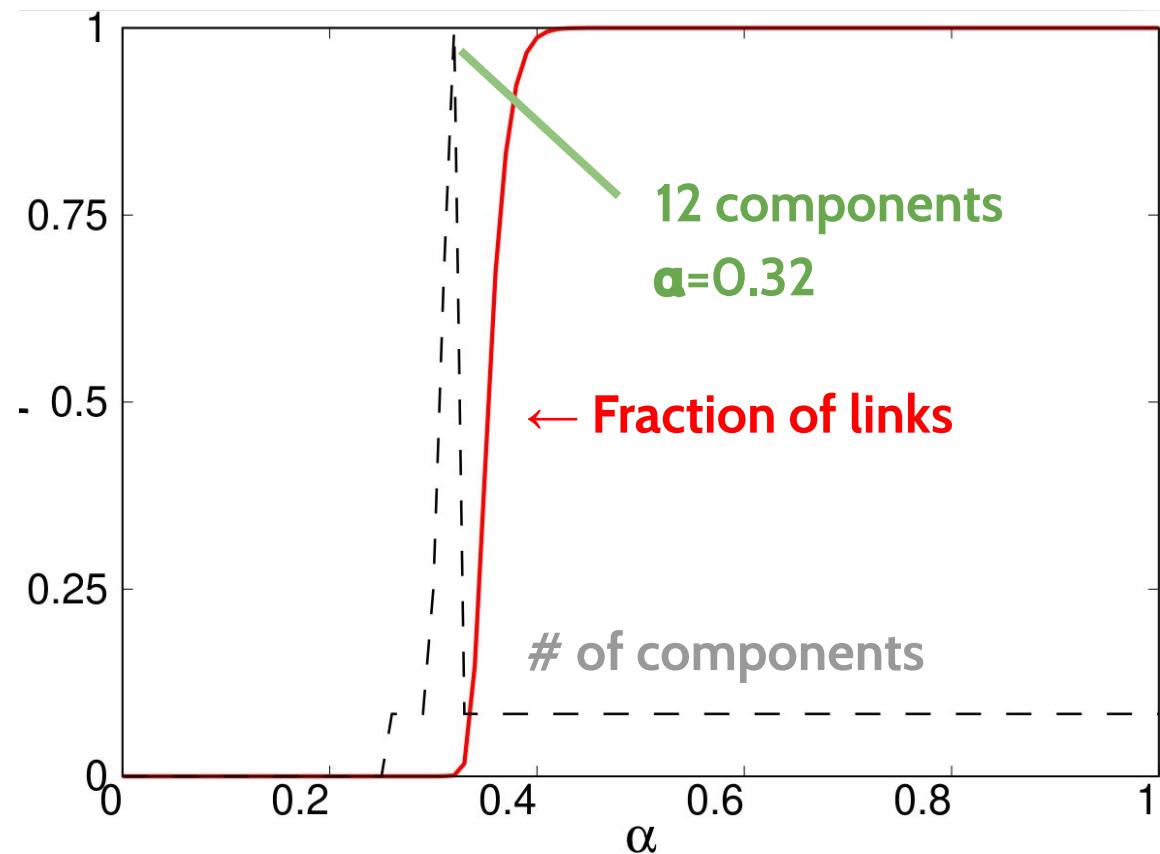


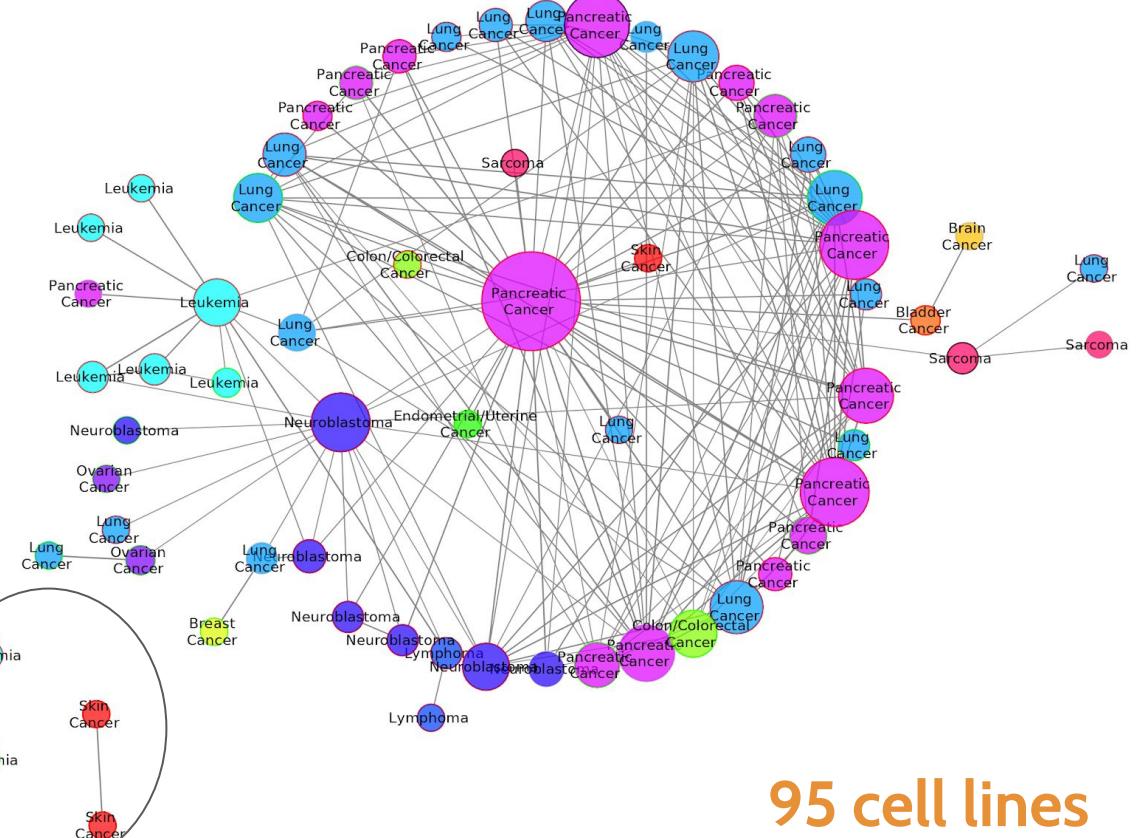
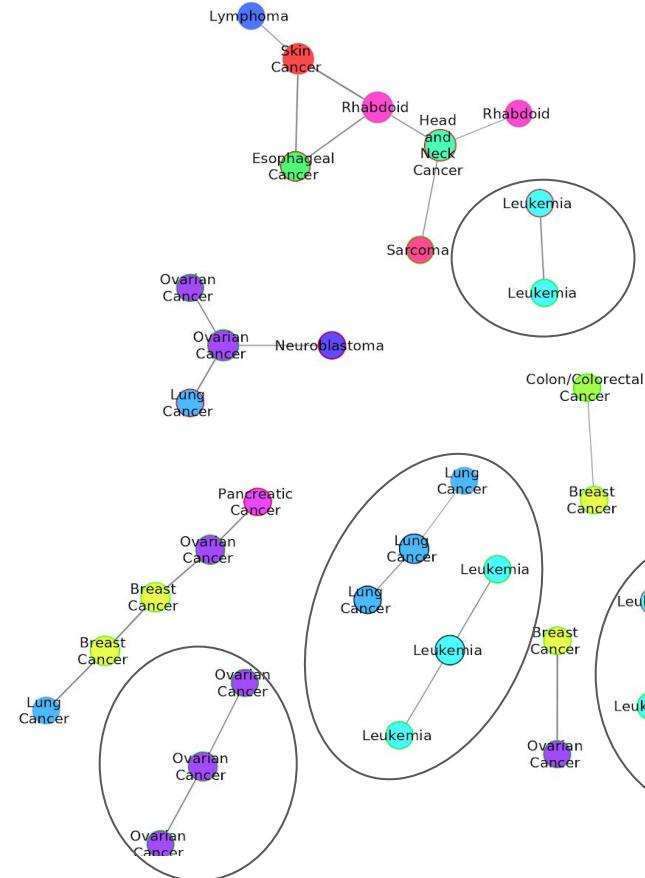
Edge Thresholding



Edge Thresholding

Link weight
thresholding
on the cell line rank
correlation matrix





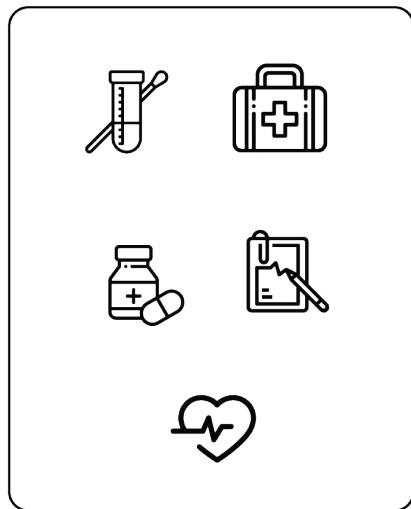
95 cell lines
213 edges

**Ok, we have a few ideas about
measuring ‘interesting-ness’...**

**But what about when we have
different kinds of variables?**



+



	Test Tube	First Aid Kit	Pulse	Clipboard	Medicine
Cloudy Hair	Orange	Cyan	Yellow	Red	Blue
Dark Hair	Red	Green	Orange	Blue	Blue
Light Hair	Orange	Cyan	Orange	Red	Blue
Glasses	Yellow	Green	Orange	Blue	Blue
Bald	Orange	Green	Red	Red	Blue
Curly Hair	Yellow	Cyan	Orange	Red	Blue

We have lots of different data types...



Cardinal



Quantitative

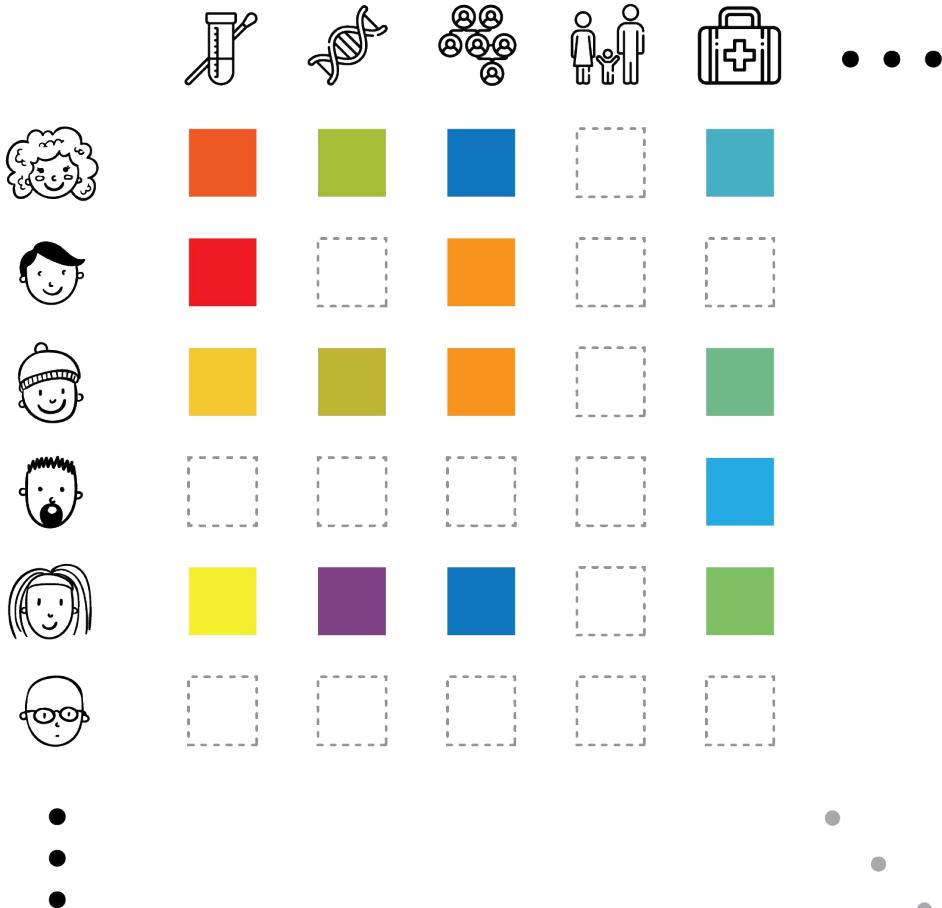


Binary

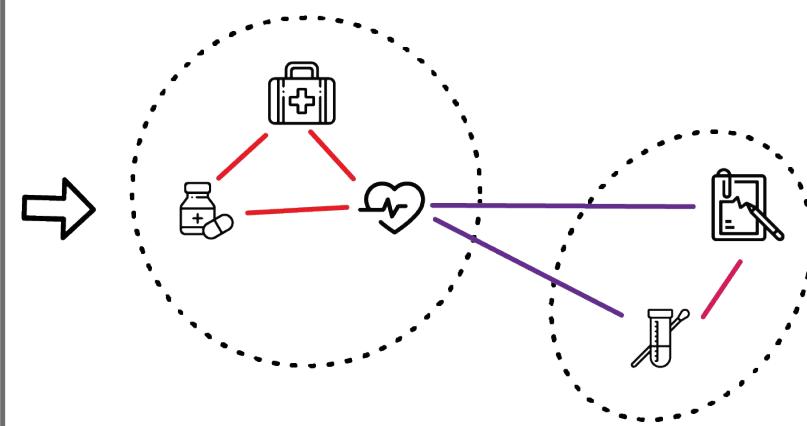
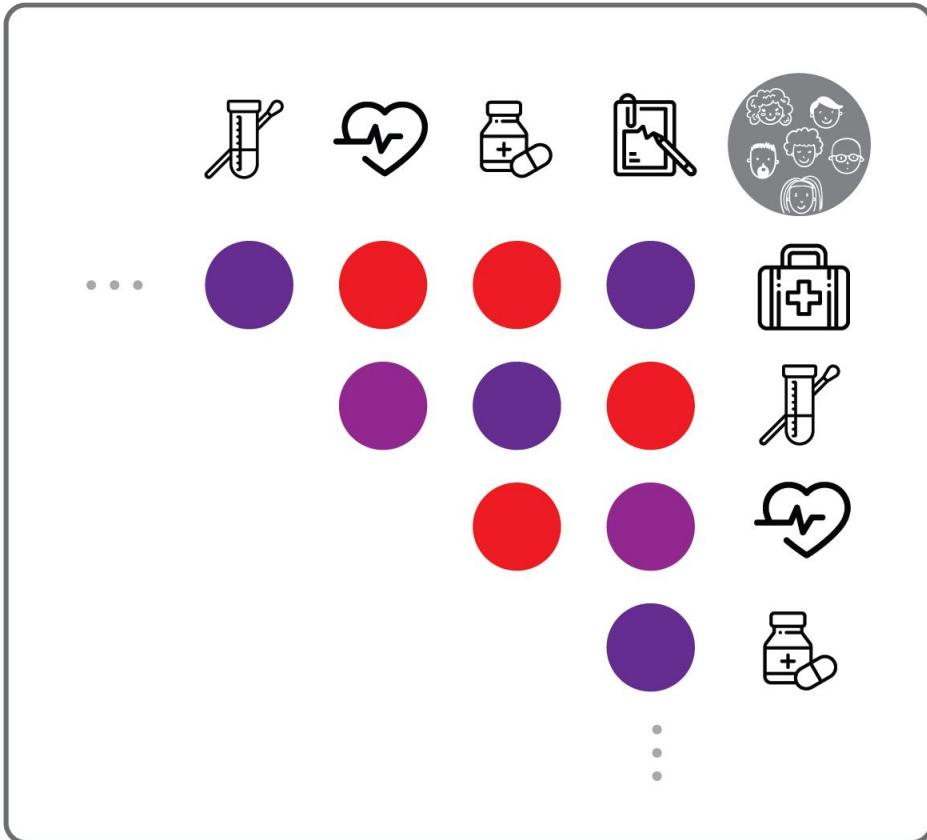


Categorical

**...and we have
a lot of
missing data**

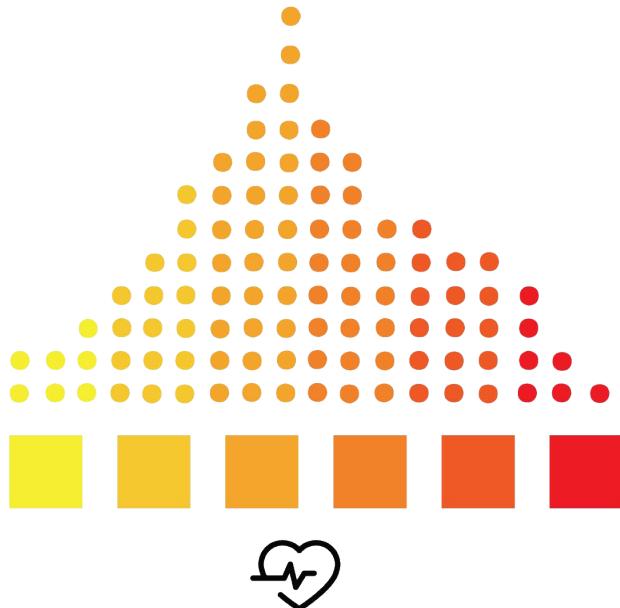




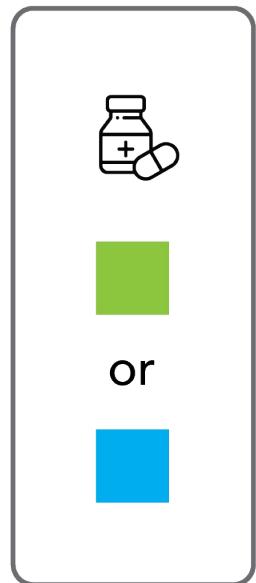


We do this using **mutual information** :

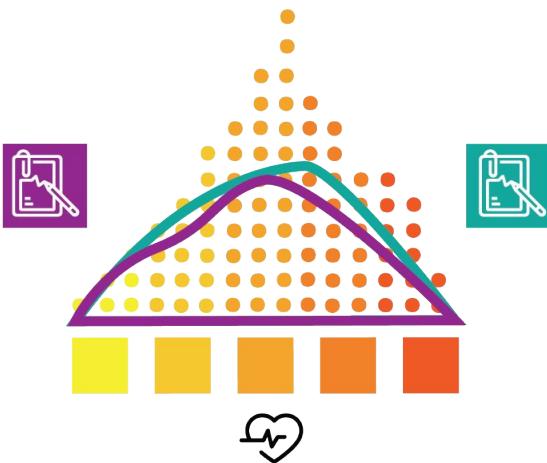
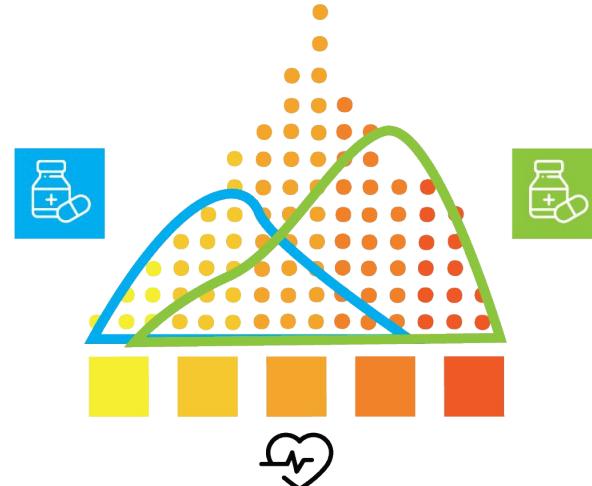
a measure of how much we learn
about one feature
from another



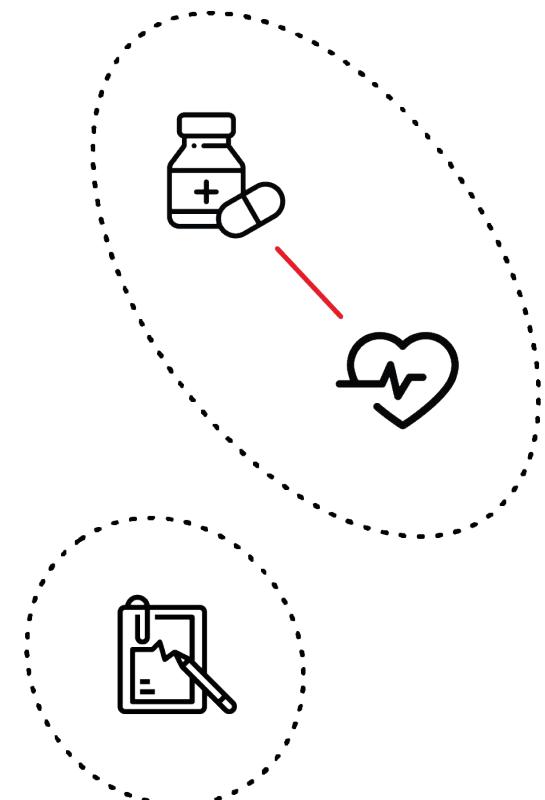
+



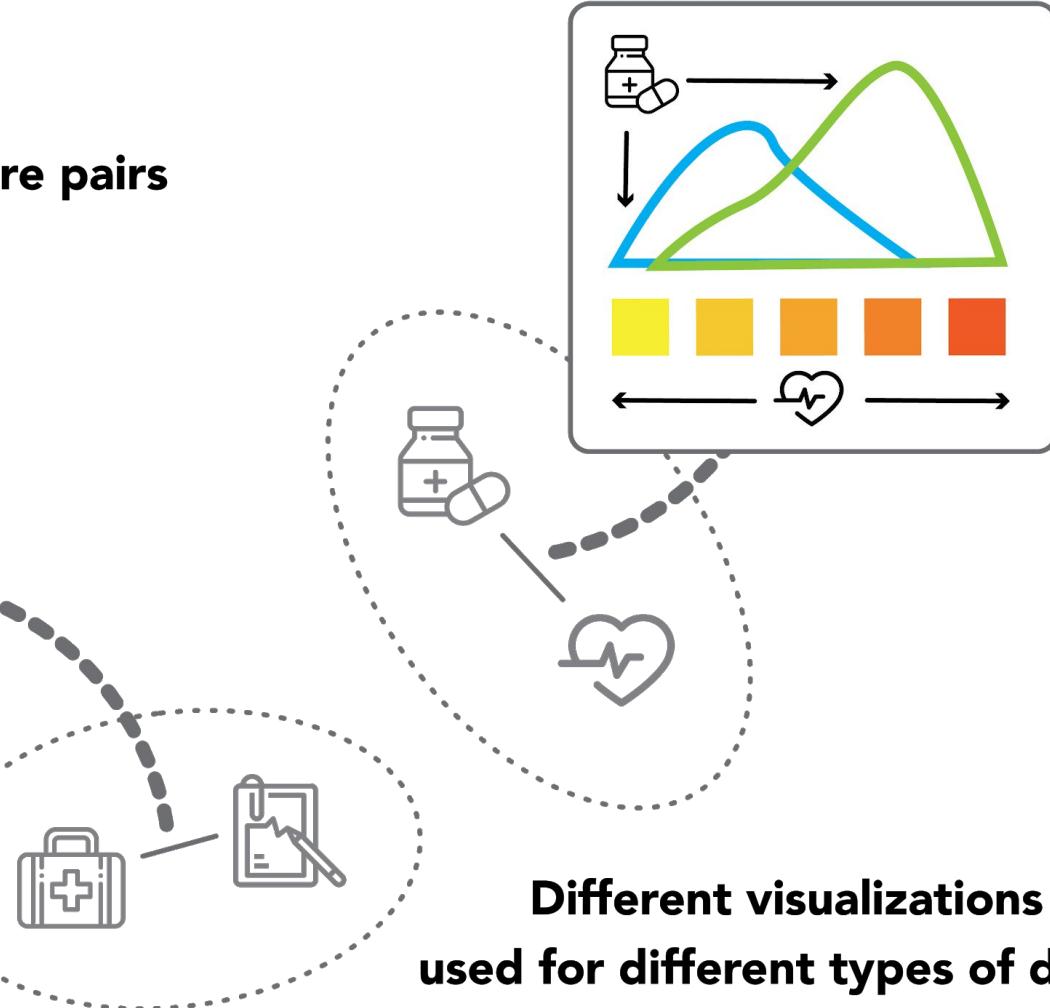
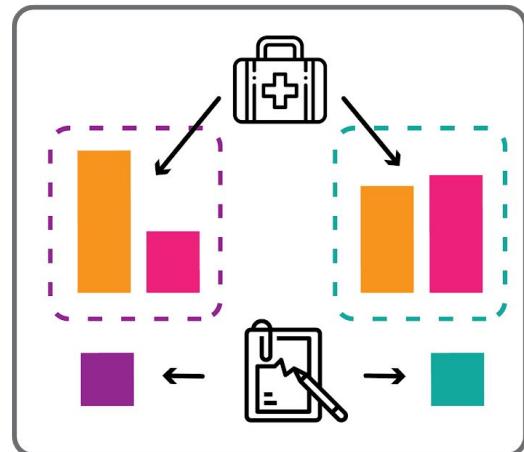
**If a feature pair has high mutual information,
we draw an edge
between them
in the graph**



**Feature pairs
with low mutual
information are
not connected**

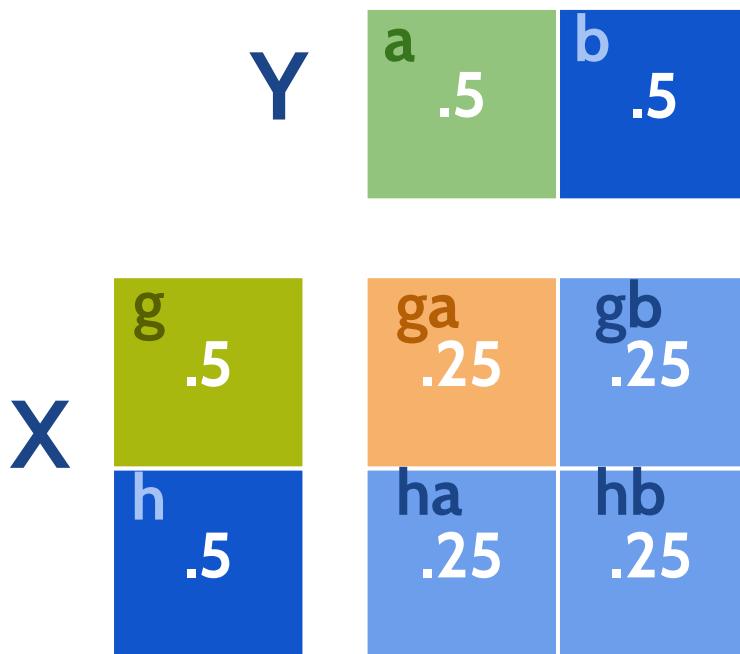


**This feature network
suggests specific feature pairs
to visualize**

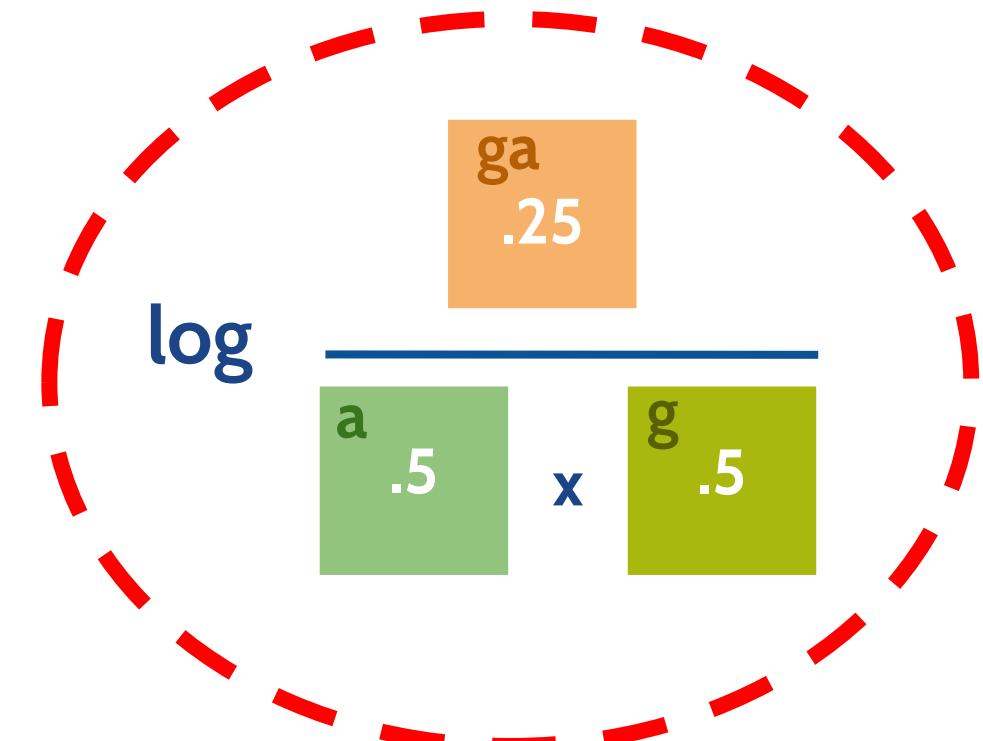


**Different visualizations are
used for different types of data**

Pointwise Mutual Information



Discrete vs. Discrete



Discrete Mutual Information

$$\sum_{X=i} \sum_{Y=j}$$

$$ij \\ .25$$

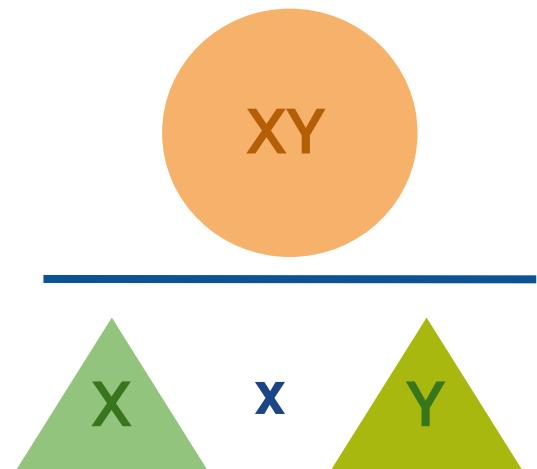
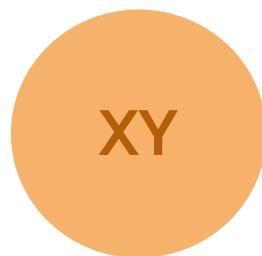
x log

$$\frac{ij \\ .25}{i \\ .5 \times j \\ .5}$$

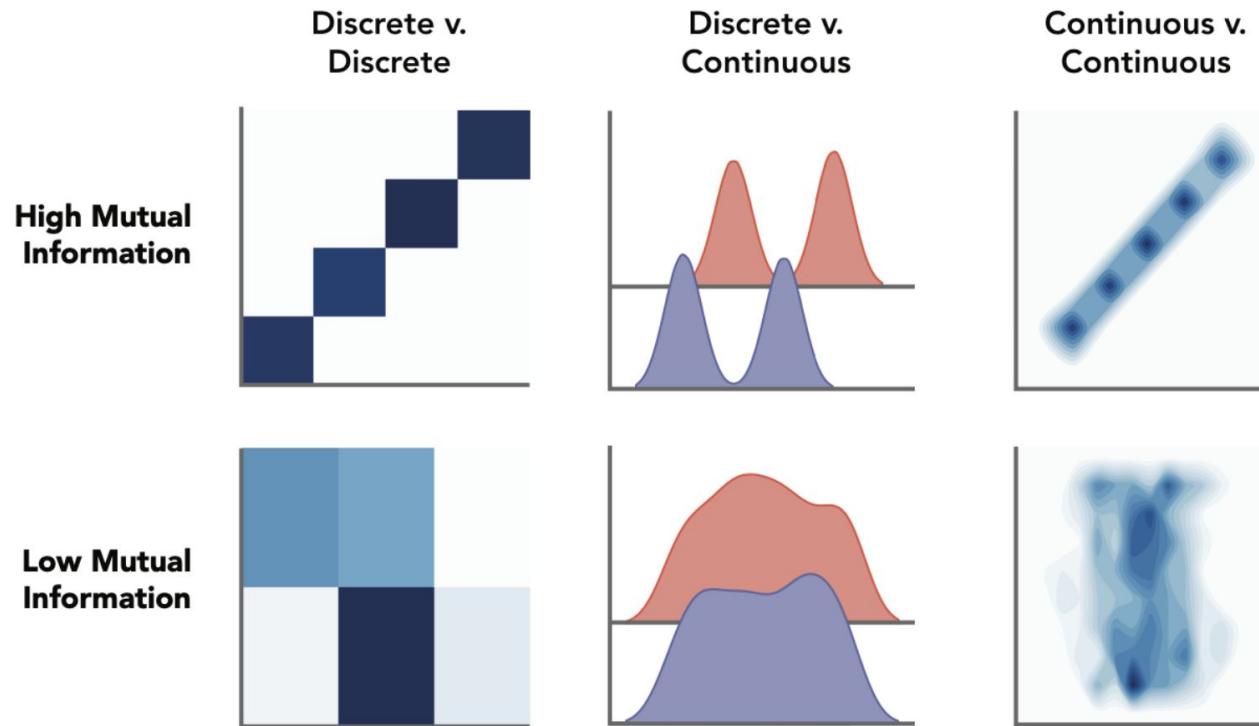
Continuous Mutual Information



$$\int_{x: i_{\min} - i_{\max}} \int_{y: j_{\min} - j_{\max}} x \log$$



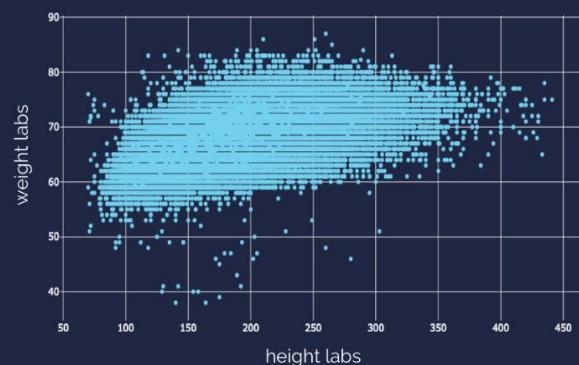
Mutual Information



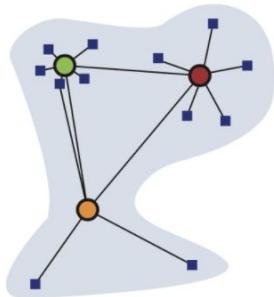


Edge Details

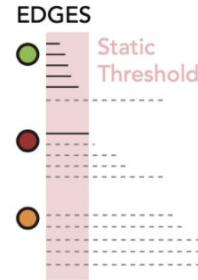
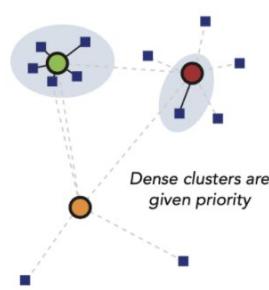
HEIGHT LABS VS WEIGHT LABS



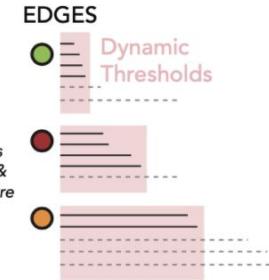
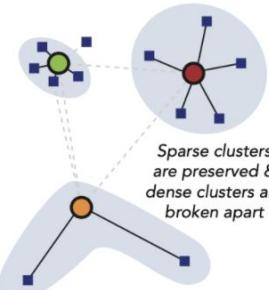
Consider this network:



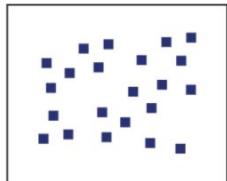
Static thresholding:
Edges below a certain weight are cut



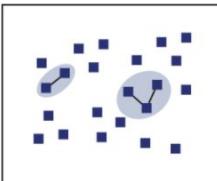
Backbone sparsification:
Edge importance determined by local significance



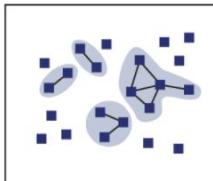
$a_1: C = 0$



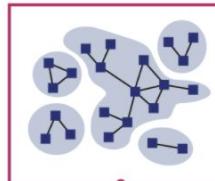
$a_2: C = 2$



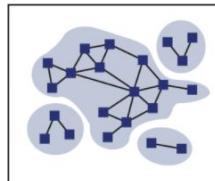
$a_3: C = 4$



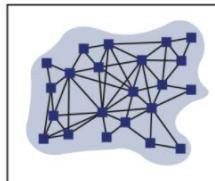
$a_3: C = 5$



$a_4: C = 4$

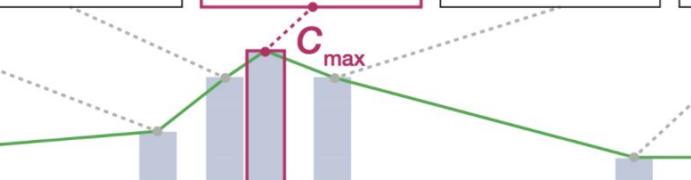


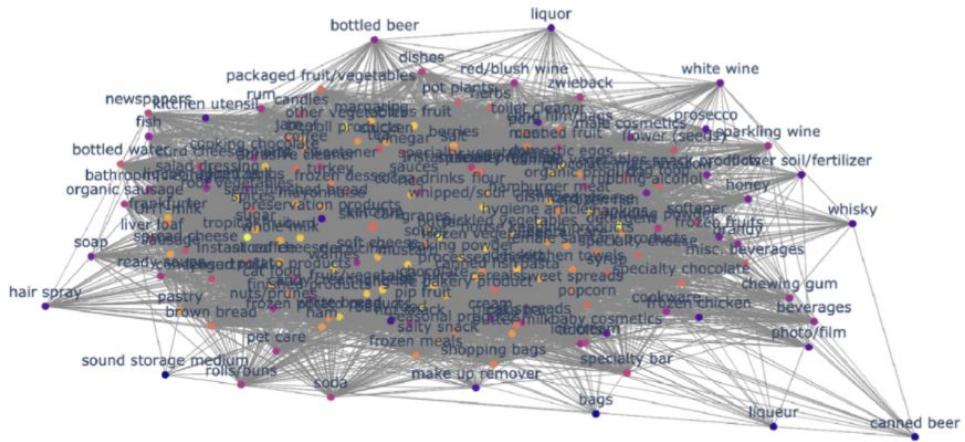
$a_5: C = 1$



Number of Connected Components (C)

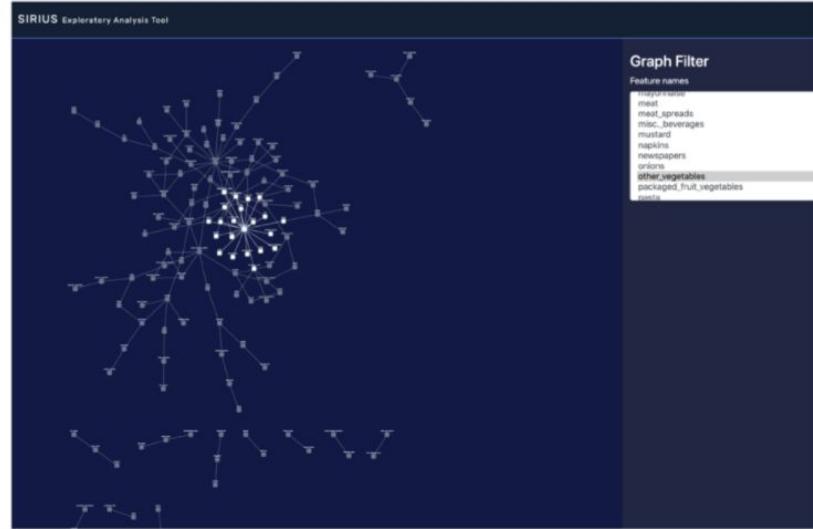
Edge Weight Significance Threshold (α)





Association Network

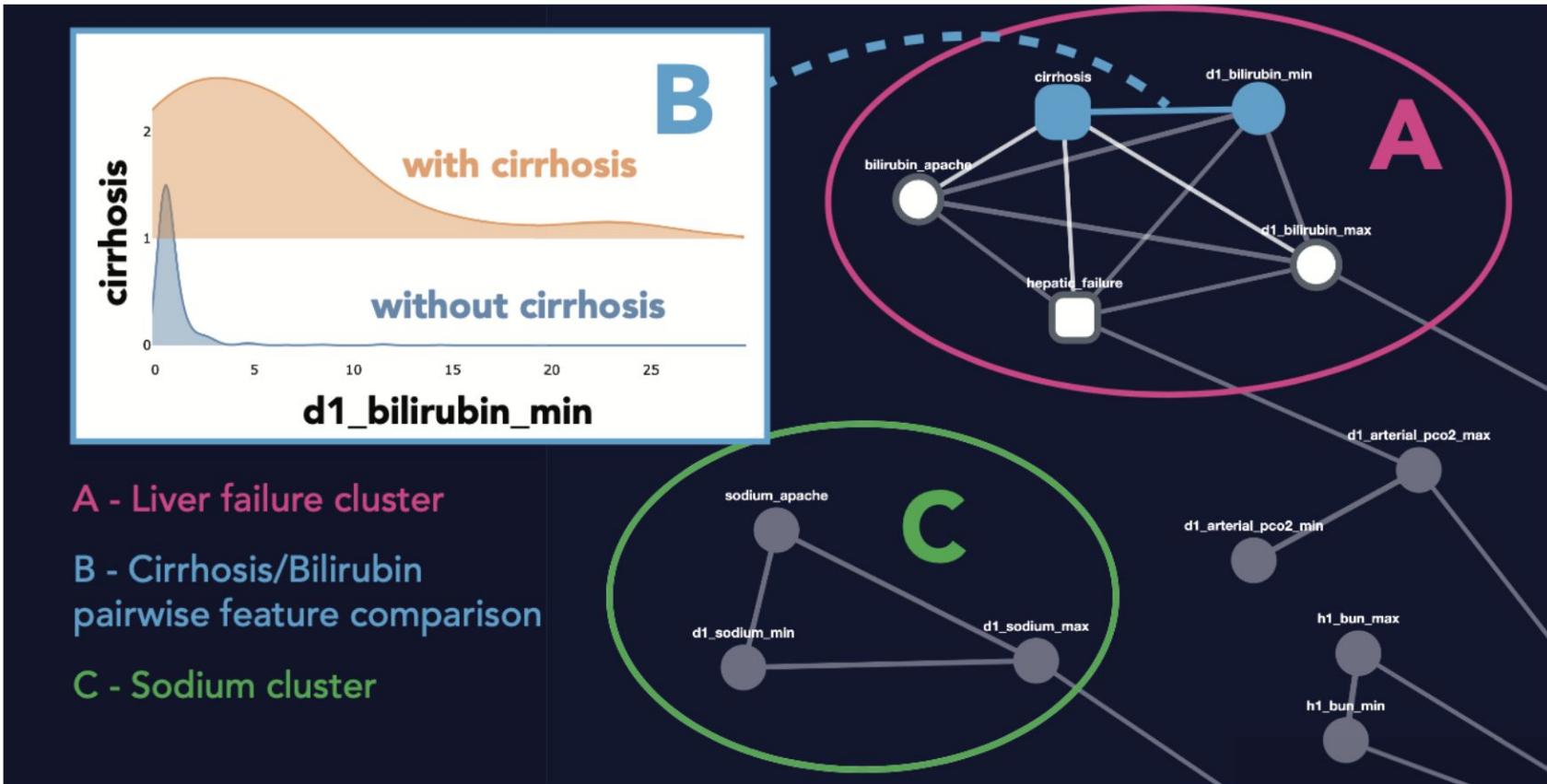
*Lift threshold at static value,
such that 50% of edges remain*



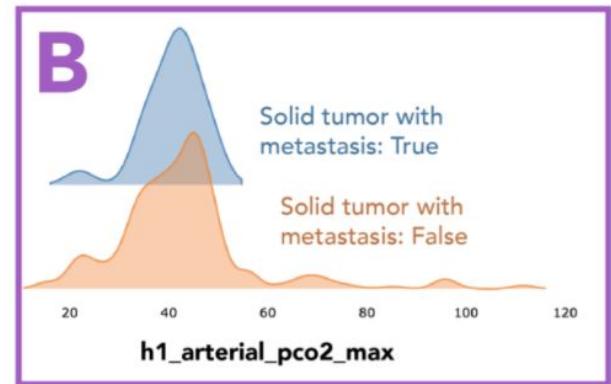
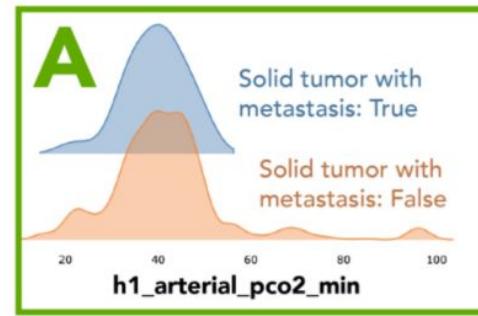
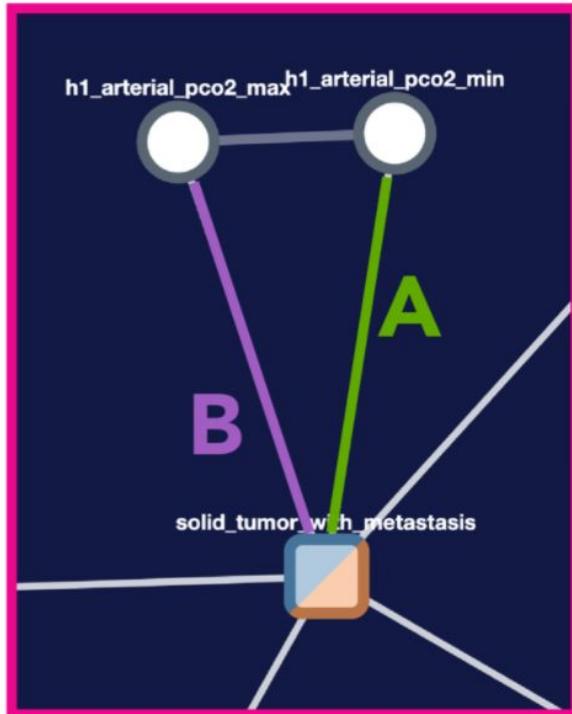
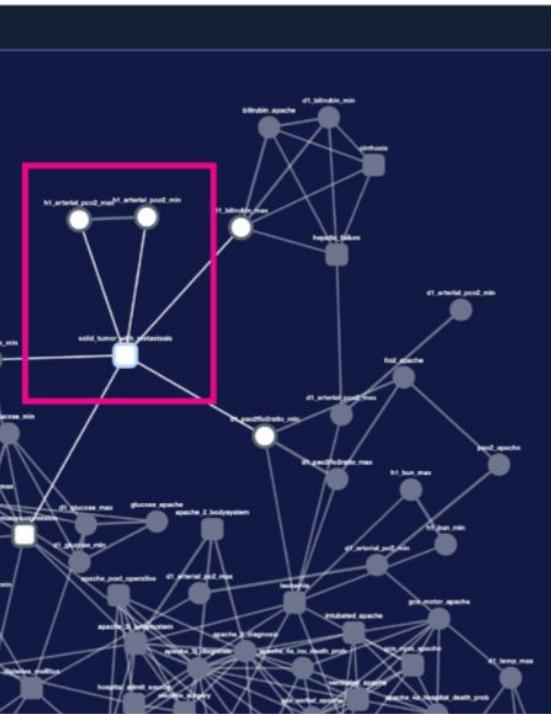
Sirius Mutual Information Network

*Mutual information threshold dynamically
using alpha value from backbone sparsification*

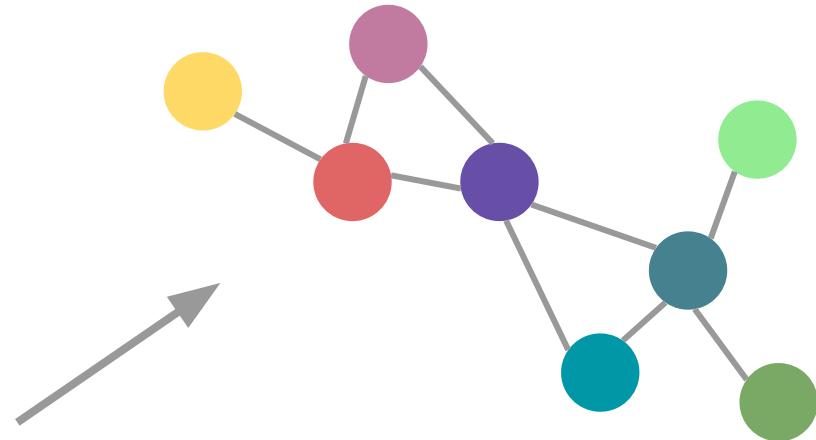
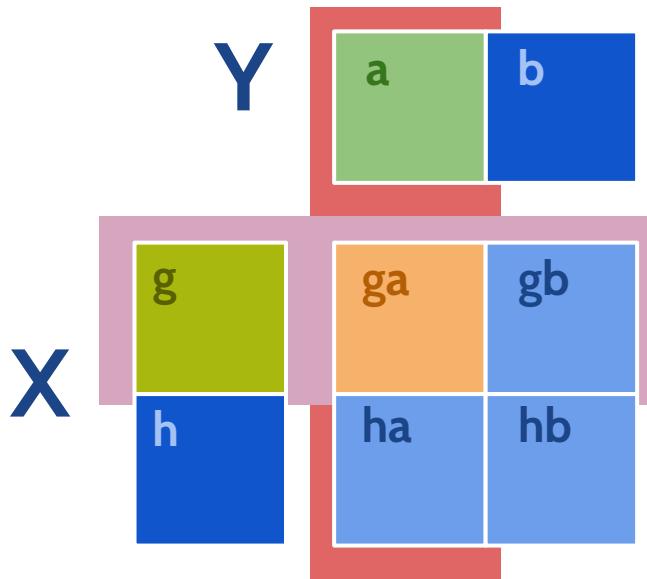
Returning to Pairwise Comparisons



Returning to Pairwise Comparisons



Thanks!



twitter:
@artistjaneadams