

# Natural Language Processing

Lecture 11

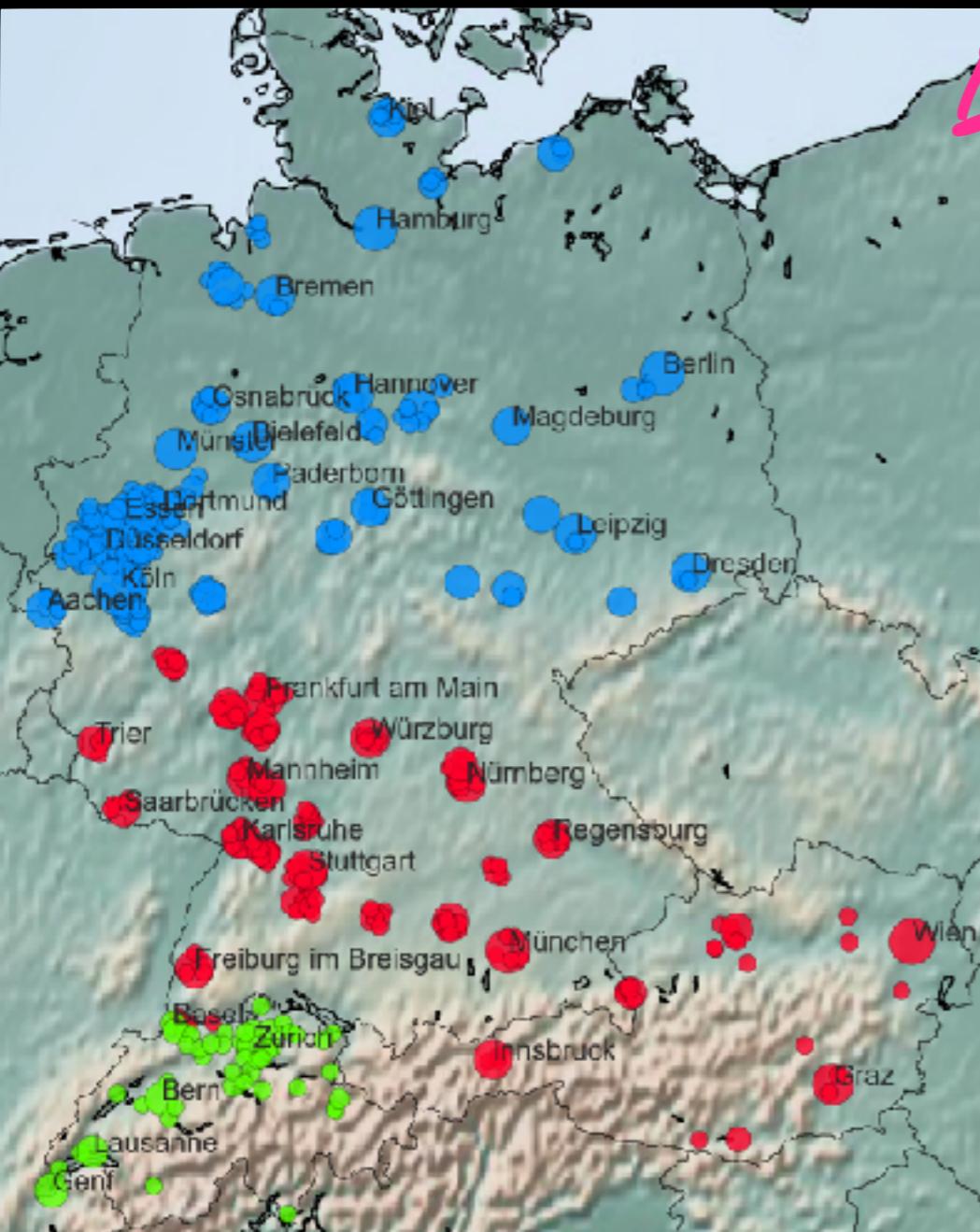
Dirk Hovy

[dirk.hovy@unibocconi.it](mailto:dirk.hovy@unibocconi.it)

 @dirk\_hovy

Bocconi

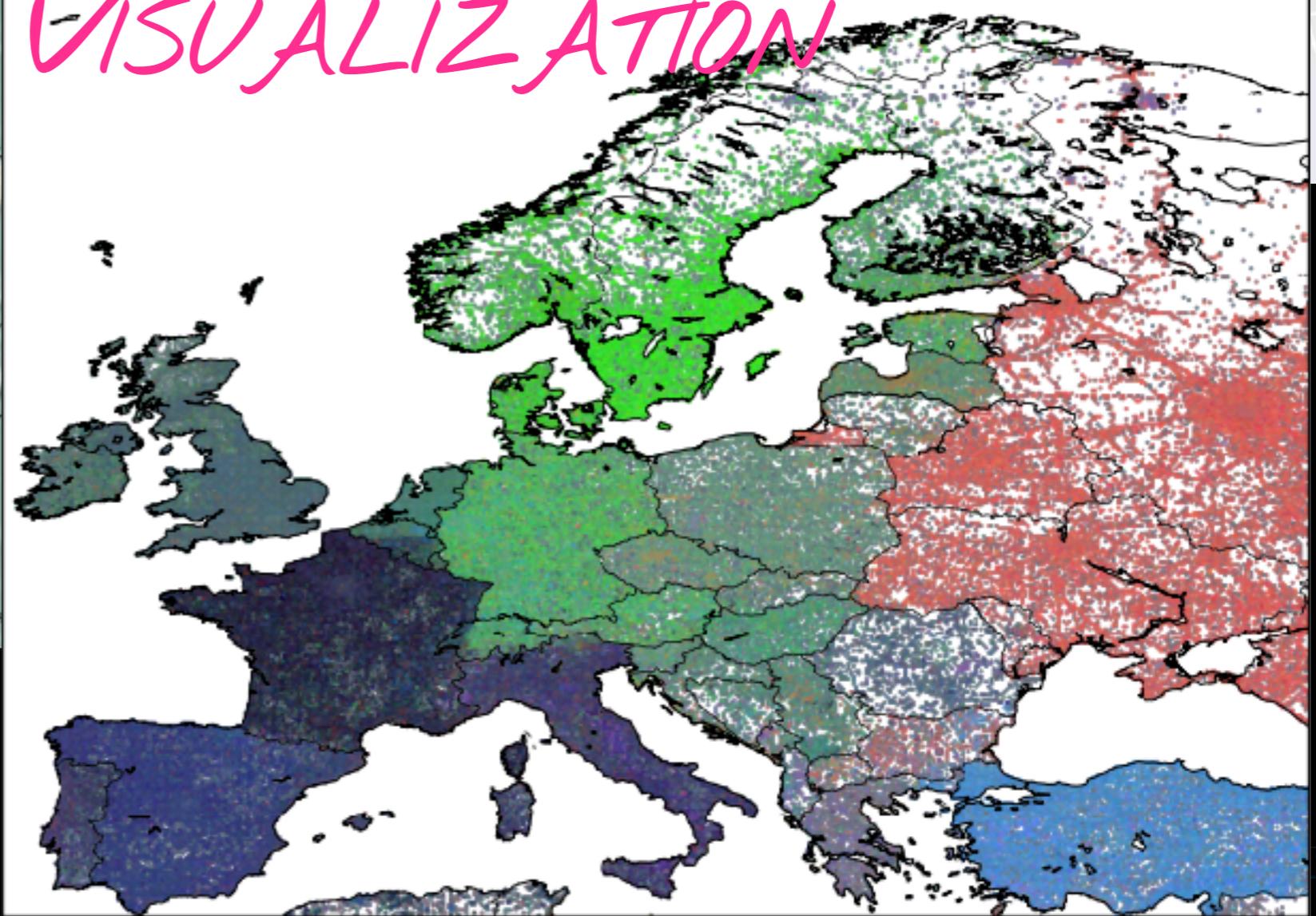
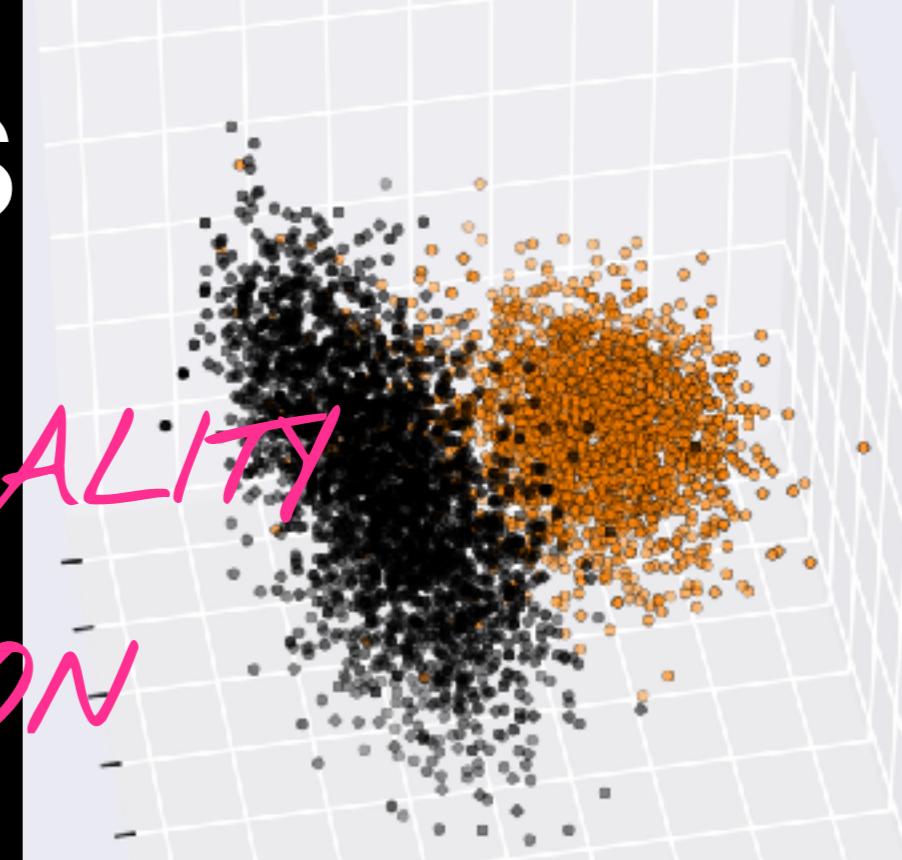
# Examples



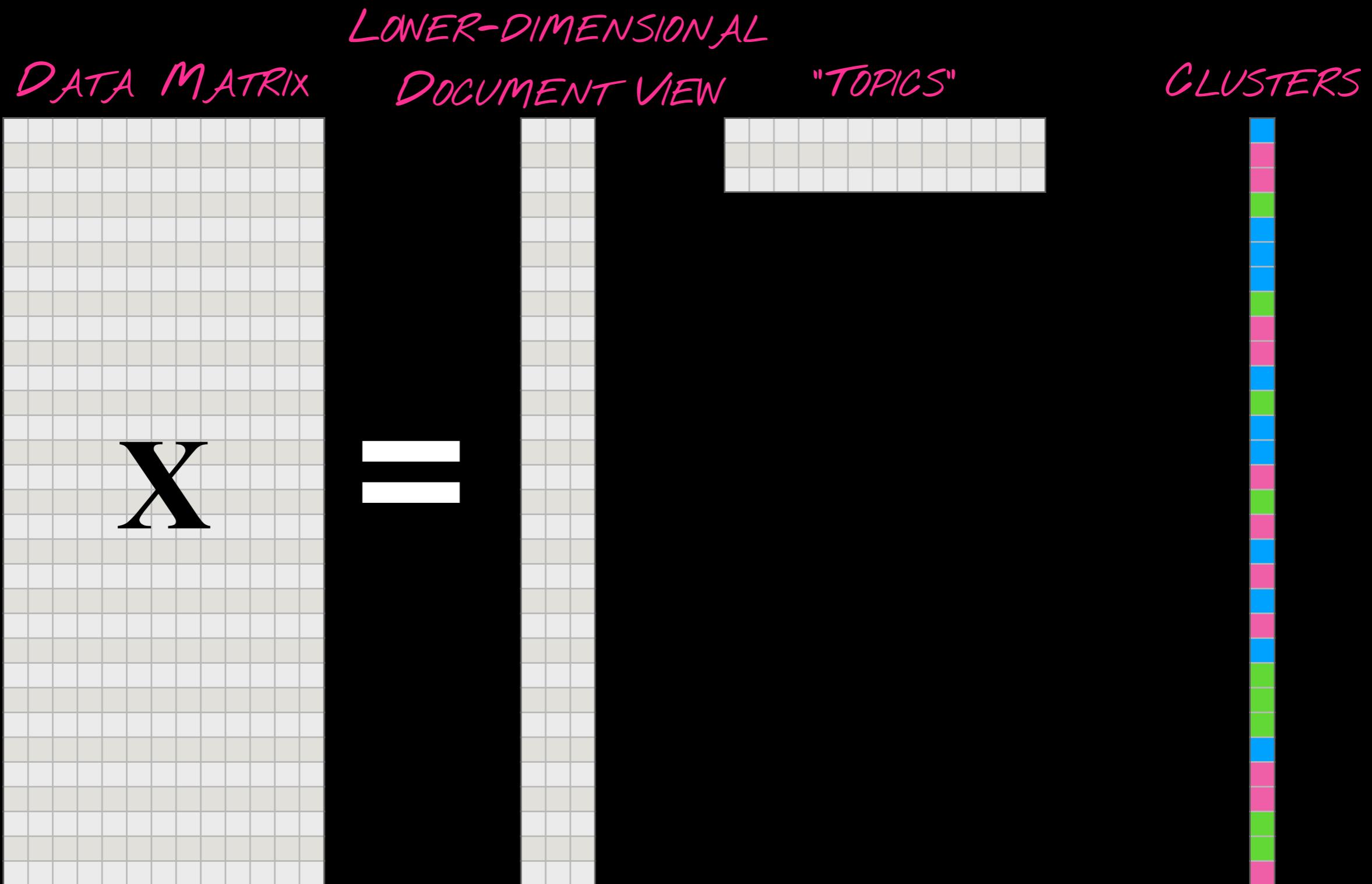
CLUSTERING

DIMENSIONALITY  
REDUCTION

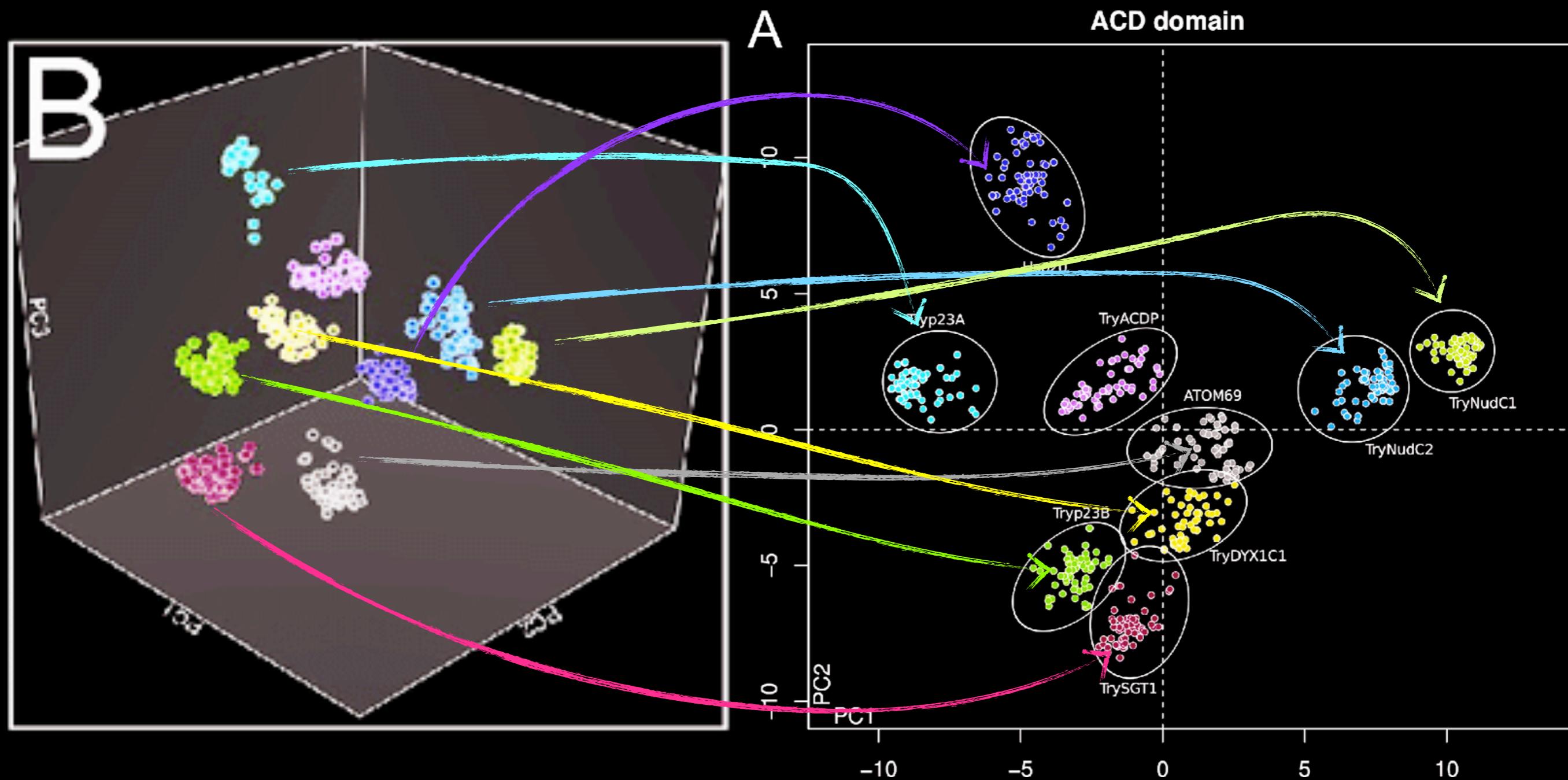
VISUALIZATION



# Latent Dimensions



# Latent Dimensions

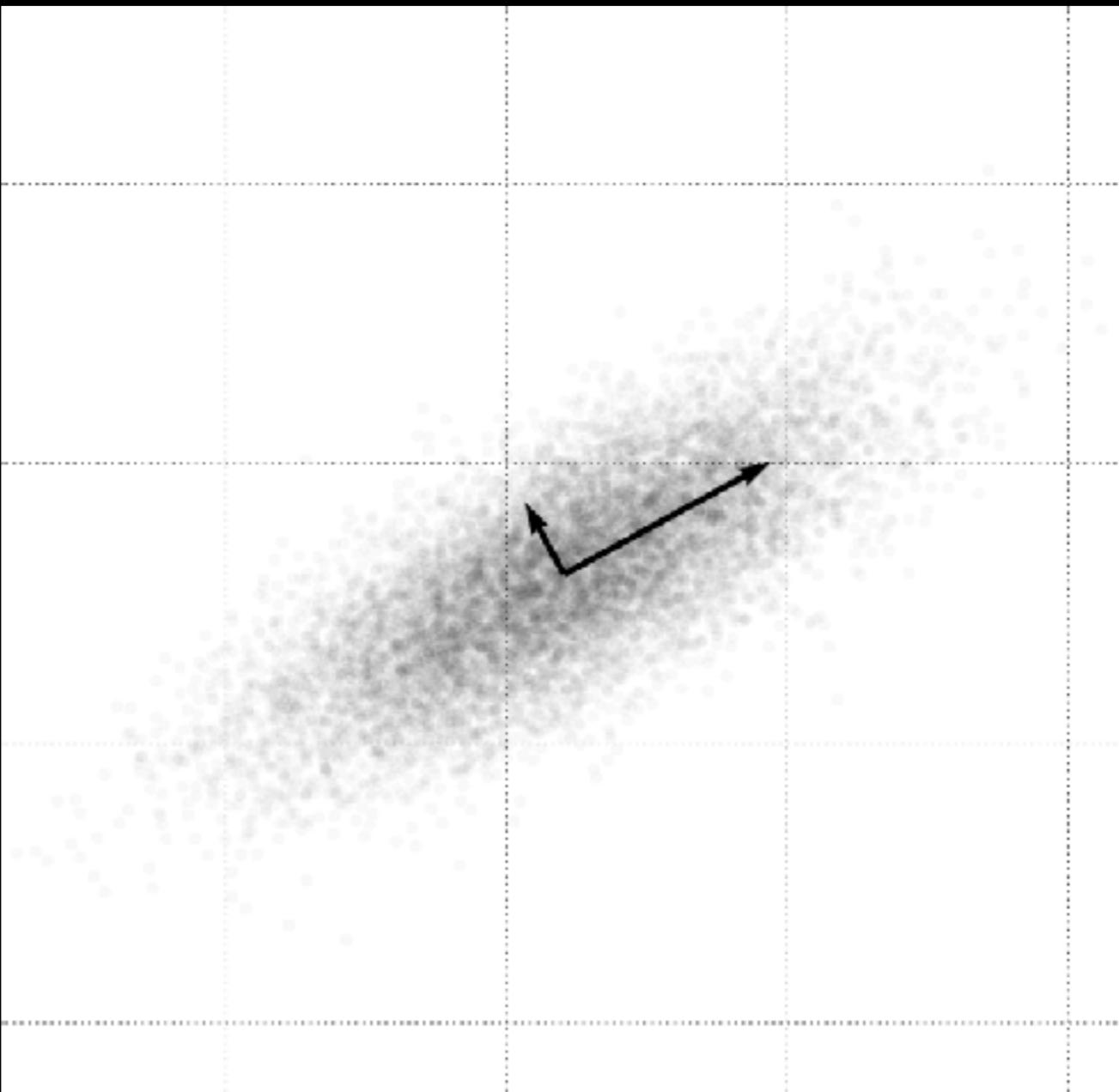


# Goals for Today

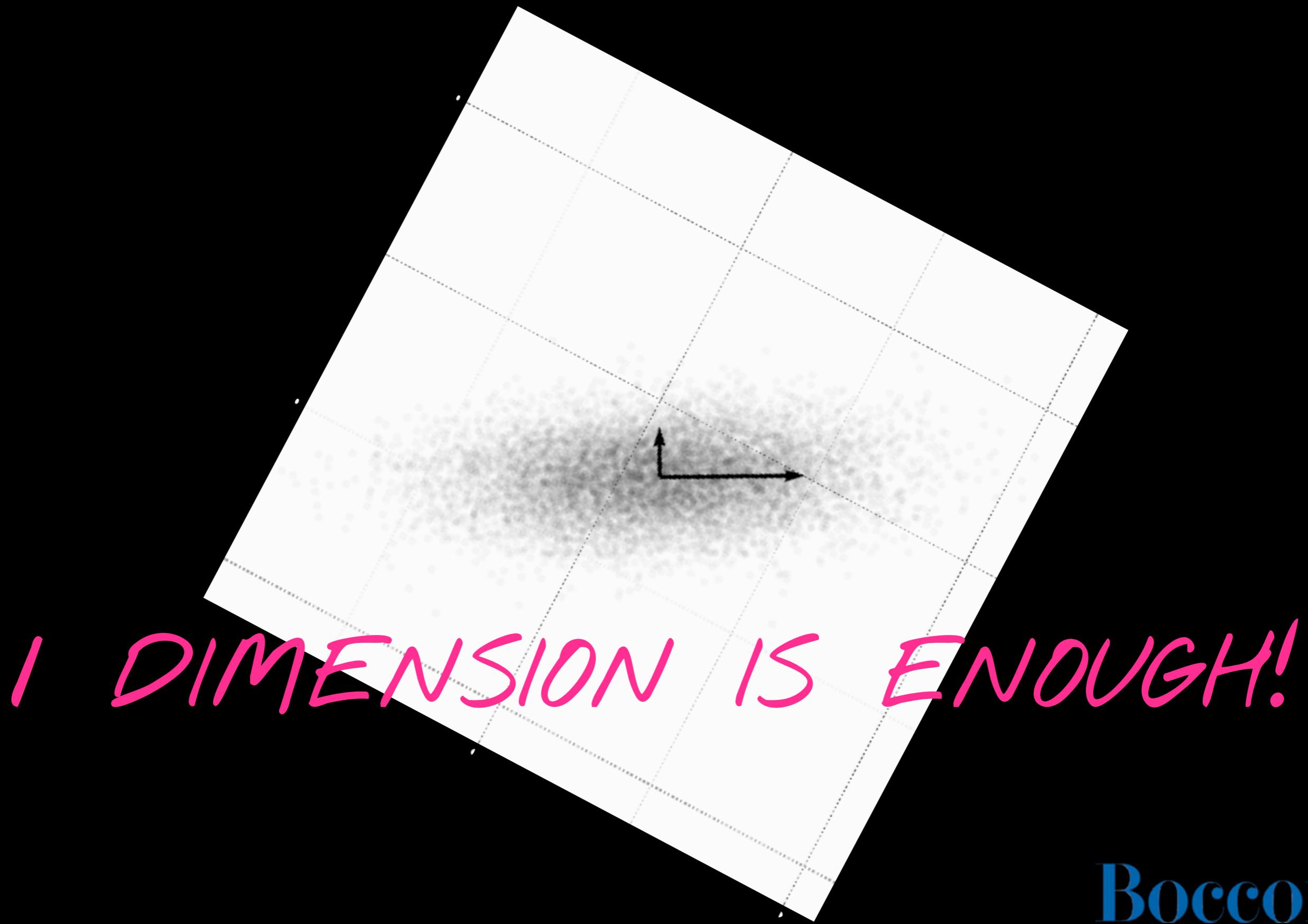
- Learn about **matrix factorization** and its use for **semantic similarity** and **visualization**
- Learn about **k-means** and **agglomerative clustering**
- Learn about **evaluation** criteria

# Matrix Factorization

# Singular Value Decomposition



# Singular Value Decomposition



*1 DIMENSION IS ENOUGH!*

# Singular Value Decomposition

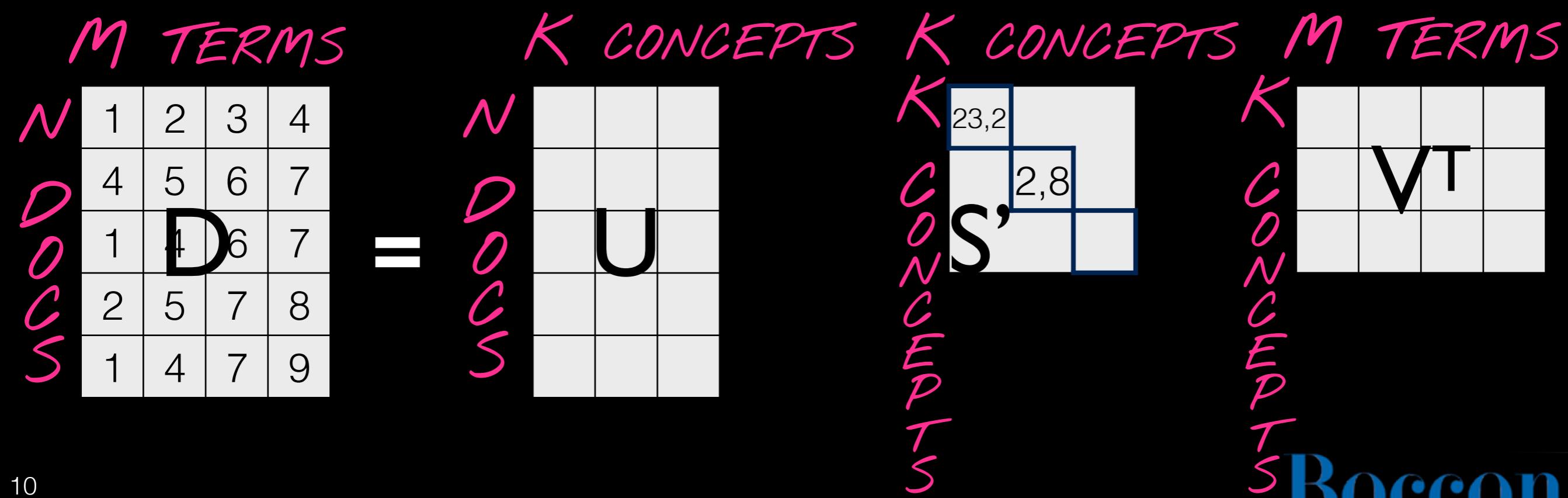
- “principal component analysis”: discover the dimensions that matter
- idea: matrix is made up of few hidden dimensions
- Dimensions correspond to **documents**, **terms**, and **latent concepts**

$$\begin{matrix} M \text{ TERMS} \\ N \\ D \\ O \\ C \\ S \end{matrix} \quad \begin{matrix} K \text{ CONCEPTS} \\ N \\ D \\ O \\ C \\ S \end{matrix} \quad = \quad \begin{matrix} K \text{ CONCEPTS} \\ K \\ C \\ O \\ N \\ C \\ E \\ P \\ T \\ S \end{matrix} \quad \begin{matrix} M \text{ TERMS} \\ K \\ C \\ O \\ N \\ C \\ E \\ P \\ T \\ S \end{matrix}$$

The diagram illustrates the Singular Value Decomposition (SVD) of a document-term matrix  $D$ . The matrix  $D$  has dimensions  $M \times N$  (5 terms by 4 documents). It is decomposed into three matrices:  $U$  (orthogonal matrix of size  $N \times N$ ),  $S$  (diagonal matrix of size  $K \times K$  containing singular values), and  $V^T$  (orthogonal matrix of size  $M \times M$  containing latent concepts). The matrix  $S$  is shown with its top-left element highlighted as 23,2.

# Singular Value Decomposition

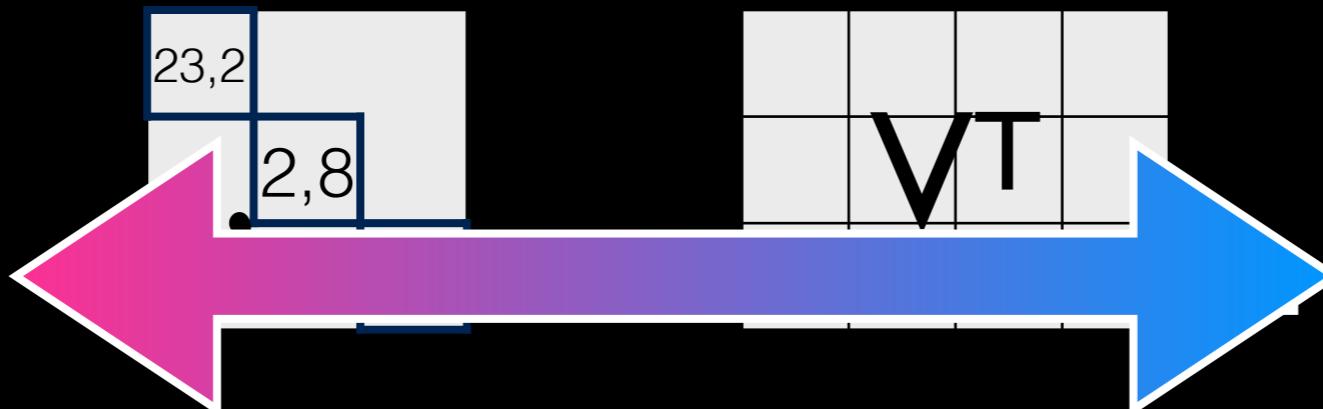
- reduce principal components/concepts to smaller number



# Singular Value Decomposition

- reconstruct original matrix in new concept space:  
**Latent Semantic Analysis**

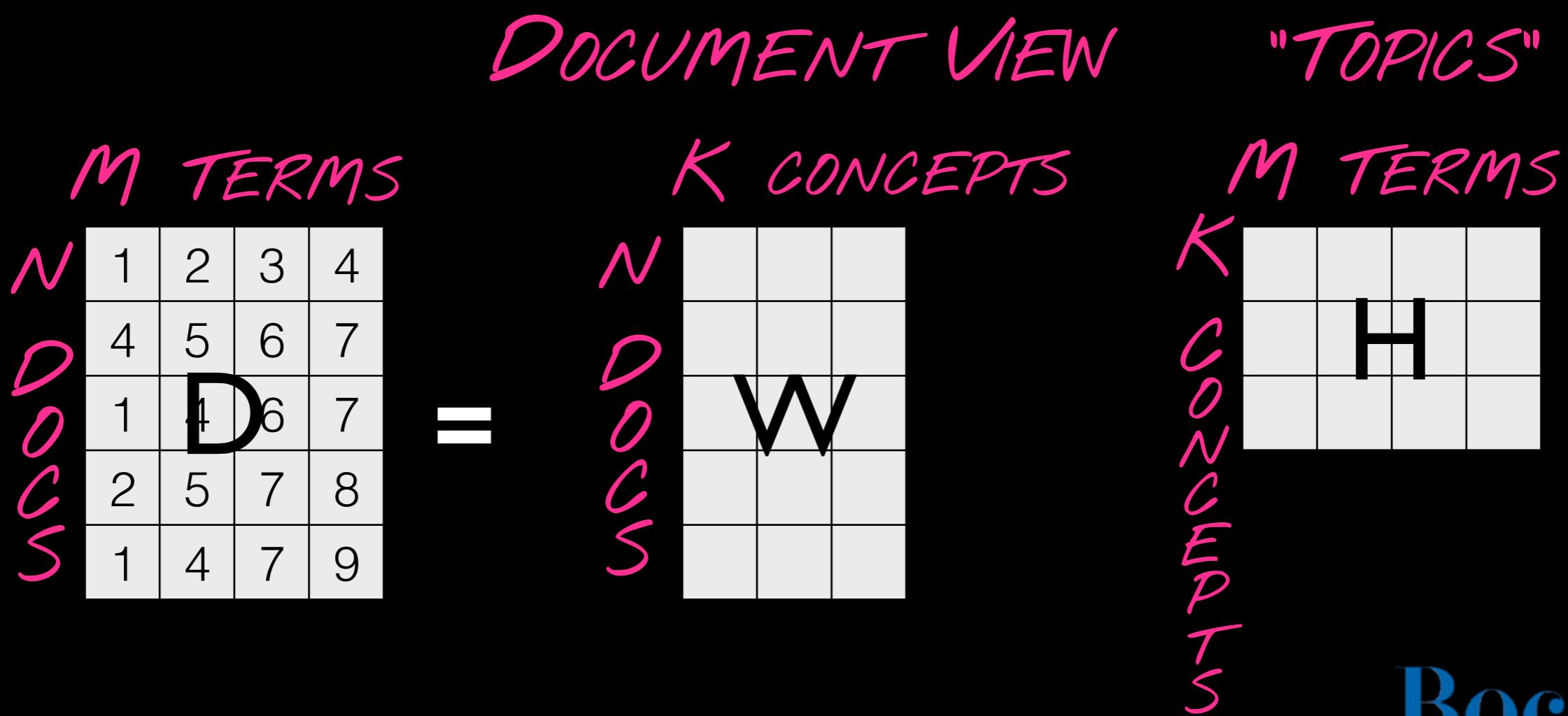
1	2	3	4
4	5	6	7
1	4	6	7
2	5	7	8
1	4	7	9



0,9	2,1	3,2	3,8
3,9	5,1	6	6,9
1,1	3,8	4,9	7,2
2,2	4,7	6,9	8,2
0,8	4,3	7,1	8,8

# Non-negative Matrix Factorization

- Use only positive values
- Find approximation of two components



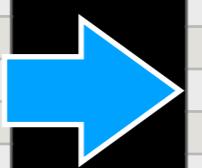
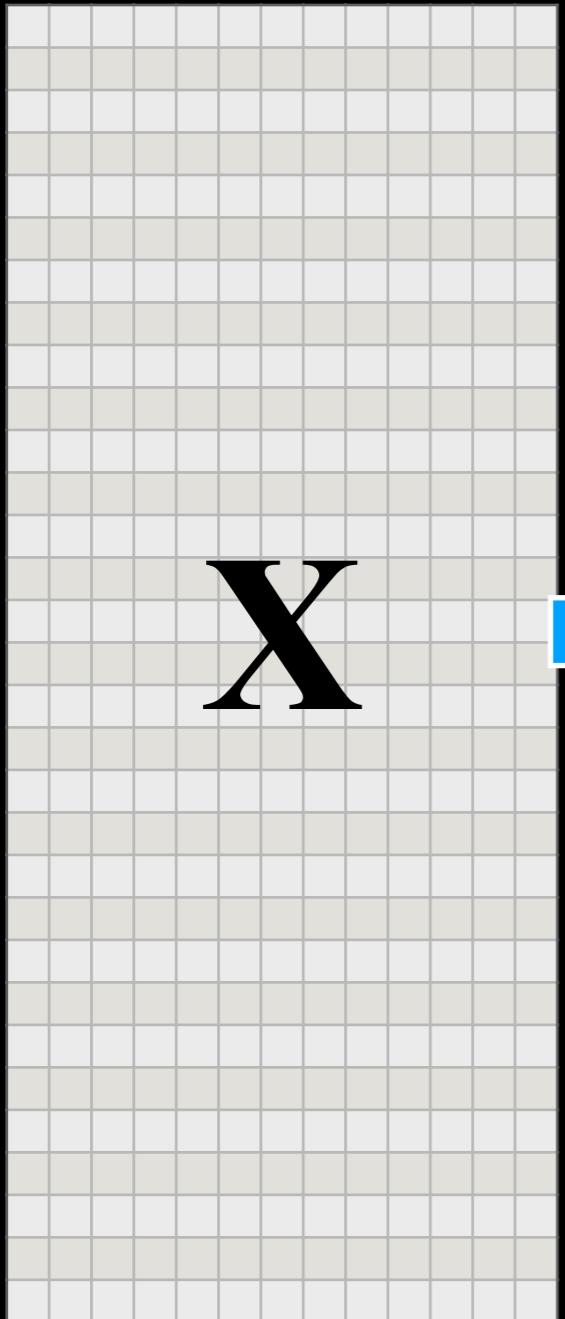
# Comparison

	SVD	NMF
<b>Negative values (embeddings) as input?</b>	yes	no
<b>#components</b>	$3: U, S, V$	$2: W, H$
<b>document view?</b>	yes: $U$	yes: $W$
<b>term view?</b>	yes: $V$	yes: $H$
<b>strength ranking?</b>	yes: $S$	no
<b>exact?</b>	yes	no
<b>"topic" quality</b>	mixed	better
<b>sparsity</b>	low	medium

# Yes, but: What is it Good for?

- Find latent **topic** dimensions (alternative: LDA)
- Find **word similarity** in latent space (alternative: Word2Vec)
- Find **document similarity** in latent space (alternative: Doc2Vec)
- Reduce dimensionality for **visualization**

# Latent Word Dimension Topics



EVT

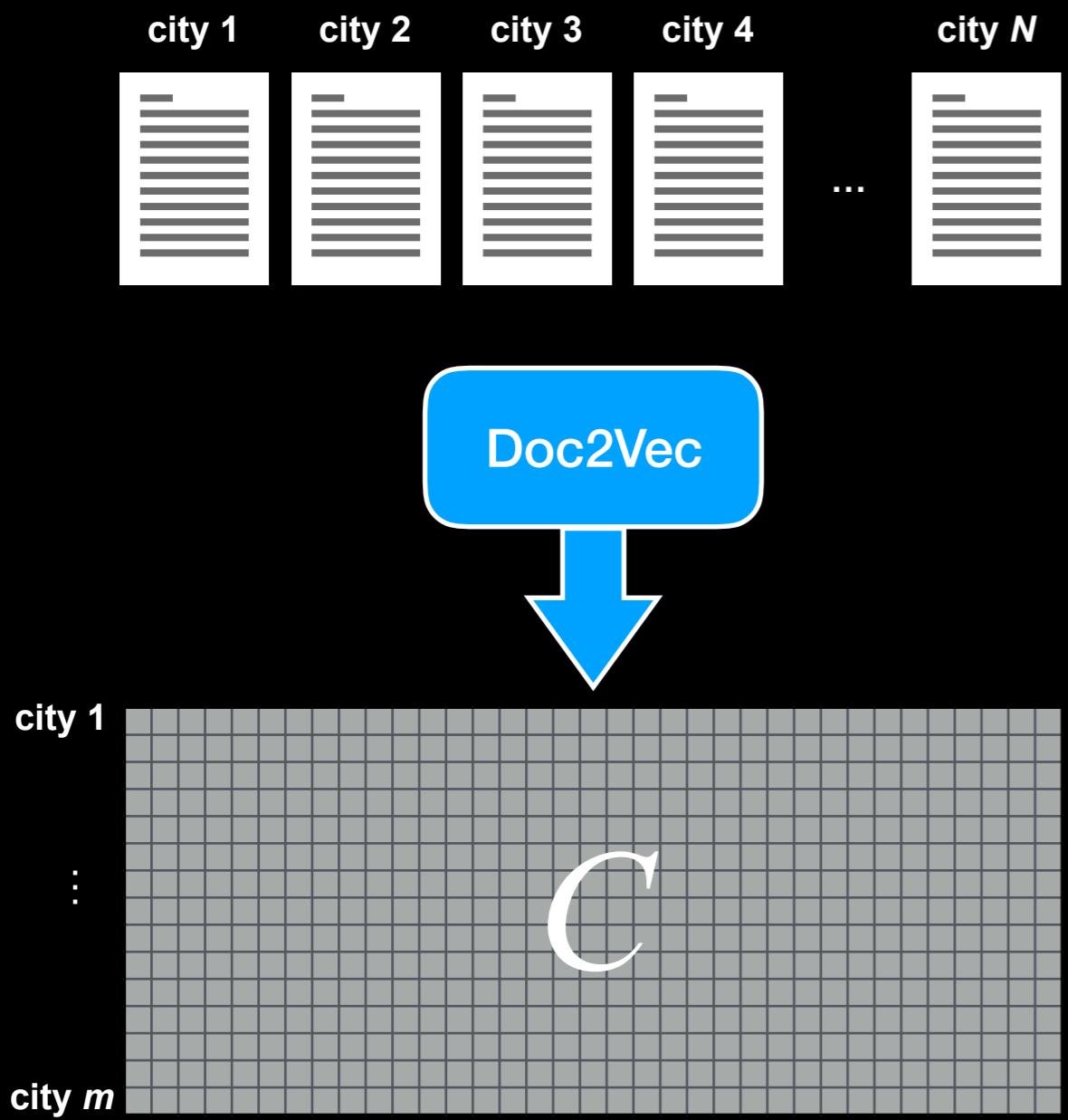
NMF(10)

FOR MOBY DICK

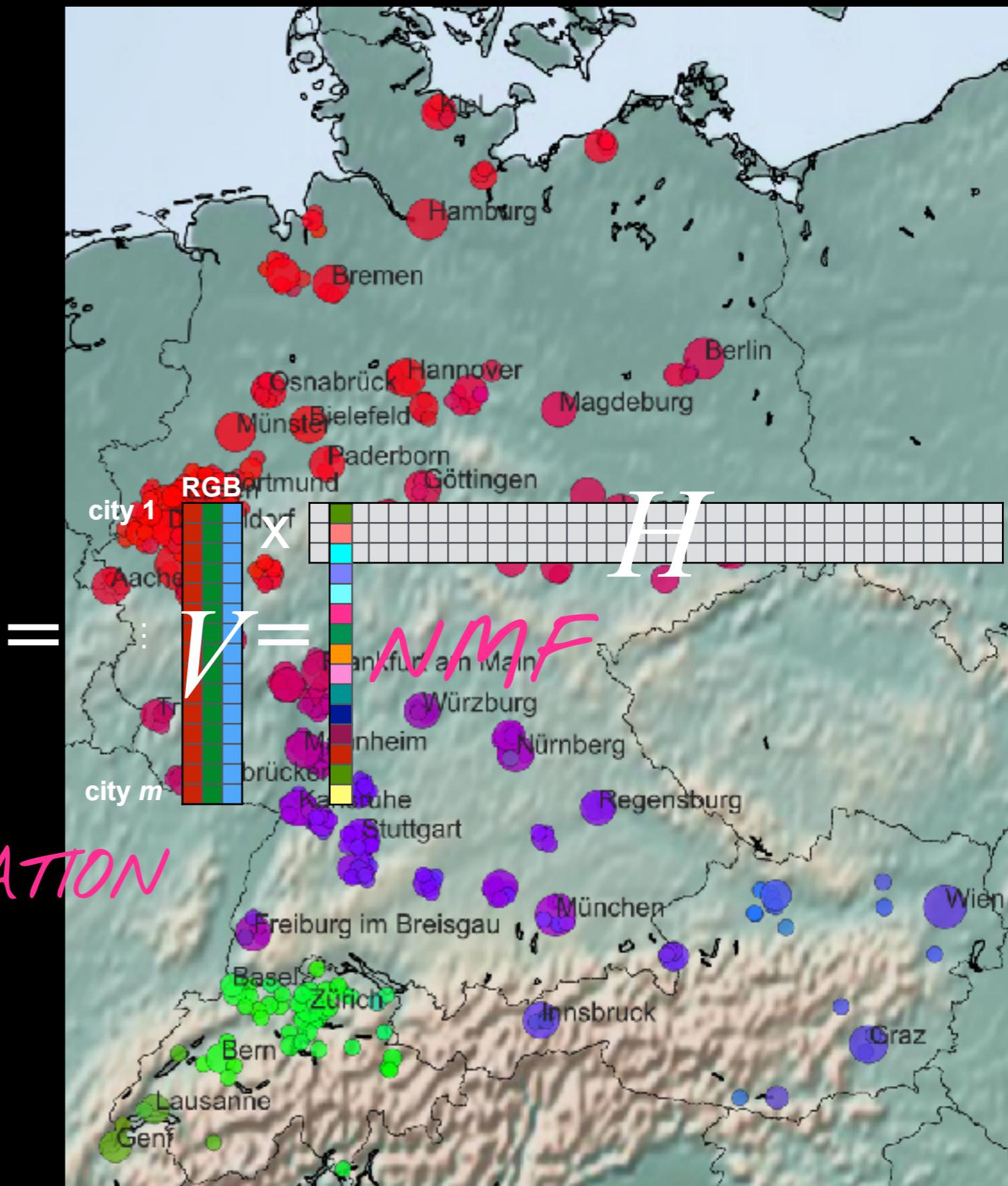
AHAB	ahab, captain, cried, captain ahab, cried ahab
STRUCTURE	chapter, folio, octavo, ii, iii
???	like, ship, sea, time, way
MEN	man, old, old man, look, young man
???	oh, life, starbuck, sweet, god
CHARACTERS	said, stubb, queequeg, don, starbuck
???	sir, aye, let, shall, think
OLD-TIMEY	thou, thee, thy, st, god
WHALES	whale, sperm, sperm whale, white, white whale
???	ye, look, say, ye ye, men

# Dimensionality Reduction for Visualizations

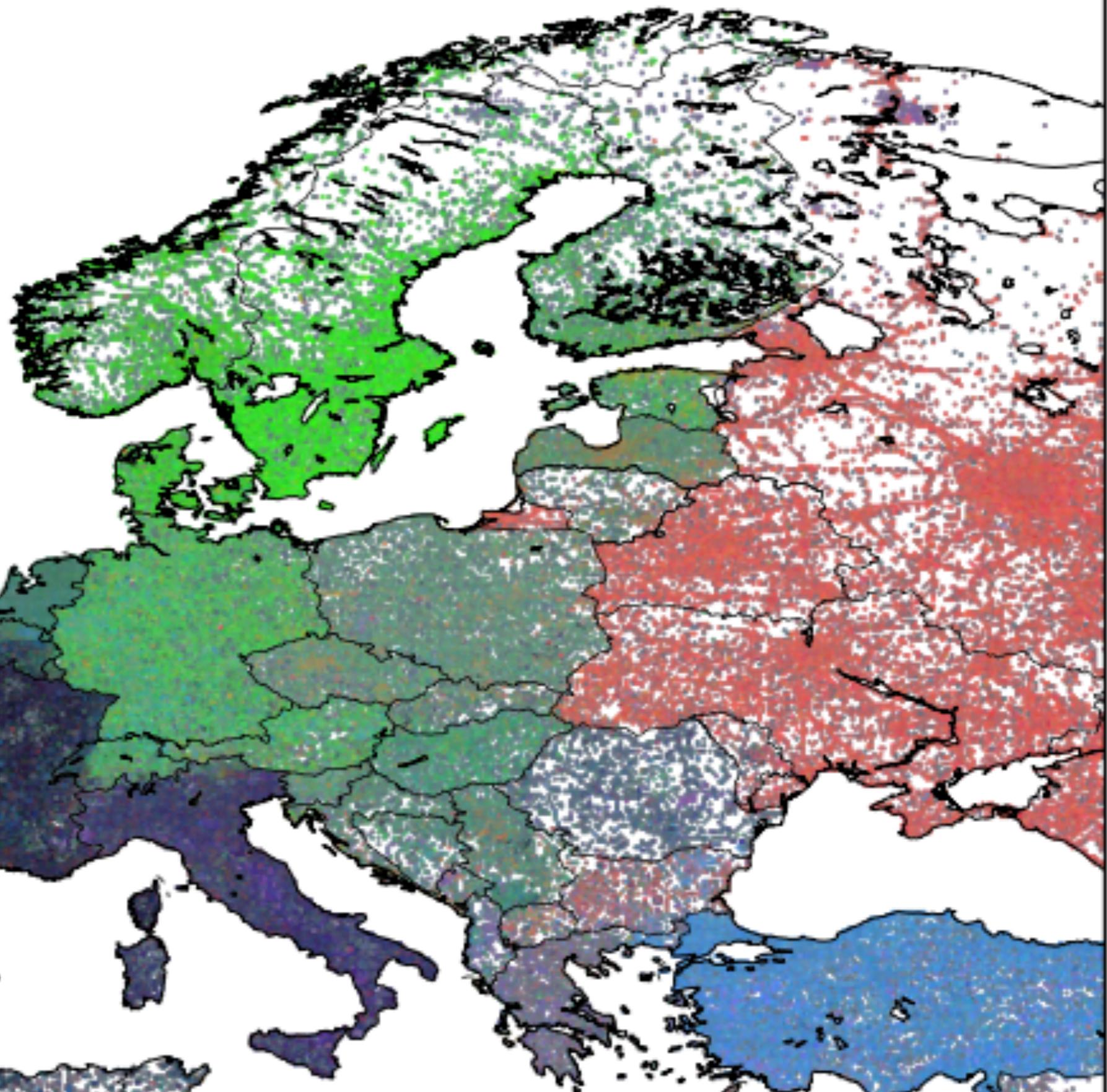
# Dimensions as RGB



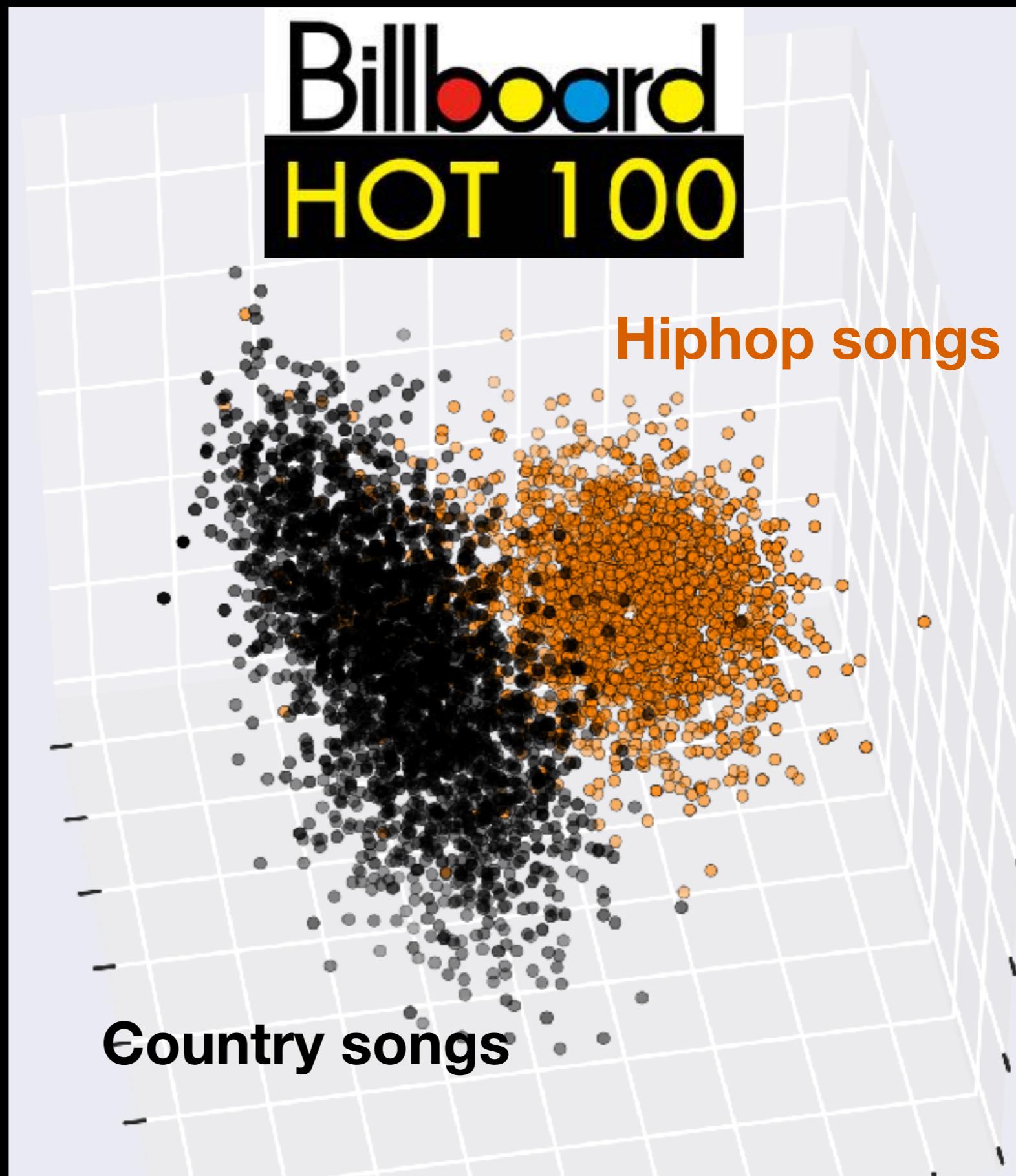
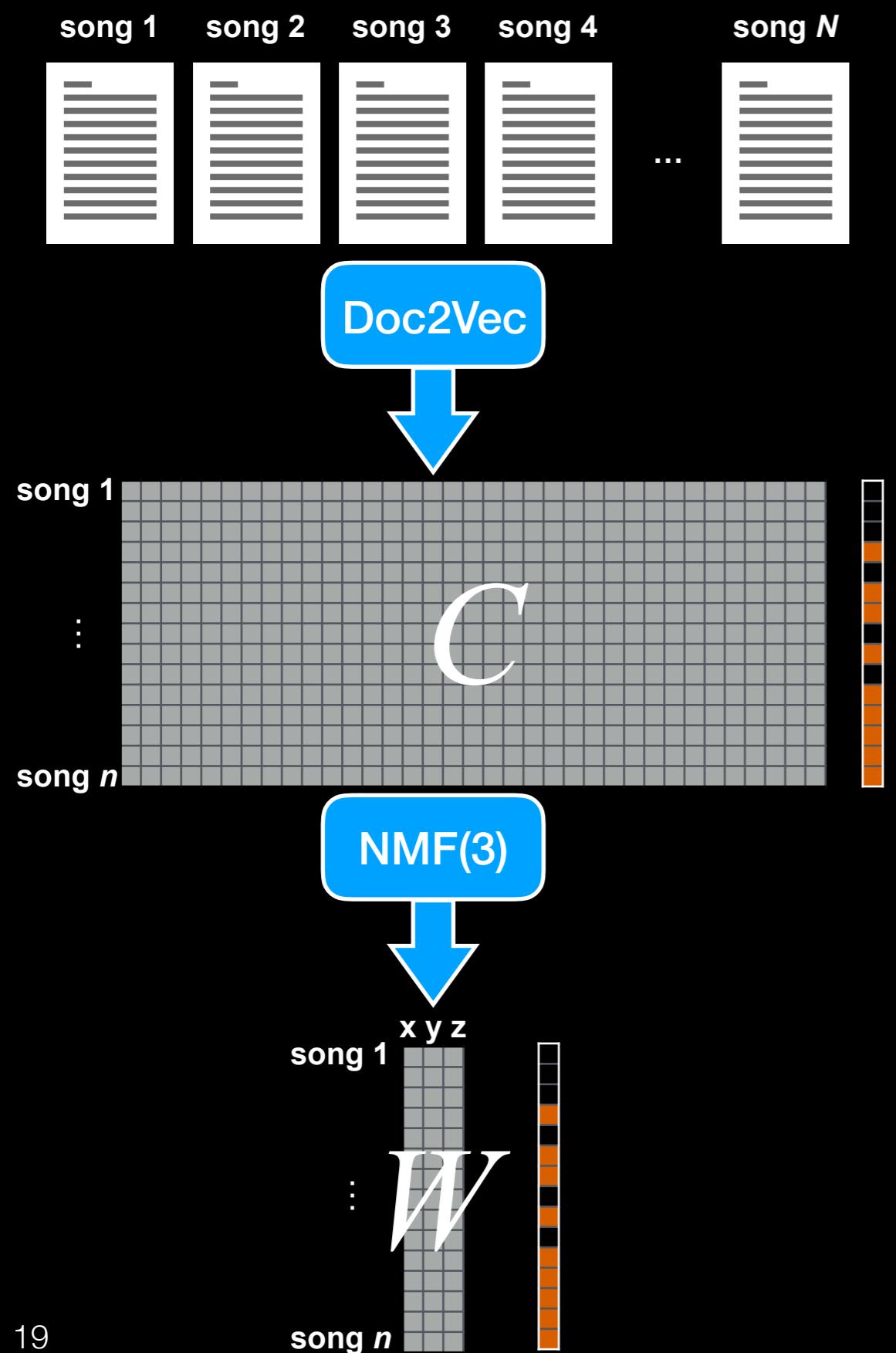
*DENSE REPRESENTATION*



# Dimensions as RGB

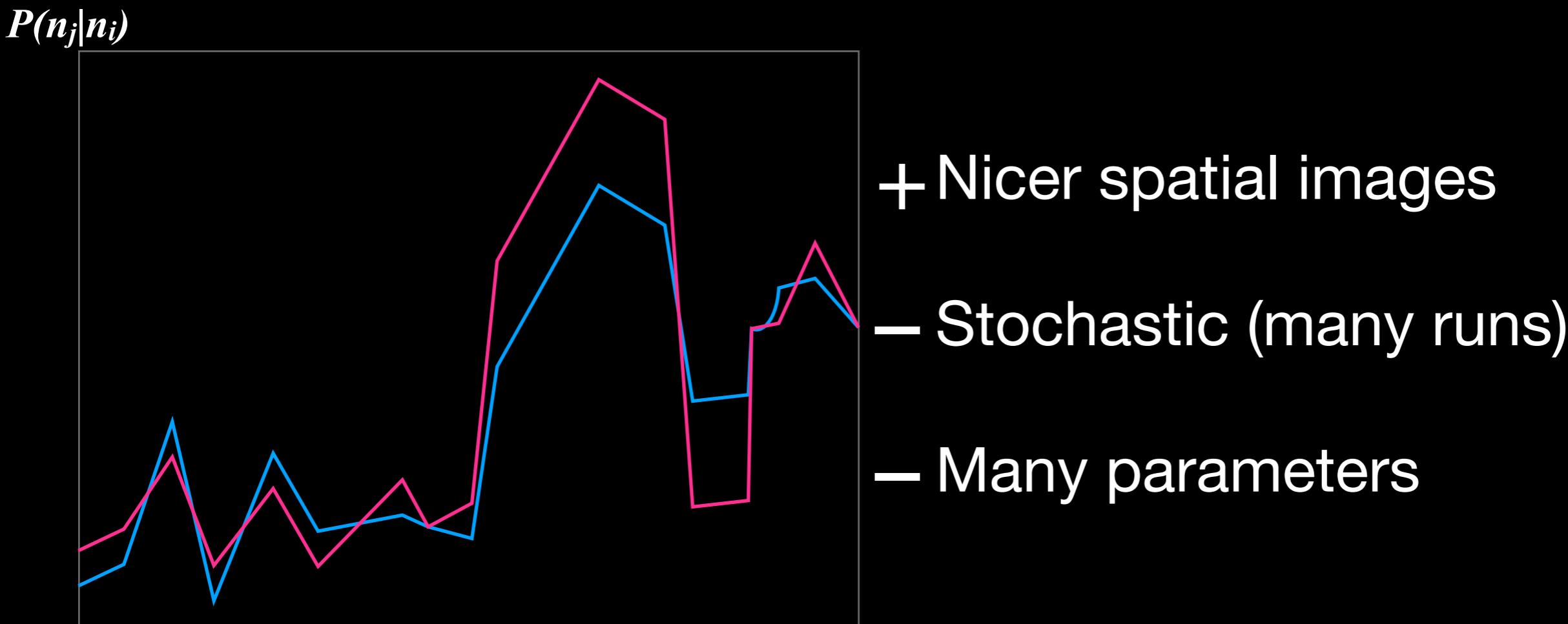


# Dimensions as Coordinates



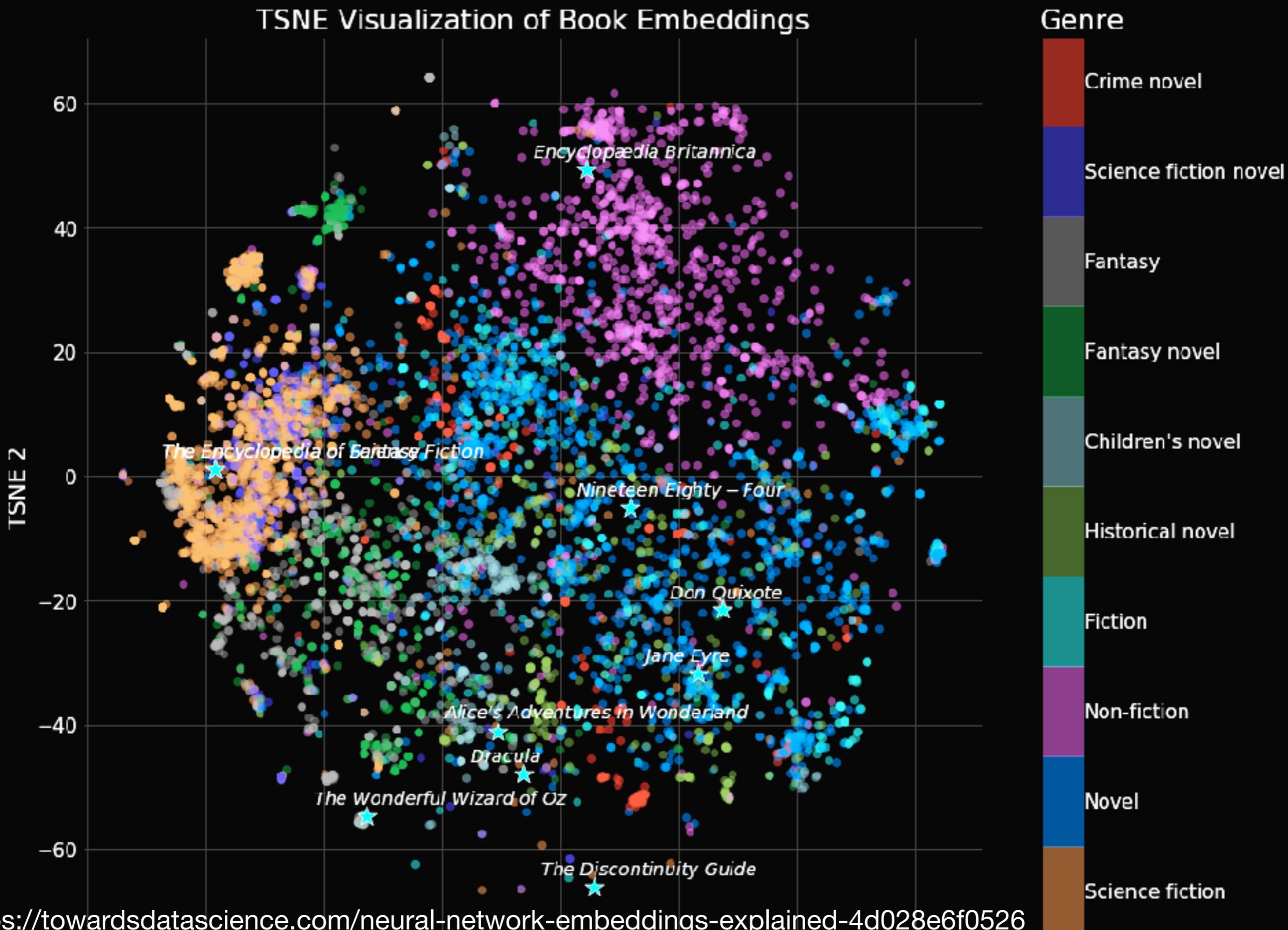
# t-SNE

- Map/preserve neighborhood structure of high-dimensional space in lower-dimensional space
- Minimize difference between probability distributions over neighbors in both dimensions



# t-SNE

TSNE Visualization of Book Embeddings



# Data Visualization

# How Vizs Help

- Goal: get a sense of the data
- understand the type of distribution
- spot outliers
- see (and then test) possible correlations
- ***always*** visualize your data!

# Communicate Insights & Patterns

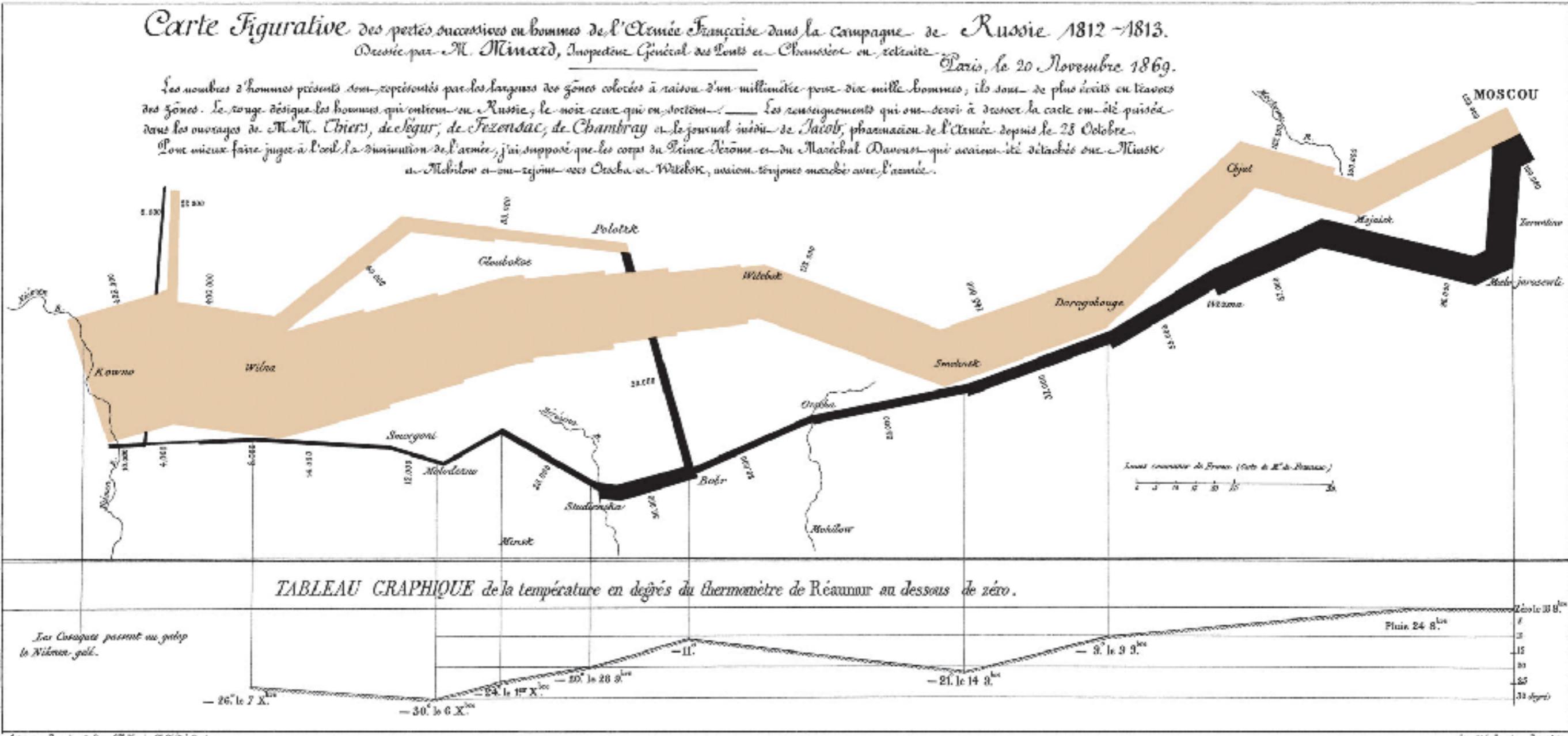
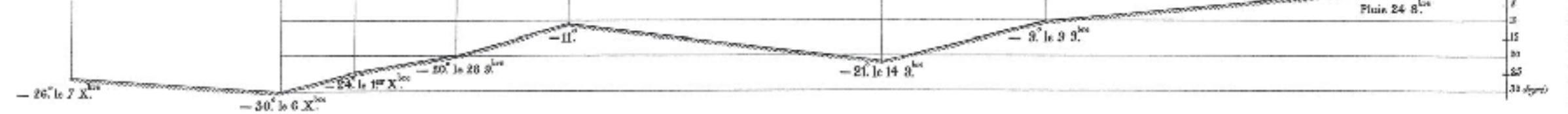


TABLEAU GRAPHIQUE de la température en degrés du thermomètre de Réaumur au dessous de zéro.

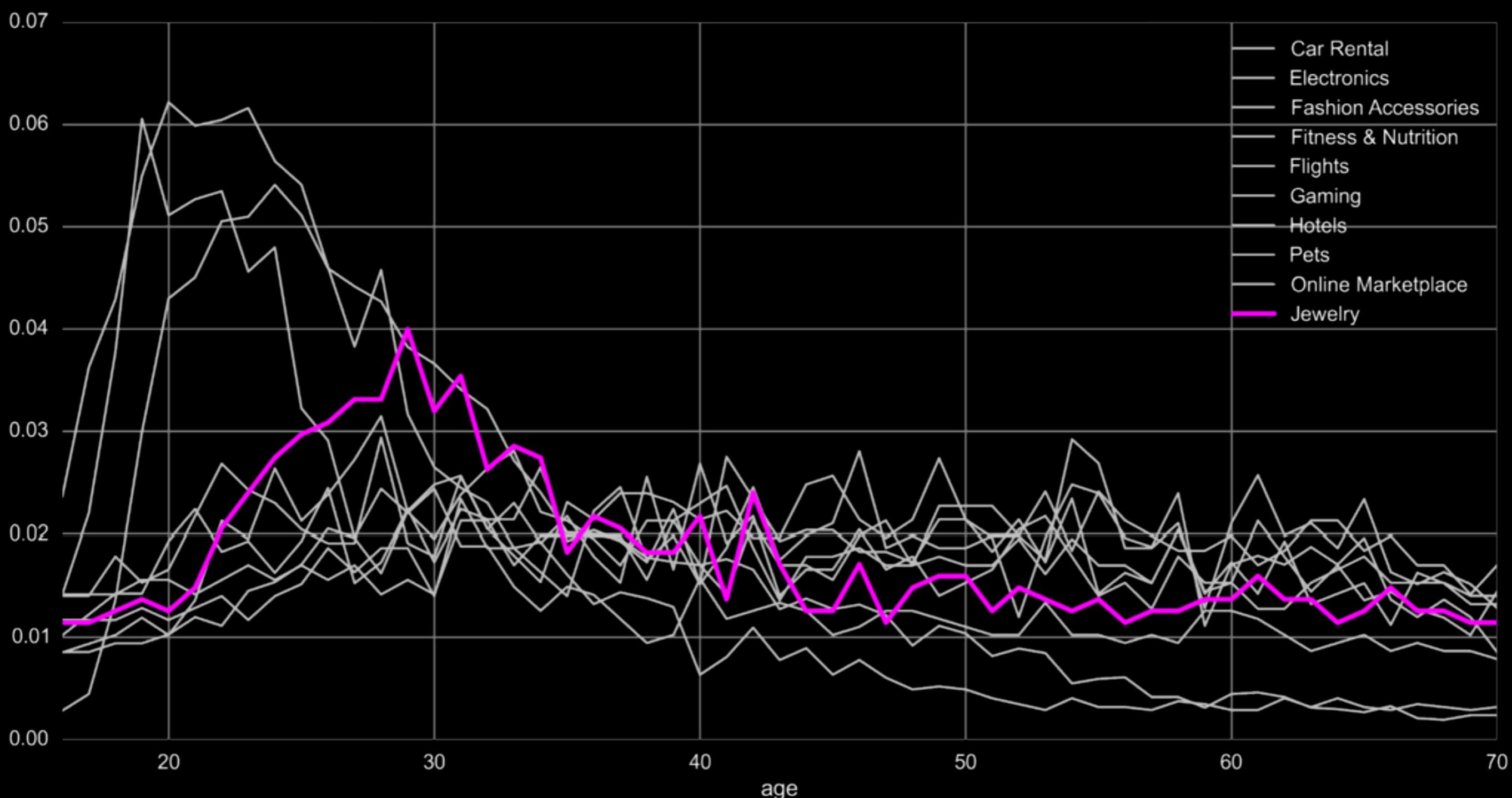


Charles Joseph Minard: Napoleon losses 1812

Bocconi

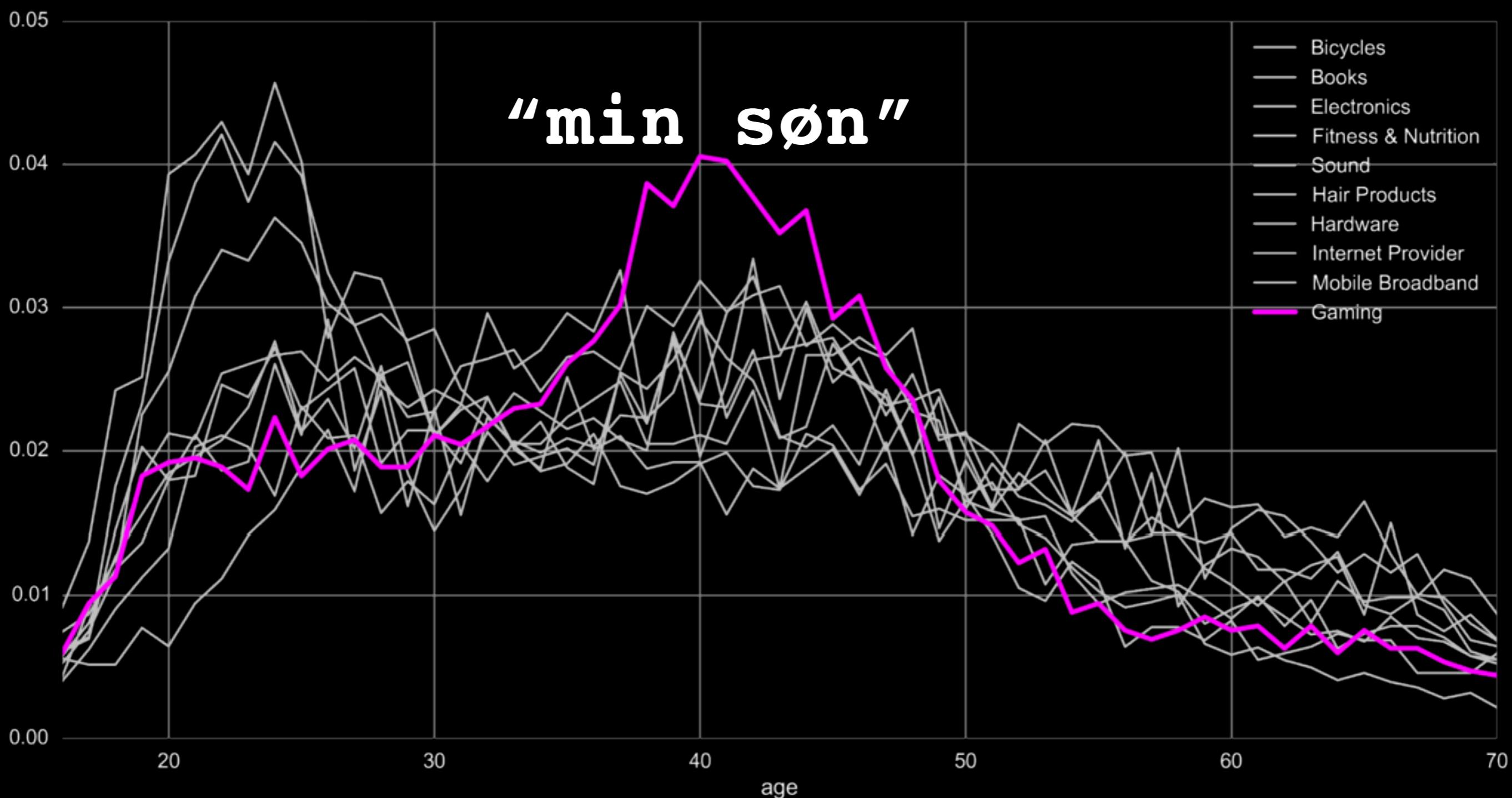
# Workhorses: Line Charts

Top 10 Review Categories for American Men

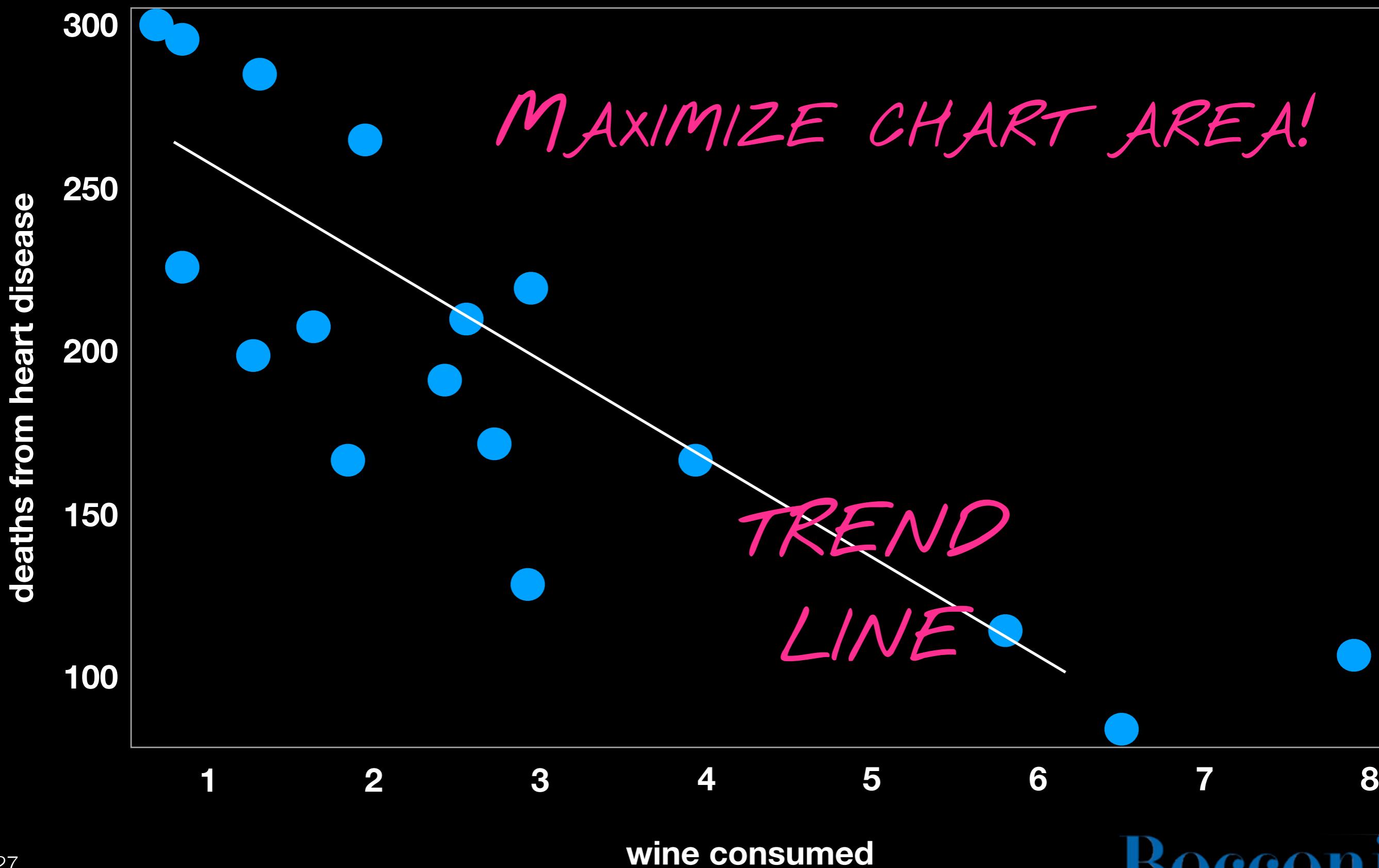


# Caveat: False Leads!

Top 10 Review Categories for Danish Women



# Workhorses: Scatter Plot



# Workhorses: Bar Charts

Cherry Pie

Apple Pie

Meat Pie

Shepherds Pie

*LEGEND!*

60

55

45

30

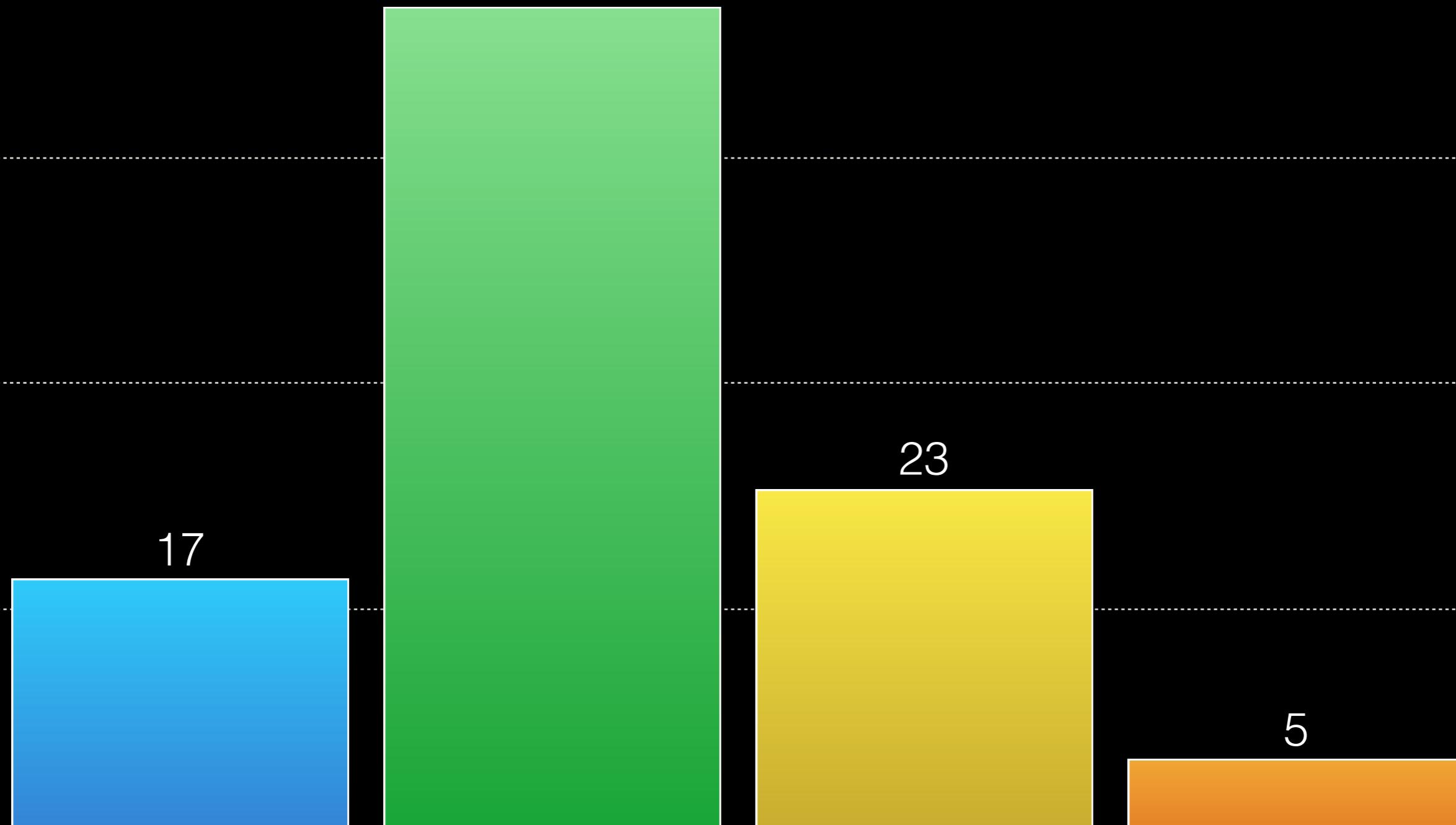
15

17

23

5

Pie Types



Bocconi

# Workhorses: Bar Charts

■ Filling

■ Dough

60

45

30

15

0

17

38

10

13

5

12

2

3

EACH ELEMENT HAS  
COMPONENTS

Cherry Pie

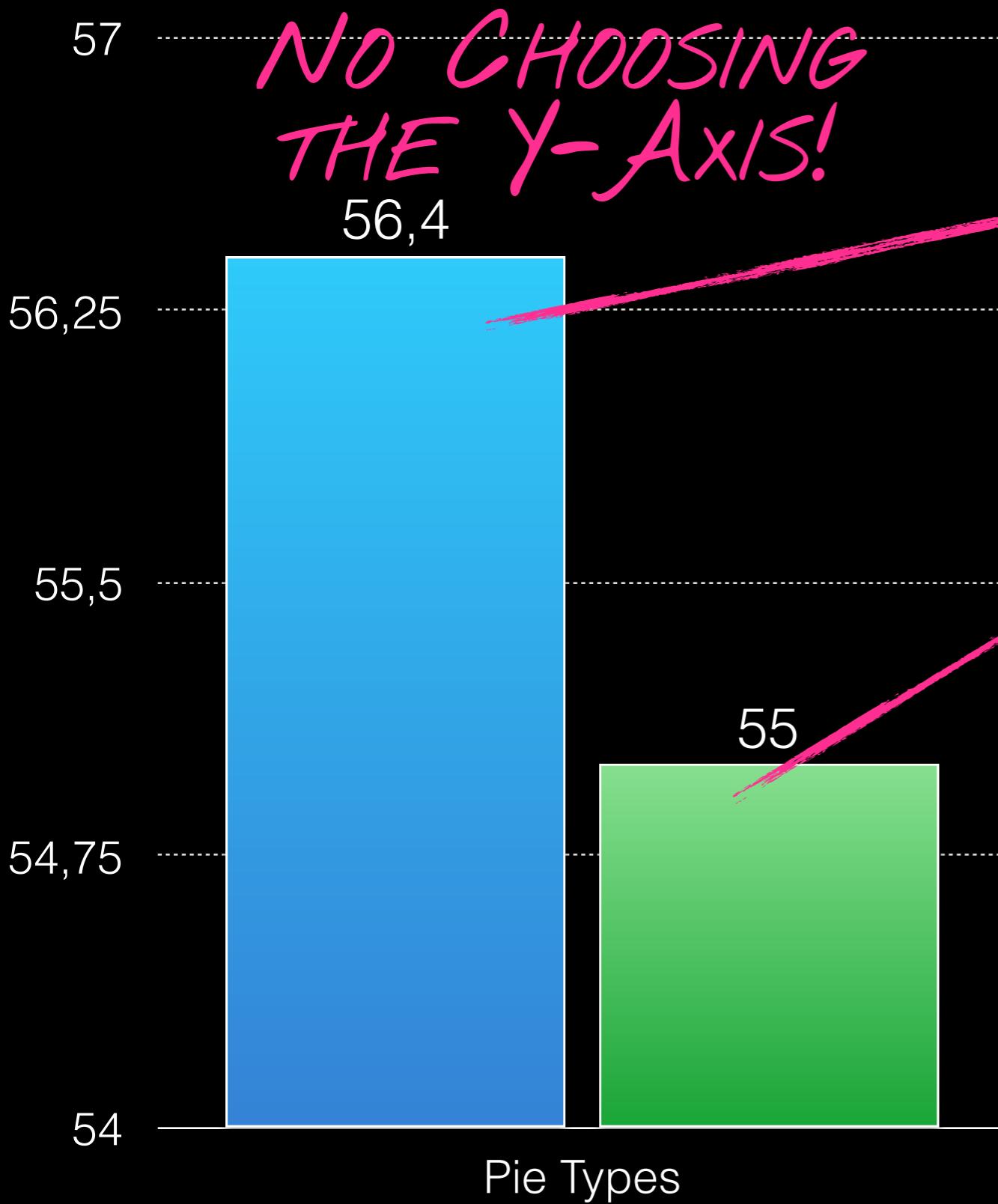
Apple Pie

Meat Pie

Bocconi

# Chart No-nos!

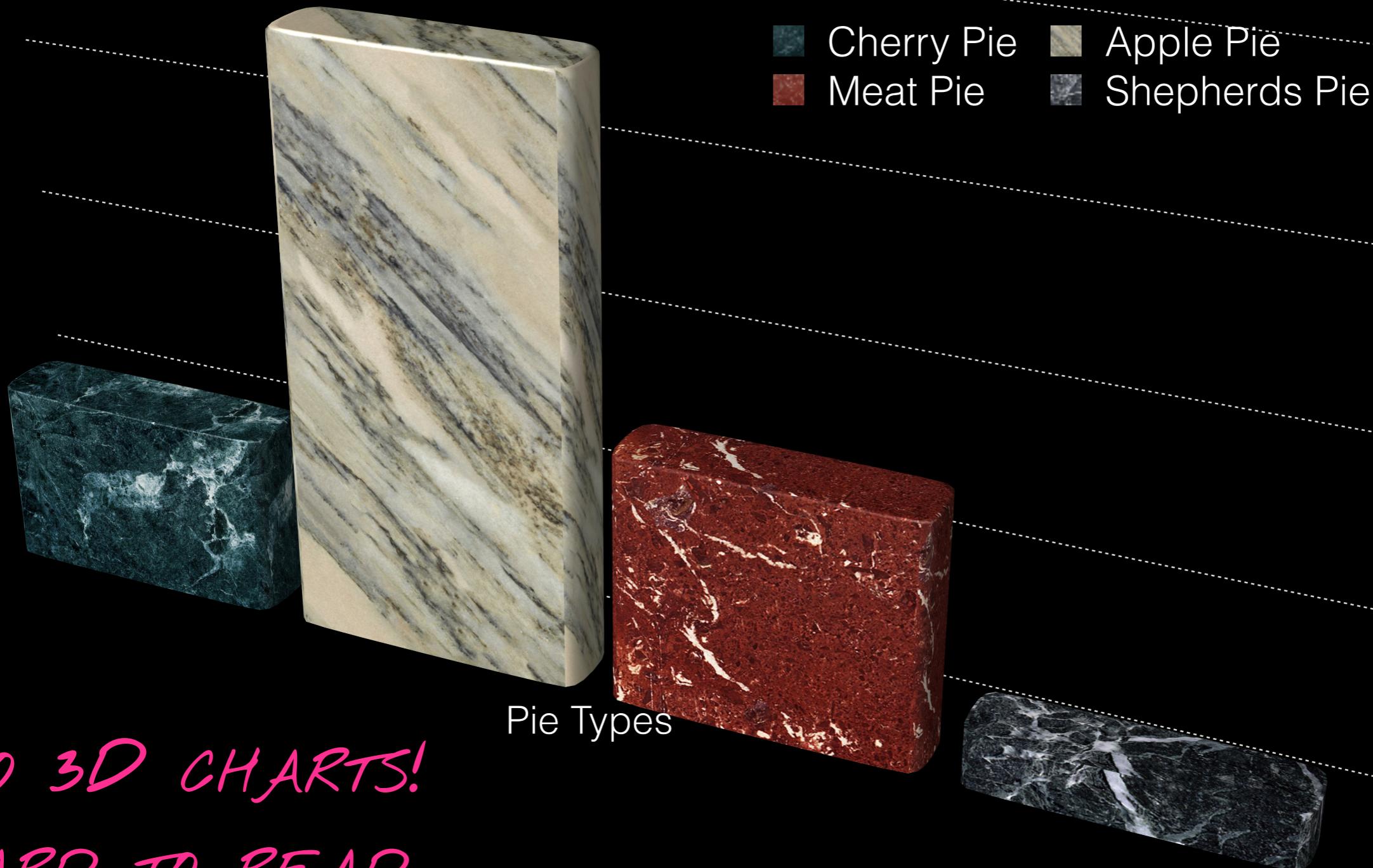
Candidate 1 Candidate 2



Candidate 1 Candidate 2



# Chart No-nos!



NO 3D CHARTS!

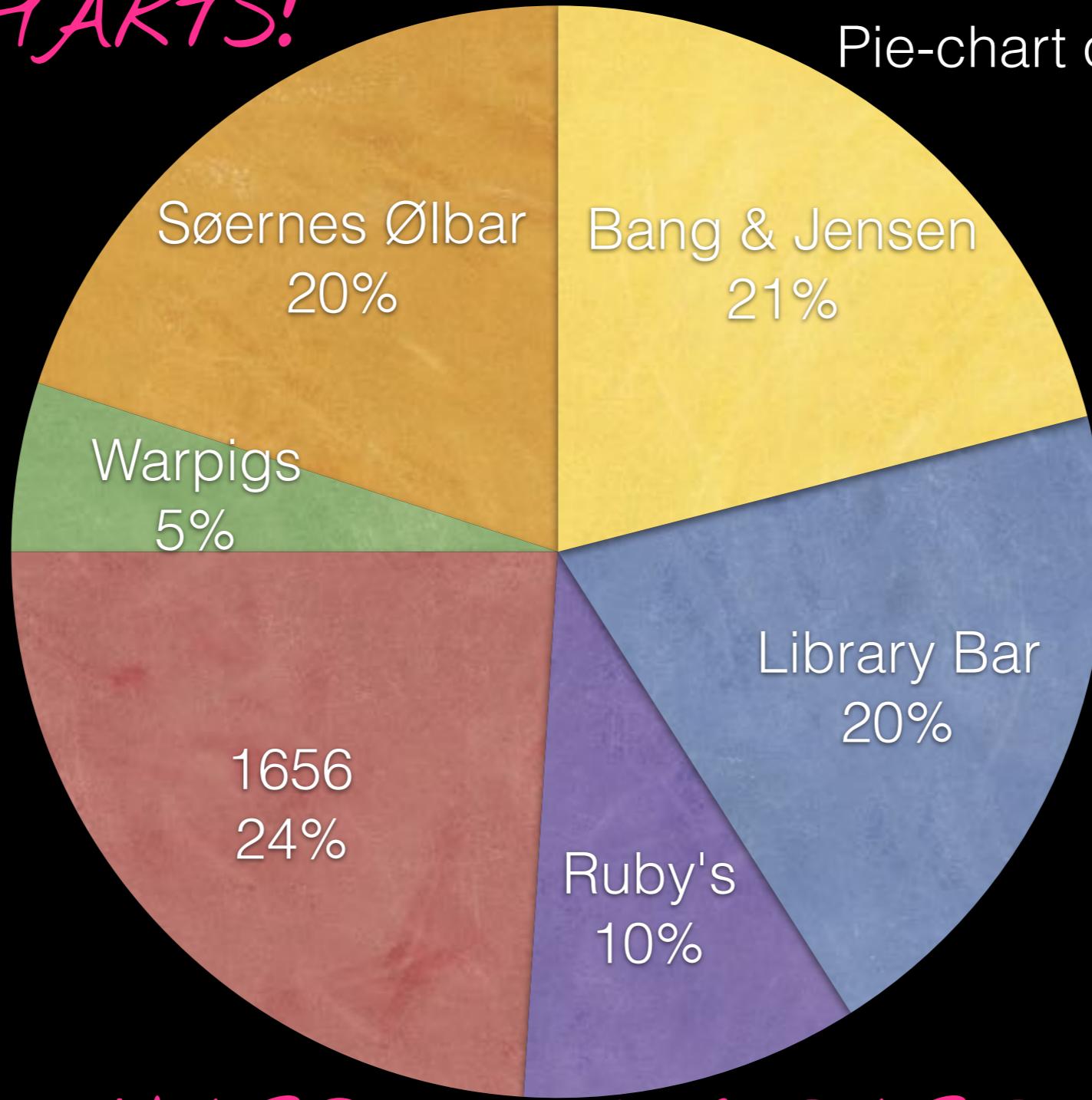
HARD TO READ

UNNECESSARY INFORMATION (DEPTH, TEXTURE)

# More No-nos!

NO PIE CHARTS!

Pie-chart of my favorite bars



HARD TO COMPARE

AREA IS IRREGULAR SHAPE

# The ONLY Exception!



- Sky
- Sunny side of pyramid
- Shady side of pyramid

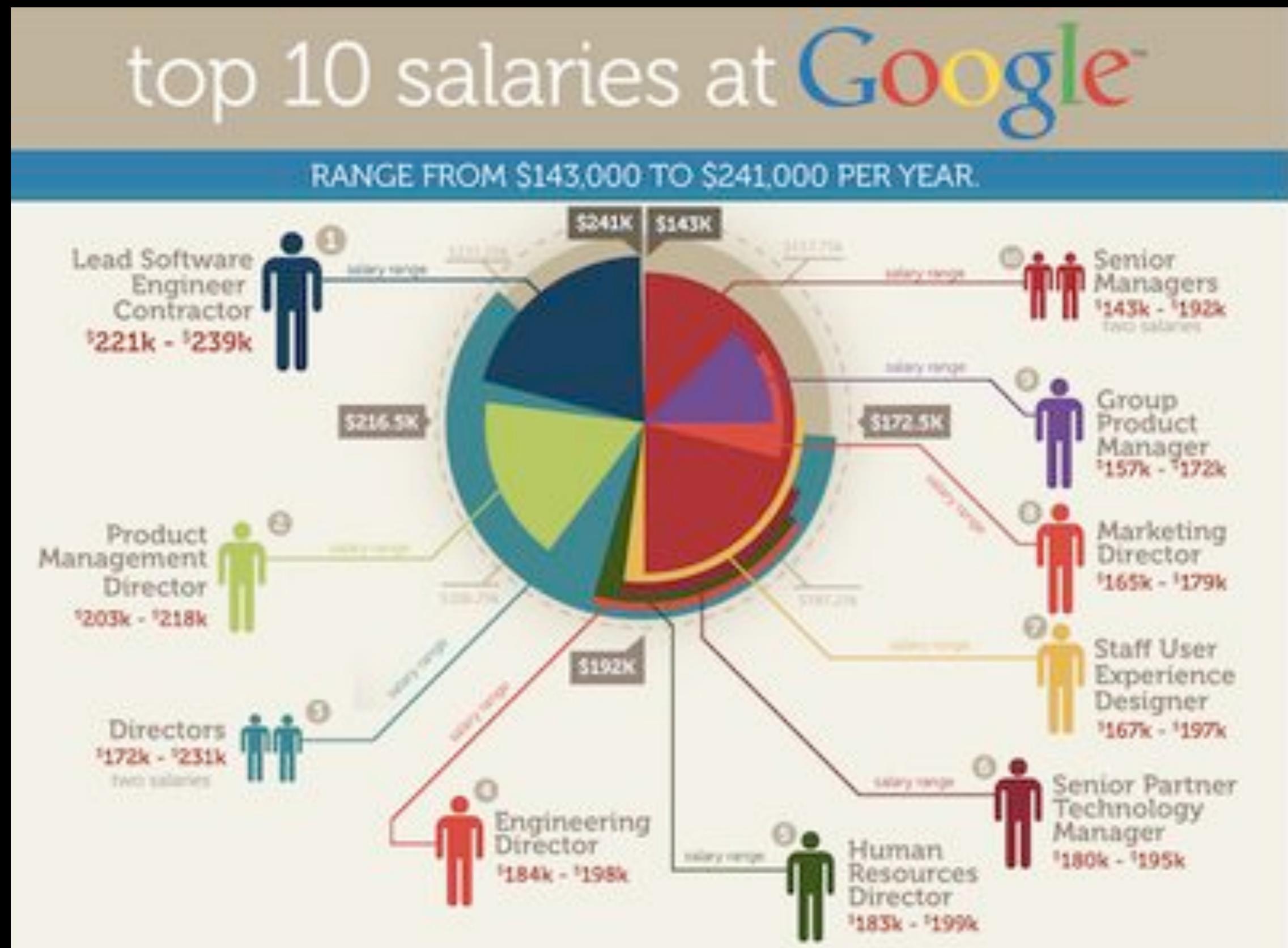
**And... Just don't...**



NO DISCERNIBLE INFORMATION VALUE  
LONG WORDS GET MORE SPACE  
UNSCIENTIFIC

# Improvements

[http://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](http://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html)



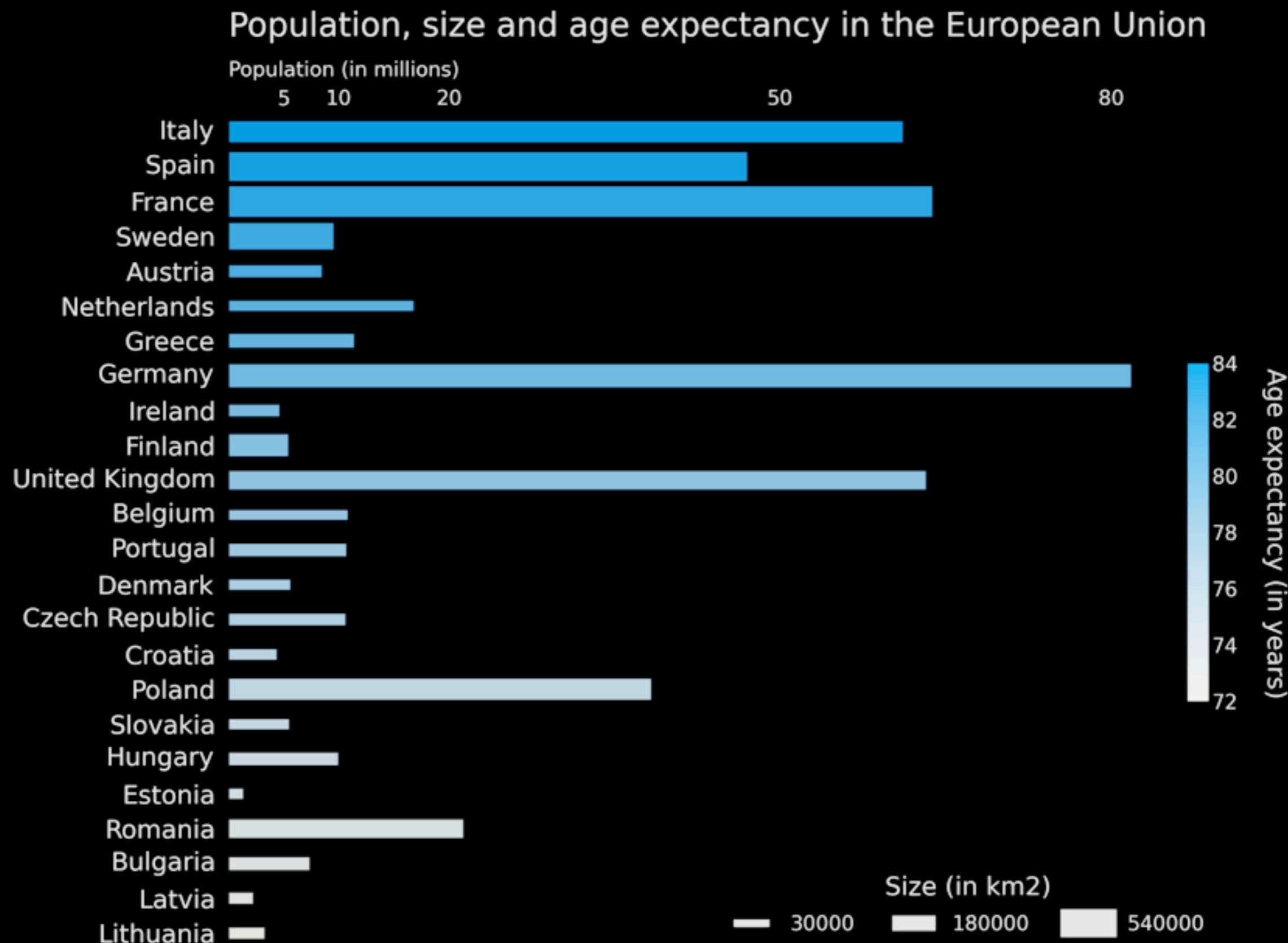
# Improvements

[http://junkcharts.typepad.com/junk\\_charts/2011/10/the-massive-burden-of-pie-charts.html](http://junkcharts.typepad.com/junk_charts/2011/10/the-massive-burden-of-pie-charts.html)



# Improvements

<https://datasciencecelab.wordpress.com/2013/12/21/beautiful-plots-with-pandas-and-matplotlib>



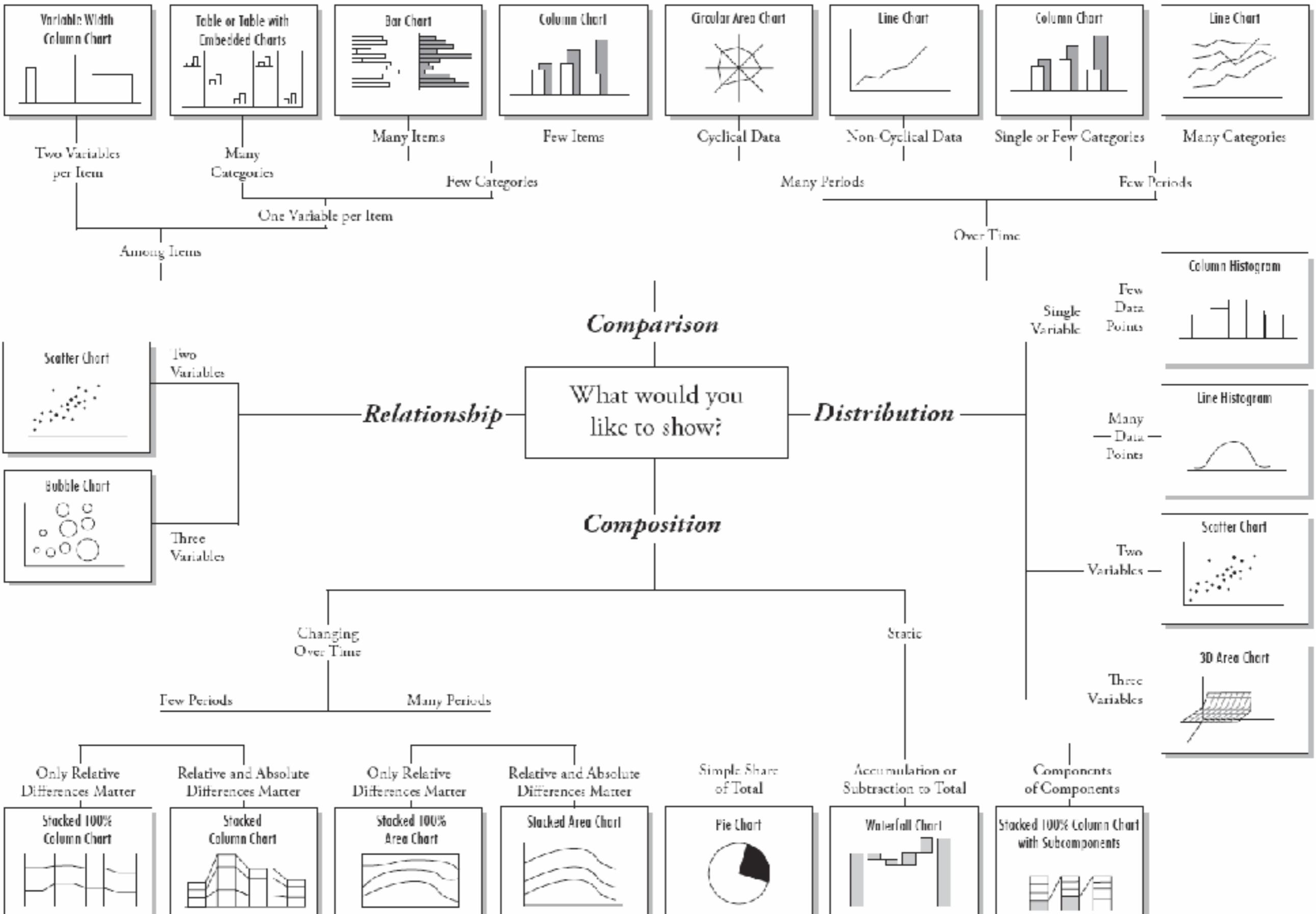
oni

# Dark Horses: Heatmaps

New COVID19 cases per 100k by week and age group

	Neu gemeldete COVID19-Fälle pro 100.000 Personen in Deutschland nach Meldeweche und Altergruppe																
	W20	W21	W22	W23	W24	W25	W26	W27	W28	W29	W30	W31	W32	W33	W34	W35	W36
0-4	3.0	3.1	2.2	2.7	2.8	4.0	3.4	3.5	3.7	3.5	4.0	4.0	4.7	6.9	7.1	5.6	5.2
5-9	3.0	2.3	3.0	2.7	3.3	4.7	3.7	3.3	3.1	3.7	4.1	5.3	7.6	10.9	10.1	9.3	7.4
10-14	3.2	3.0	3.2	2.8	3.2	5.2	4.0	2.8	2.7	4.1	5.2	6.8	9.8	14.7	15.5	12.9	9.9
15-19	4.7	4.0	4.2	4.2	2.9	6.3	4.3	3.6	3.9	5.0	6.5	9.5	14.2	19.2	22.8	21.5	17.9
20-24	9.0	6.8	6.4	5.0	5.4	8.1	5.5	5.6	5.1	6.5	8.0	11.3	15.9	25.3	33.3	28.6	23.9
25-29	7.9	7.1	6.4	4.5	4.9	8.0	6.5	6.4	4.6	6.2	8.3	10.2	12.3	16.8	24.0	24.4	20.9
30-34	7.0	5.7	4.8	3.6	4.0	7.2	5.9	4.6	4.6	5.3	6.7	7.6	9.1	12.4	16.8	15.1	15.0
35-39	6.7	5.3	4.6	2.9	3.6	7.0	5.6	4.2	3.8	4.8	6.7	7.2	9.7	12.1	14.2	13.3	10.8
40-44	7.1	4.8	4.2	3.4	3.1	8.8	5.4	4.2	3.8	5.0	6.3	8.5	9.6	12.5	14.1	12.9	11.1
45-49	6.5	5.1	4.3	3.2	2.9	7.4	5.5	4.4	3.6	4.6	6.6	7.7	9.0	11.4	13.5	11.9	10.6
50-54	5.8	3.9	3.7	2.3	2.3	4.8	4.1	2.8	2.6	3.2	4.5	5.4	6.4	7.2	8.6	8.2	7.4
55-59	5.2	3.9	3.5	2.1	2.0	2.7	3.0	2.3	1.7	2.5	3.6	4.2	4.5	4.8	6.3	5.1	5.1
60-64	4.1	3.1	2.9	1.9	1.5	2.1	1.9	1.9	1.6	1.9	2.6	2.6	3.5	3.3	4.0	4.2	3.8
65-69	3.0	2.5	2.2	1.2	1.1	1.4	1.3	1.3	1.3	1.4	1.7	1.7	2.5	2.2	2.4	2.7	2.2
70-74	3.2	2.3	2.0	1.2	1.0	1.4	1.0	1.0	1.4	1.5	1.6	2.3	2.3	2.0	2.4	2.3	2.4
75-79	5.1	3.0	2.4	1.7	1.3	1.9	1.6	1.2	1.2	1.2	1.6	2.0	1.6	2.0	1.8	1.9	2.0
80+	9.6	6.5	4.6	2.9	2.8	2.4	1.9	1.7	1.2	1.5	1.7	2.9	2.4	2.8	1.8	1.9	1.6

# Chart Suggestions—A Thought-Starter



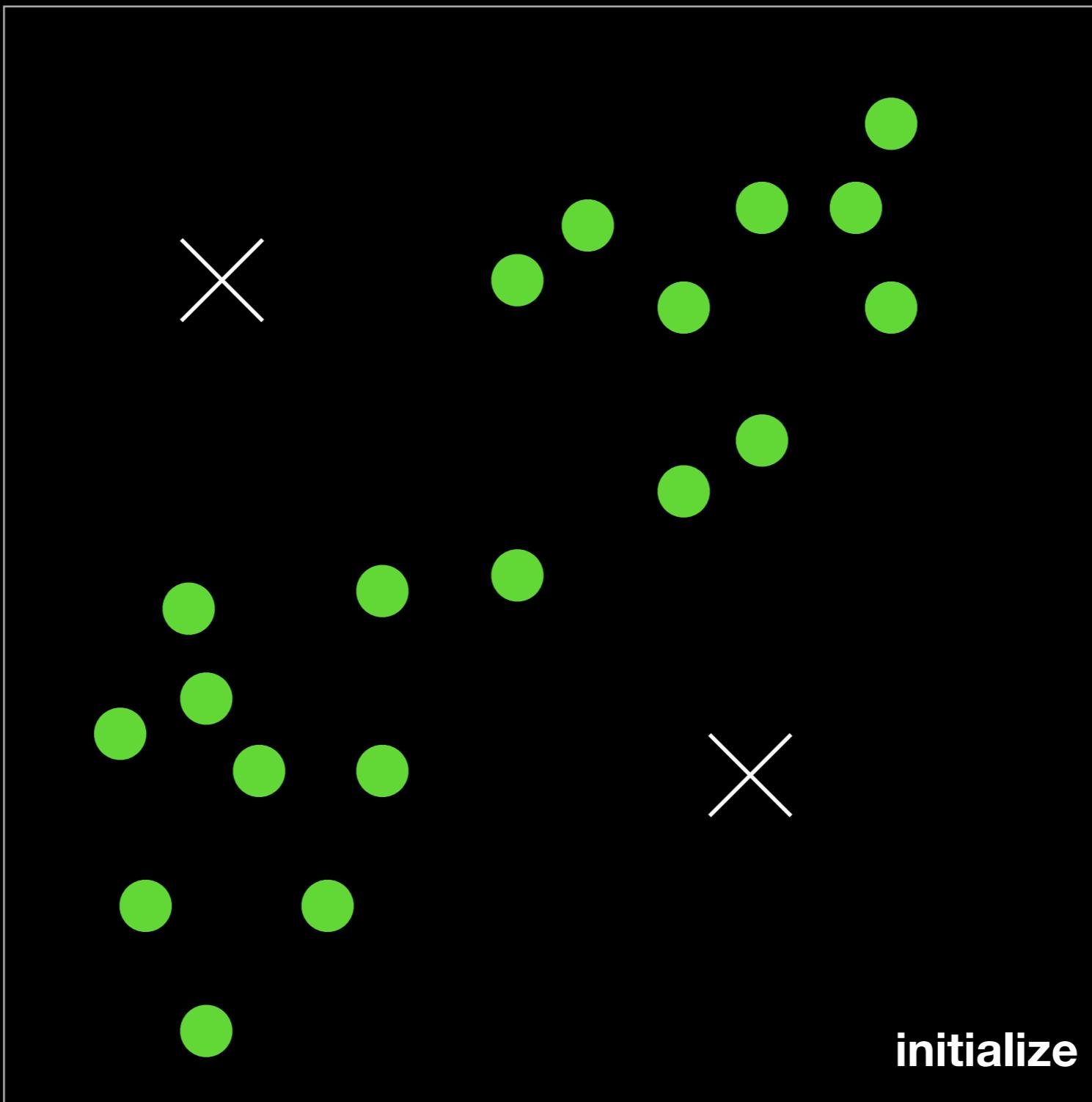
# Reading

- <http://datavisualization.ch/>
- <http://www.princeton.edu/~ina/infographics/index.html>
- <http://www.smashingmagazine.com/2007/08/02/data-visualization-modern-approaches/>
- <http://abeautifulwww.com/2007/05/20/visualizing-the-power-struggle-in-wikipedia/>
- [http://www.designlabelblog.com/2008/03/data-visualization-and-infographics\\_30.html](http://www.designlabelblog.com/2008/03/data-visualization-and-infographics_30.html)
- [http://www.swiss-miss.com/data\\_visualization/page/2](http://www.swiss-miss.com/data_visualization/page/2)
- <http://www.visualcomplexity.com/vc/>
- <http://flowingdata.com/2008/03/19/21-ways-to-visualize-and-explore-your-email-inbox/>
- <http://www.math.yorku.ca/SCS/Gallery/>
- <http://christopherbaker.net/projects/mymap/>
- <http://mashable.com/2007/05/15/16-awesome-data-visualization-tools/>
- <http://www.creativesynthesis.net/blog/2007/05/16/chi-review-qualities-of-perceived-aesthetic-in-data-visualization/>

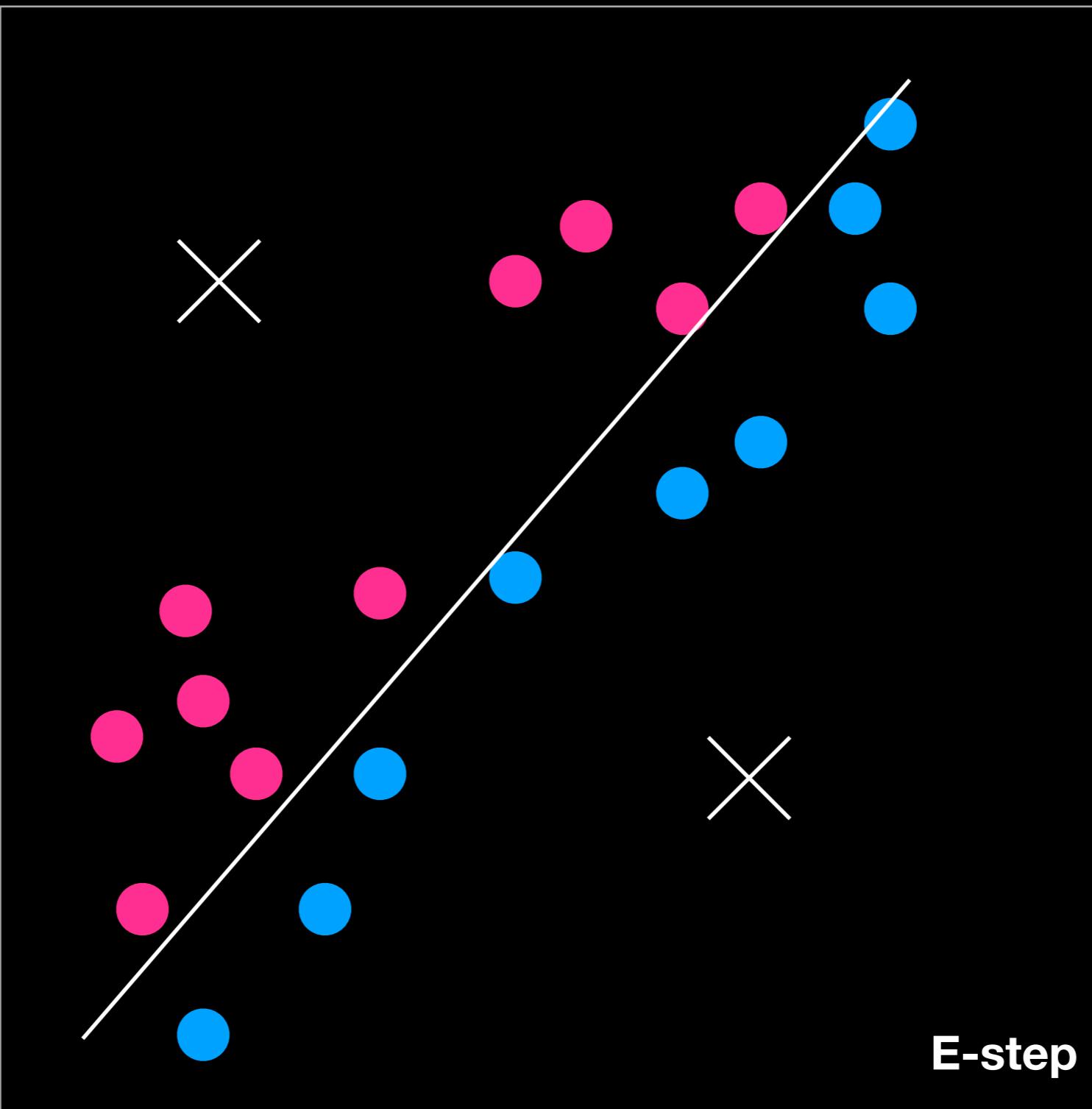
# Clustering

# $k$ -Means Clustering

# $k$ -Means

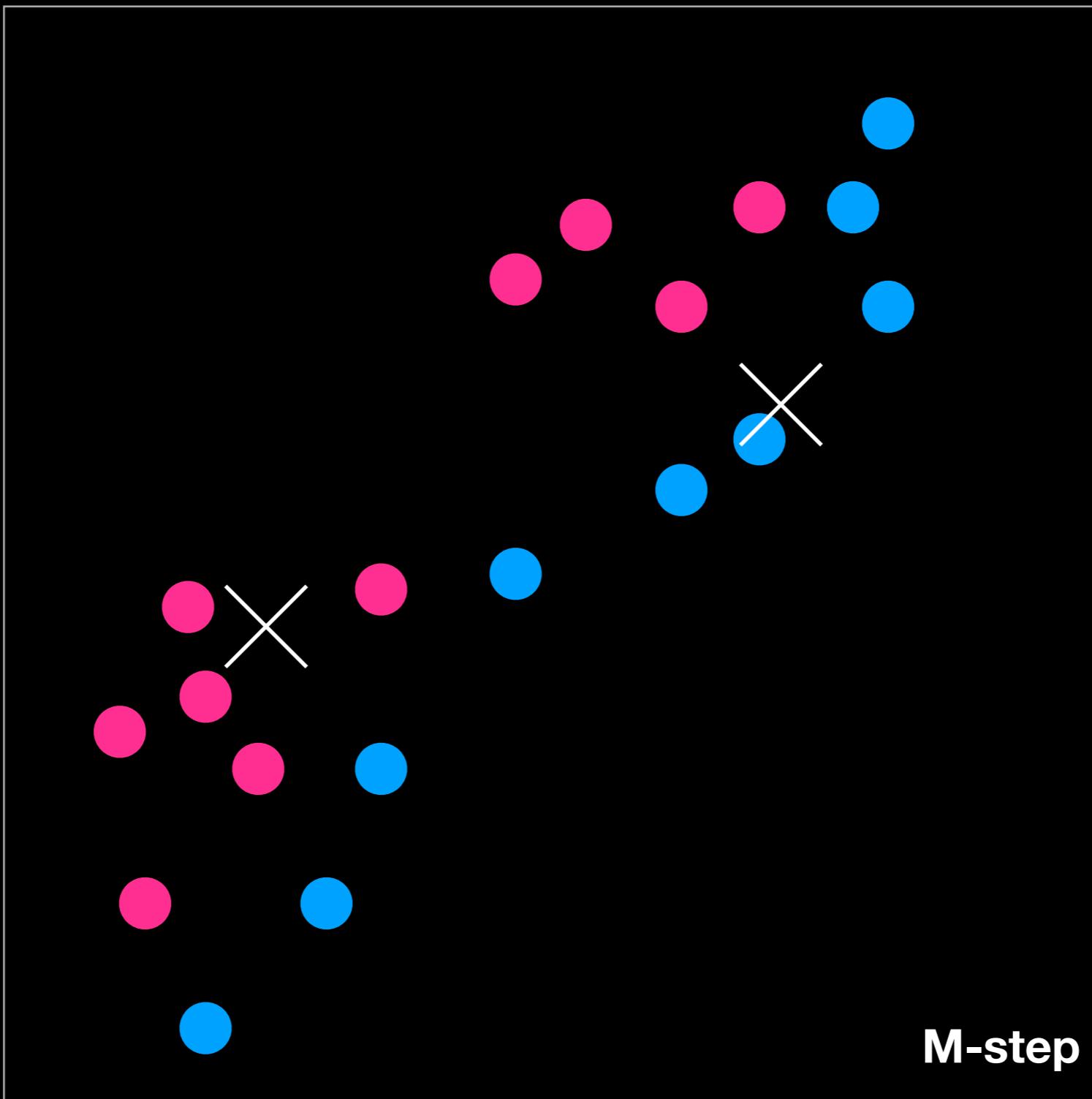


# $k$ -Means



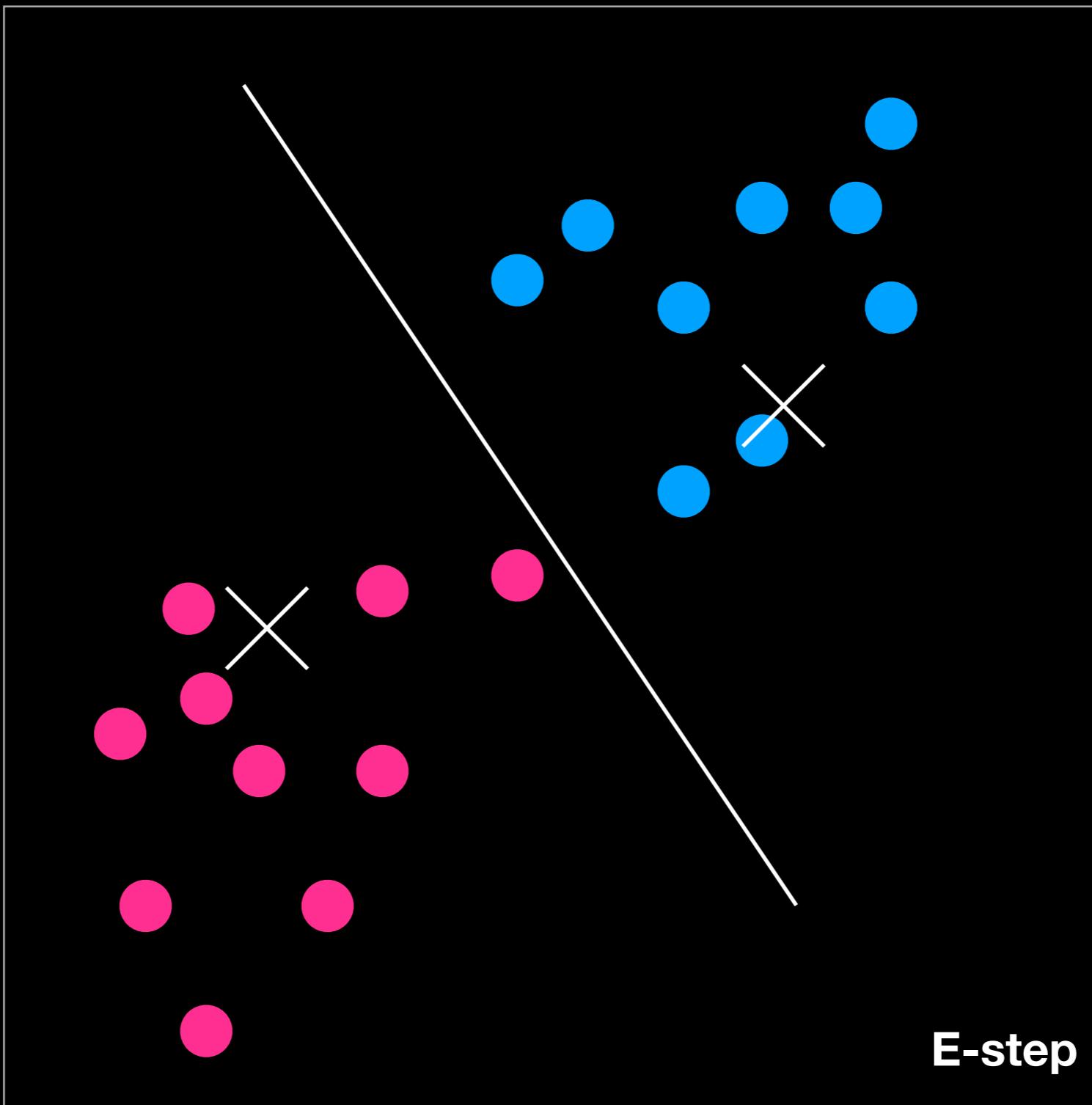
*ASSIGN POINTS TO CENTROIDS*

# $k$ -Means

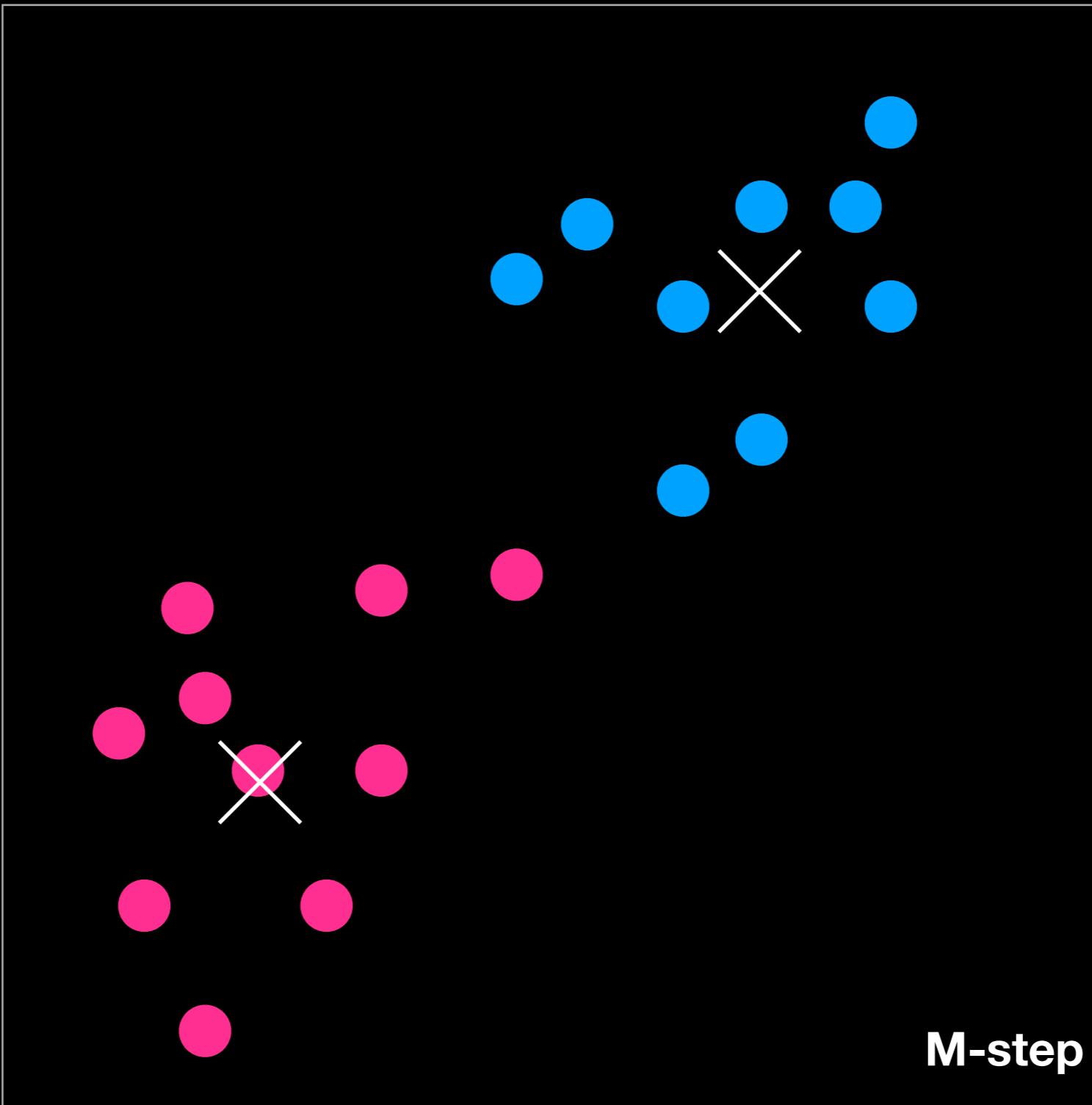


*RECOMPUTE CENTROIDS*

# $k$ -Means

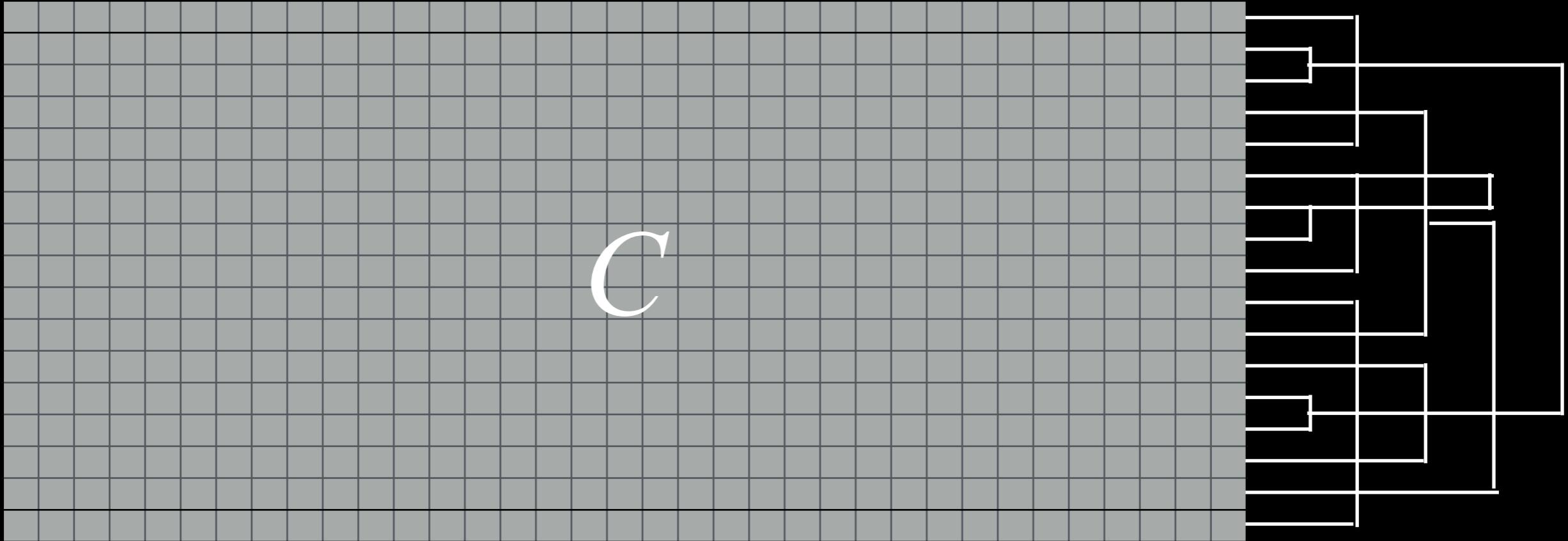


# $k$ -Means

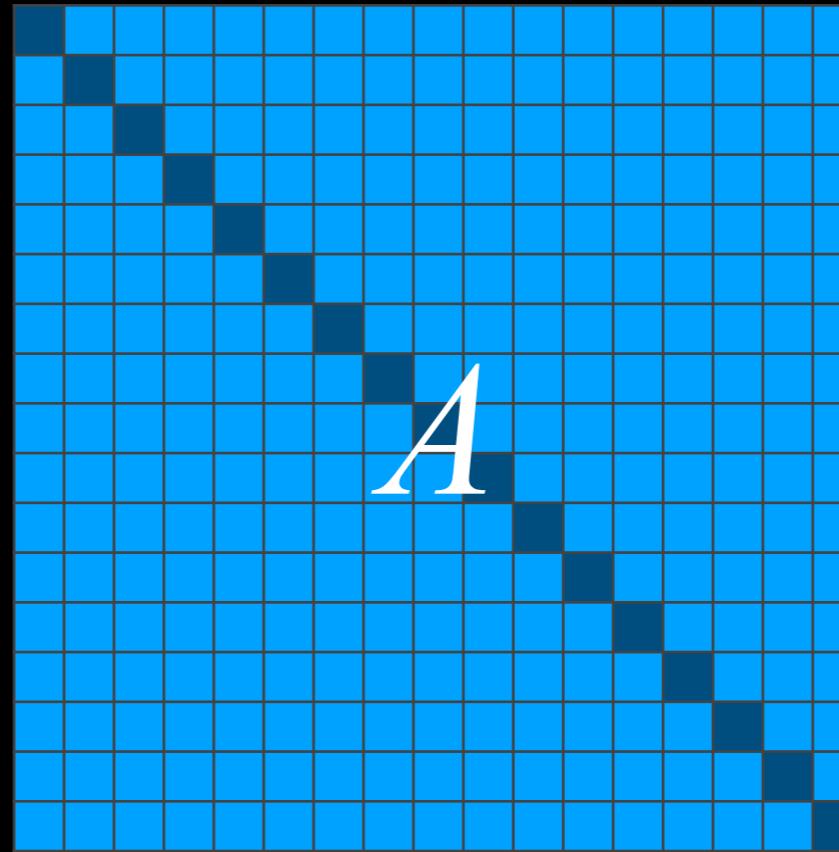


# Agglomerative Clustering

# Building Up

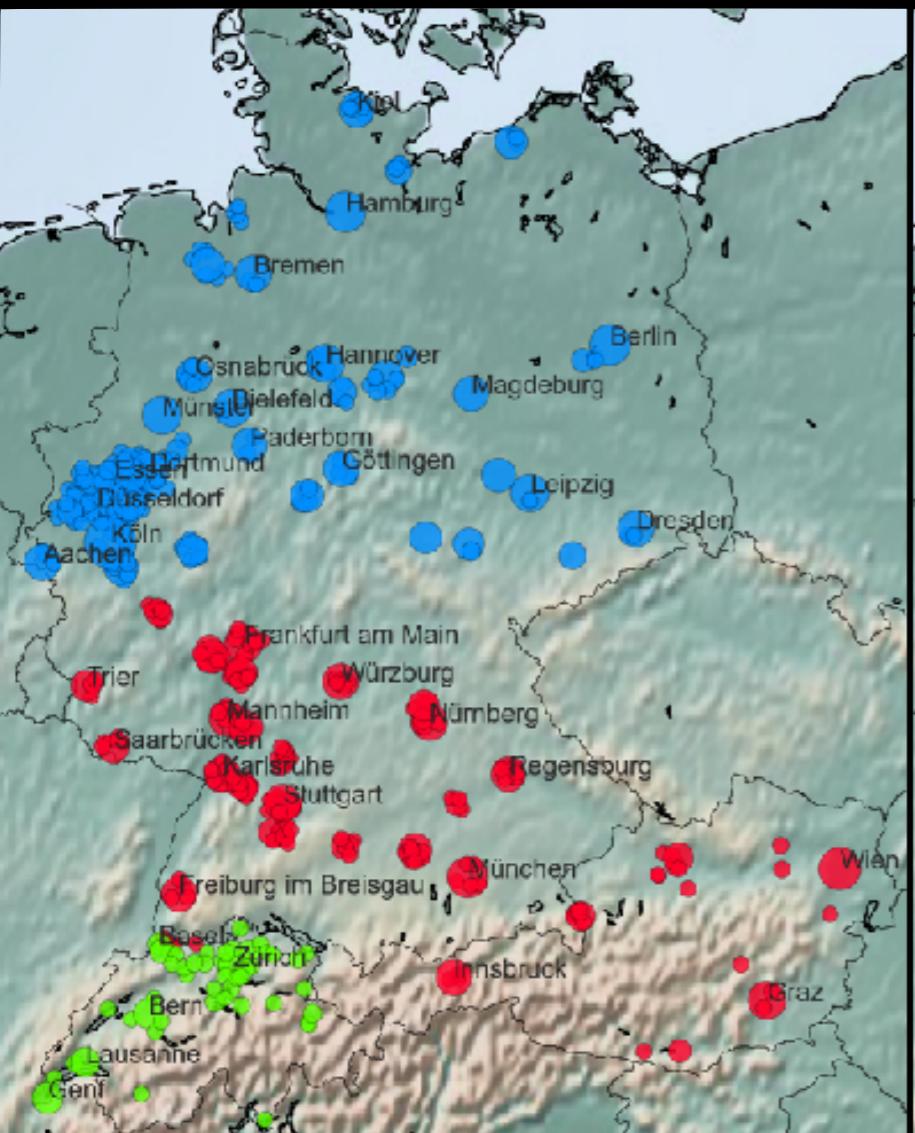


*ADJACENCY  
MATRIX*

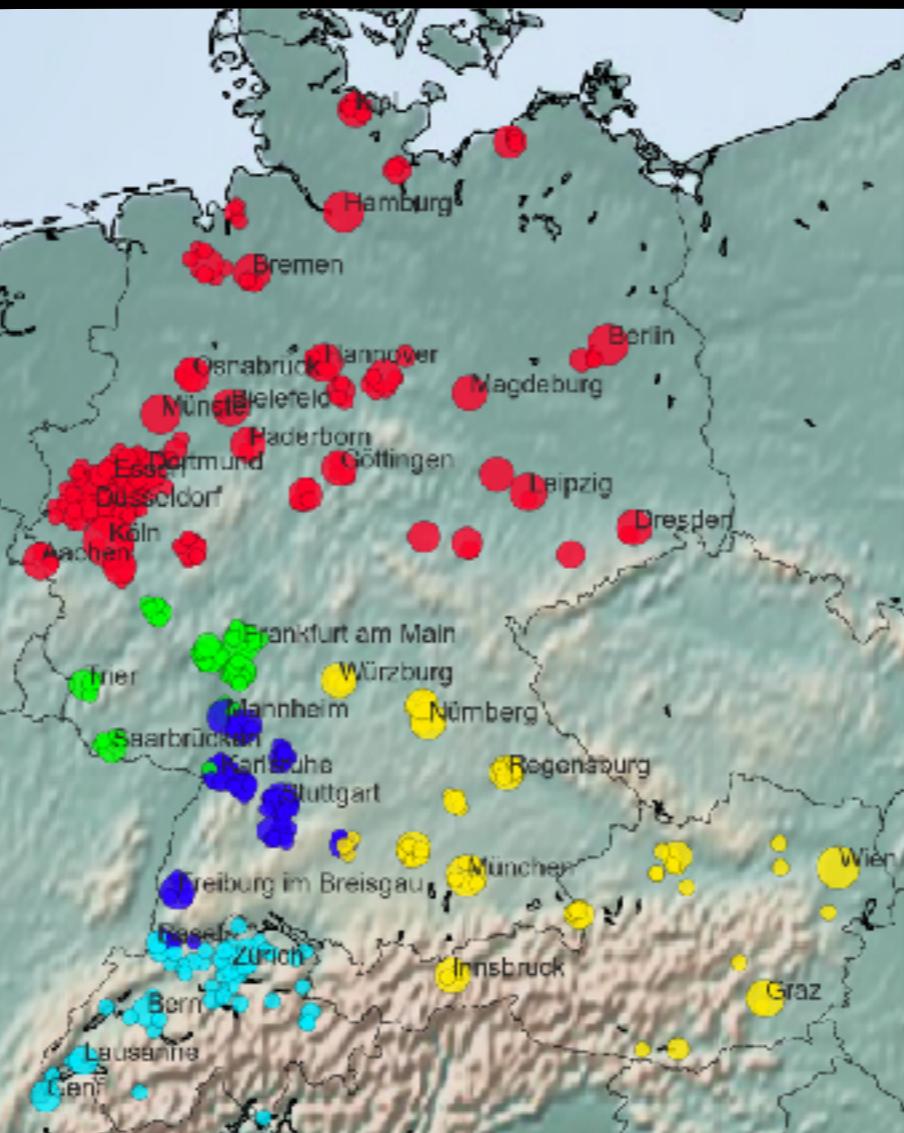


# Dialect Clusters

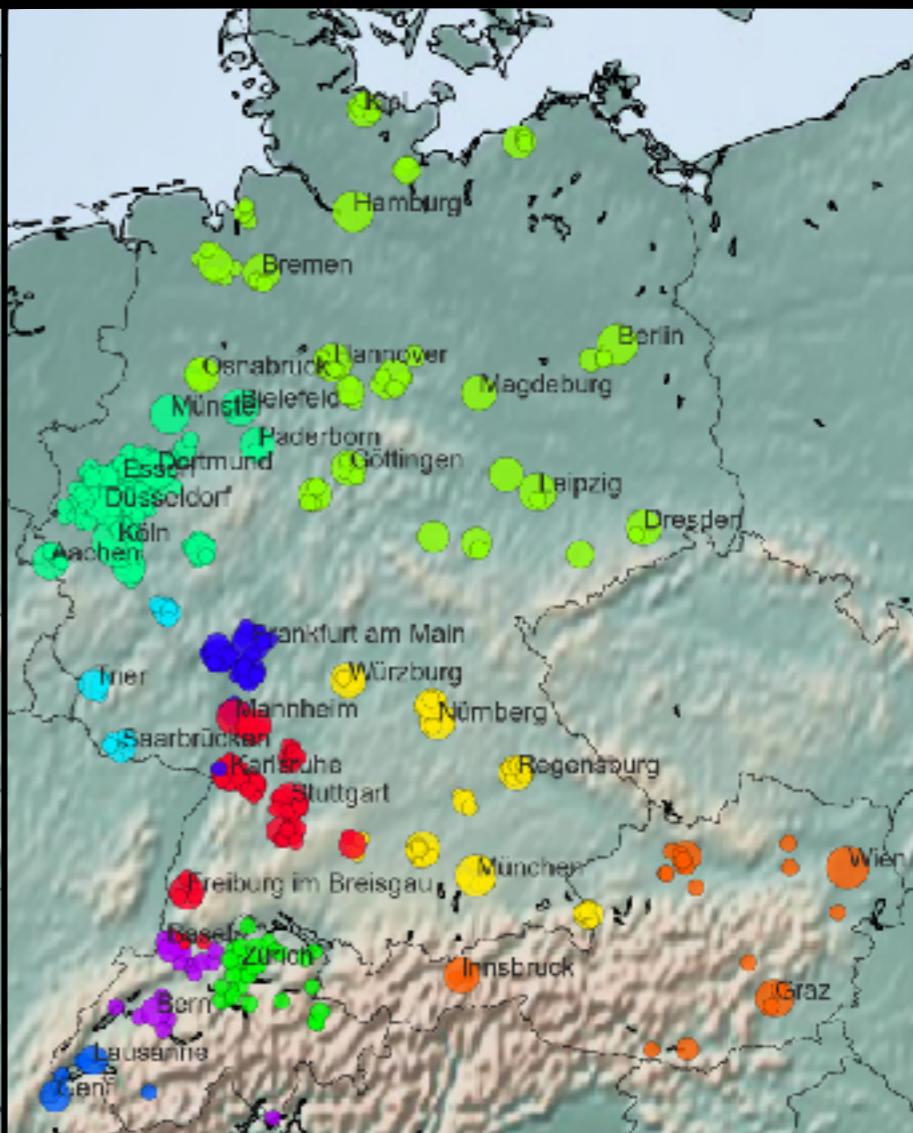
3



5

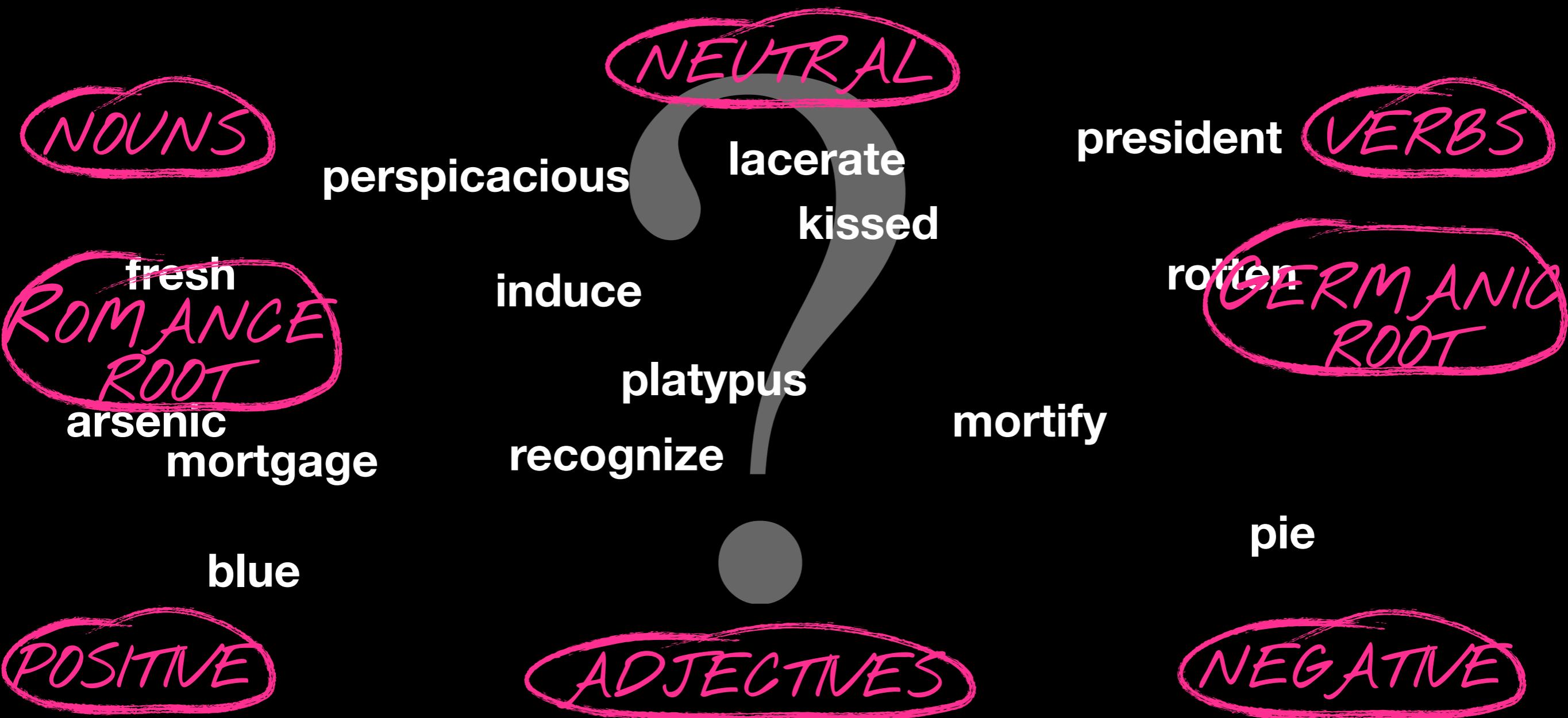


10



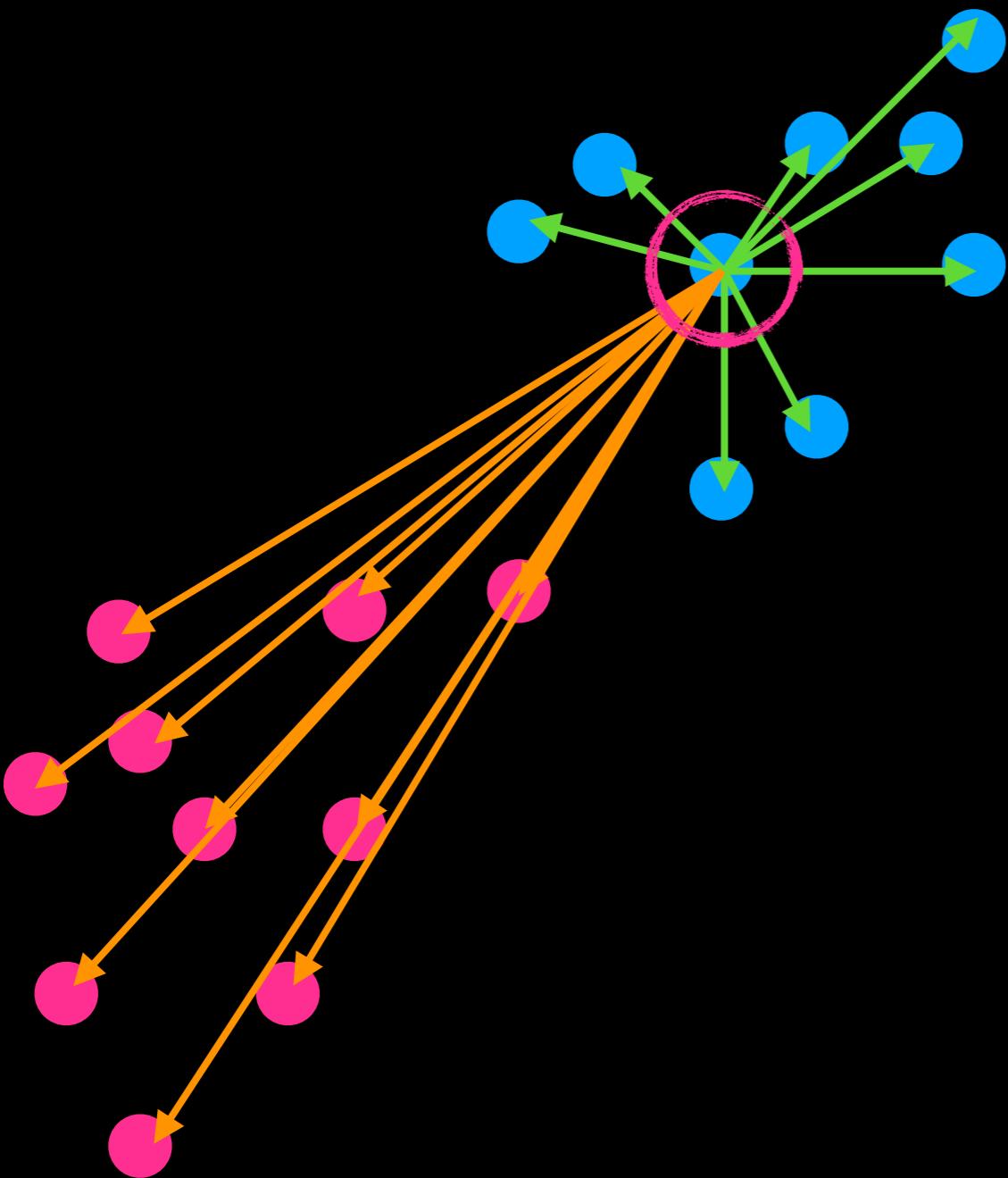
# Evaluating Clusters

# Making Sense of Clusters



# How Many Clusters?

## Silhouette Score



$a = \text{mean intra-cluster distance}$

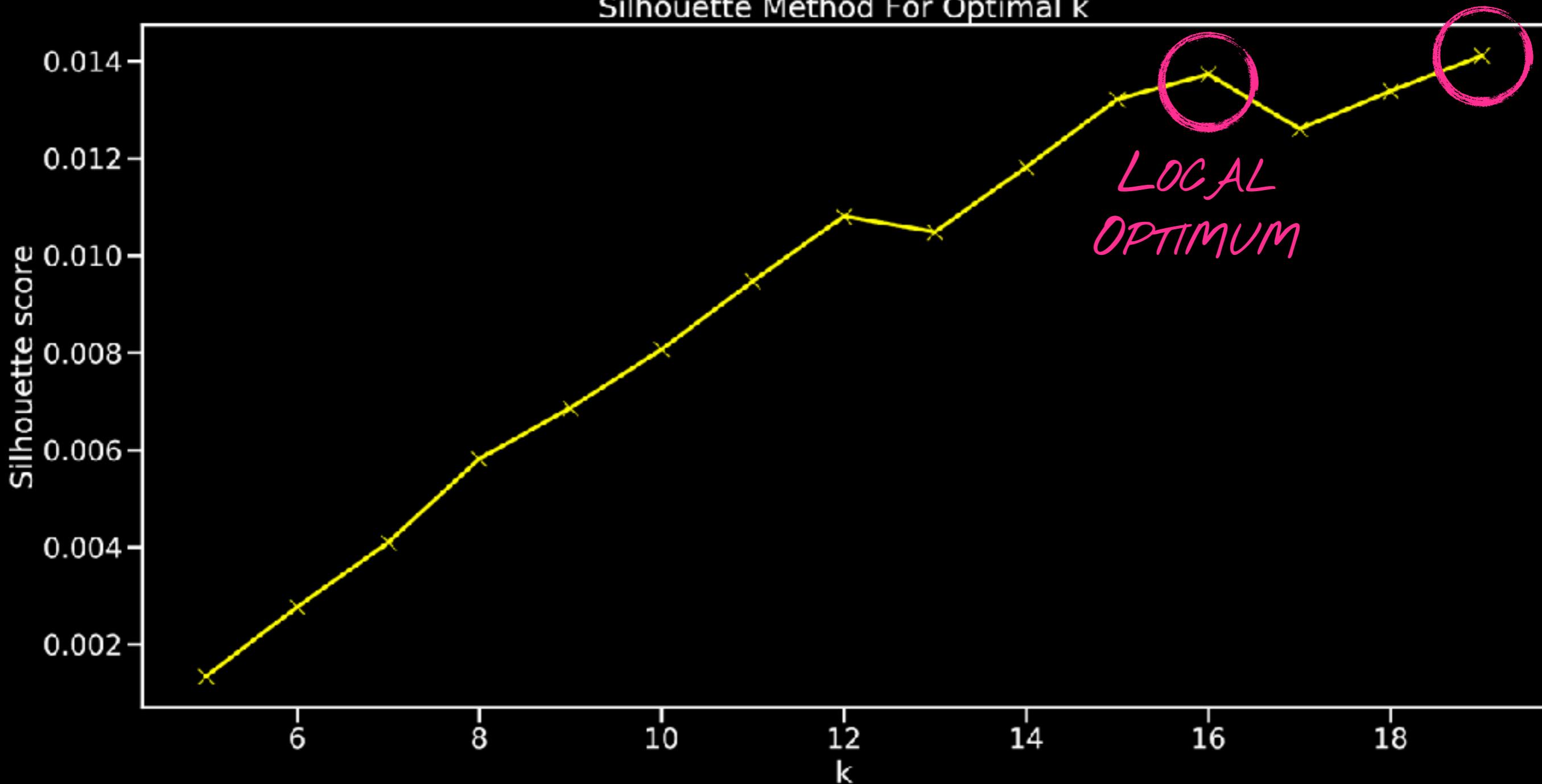
$$S = \frac{(b - a)}{\max(a, b)}$$

$b = \text{mean dist. nearest cluster}$

# Silhouette Scores

*DEPENDS ON PATIENCE/COMPUTE POWER*

Silhouette Method For Optimal k



# Supervised Evaluation Metrics

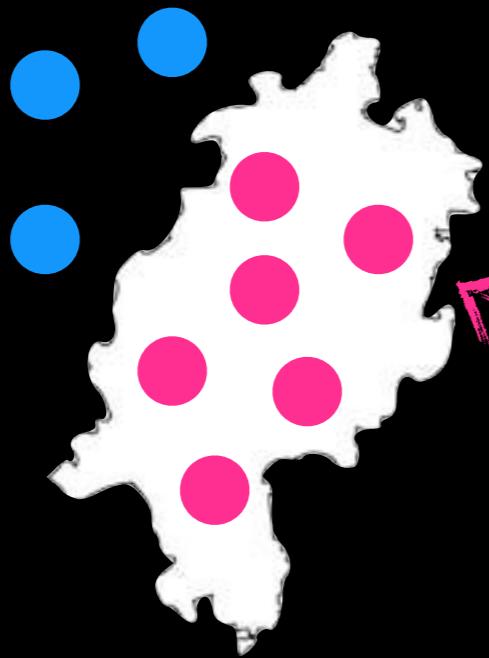
## Homogeneity

cluster has only 1 gold label

## Completeness

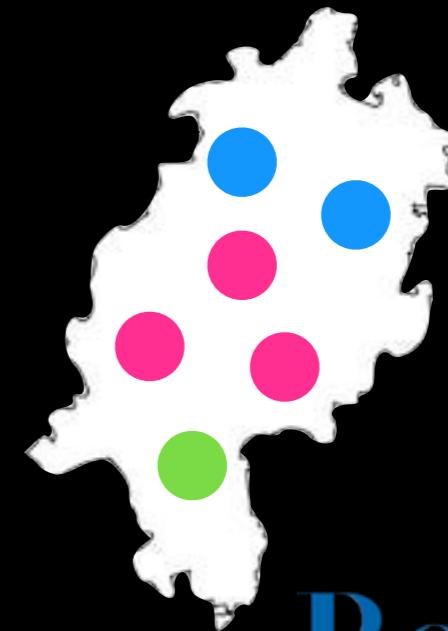
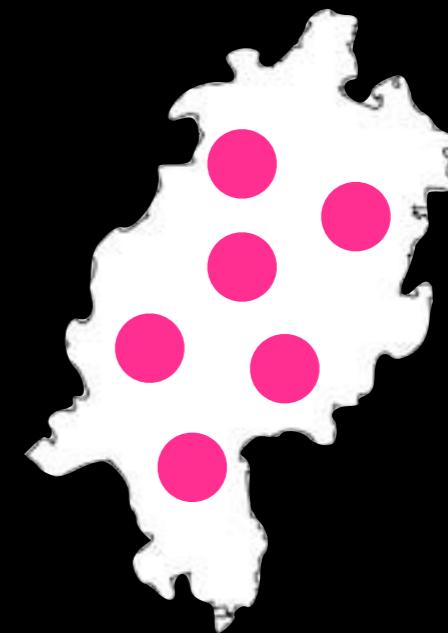
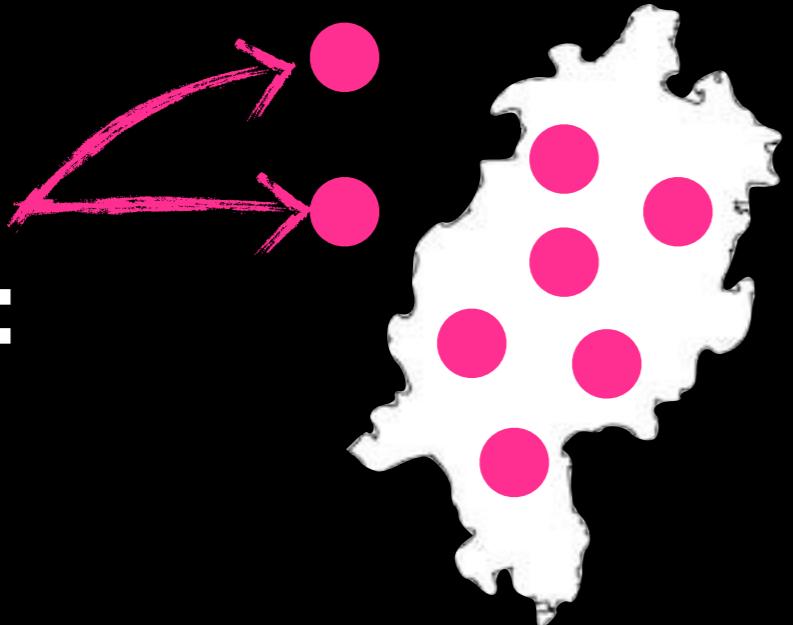
gold label has only 1 cluster

Good:



GOLD LABEL  
(REGION)

Bad:



# Comparison

	<b><i>k-means</i></b>	<b>Agg</b>
<b>scalable</b>	yes	no (up to ~20k)
<b>repeatable result</b>	no	yes
<b>include external info</b>	no	yes
<b>Good on dense clusters?</b>	no	yes

# Wrapping Up

# When to Use What

	Discrete Features	Embeddings
Latent topics	NMF	<i>Not applicable</i>
RGB translation	NMF	SVD + scaling
Plotting	SVD	t-SNE
Clustering	<i>Reduce dimensions</i>	<i>Use as-is</i>

# Take-Home Points

- **Matrix factorization** assumes latent concept dimensions
  - Can be used for semantic similarity (**LSA**)
  - Reduced components can be **visualized** in **graphs** or as **RGB** colors
- **Clusters** can group input in new ways
- Trade-off between speed and interpretability