**Carnegie Mellon University**

**HeinzCollege**

# 95-865 Unstructured Data Analytics

Recitation: More on PCA & manifold learning

Slides by George H. Chen

# 2D PCA Plots

Demo

# t-SNE Interpretation

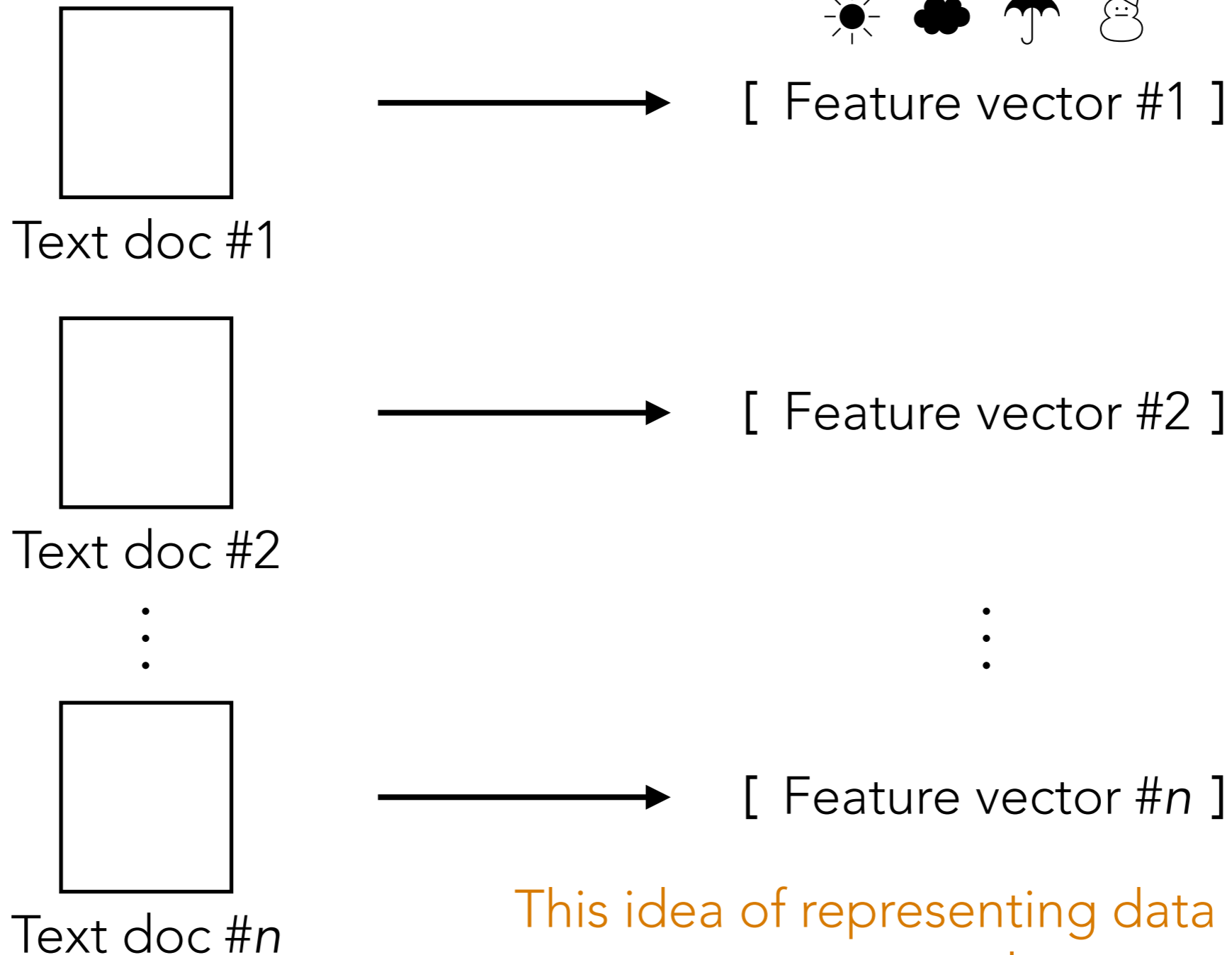https://distill.pub/2016/misread-tsne/

# Dimensionality Reduction for Visualization

- There are *many* methods (I've posted a link on the course webpage to a scikit-learn example using ~10 methods)

- PCA is very well-understood; the new axes can be interpreted

- Nonlinear dimensionality reduction (manifold learning): new axes may not really be all that interpretable

- PCA is good to try first (look at plot & explained variance ratios)
  - If PCA works poorly, then t-SNE could be a good 2nd thing to try

- If you have good reason to believe that only certain features matter, of course you could restrict your analysis to those!

- t-SNE can be annoying to use but is still very popular
  - Promising recently developed alternative: PaCMAP (Wang et al 2021) accounts for local and global structure simultaneously and also uses "mid-near" neighbors of points — link on course webpage

# Let's look at images
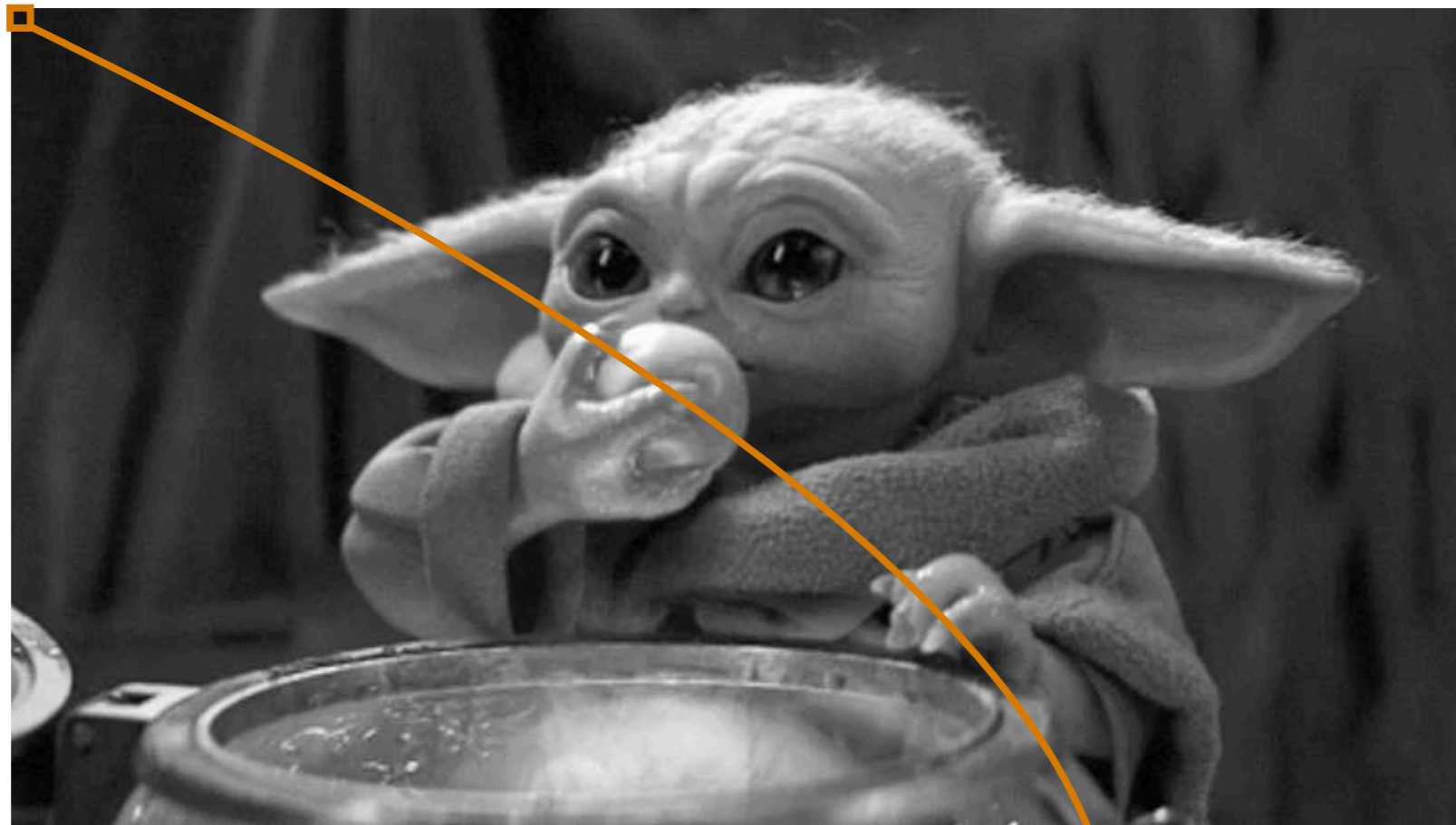
# (Flashback) Multiple Documents

Choose a common vocabulary to use across all documents

[ Feature vector #1 ]

Text doc #1

[ Feature vector #2 ]

Text doc #2

⋮

⋮

[ Feature vector #n ]

Text doc #n

This idea of representing data as feature vectors is very general — not just for text!

# Example: Representing an Image
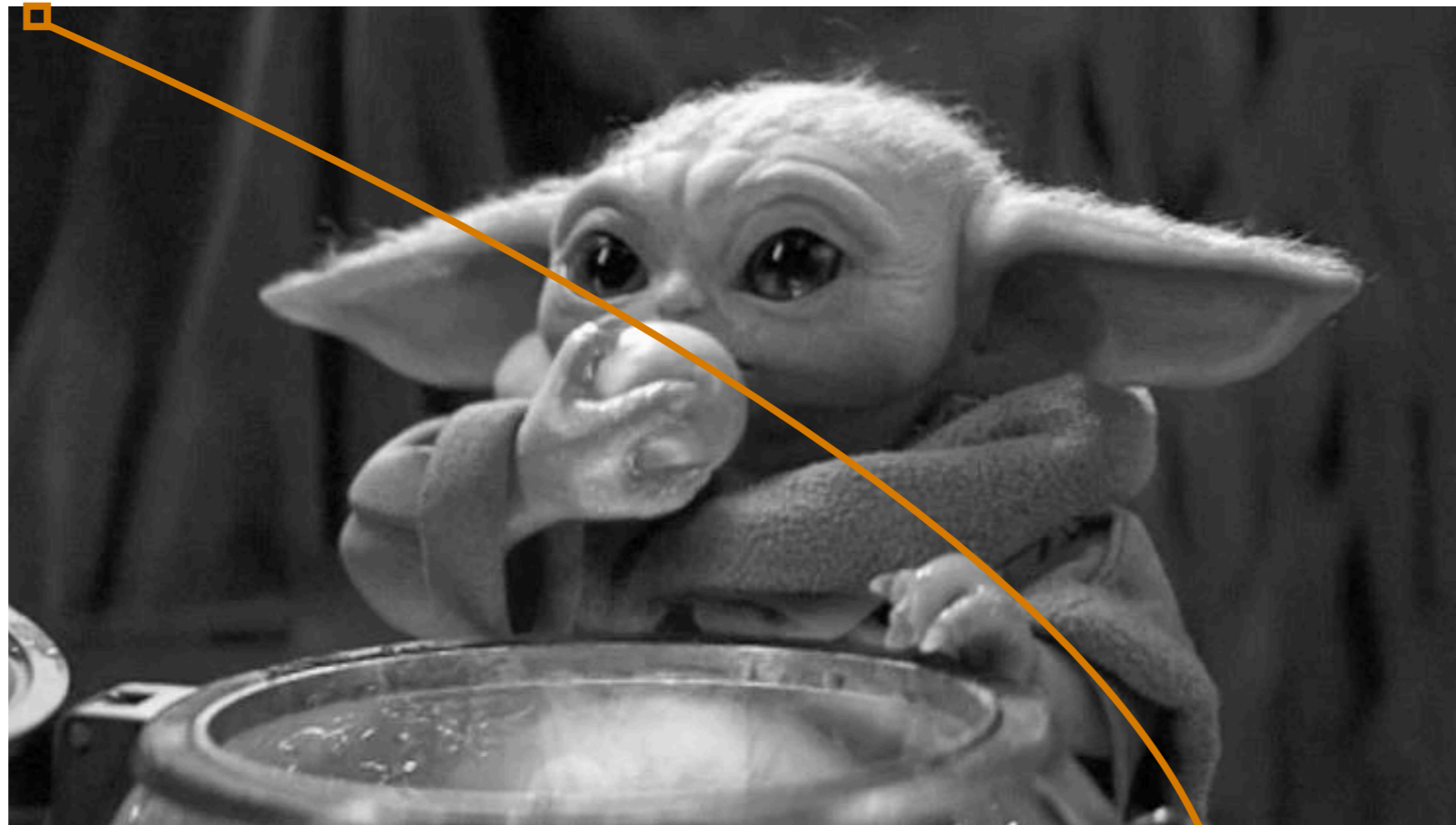
Go row by row and look at pixel values



1: black
0: white

[ 1                    ]

*Image source: The Mandalorian*

# Example: Representing an Image

Go row by row and look at pixel values



1: black
0: white

[ 1   0.9                          ]

*Image source: The Mandalorian*

# Example: Representing an Image

Go row by row and look at pixel values



1: black
0: white

[ 1   0.9   $\cdots$   0.1          ]

*Image source: The Mandalorian*

# Example: Representing an Image

Go row by row and look at pixel values



1: black
0: white

[ 1   0.9  ⋯  0.1  ⋯  0.9 ]

# dimensions = image width × image height

*Image source: The Mandalorian*   Very high dimensional!

# Terminology Remark

$$[ \ 1 \quad 0.9 \quad \cdots \quad 0.1 \quad \cdots \quad 0.9 \ ]$$

⚠️ We use "dimension" to means two different things:

- number of axes we can index into for a table/array (e.g., 2D means there are rows & columns)

# dimensions = 1

- total number of entries in the table/array

# dimensions = image height × image width

# Dimensionality Reduction for Images

Demo