# Natural Language Processing

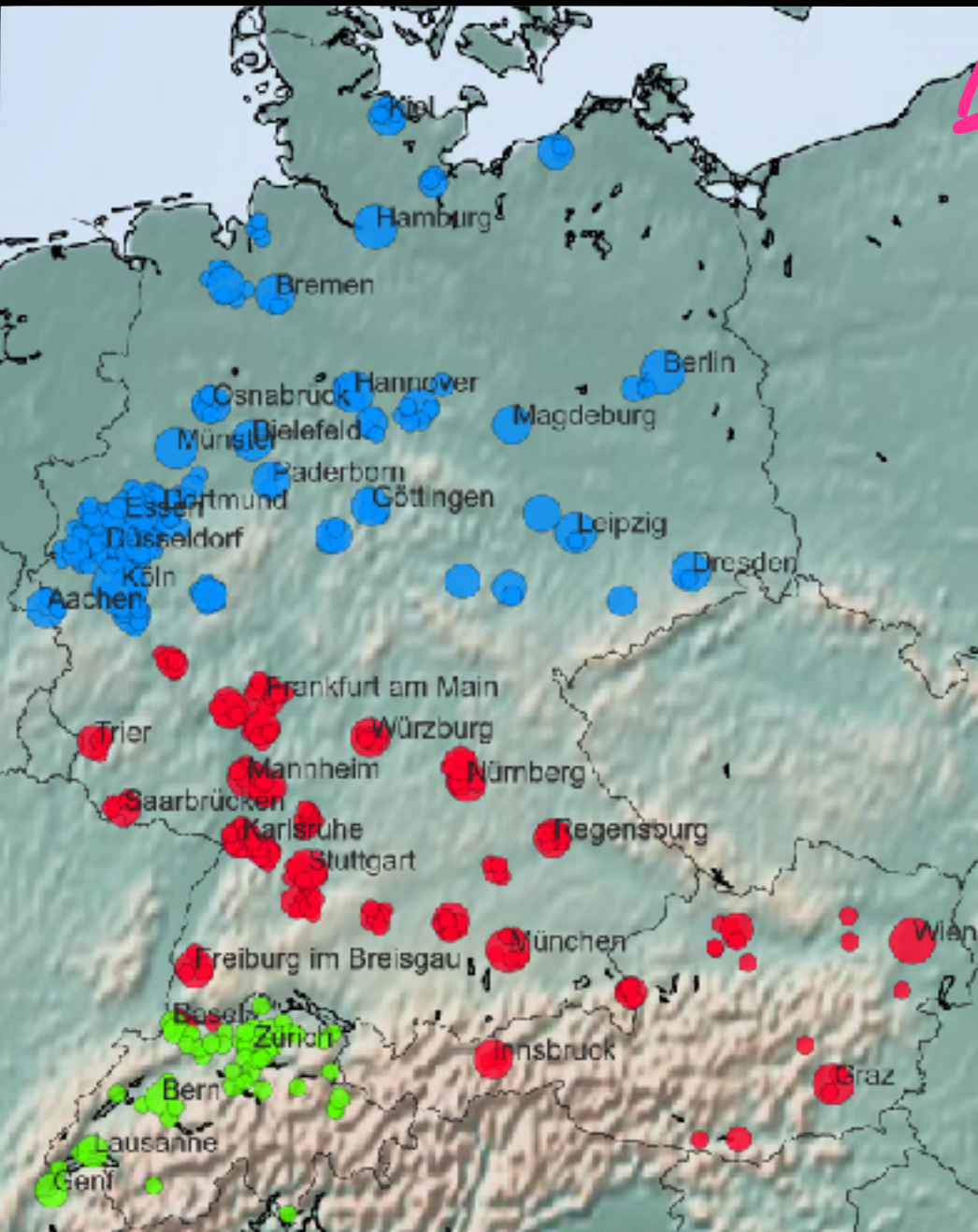**Lecture 11**

Dirk Hovy

<u>dirk.hovy@unibocconi.it</u>
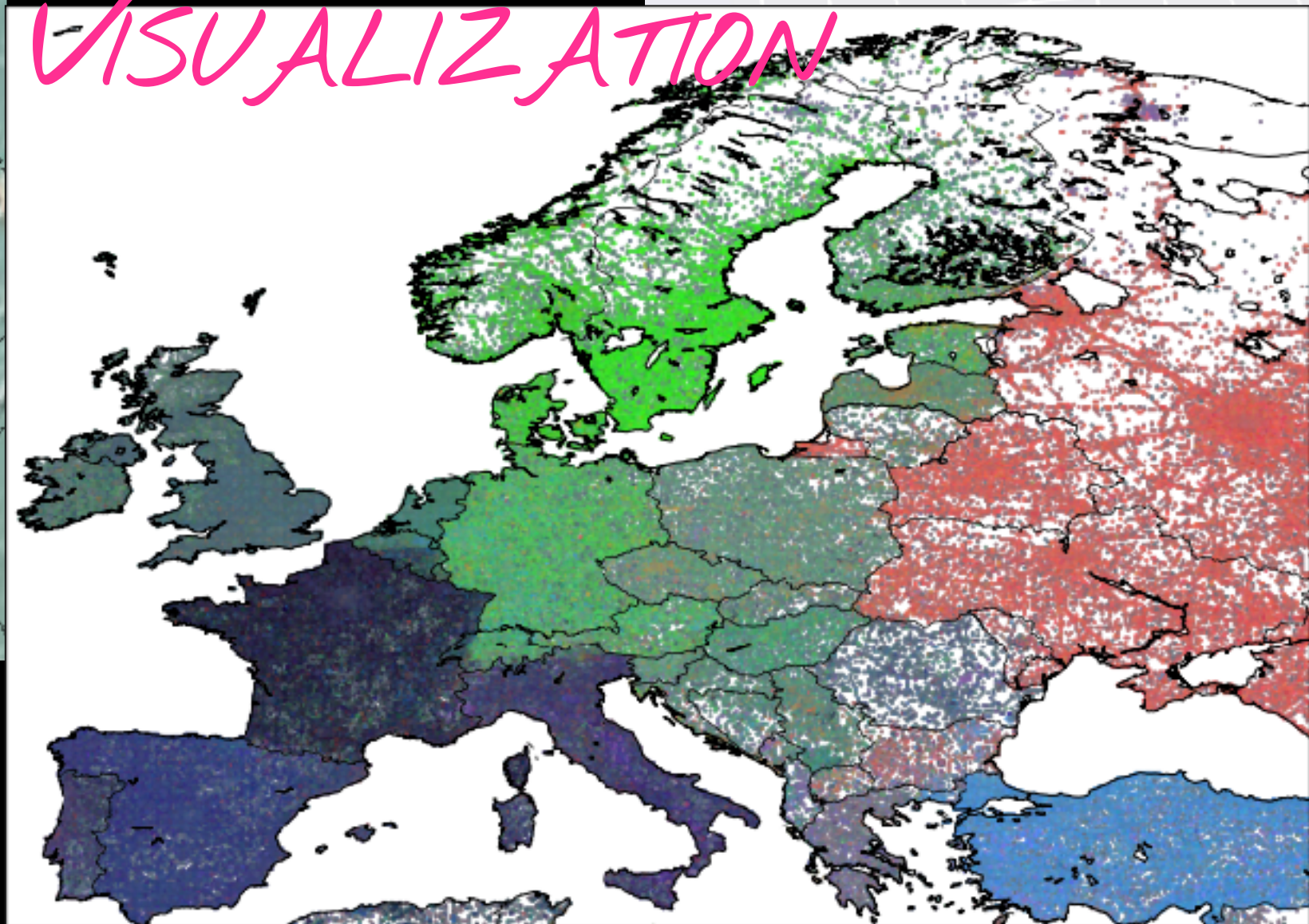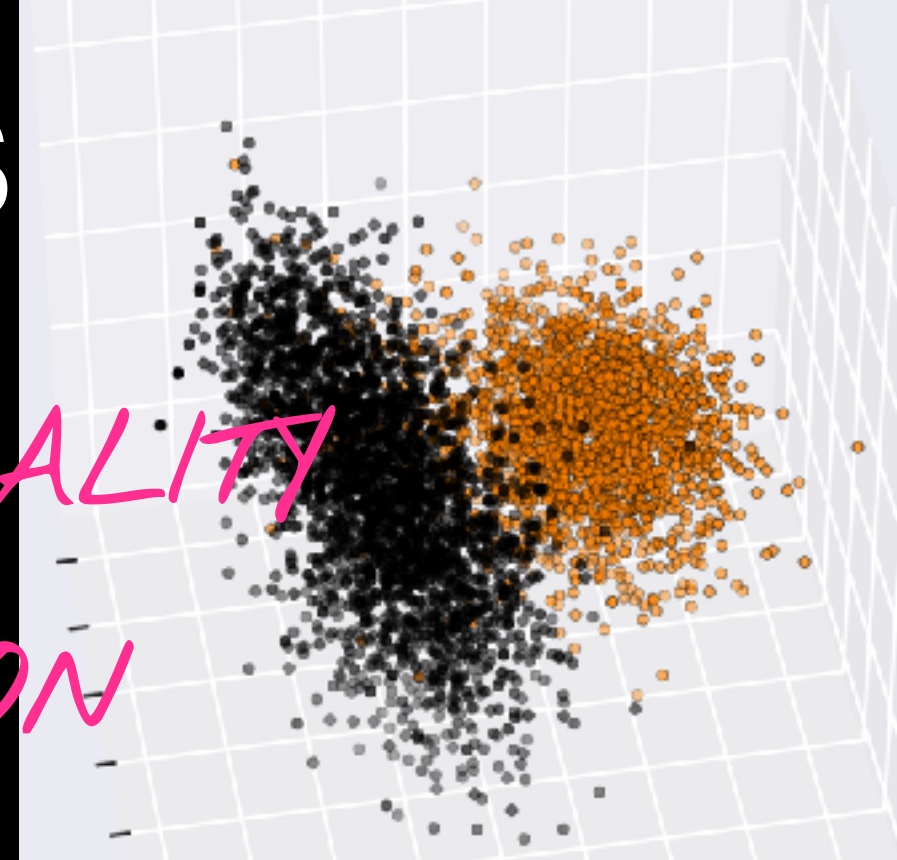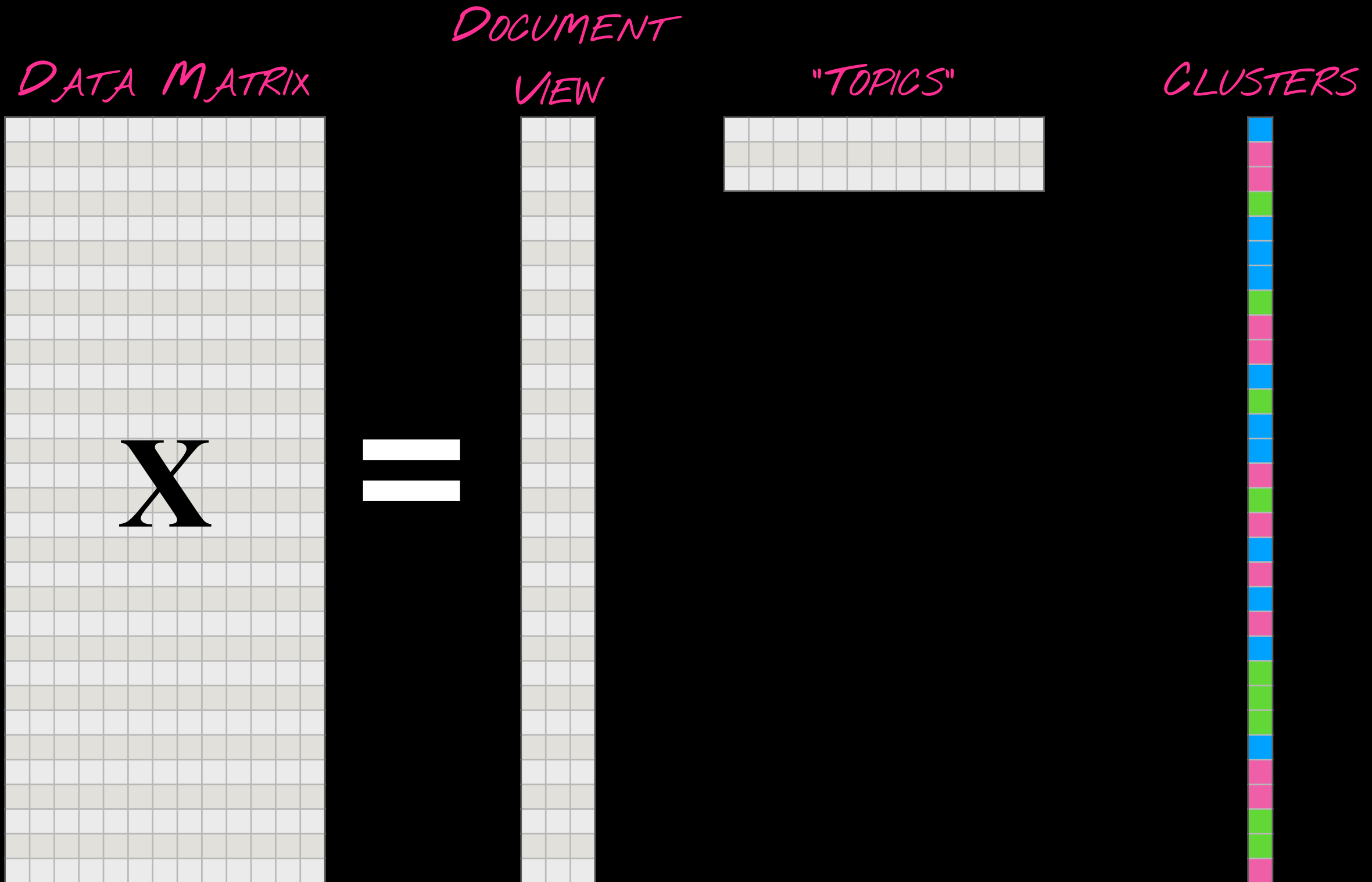
@dirk_hovy

# Examples



DIMENSIONALITY REDUCTION
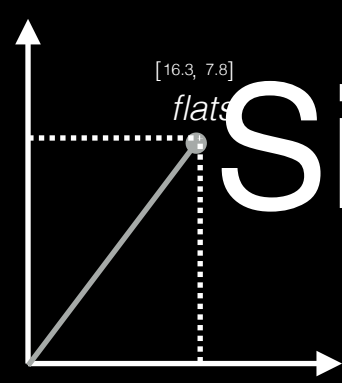
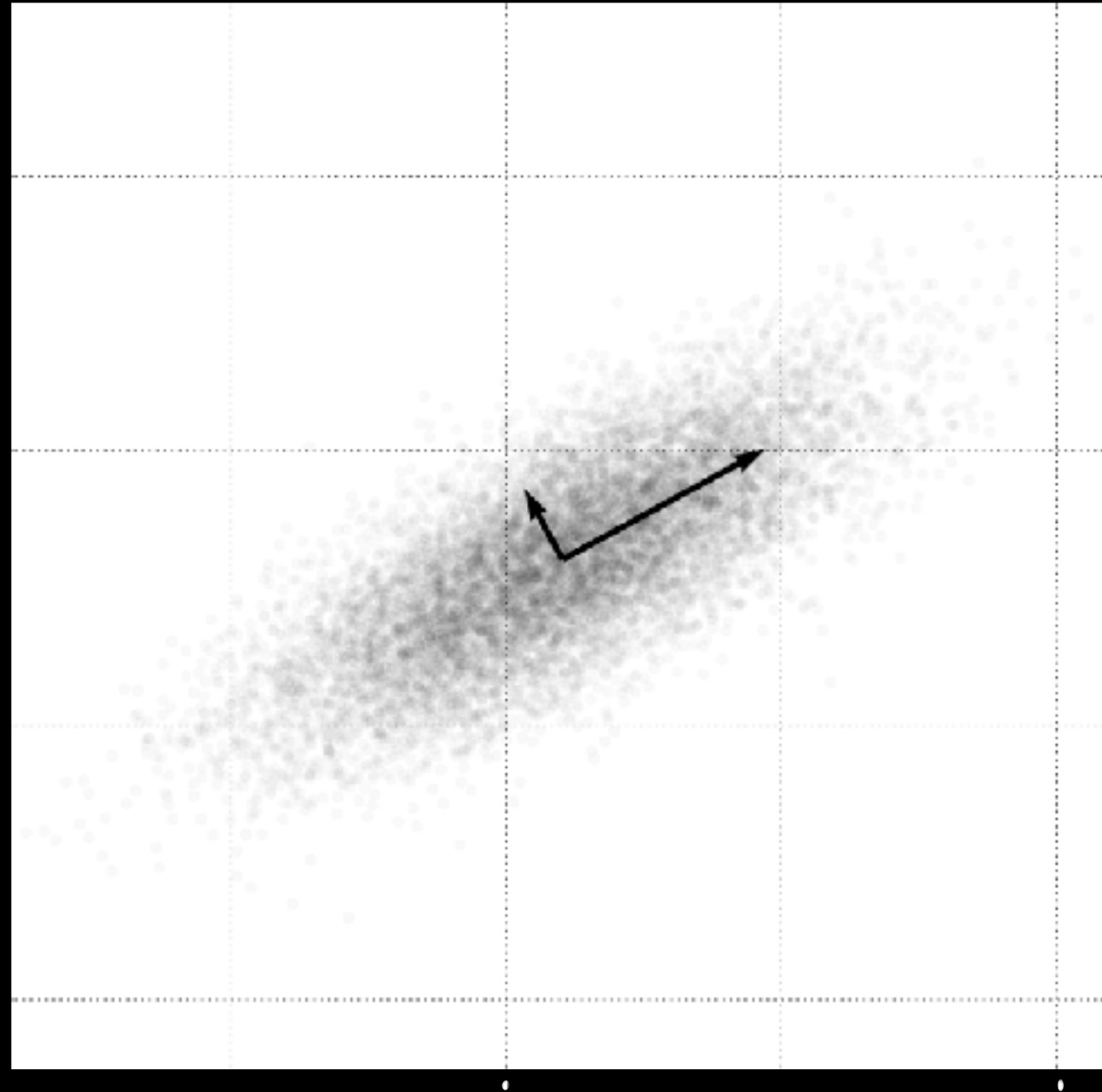VISUALIZATION

CLUSTERING

2

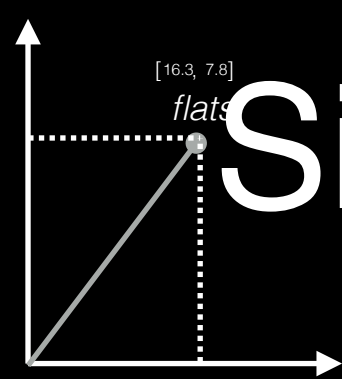# Latent Dimensions

# Goals for Today

- Learn about **matrix factorization** and its use for **semantic similarity** and **visualization**

- Learn about *k*-**means** and **agglomerative clustering**

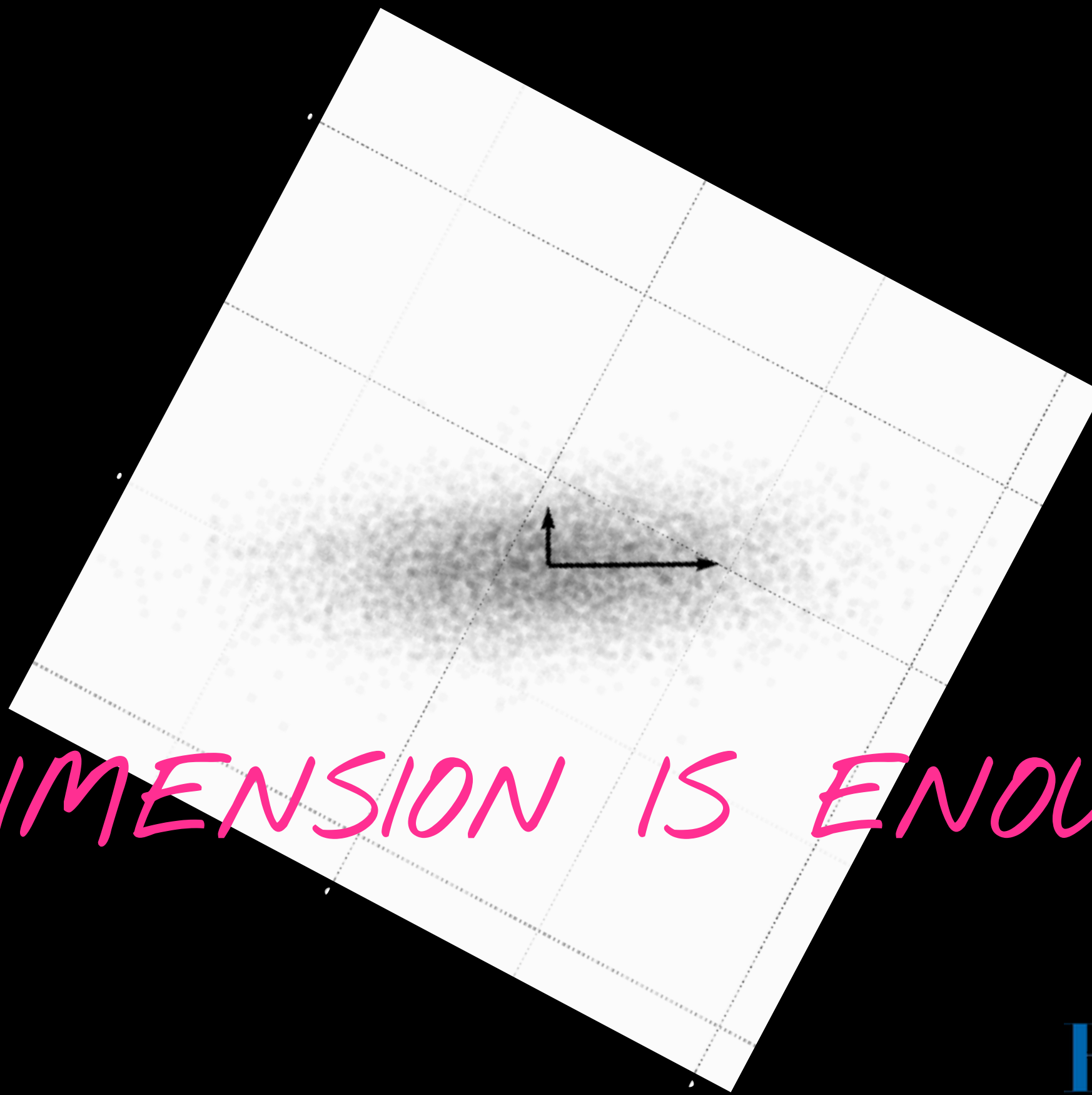- Learn about **evaluation** criteria

Bocconi

# Matrix Factorization

1 DIMENSION IS ENOUGH!

Bocconi

# 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**

**M TERMS**

**N DOCS**

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

D

=

**K CONCEPTS**

**N DOCS**

U

**K CONCEPTS**

**K CONCEPTS**

| 23,2 | | |
| | 2,8 | |
| | | 0,8 |

S

**M TERMS**

**K CONCEPTS**

V$^T$

Bocconi

# Singular Value Decomposition

- reduce principal components/concepts to smaller number

# Singular Value Decomposition

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

# Non-negative Matrix Factorization

- Use only positive values

- Find approximation of two components

DOCUMENT VIEW         "TOPICS"

M TERMS              K CONCEPTS              M TERMS

| N | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| D | 4 | 5 | 6 | 7 |
| O | 1 | 4 6 | 7 |
| C | 2 | 5 | 7 | 8 |
| S | 1 | 4 | 7 | 9 |

D = W · H

11

# Comparison

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

Bocconi

# Yes, but: What is it Good for?

- Find top words for each latent dimension (≈ "topics")

- Find word similarity in latent space (alternative: Word2Vec)

- Find document similarity in latent space (alternative: Doc2Vec)

- Reduce dimensionality for visualization

Bocconi

# Latent Word Dimensions



X ➡ H/V$^T$

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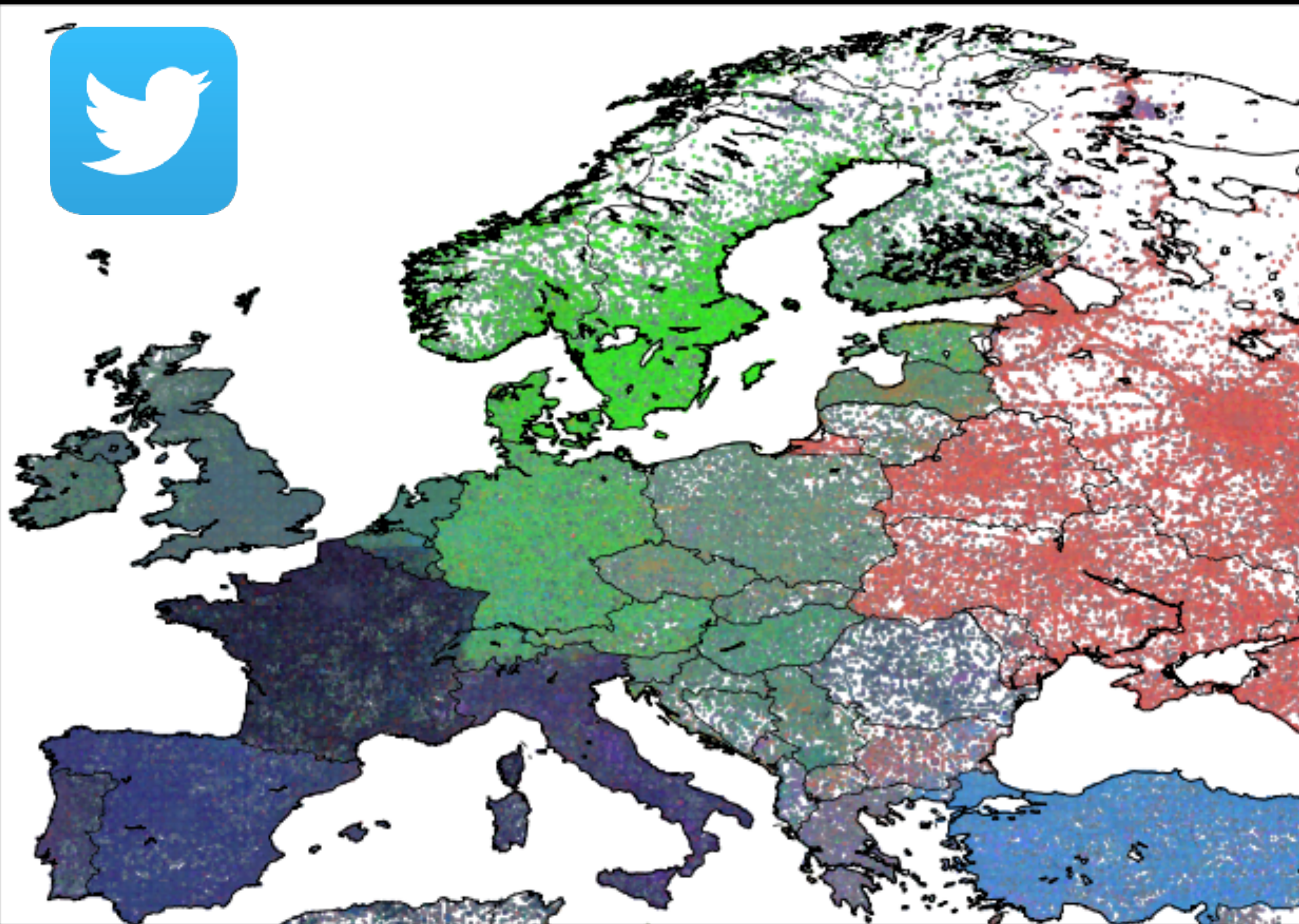sir, aye, let, shall, think
thou, thee, thy, st, god
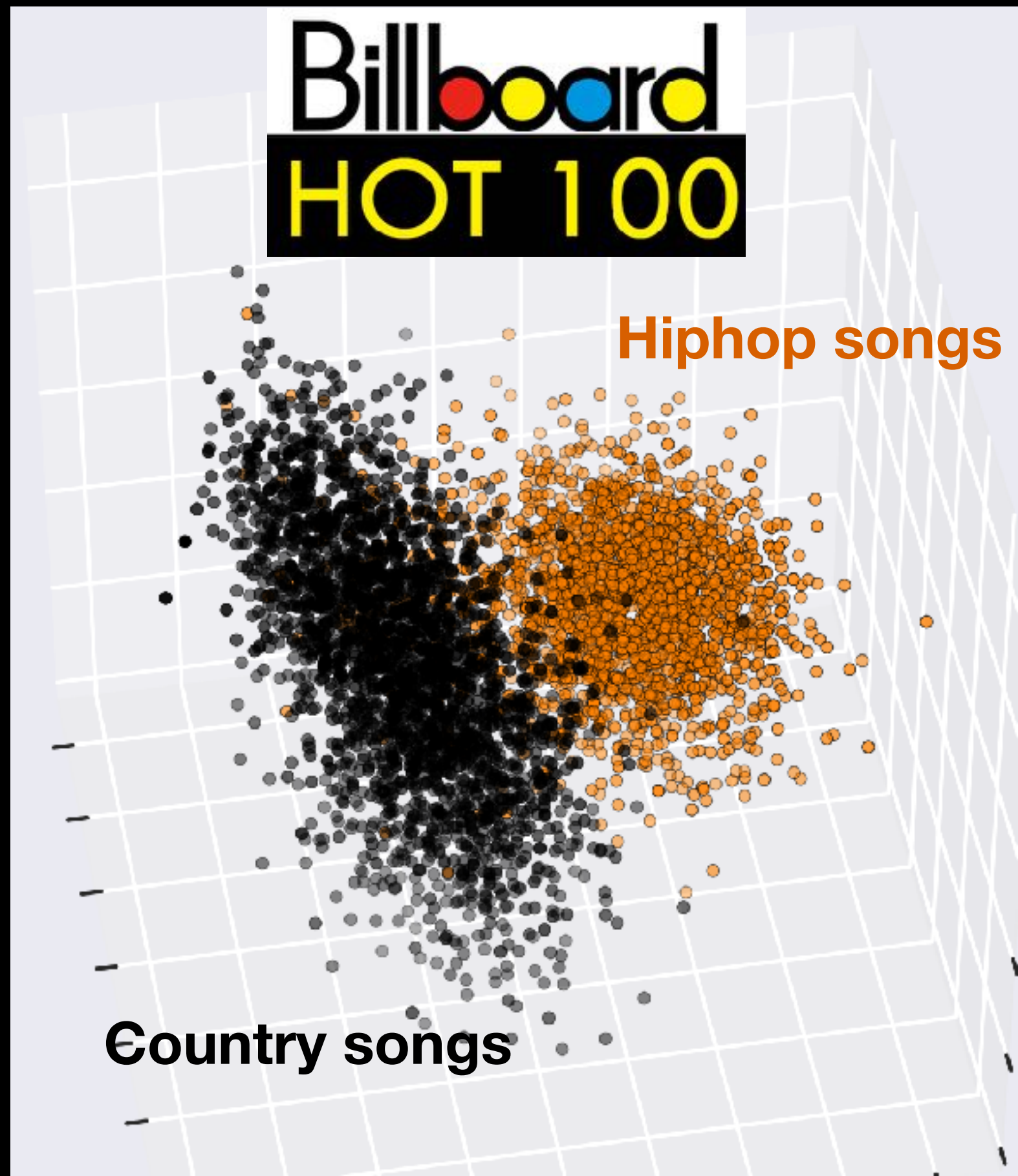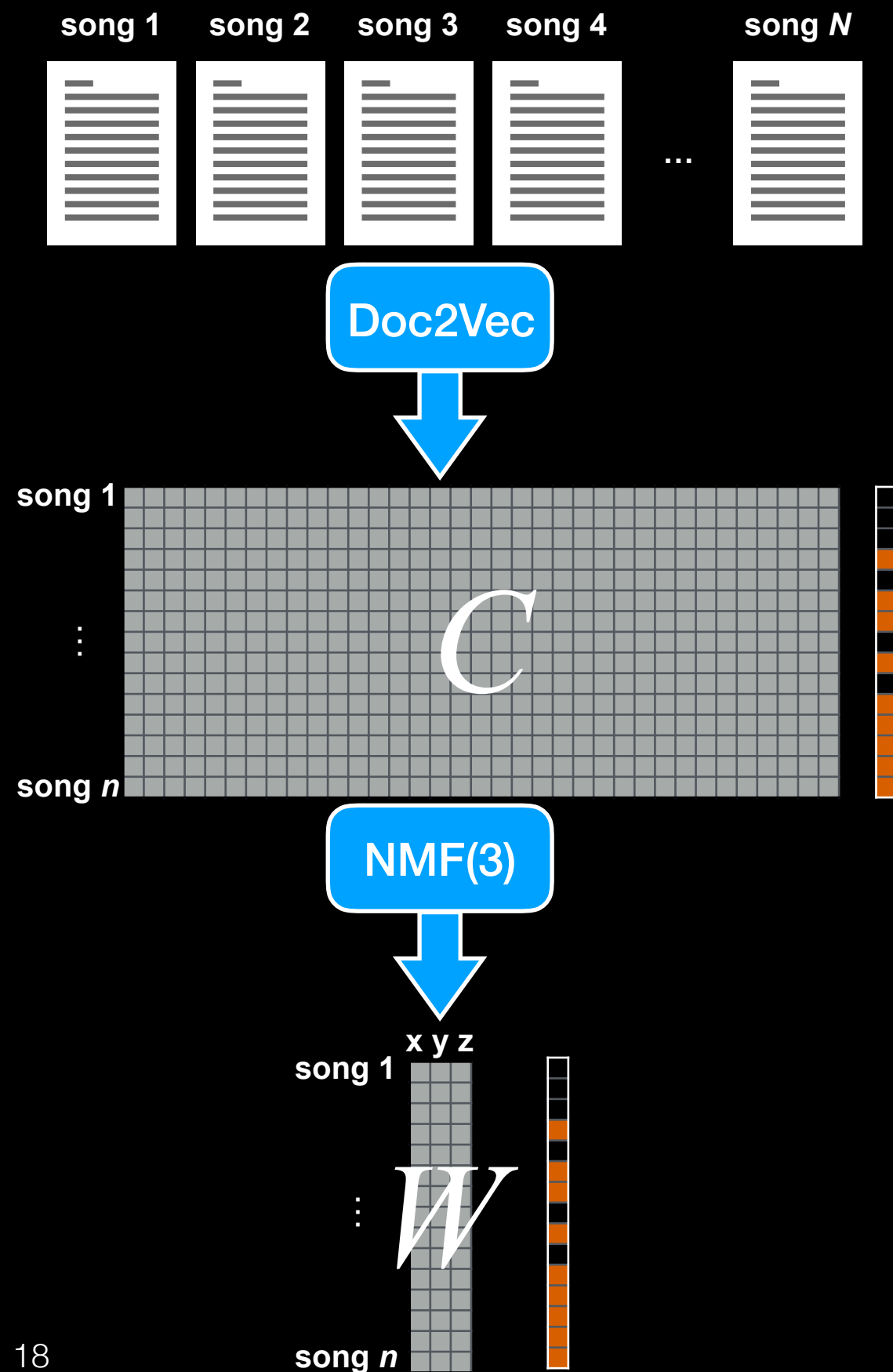whale, sperm, sperm whale, white, white whale
ye, look, say, ye ye, men

Bocconi

# Dimensionality Reduction for Visualizations

# Dimensions as RGB

city 1   city 2   city 3   city 4   city *N*

...

Doc2Vec

city 1

$C$

$\vdots$

city *m*

*DENSE REPRESENTATION*

=

city 1

$\vdots$

city *m*

RGB

$V =$

x

$H$

*NMF*

16

# Dimensions as RGB

# Dimensions as Position

song 1  song 2  song 3  song 4  song *N*

...

**Doc2Vec**

song 1

$C$

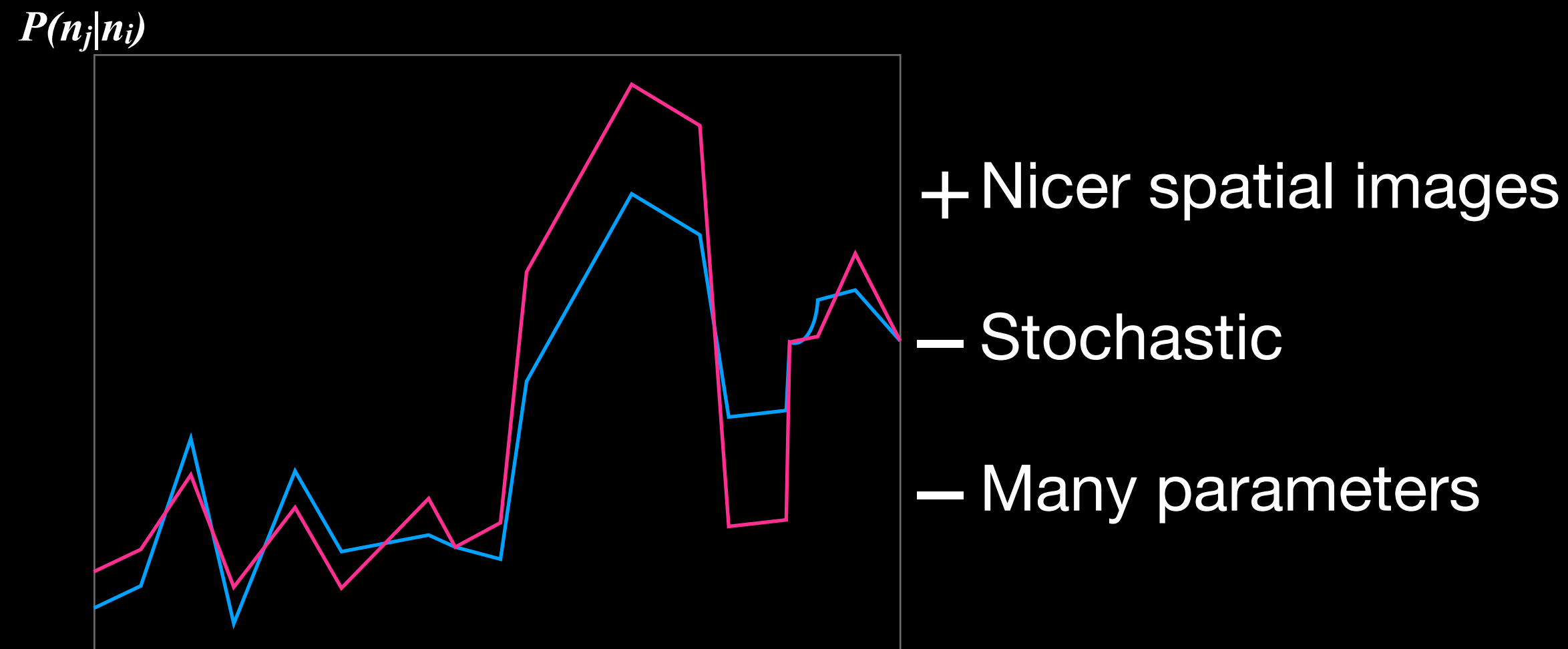song *n*

**NMF(3)**

x y z

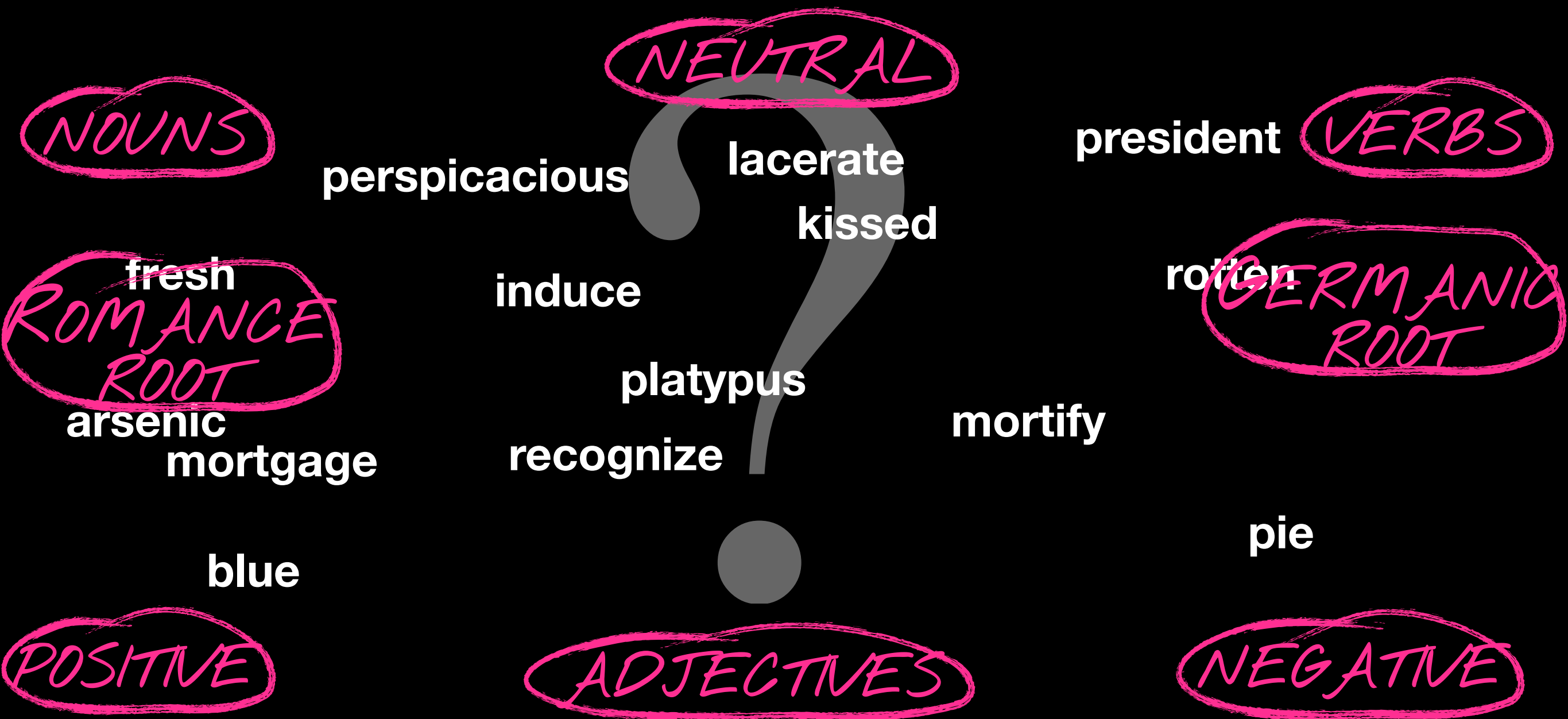song 1

$W$

song *n*

Hiphop songs

Country songs

# t-SNE

- Map/preserve neighborhood structure of high-dimensional space in lower-dimensional space

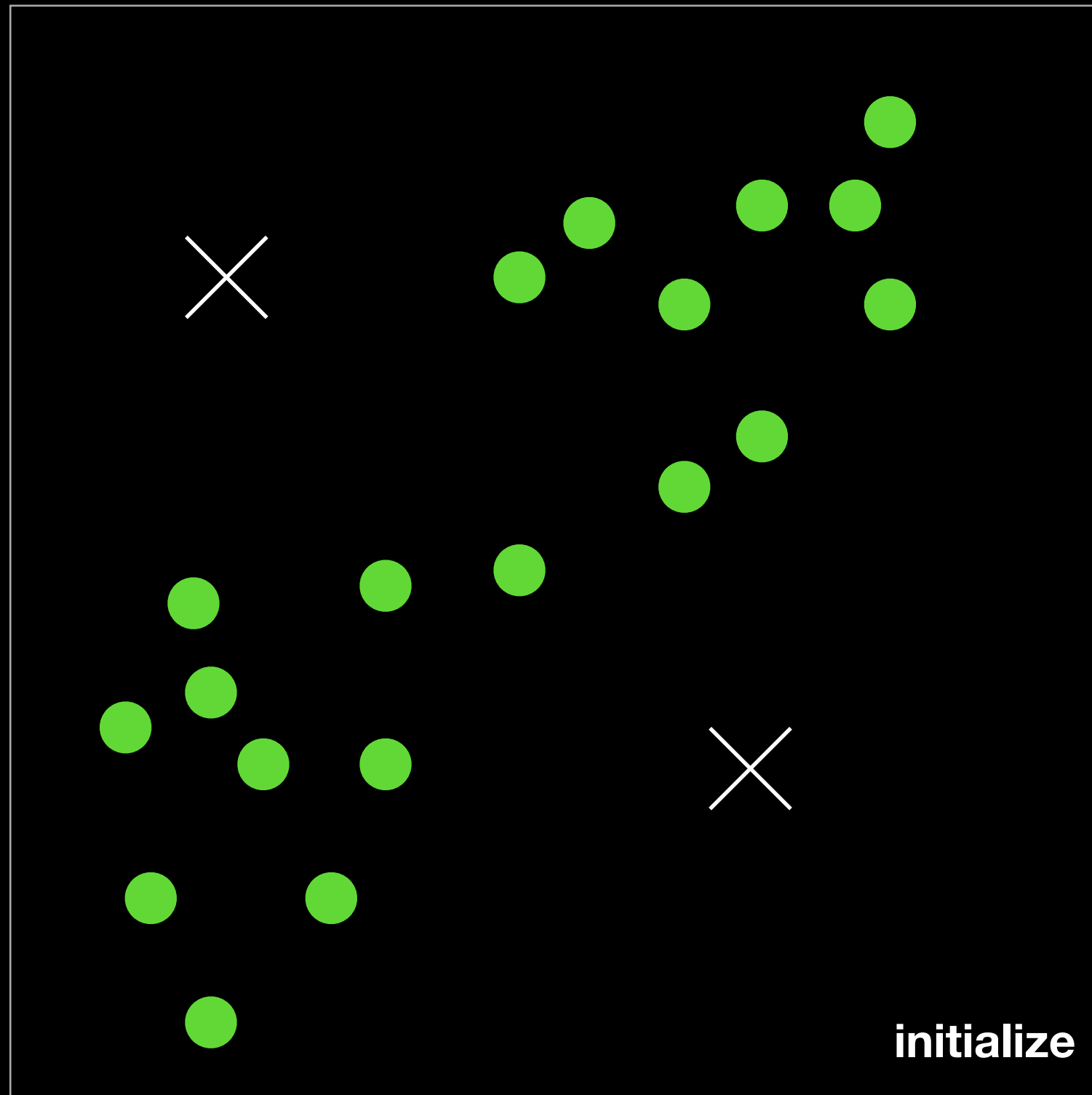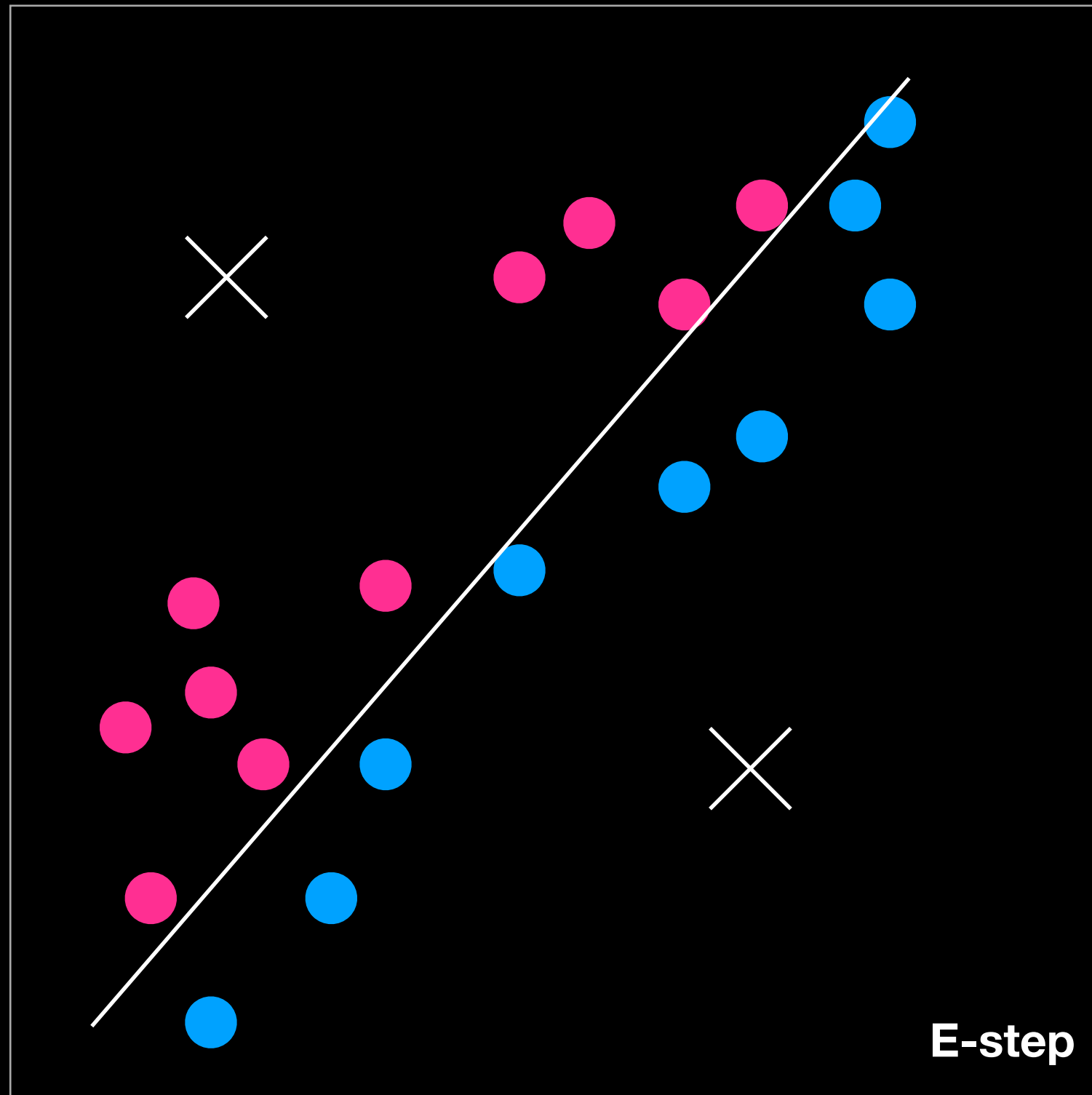- Minimize difference between probability distributions over neighbors in both dimensions

$P(n_j|n_i)$



+ Nicer spatial images

— Stochastic

— Many parameters

Bocconi

# Clustering

# Making Sense of Clusters

NEUTRAL

NOUNS

perspicacious

lacerate

president VERBS

kissed

fresh

ROMANCE ROOT

induce

rotten GERMANIC ROOT

arsenic

platypus

mortify

mortgage

recognize

pie

blue

POSITIVE

ADJECTIVES

NEGATIVE

# *k*-Means Clustering

# *k*-Means



initialize

23

# *k*-Means



E-step

ASSIGN POINTS TO CENTROIDS

# *k*-Means



M-step

RECOMPUTE CENTROIDS

# *k*-Means



E-step

26

Bocconi

# *k*-Means



M-step

27

# Agglomerative Clustering

# Building Up

*C*

*ADJACENCY*

*MATRIX*

*A*

**Bocconi**
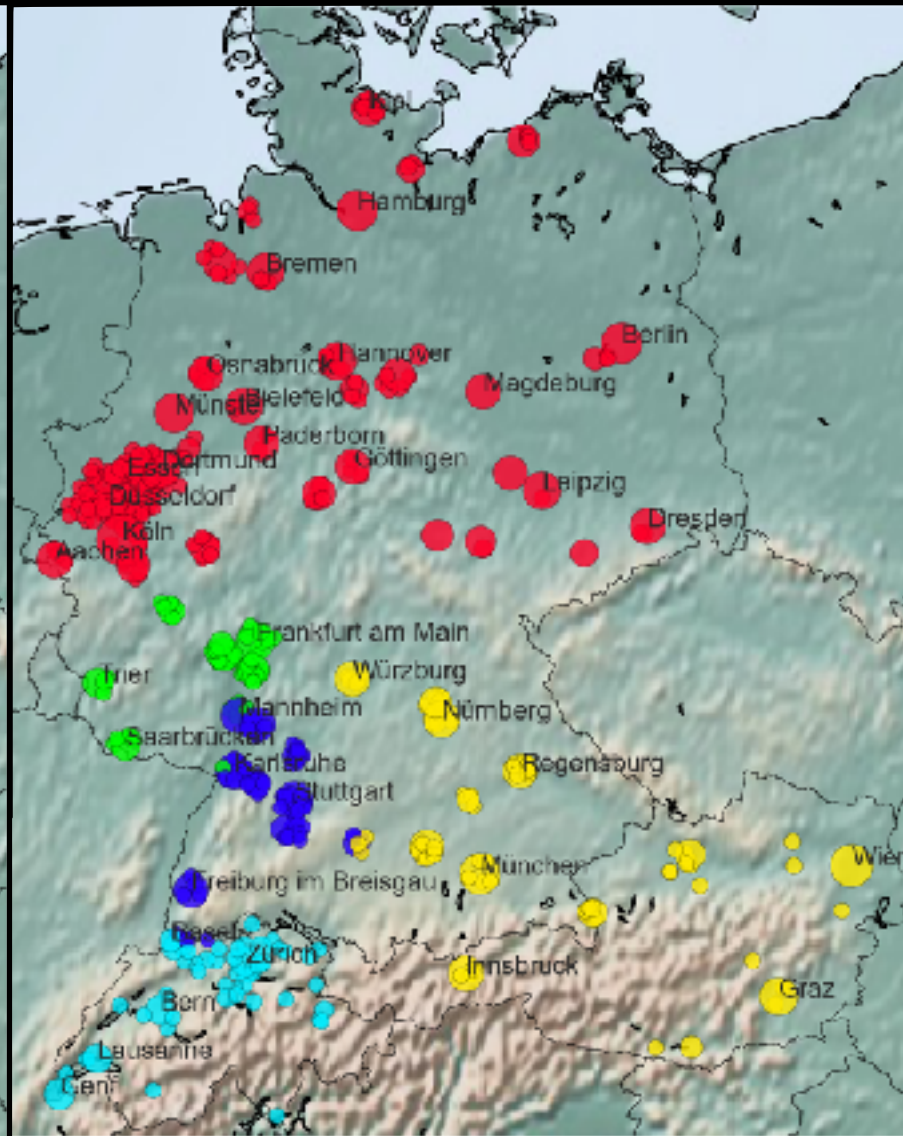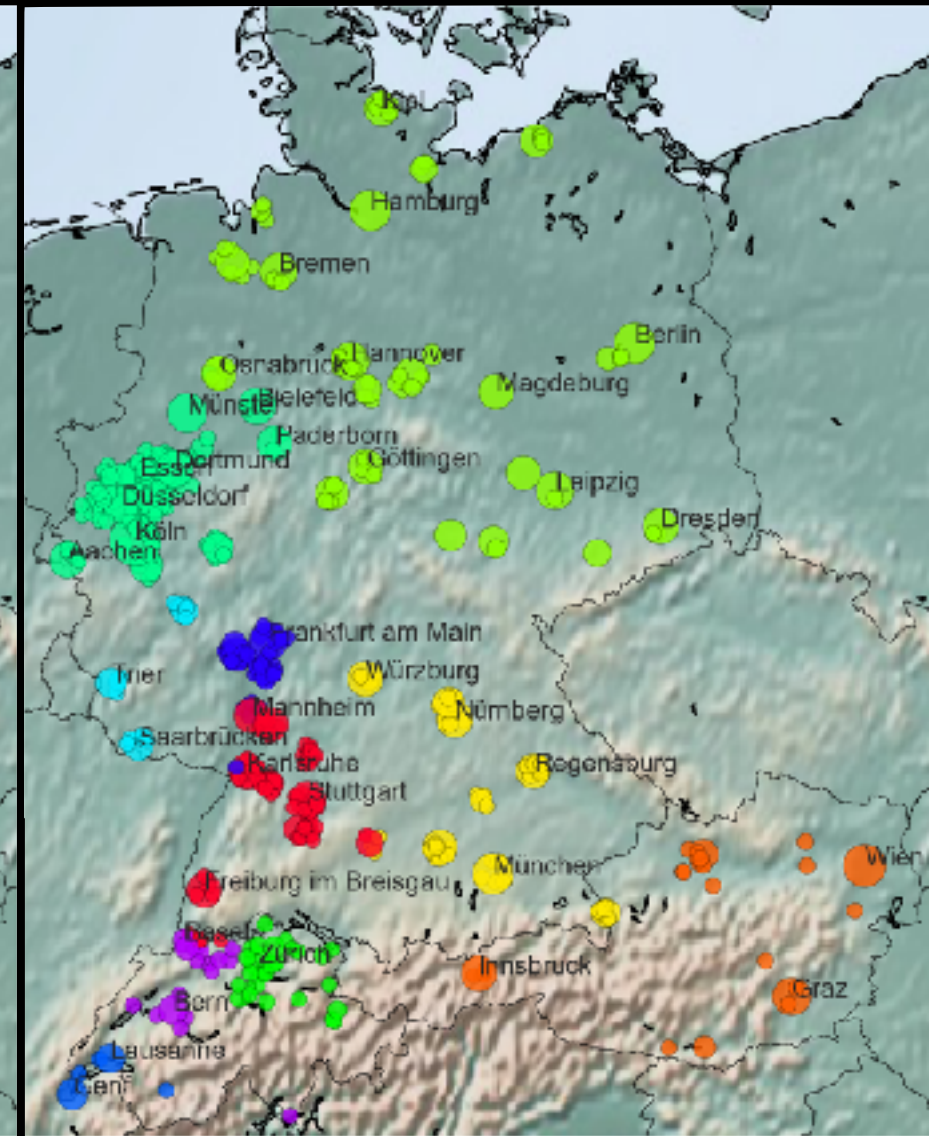
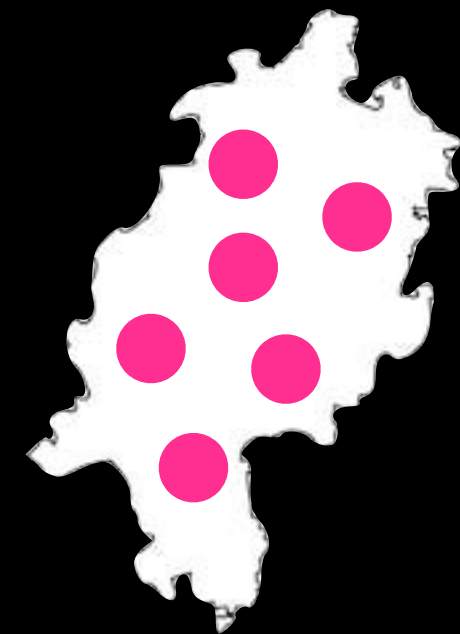# Dialect Clusters



3

5

10

# Evaluation Metrics

**Homogeneity**
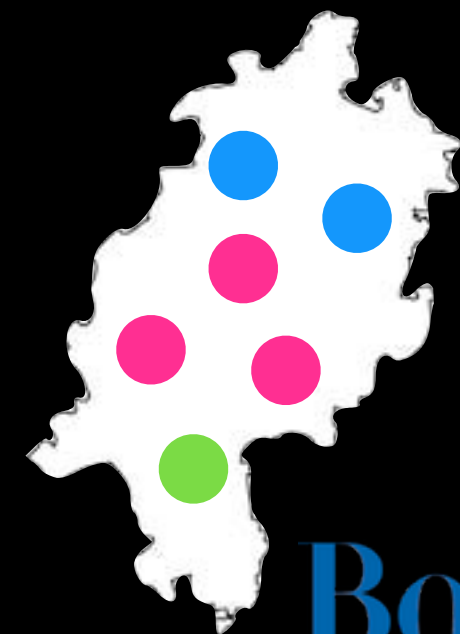cluster has only 1 gold label
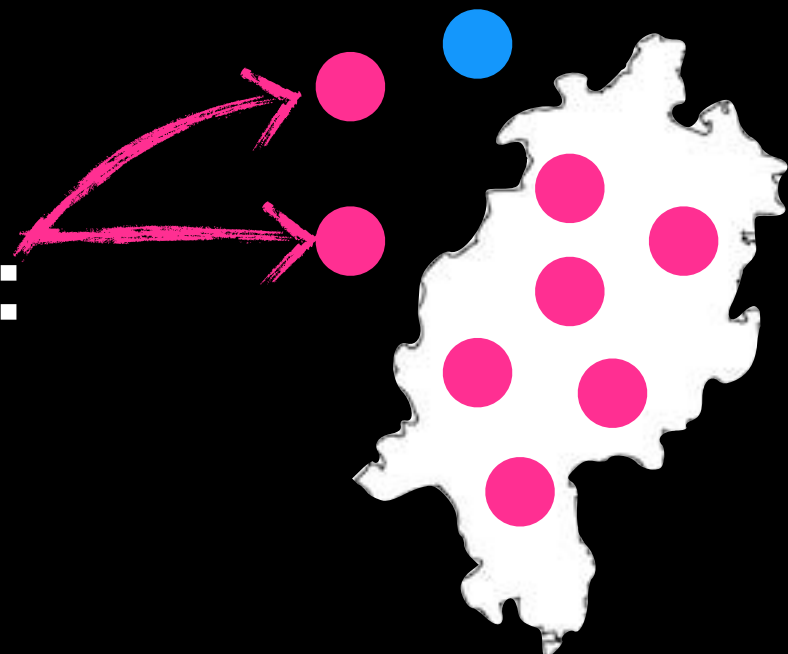
**Completeness**
gold label has only 1 cluster

**Good:**

**Bad:**

*GOLD LABEL (REGION)*

# Comparison

| | *k*-means | Agg |
|---|---|---|
| **scalable** | yes | no (up to ~20k) |
| **repeatable result** | no | yes |
| **include external info** | no | yes |
| **Good on dense clusters?** | no | yes |

# Wrapping Up

# When to Use What

| | Discrete Features | Embeddings |
|---|---|---|
| **Latent topics** | NMF | *Not applicable* |
| **RGB translation** | NMF | SVD + scaling |
| **Plotting** | SVD | t-SNE |
| **Clustering** | *Reduce dimensions* | *Use as-is* |

**Bocconi**

# 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

Bocconi