# Dimensionality Reduction

Principal Component Analysis (PCA)

t-Stochastic Neighbor Embedding (t-SNE)

Uniform Manifold Approximation and Projection (UMAP)

# Quick review

# Quick review

We've learned several useful methods already.
What sorts of things are we now good at?

# Quick review

We've learned several useful methods already. What sorts of things are we now good at?

Can I always use these? If not, what are the requirements in order to run these methods?

COSMIC DAWN CENTER

DAWN

# Quick review

We've learned several useful methods already. What sorts of things are we now good at?

Can I always use these? If not, what are the requirements in order to run these methods?
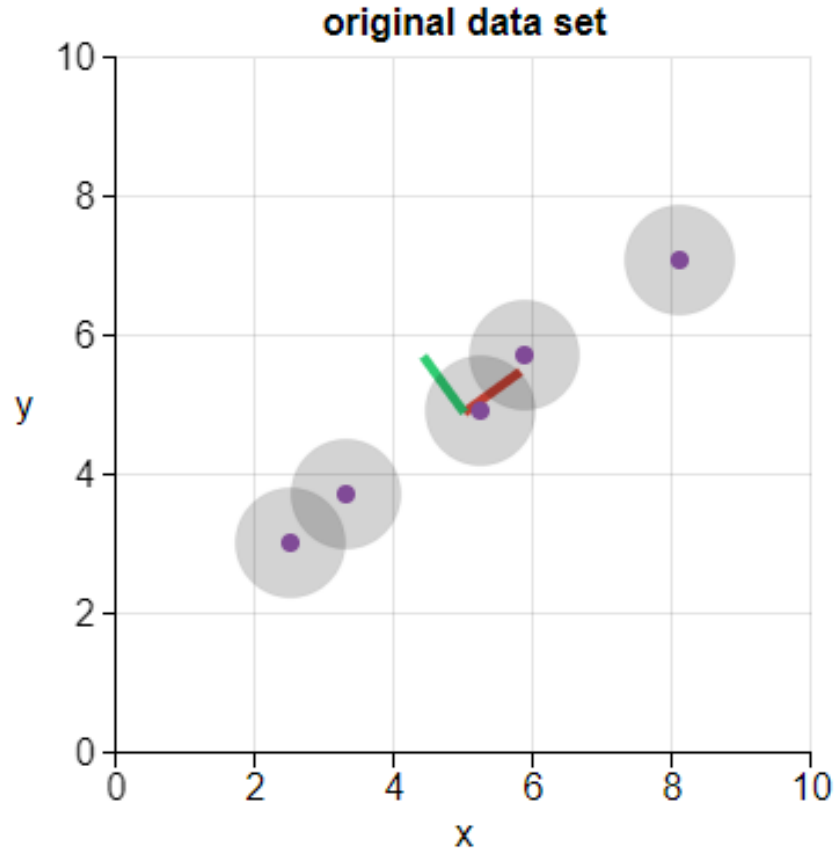
What should we do if we don't have labels?

COSMIC DAWN CENTER

DAWN

# What should we do if we don't have labels?

# Principal Component Analysis (PCA)

https://setosa.io/ev/principal-component-analysis/

# Principal Component Analysis (PCA)

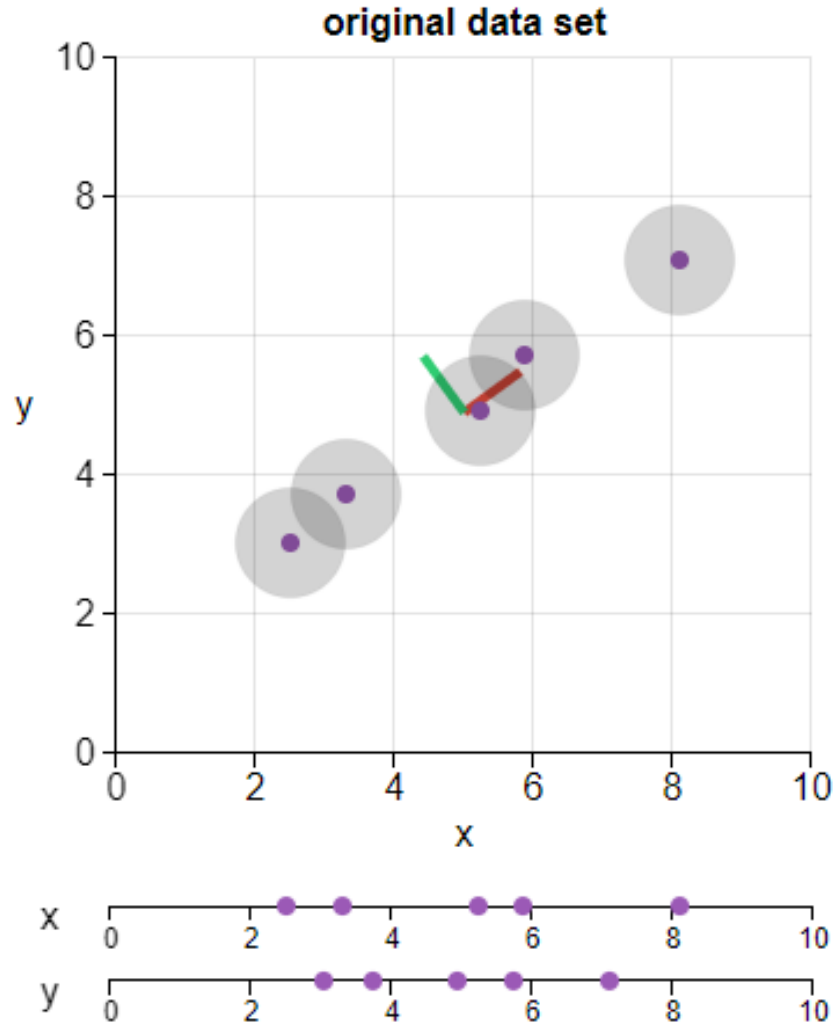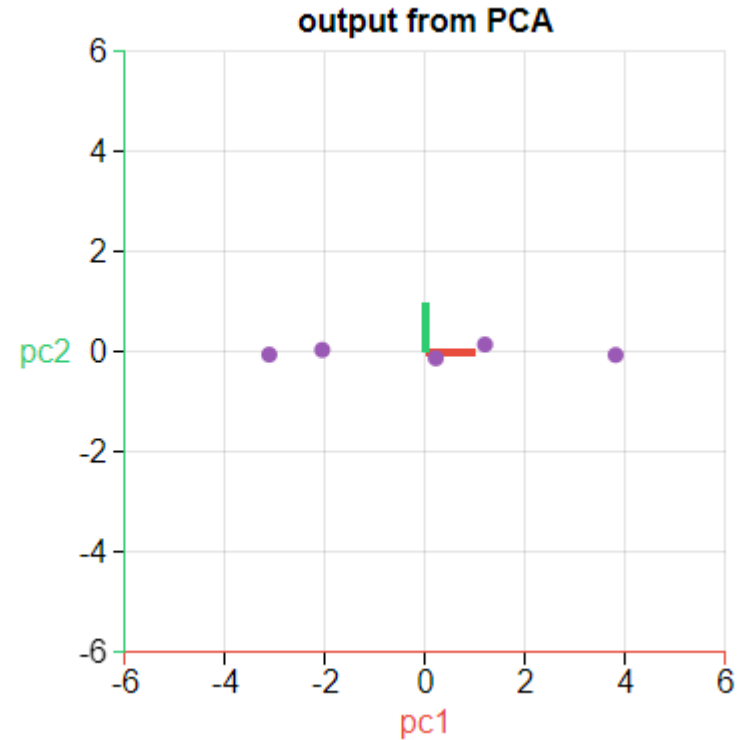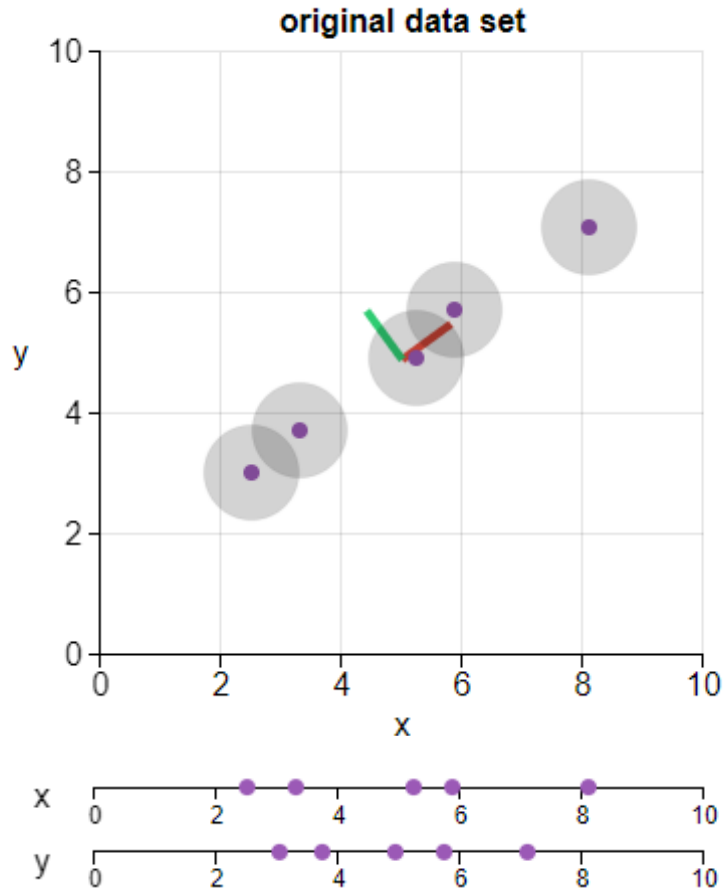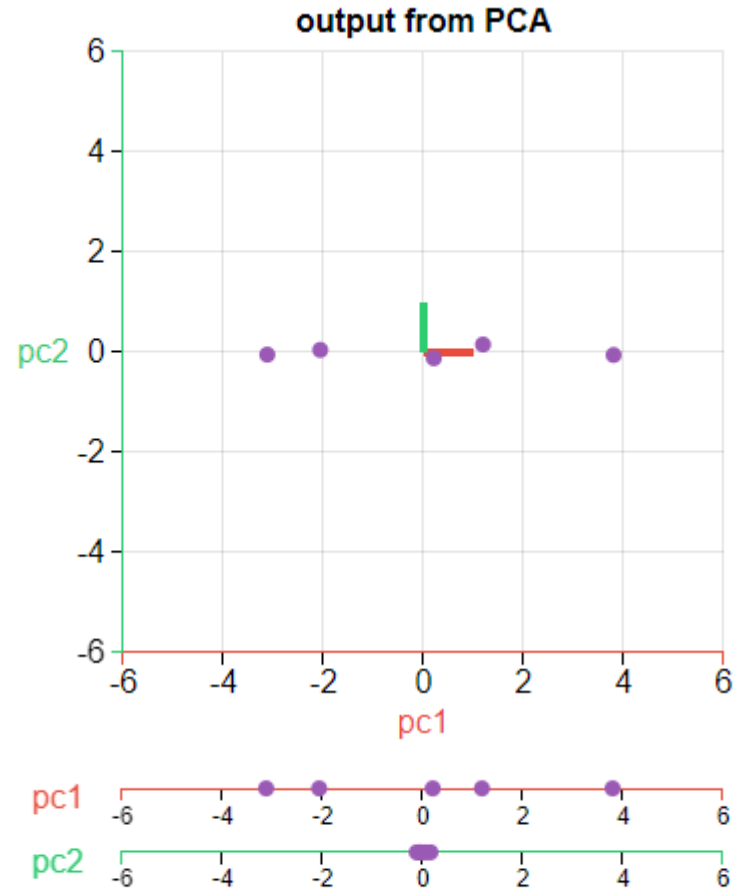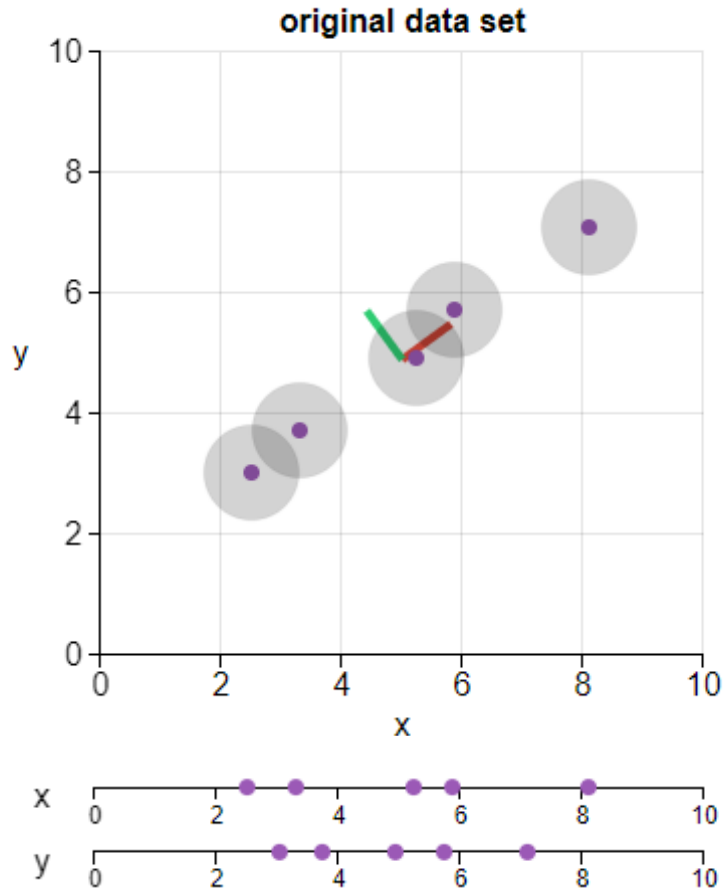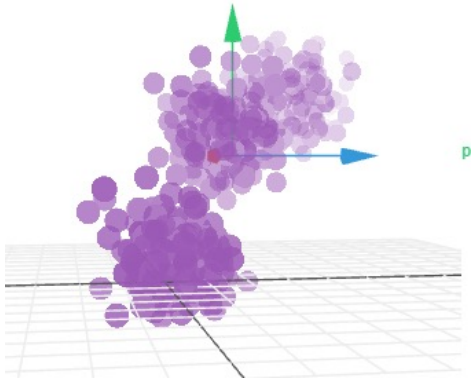https://setosa.io/ev/principal-component-analysis/

# Principal Component Analysis (PCA)

# Principal Component Analysis (PCA)

# PCA in 3D

COSMIC DAWN CENTER

DAWN

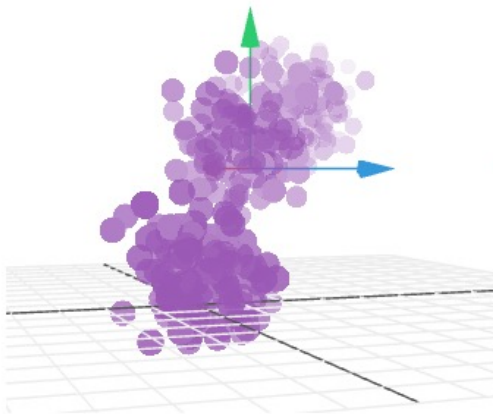Wednesday May 3, 2023

# PCA in 3D

# PCA in 3D

# PCA in 3D

COSMIC DAWN CENTER

DAWN

# PCA in 3D

COSMIC DAWN CENTER
DAWN

# PCA in 3D

COSMIC DAWN CENTER

DAWN

# PCA in 17D

| | England | N Ireland | Scotland | Wales |
|---|---|---|---|---|
| Alcoholic drinks | 375 | 135 | 458 | 475 |
| Beverages | 57 | 47 | 53 | 73 |
| Carcase meat | 245 | 267 | 242 | 227 |
| Cereals | 1472 | 1494 | 1462 | 1582 |
| Cheese | 105 | 66 | 103 | 103 |
| Confectionery | 54 | 41 | 62 | 64 |
| Fats and oils | 193 | 209 | 184 | 235 |
| Fish | 147 | 93 | 122 | 160 |
| Fresh fruit | 1102 | 674 | 957 | 1137 |
| Fresh potatoes | 720 | 1033 | 566 | 874 |
| Fresh Veg | 253 | 143 | 171 | 265 |
| Other meat | 685 | 586 | 750 | 803 |
| Other Veg | 488 | 355 | 418 | 570 |
| Processed potatoes | 198 | 187 | 220 | 203 |
| Processed Veg | 360 | 334 | 337 | 365 |
| Soft drinks | 1374 | 1506 | 1572 | 1256 |
| Sugars | 156 | 139 | 147 | 175 |

# PCA in 17D

https://setosa.io/ev/principal-component-analysis/

| | England | N Ireland | Scotland | Wales |
|---|---|---|---|---|
| Alcoholic drinks | 375 | 135 | 458 | 475 |
| Beverages | 57 | 47 | 53 | 73 |
| Carcase meat | 245 | 267 | 242 | 227 |
| Cereals | 1472 | 1494 | 1462 | 1582 |
| Cheese | 105 | 66 | 103 | 103 |
| Confectionery | 54 | 41 | 62 | 64 |
| Fats and oils | 193 | 209 | 184 | 235 |
| Fish | 147 | 93 | 122 | 160 |
| Fresh fruit | 1102 | 674 | 957 | 1137 |
| Fresh potatoes | 720 | 1033 | 566 | 874 |
| Fresh Veg | 253 | 143 | 171 | 265 |
| Other meat | 685 | 586 | 750 | 803 |
| Other Veg | 488 | 355 | 418 | 570 |
| Processed potatoes | 198 | 187 | 220 | 203 |
| Processed Veg | 360 | 334 | 337 | 365 |
| Soft drinks | 1374 | 1506 | 1572 | 1256 |
| Sugars | 156 | 139 | 147 | 175 |

COSMIC DAWN CENTER
DAWN

# OK, so how can we find the right basis?

1.  Standardization

# OK, so how can we find the right basis?

1.  Standardization

2.  Compute covariance matrix

$$
\begin{bmatrix}
Cov(x,x) & Cov(x,y) & Cov(x,z) \\
Cov(y,x) & Cov(y,y) & Cov(y,z) \\
Cov(z,x) & Cov(z,y) & Cov(z,z)
\end{bmatrix}
$$

# OK, so how can we find the right basis?

1. Standardization

2. Compute covariance matrix

$$
\left[
\begin{array}{ccc}
Cov(x,x) & Cov(x,y) & Cov(x,z) \\
Cov(y,x) & Cov(y,y) & Cov(y,z) \\
Cov(z,x) & Cov(z,y) & Cov(z,z)
\end{array}
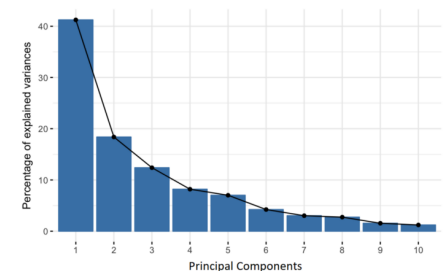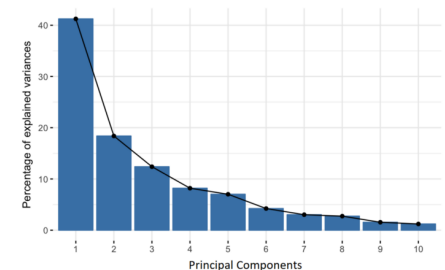\right]
$$

3. Compute eigenvectors and eigenvalues

# OK, so how can we find the right basis?

1. Standardization

2. Compute covariance matrix

$$\begin{bmatrix} Cov(x,x) & Cov(x,y) & Cov(x,z) \\ Cov(y,x) & Cov(y,y) & Cov(y,z) \\ Cov(z,x) & Cov(z,y) & Cov(z,z) \end{bmatrix}$$

3. Compute eigenvectors and eigenvalues
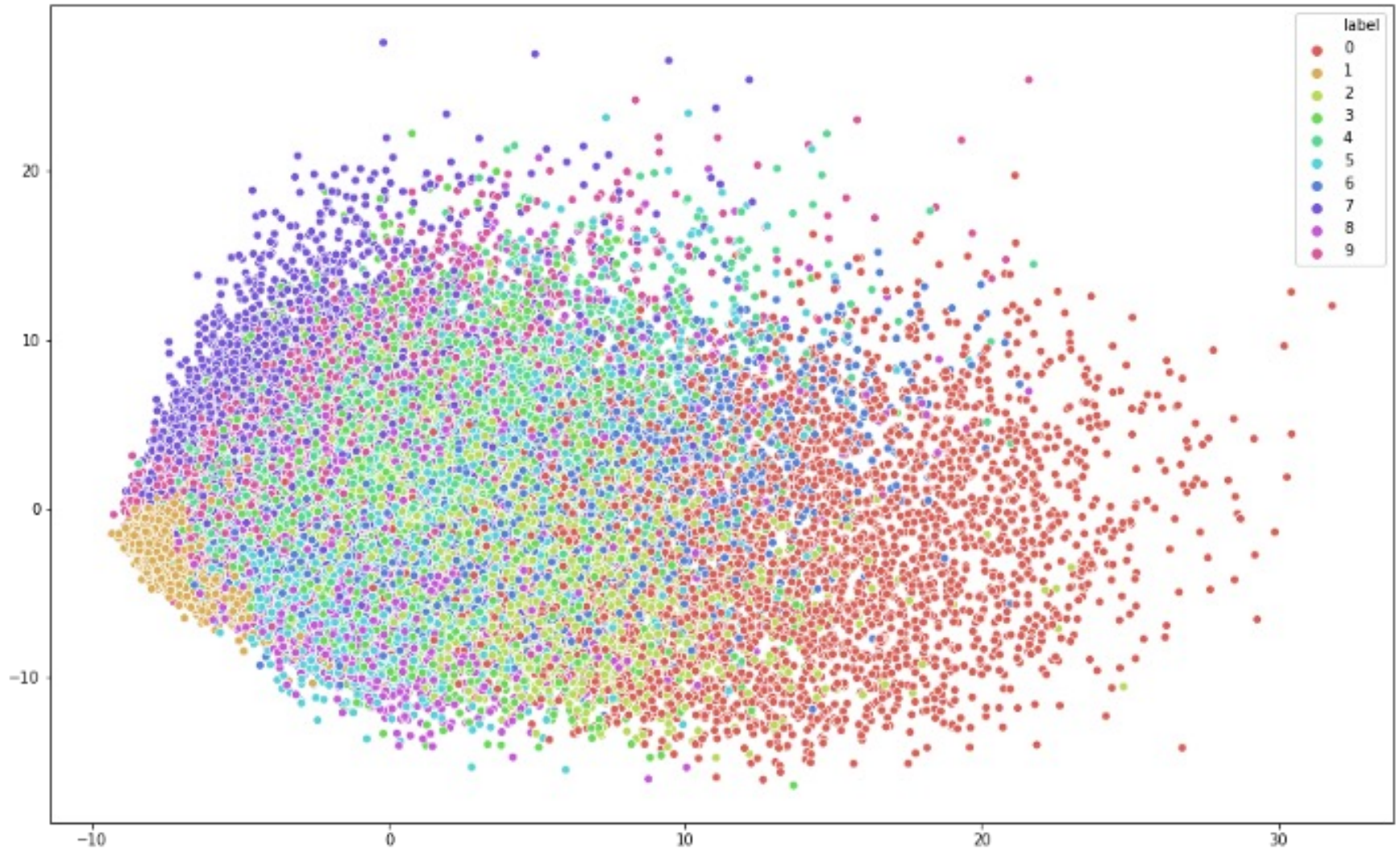
4. Discard vectors that are not important enough

COSMIC DAWN CENTER

DAWN

# Example: Handwritten Digits

MNIST dataset

# Group "similar" things together

Principal Component Analysis

# Example: Handwritten Digits

MNIST dataset

# Some things aren't linear!

Wikimedia Commons



| Method | Variance unexplaine |
|--------|---------------------|
| PCA | 23.23% |
| SOM | 6.86% |

COSMIC DAWN CENTER
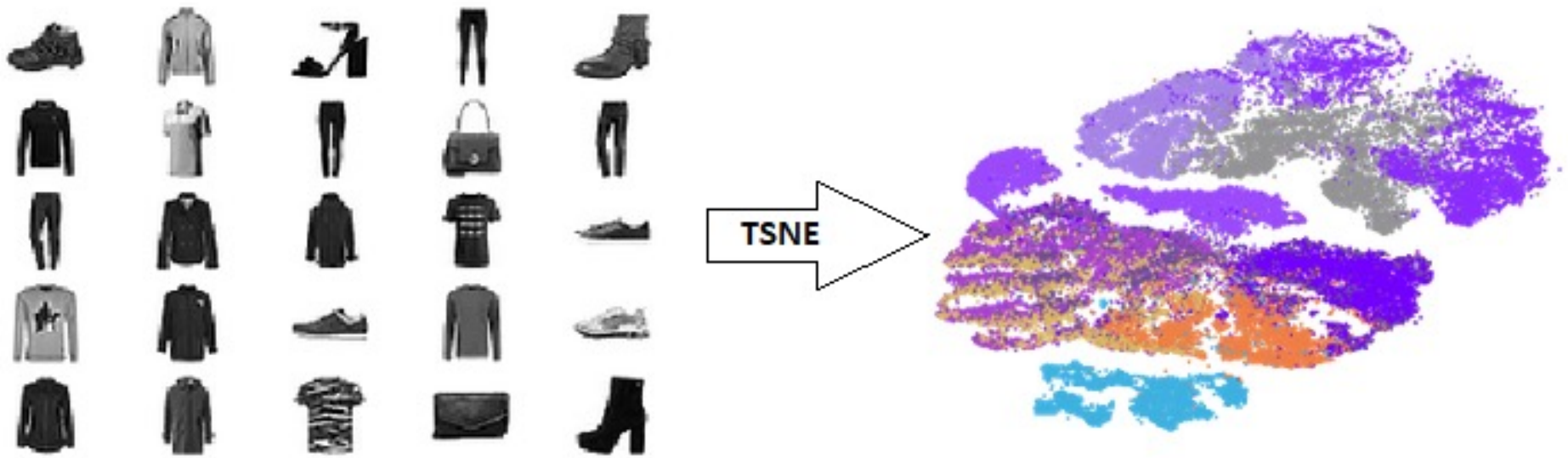DAWN

# Group "similar" things together

# Group "similar" things together

"Fashion MNIST" datasets, t-SNE

# Group "similar" things together

Wang et al. 2020

# Example: Separate Short and Long GRBs

R. Mallozzi, updated Aug 2018 at https://gammaray.msfc.nasa.gov/batse/grb/duration/



BATSE 4B Catalog

# t-SNE map for *Swift* light curves

Kragh Jespersen et al. 2020

# t-SNE map for *Swift* light curves

Kraah Jespersen et al. 2020

COSMIC DAWN CENTER

DAWN

# t-SNE map for *Swift* light curves

Kragh Jespersen et al. 2020

# Embeddings, colored by duration

# Embeddings, colored by duration

# Hardness distribution

Kragh Jespersen et al. 2020

COSMIC DAWN CENTER

DAWN

# Redshift distribution

Kragh Jespersen et al. 2020

# Possible subgroupings?

Kragh Jespersen et al. 2020



**Legend:**
- ▼ SNs
- ▼ Unclear SNs
- ◆ Kns
- ◆ Unclear Kns

# Structure is durable, not location!

Steinhardt, Mann, Rusakov, and Kragh Jespersen 2023

# Objects can "jump"

Steinhardt, Mann, Rusakov, and Kragh Jespersen 2023

# Example: Photometry

V  i  z  J  H

Bouwens et al. 2006

z~7.4

z~6.8

z~6.8

z~6.8

no detection        detection

COSMIC DAWN

DAW

3, 2023

# Example: Galaxy spectra

SDSS/Galaxy Zoo



Survey: *sdss* Program: *legacy* Target: *GALAXY ROSAT_D ROSAT_E*
RA=25.65806, Dec=−1.22998, Plate=401, Fiber=125, MJD=51788
*z*=0.04263±0.00002 Class=GALAXY AGN
No warnings.

Wednesday May 3, 2023

# Spectra have similar features!

Zaroubi et al. 2013

# Maybe limited information is enough

Zaroubi et al. 2013

# Our goal:



| | |
|---|---|
| Filters | |
| V i z J H | |

z~7.4
z~6.8
z~6.8
z~6.8

no detection          detection

→

"V band"
"i band"
"z band"
"J band"
"H band"

etc.

→

(Stellar) Mass
Star formation rate
Star formation history
Distance/Redshift
Age

etc.

COSMIC DAWN CENTER
DAWN

# Two Fundamental Assumptions for Photometry

# Two Fundamental Assumptions for Photometry

1.  If an object is sufficiently well-measured, there is a surjective (one-to-one or many-to-one, but not one-to-many) mapping from photometric fluxes to astrophysical properties

# Our goal:



"V band"
"i band"
"z band"
"J band"
"H band"

etc.

(Stellar) Mass
Star formation rate
Star formation history
Distance/Redshift
Age

etc.

COSMIC DAWN CENTER
DAWN

# Two Fundamental Assumptions for Photometry

1. If an object is sufficiently well-measured, there is a surjective (one-to-one or many-to-one, but not one-to-many) mapping from photometric fluxes to astrophysical properties.

2. Objects with sufficiently similar photometry should be mapped to similar astrophysical properties.

# Our goal:



”V band”
”i band”
”z band”
”J band”
”H band”

etc.

(Stellar) Mass
Star formation rate
Star formation history
Distance/Redshift
Age

etc.

COSMIC DAWN CENTER
DAWN

# Our goal:



"V band"
"i band"
"z band"
"J band"
"H band"

etc.

Choose a few good
examples

(Stellar) Mass
Star formation rate
Star formation history
Distance/Redshift
Age

etc.

COSMIC DAWN CENTER
DAWN

# Approach 1: Color space map using two colors (three bands)

COSMOS2015 catalog, objects at z≈1

# Approach 1: Color space map using two colors (three bands)

COSMOS2015 catalog, objects at z≈1

# Approach 2: Color space map using all bands

# ~~Two~~ Three Fundamental Assumptions for Photometry

1. If an object is sufficiently well-measured, there is a surjective (one-to-one or many-to-one, but not one-to-many) mapping from photometric fluxes to astrophysical properties.

2. Objects with sufficiently similar photometry should be mapped to similar astrophysical properties.

3. We can map objects from the full, n-dimensional space with all bands to a smaller one with many neighbors, and the other two assumptions will continue to hold.
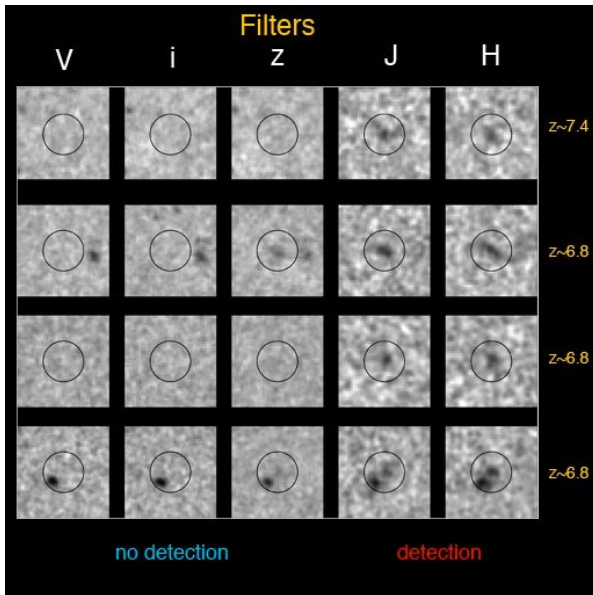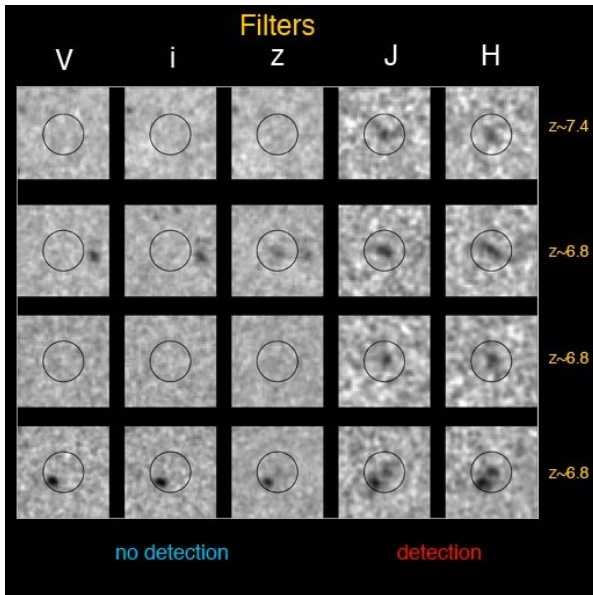
# ~~Two~~ Three Fundamental Assumptions for Photometry

1. If an object is sufficiently well-measured, there is a surjective (one-to-one or many-to-one, but not one-to-many) mapping from photometric fluxes to astrophysical properties.

2. Objects with sufficiently similar photometry should be mapped to similar astrophysical properties.

3. We can map objects from the full, n-dimensional space with all bands to a smaller one with many neighbors, and the other two assumptions will continue to hold.

**Can we somehow decide what information is "important" even without labels?**

Wednesday May 3, 2023

# Approach 2: First, make a t-SNE map reducing to two dimensions

COSMOS2015 catalog, objects at z≈1

# Approach 2: Similar galaxies are nearby

COSMOS2015 catalog, objects at z≈1

COSMIC DAWN CENTER
DAWN

# Approach 2: Similar galaxies are nearby

**WARNING: positions are neither fixed nor meaningful. Topology is meaningful.**

# Approach 2: Arranging by photometry also calculates other useful things!

COSMOS2015 catalog, objects at z≈1



$Log_{10}(M_*)$

9.5   10.0   10.5   11.0   11.5

$S/N_{MIPS} < 5$
$S/N_{MIPS} > 5$

COSMIC DAWN CENTER
DAWN

# projector.tensorflow.org

**Embedding Projector**                                                    ⑦ 🐛

## DATA

5 tensors found
Word2Vec 10K                                    ▼

**Label by**          **Color by**
word         ▼        No color map          ▼

**Edit by**
word         ▼        Tag selection as

[ Load ]  [ Publish ]  [ Download ]  [ Label ]

☑ Sphereize data  ❓

Checkpoint: Demo datasets

Metadata:    oss_data/word2vec_10000_200d_
             labels.tsv

---

**UMAP**   T-SNE   PCA   CUSTOM

Dimension          2D  ⬤  3D

Neighbors ❓    ●————————  15

[ Run ]

For faster results, the data will be sampled
down to 5,000 points.

📄 Learn more about UMAP.

---

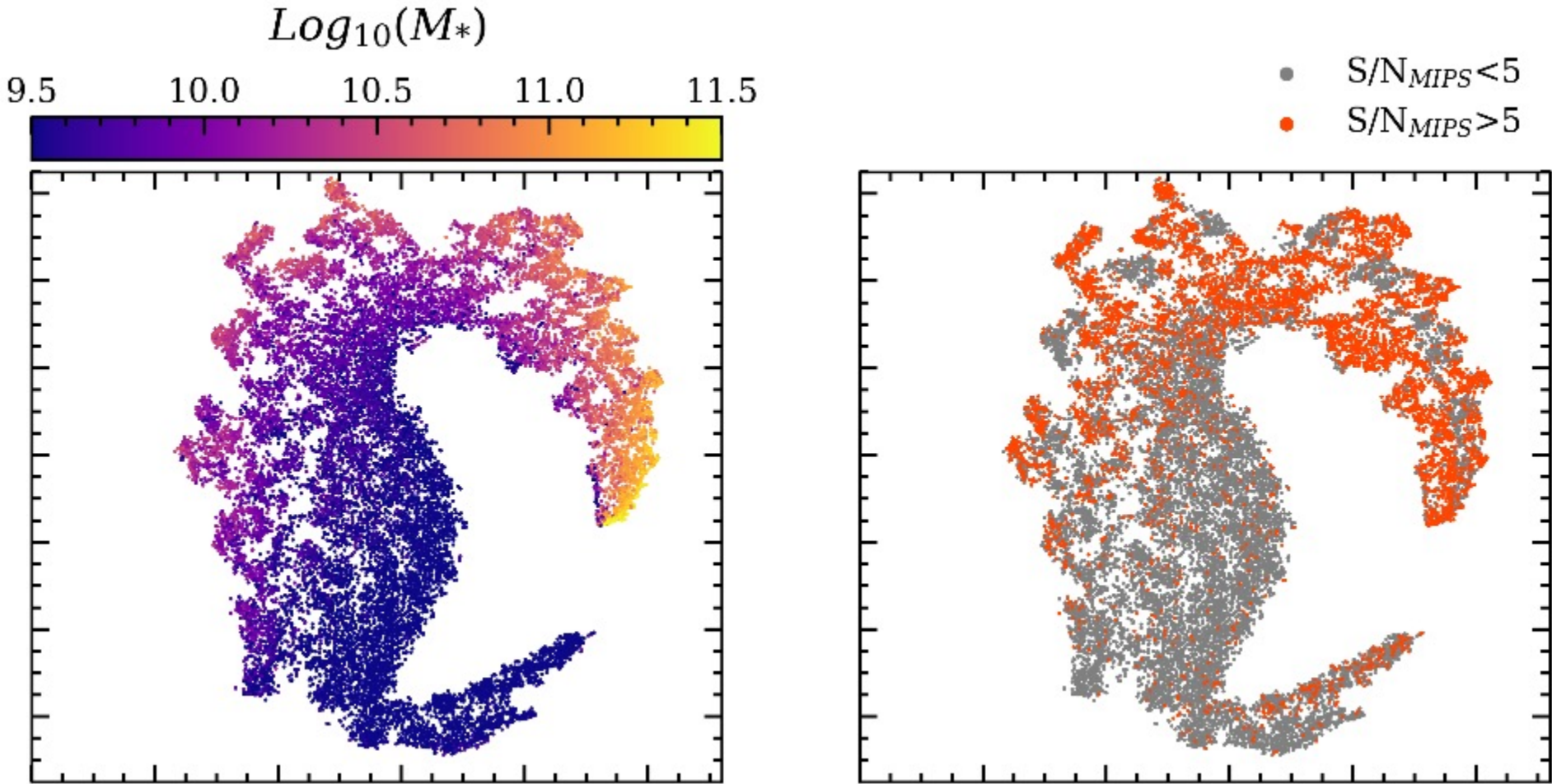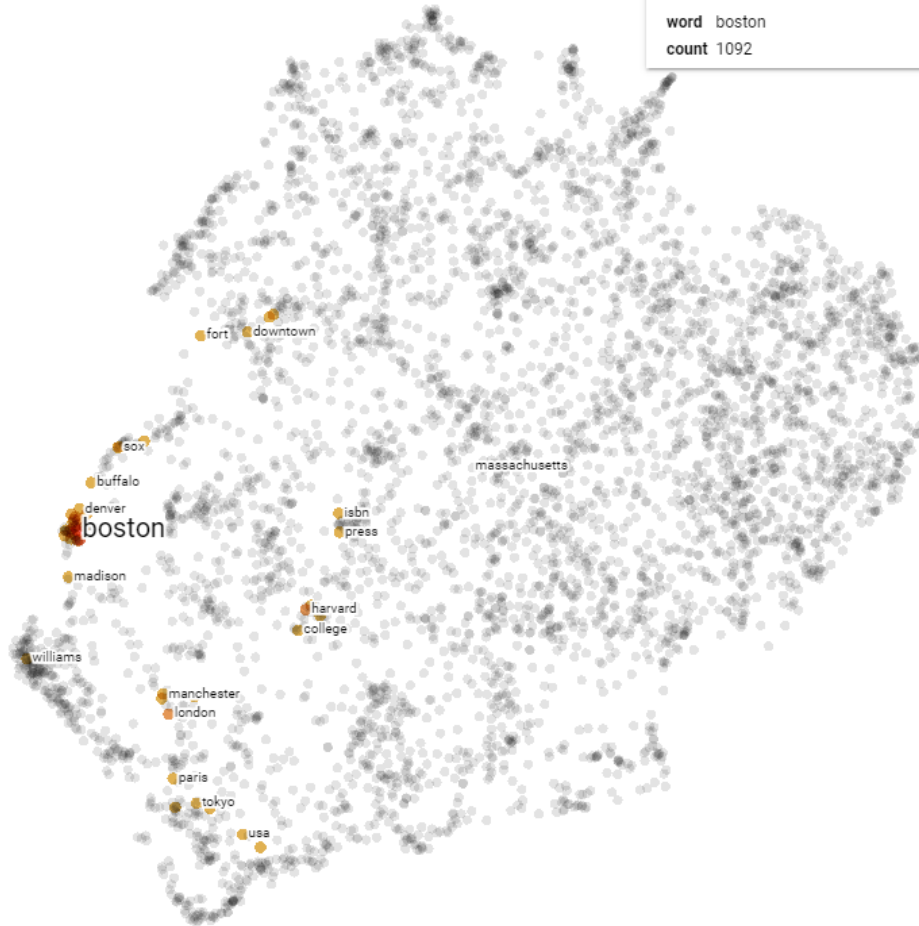▫ 🌓 ● 🅰 | Points: 10000 | Dimension: 200 | Selected 101 points

[ Show All Data ]  [ Isolate 101 points ]  [ Clear selection ]

⌂

❓

boston                                          ⌃
  **word**   boston
  **count**  1092

fort  downtown

sox
buffalo
denver
**boston**            massachusetts
madison

              isbn
              press

                 harvard
                 college

williams

     manchester
     london

   paris
   tokyo
       usa

Search          [ .* ]        by  word  ▼

neighbors ❓  ●————————  100

distance        COSINE  EUCLIDEAN

Nearest points in the original space:

| | |
|---|---|
| chicago | 0.406 |
| massachusetts | 0.406 |
| philadelphia | 0.413 |
| atlanta | 0.497 |
| harvard | 0.502 |
| london | 0.508 |
| illinois | 0.512 |
| baltimore | 0.514 |
| maryland | 0.546 |
| york | 0.551 |
| cincinnati | 0.554 |
| toronto | 0.556 |
| seattle | 0.569 |
| brooklyn | 0.578 |
| miami | 0.581 |
| cambridge | 0.583 |
| pennsylvania | 0.584 |
| pittsburgh | 0.587 |
| detroit | 0.593 |
| california | 0.595 |
| sox | 0.607 |
| kansas | 0.608 |

BOOKMARKS (0) ❓