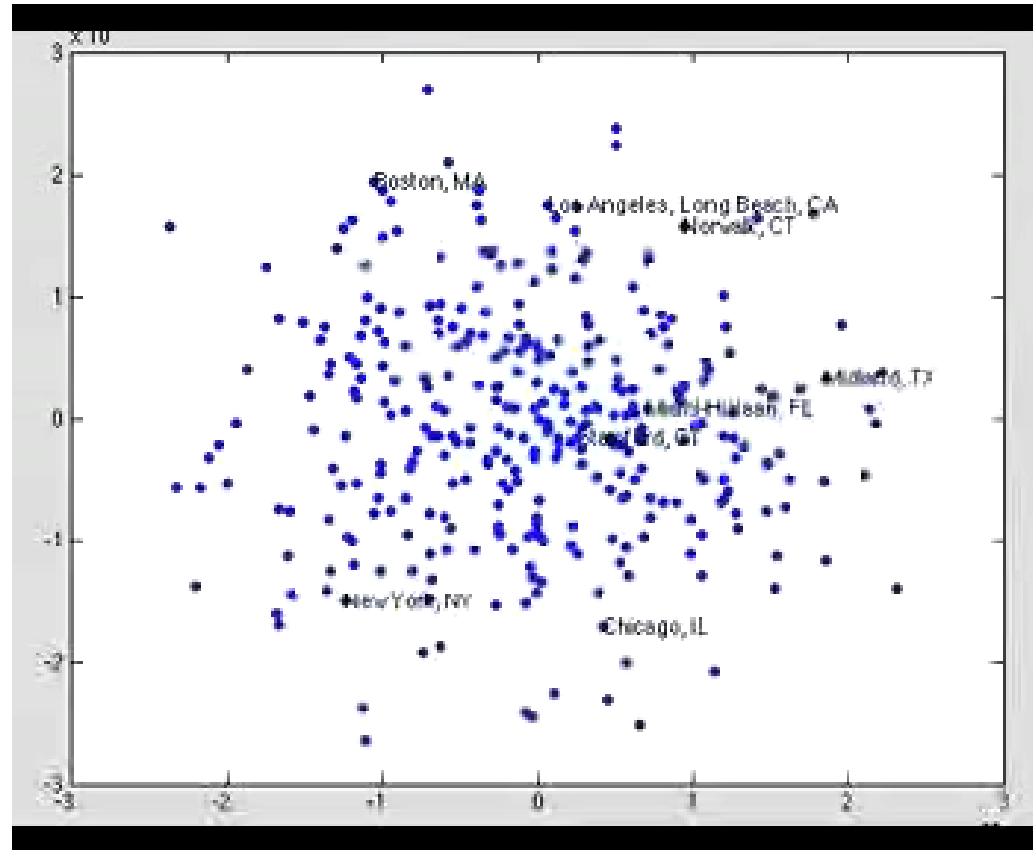# CPSC 340:
# Machine Learning and Data Mining

Multi-Dimensional Scaling

Summer 2021

# Admin

- Assignment 6 out, due Friday 11:55pm
- Today is final exam coverage cut-off
- Final exam is next Wednesday (June 23)
  - Prep materials go up soon

- Course evaluation is open.
  - Please give me an honest feedback! How did I do?

# Last Time: Multi-Dimensional Scaling



$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \| z_i - z_j \| - \| x_i - x_j \| \right)^2$$

"Try to make scatterplot distances match high-dimensional distance"

sum over pairs of examples

distance in scatterplot

Distance between points in original 'd' dimensions

# Multi-Dimensional Scaling

- **Multi-dimensional scaling (MDS):**
  - Optimize the final locations of the $z_i$ values.

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \|z_i - z_j\| - \|x_i - x_j\| \right)^2$$
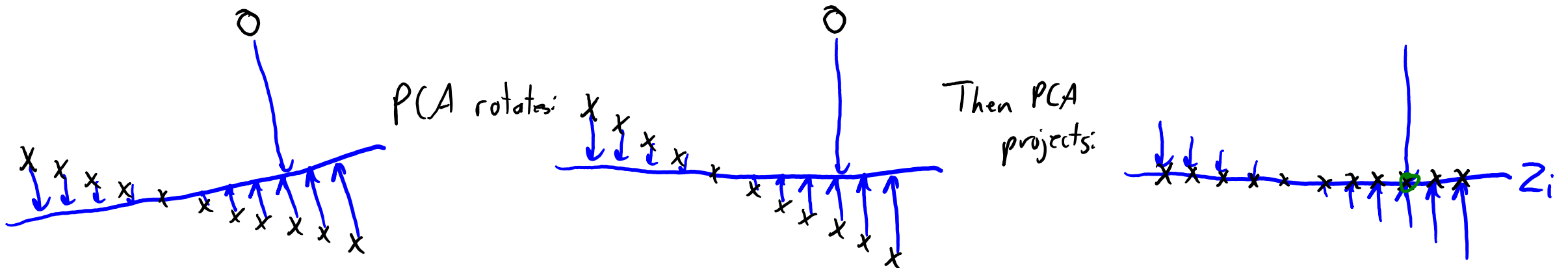
# Multi-Dimensional Scaling

- **Multi-dimensional scaling (MDS)**:
  - Optimize the final locations of the $z_i$ values.

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \|z_i - z_j\| - \|x_i - x_j\| \right)^2$$

  - Non-parametric dimensionality reduction and visualization:
    - No 'W': just trying to make $z_i$ preserve high-dimensional distances between $x_i$.



PCA rotates:    Then PCA projects:    $z_i$

# Multi-Dimensional Scaling

- **Multi-dimensional scaling (MDS):**
  - Optimize the final locations of the $z_i$ values.

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \|z_i - z_j\| - \|x_i - x_j\| \right)^2$$

  - Non-parametric dimensionality reduction and visualization:
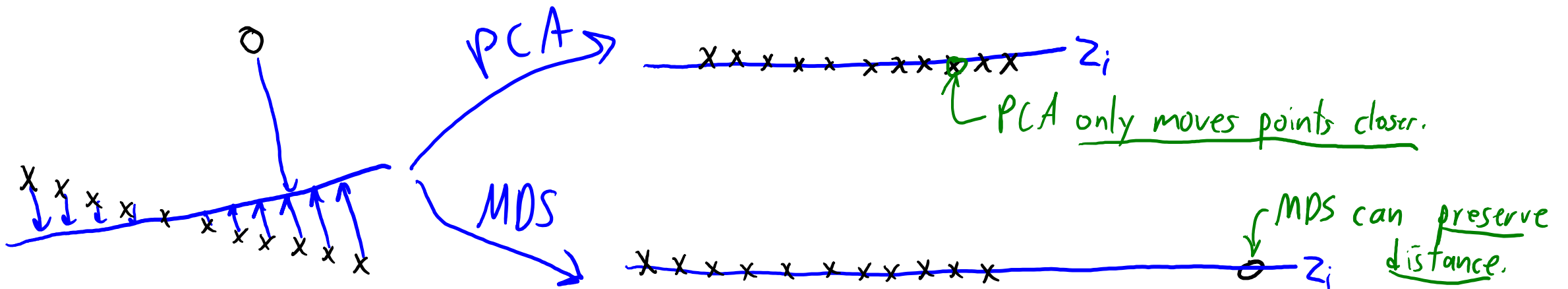    - No 'W': just trying to make $z_i$ preserve high-dimensional distances between $x_i$.



PCA

$z_i$

PCA only moves points closer.

MDS

MDS can preserve distance.

$z_i$

# Multi-Dimensional Scaling

- **Multi-dimensional scaling (MDS):**
  - Optimize the final locations of the $z_i$ values.

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \|z_i - z_j\| - \|x_i - x_j\| \right)^2$$

  - Non-parametric dimensionality reduction and visualization:
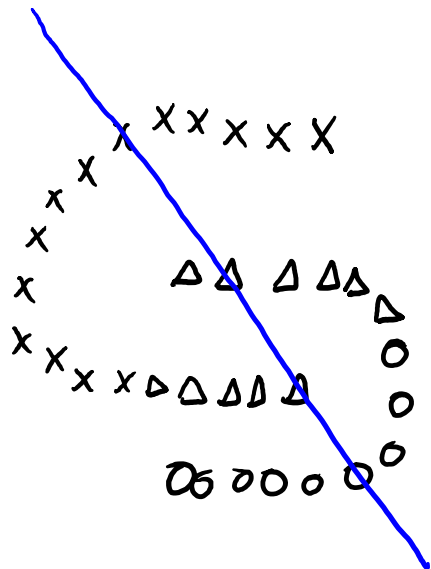    - No 'W': just trying to make $z_i$ preserve high-dimensional distances between $x_i$.
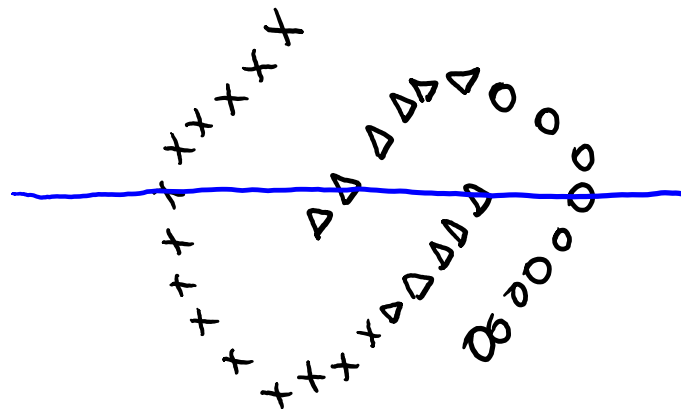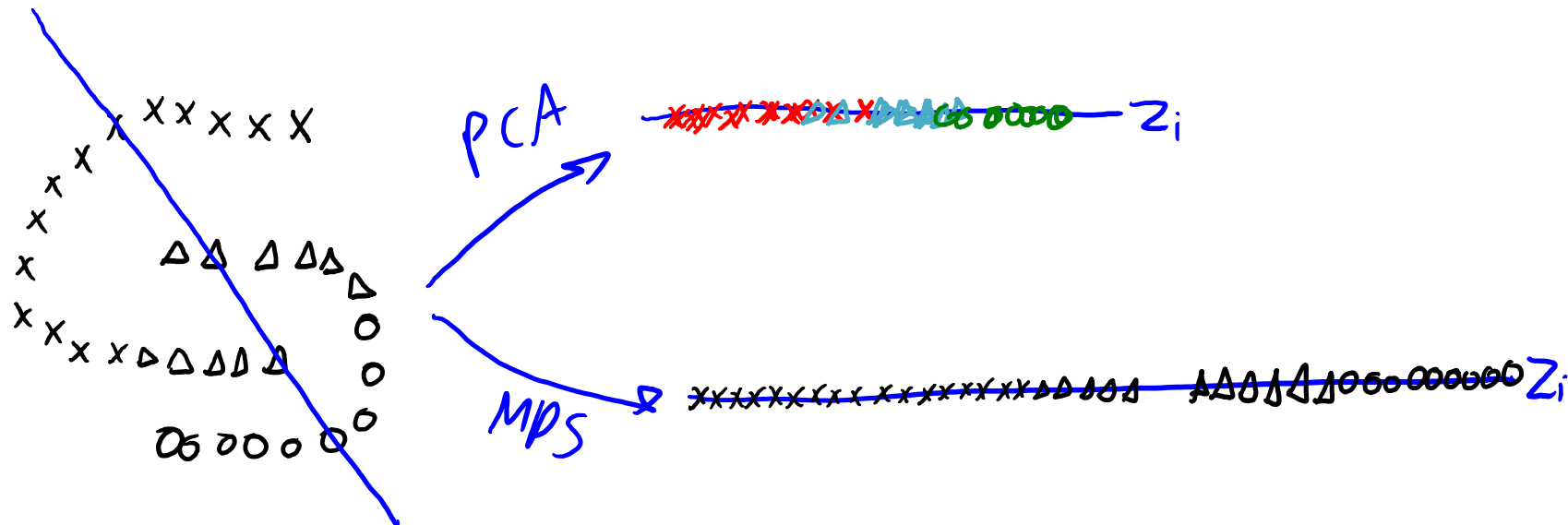
# Multi-Dimensional Scaling

- **Multi-dimensional scaling (MDS):**
  - Optimize the final locations of the $z_i$ values.

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \|z_i - z_j\| - \|x_i - x_j\| \right)^2$$

  - Non-parametric dimensionality reduction and visualization:
    - No 'W': just trying to make $z_i$ preserve high-dimensional distances between $x_i$.

# Multi-Dimensional Scaling

- **Multi-dimensional scaling (MDS):**
  - Optimize the final locations of the $z_i$ values.

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \|z_i - z_j\| - \|x_i - x_j\| \right)^2$$

- **Cannot use SVD** to compute solution:
  - Instead, do gradient descent on the $z_i$ values.
  - You "learn" a scatterplot that tries to visualize high-dimensional data.
  - Not convex and sensitive to initialization.
    - And solution is not unique due to various factors like translation and rotation.

# In This Lecture

1. Multi-Dimensional Scaling
   - Euclidean MDS
   - Sammon Mapping
   - Geodesic MDS (ISOMAP)

2. Latent Factors for Language

Coming Up Next

# EUCLIDEAN MDS VARIANTS

# Different MDS Cost Functions

- **MDS** default objective: squared difference of Euclidean norms:

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \|z_i - z_j\| - \|x_i - x_j\| \right)^2$$

> **Q**: How many distance functions are involved here?

> **Q**: Can we generalize this to other measures of distance?

# Different MDS Cost Functions

- **MDS** default objective function with general distances/similarities:

$$f(Z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} d_3(d_2(z_i, z_j) - d_1(x_i, x_j))$$

  – Functions are not necessarily the same:
    - $d_1$ := high-dimensional distance we want to match.

$$d_1 : \mathbb{R}^d \times \mathbb{R}^d \longrightarrow \mathbb{R}$$

    - $d_2$ := low-dimensional distance we can control.

$$d_2 : \mathbb{R}^2 \times \mathbb{R}^2 \longrightarrow \mathbb{R}$$

    - $d_3$ := how we compare high-/low-dimensional distances.

$$d_3 : \mathbb{R} \times \mathbb{R} \longrightarrow \mathbb{R}$$

# Different MDS Cost Functions

- **MDS** default objective function with general distances/similarities:

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} d_3\left(d_2(z_i, z_j) - d_1(x_i, x_j)\right)$$
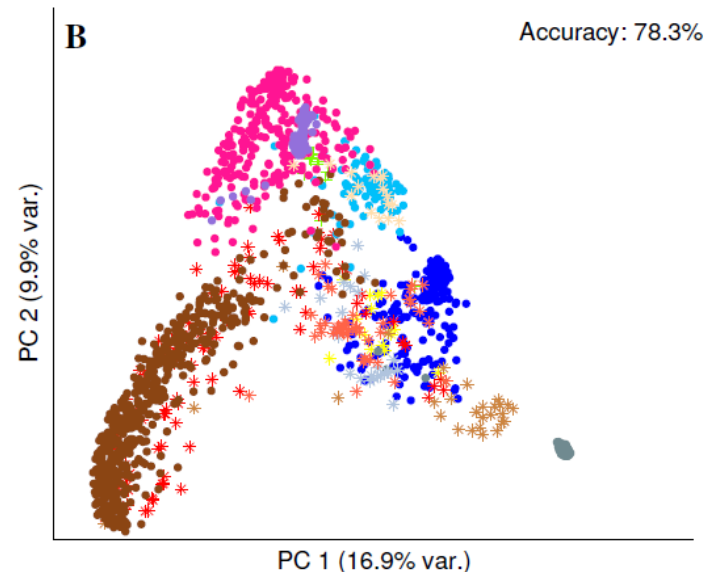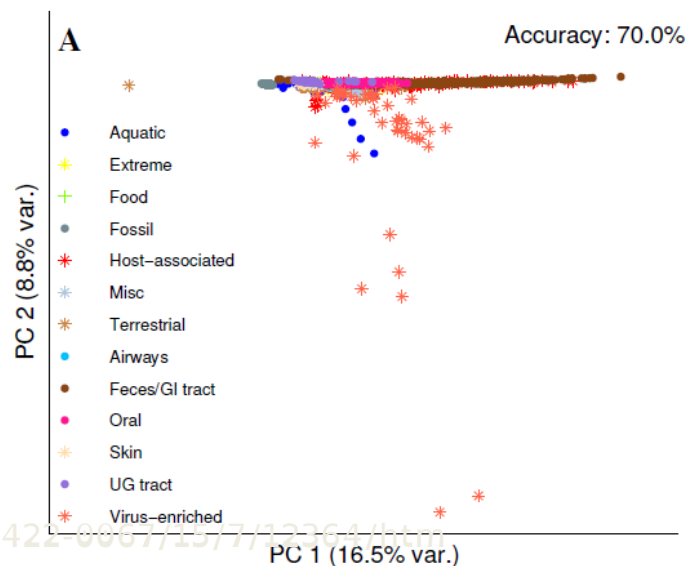
- **"Classic" MDS:**
  - $d_1(x_i, x_j) = x_i^T x_j$, $d_2(z_i, z_j) = z_i^T z_j$, $d_3(a, b) = (a - b)^2$
  - This is a factorless version of _____.
  - Not a great choice because it's _____.

# Different MDS Cost Functions

- **MDS** default objective function with general distances/similarities:

$$f(z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} d_3(d_2(z_i, z_j) - d_1(x_i, x_j))$$

- Another possibility: $d_1(x_i, x_j) = \|x_i - x_j\|_1$ and $d_2(z_i, z_j) = \|z_i - z_j\|$.
  - $z_i$ approximates high-dimensional $L_1$-norm distances.
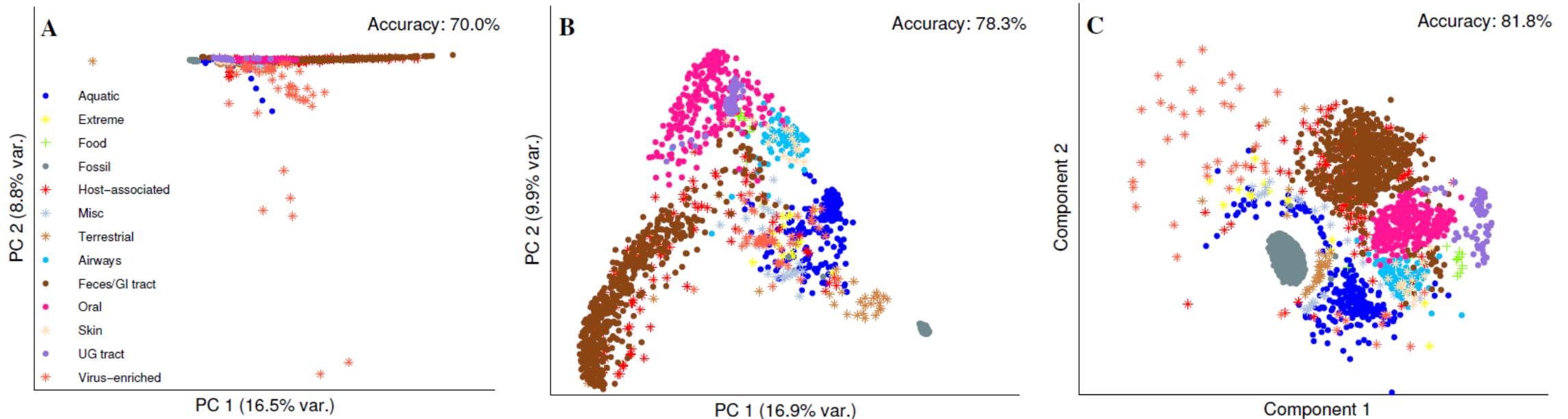
# Sammon's Mapping

- Challenge for most MDS models: they focus on _____.
  - Leads to "crowding" effect like with PCA.
- Early attempt to address this is Sammon's mapping:
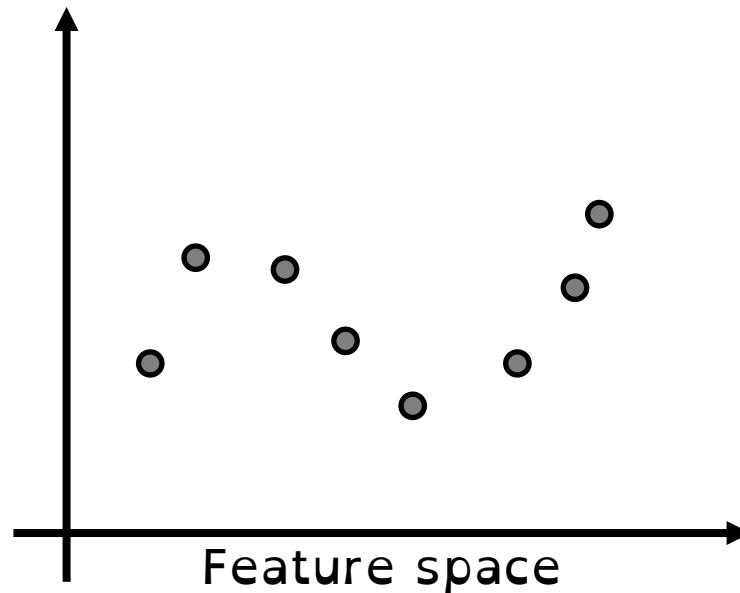  - Weighted MDS so large/small distances are more comparable.

$$f(Z) = \sum_{i=1}^{n} \sum_{j=i+1}^{n} \left( \frac{d_2(z_i, z_j) - d_1(x_i, x_j)}{d_1(x_i, x_j)} \right)^2$$

  - Denominator reduces focus on large distances.

# Sammon's Mapping

- Challenge for most MDS models: they <span style="color:red">focus on large distances</span>.
  - Leads to "crowding" effect like with PCA.
- Early attempt to address this is <span style="color:blue">Sammon's mapping</span>:
  - <span style="color:green">Weighted MDS</span> so large/small distances are more comparable.

17

Coming Up Next

# MANIFOLDS

# "Manifold"

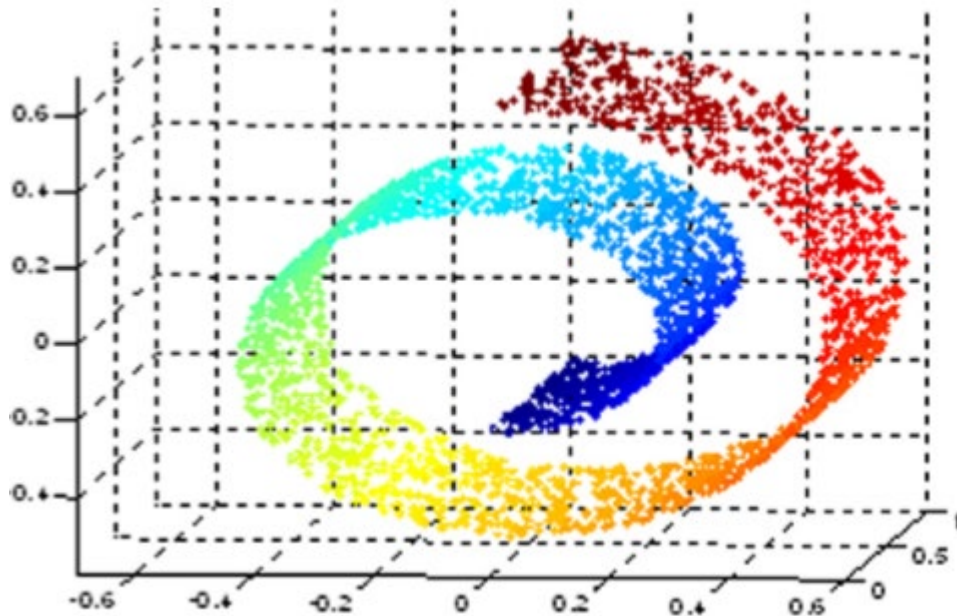- "Manifold" := non-Euclidean subspace of feature space where datapoints live



Feature space

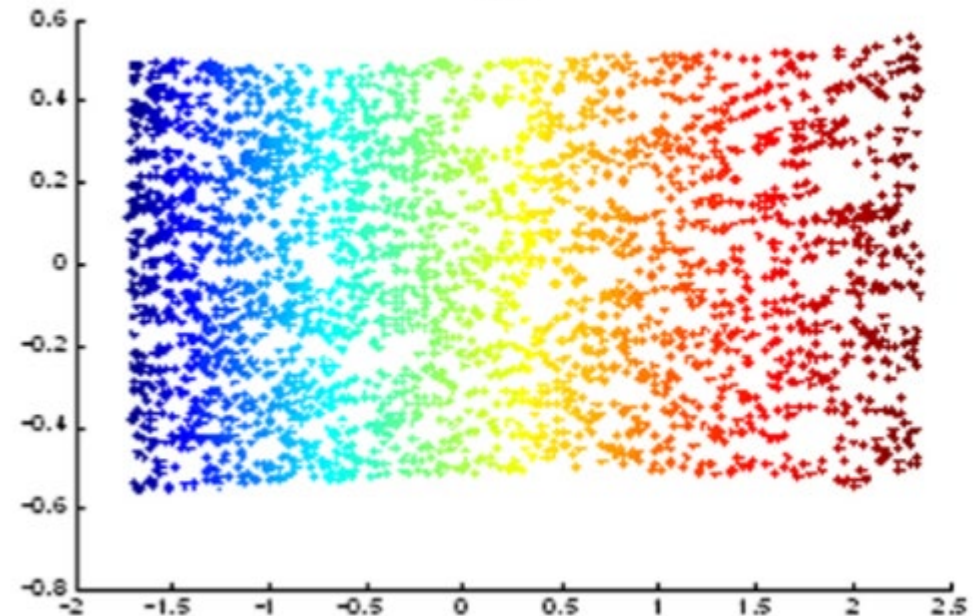- Assumption: most data live on a manifold, not a true Euclidean feature space!

# Learning Manifolds

- Consider data that lives on a low-dimensional "manifold".
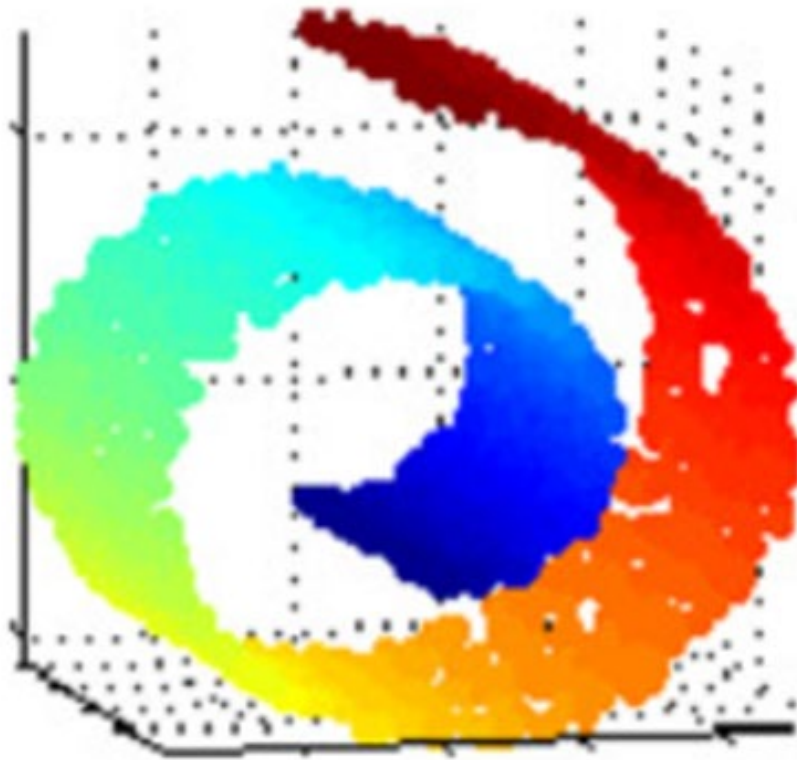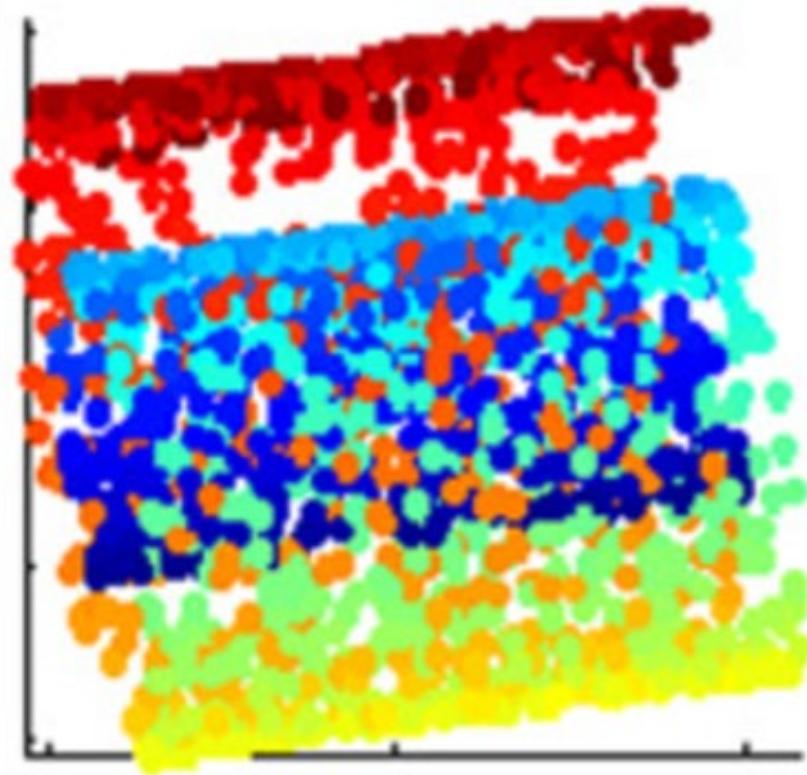- e.g. 'Swiss roll':

Original data

Two-dimensional manifold

# Learning Manifolds

- Consider data that lives on a low-dimensional "manifold".
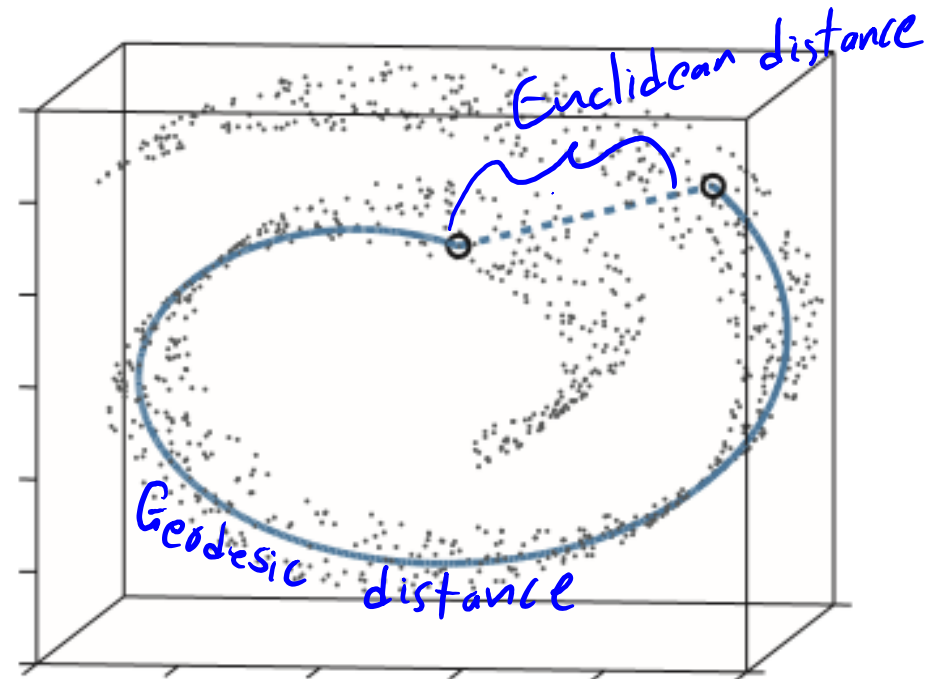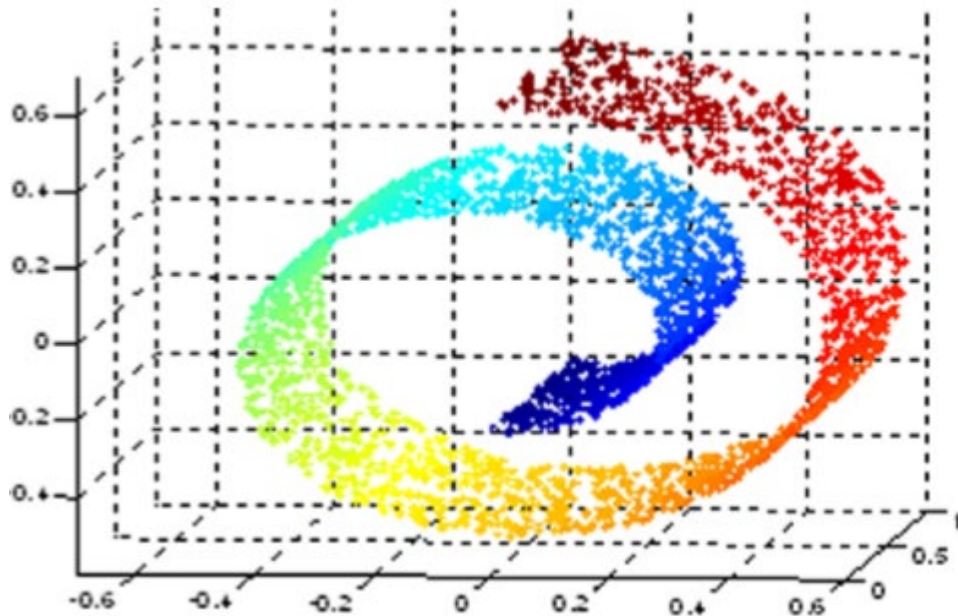  - With usual distances, PCA/MDS will not discover non-linear manifolds.



Original data

PCA

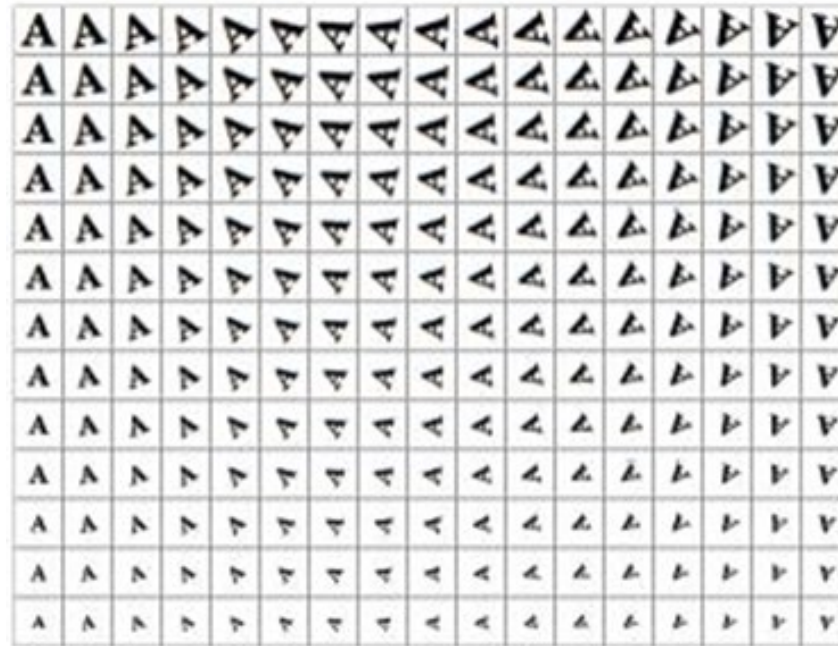# Learning Manifolds

- Consider data that lives on a low-dimensional "manifold".
  - With usual distances, PCA/MDS will not discover non-linear manifolds.
- We need geodesic distance: the _____.

# Manifolds in Image Space
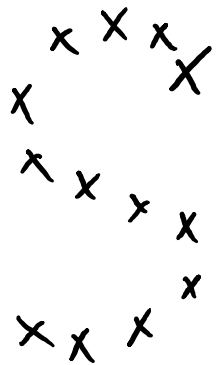
- Consider slowly-varying transformation of image:



- <span style="color:green">Images are on a manifold</span> in the high-dimensional space.
    - Euclidean distance <span style="color:red">doesn't reflect manifold structure</span>.
    - <span style="color:blue">Geodesic distance</span> is <span style="color:green">distance through space of rotations/resizings</span>.
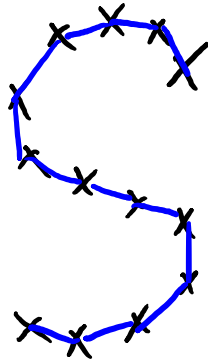
Coming Up Next

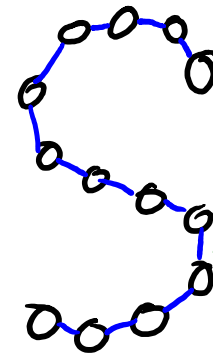# ISOMAP

# ISOMAP

- **ISOMAP** is MDS on manifolds:



find "neighbours" of each point

Represent points and neighbours as a <u>weighted</u> <u>graph</u>.

"Weight" on each edge is <u>distance</u> between points

Approximate <u>geodesic</u> distance by <u>shortest path</u> through graph.

Run MDS with these approximate geodesic distances.

ISOMAP $z_i$ values in 1D <u>or</u> 2D
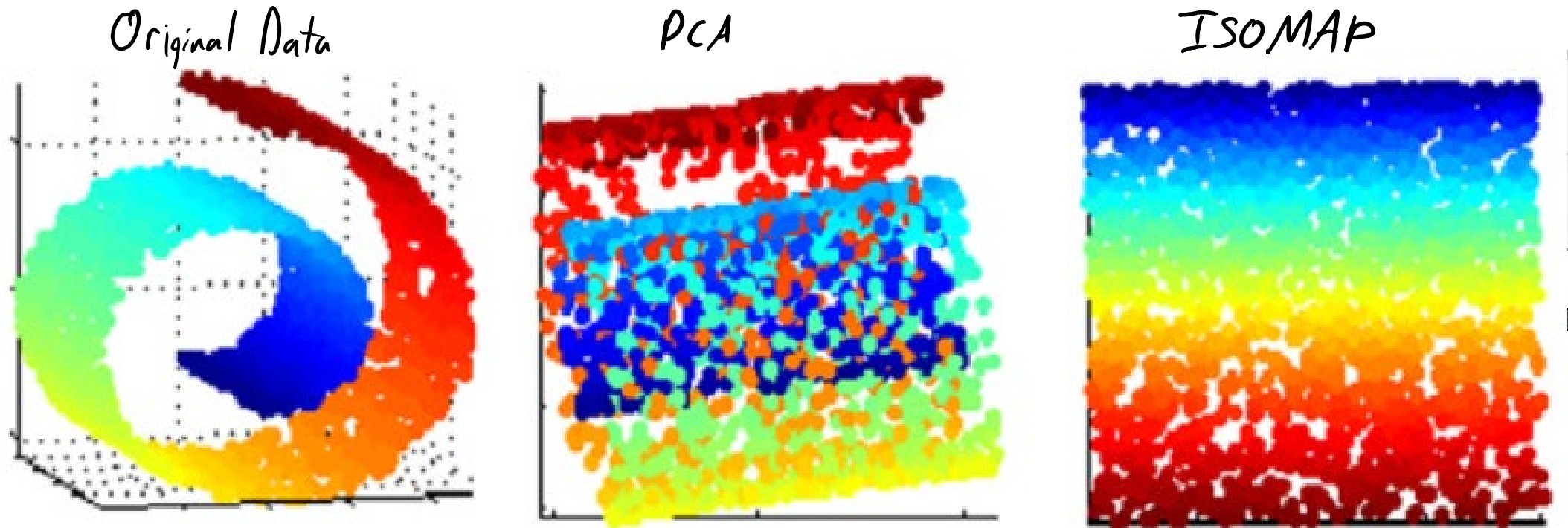
$$D = \begin{bmatrix} 0 & 1 & 2 & 3 & - - - \\ 1 & 0 & 1 & 2 & \\ 2 & 1 & 0 & 1 & \\ 3 & 2 & 1 & 0 & \\ \vdots & \vdots & \vdots & \vdots & \end{bmatrix}$$

# ISOMAP

- ISOMAP can "unwrap" the roll:
  - _____ are approximations to geodesic distances.

Original Data | PCA | ISOMAP



- Sensitive to having _____:
  - Points off of manifold and gaps in manifold cause problems.
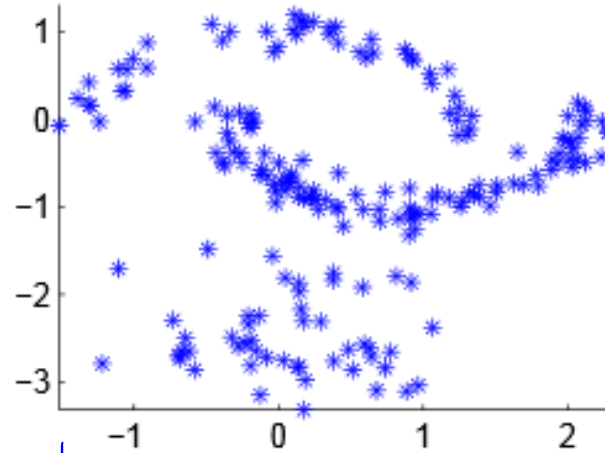
# Constructing Neighbour Graphs

- Sometimes you can define the graph/distance without features:
  - Facebook friend graph.
  - Connect YouTube videos if one video tends to follow another.

- But we can also convert from features $x_i$ to a "neighbour" graph (A6):
  - Approach 1 ("epsilon graph"): connect $x_i$ to all $x_j$ within some threshold $\varepsilon$.
    - Like we did with density-based clustering.

  - Approach 2a ("KNN graph"): connect $x_i$ to $x_j$ if:
    - $x_j$ is a KNN of $x_i$ **OR** $x_i$ is a KNN of $x_j$.

  - Approach 2b ("mutual KNN graph"): connect $x_i$ to $x_j$ if:
    - $x_j$ is a KNN of $x_i$ **AND** $x_i$ is a KNN of $x_j$.

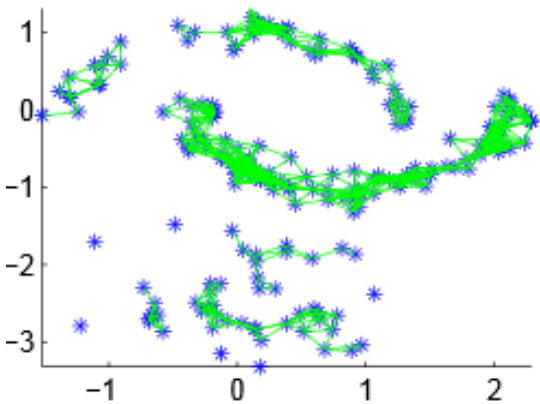# Converting from Features to Graph

add edge
if $\|x_i - x_j\| \leq 0.3$
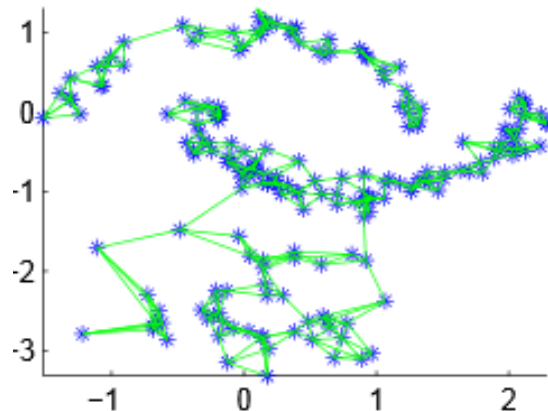
add edge if
'i' is 5-NN
of 'j'
or
'j' is
5-NN
of 'i'

add edge if
'i' and 'j'
are KNNs
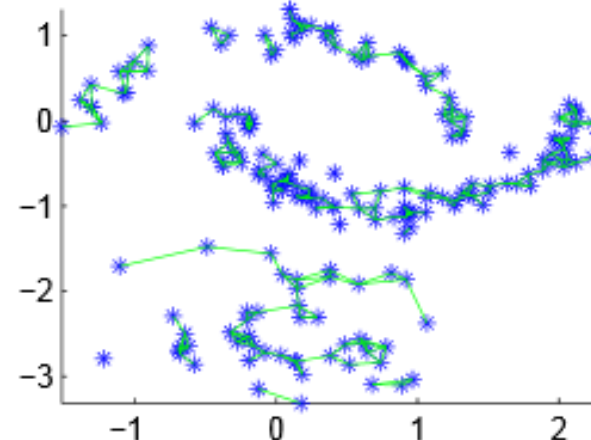of each other.

Data points

epsilon-graph, epsilon=0.3

kNN graph, k = 5

Mutual kNN graph, k = 5

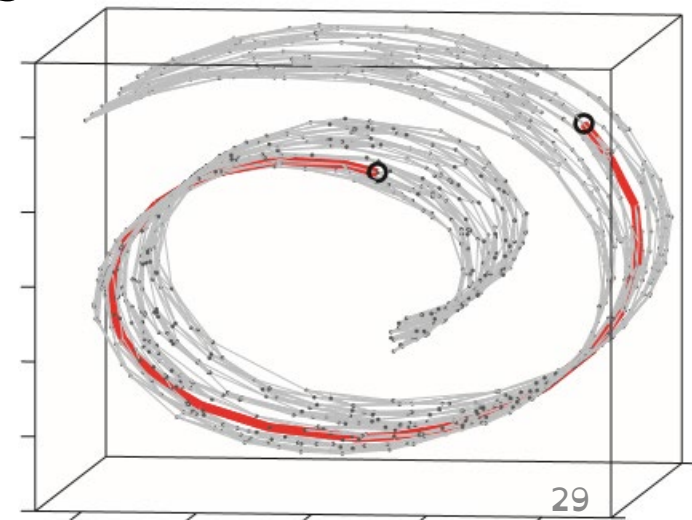# ISOMAP

- **ISOMAP** is latent-factor model for visualizing data on manifolds:
  1. Find the neighbours of each point.
     - Usually "k-nearest neighbours graph", or "epsilon graph".

  2. Compute edge weights:
     - Usually distance between neighbours.

  3. Compute weighted shortest path between all points
     - Dijkstra or other shortest path algorithm.

  4. Run MDS using these distances.

# ISOMAP on Hand Images



- Related method is "local linear embedding".

# Sammon's Map vs. ISOMAP vs. PCA



Sammon Map

ISOMAP

PCA

# Sammon's Map vs. ISOMAP vs. t-SNE



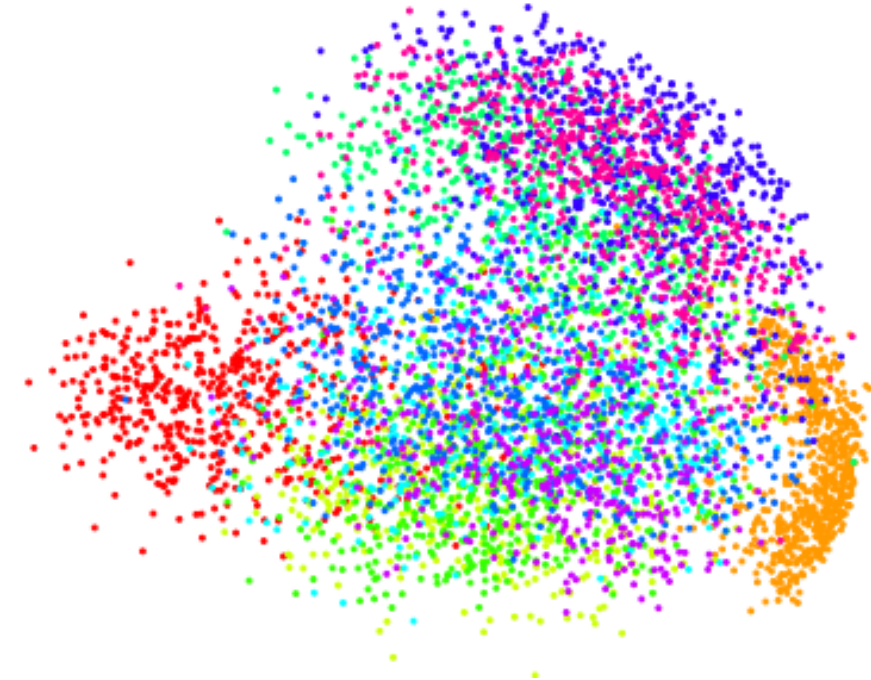Sammon Map

ISOMAP

t-SNE

# Sammon's Map vs. ISOMAP vs. t-SNE

Sammon Map

ISOMAP

t-SNE

Remember this is <u>unsupervised</u>, algorithms do <u>not</u> know the labels.

# Sammon's Map vs. ISOMAP vs. t-SNE



Legend:
- 0 (red)
- 1 (orange)
- 2 (yellow)
- 3 (green)
- 4 (light green)
- 5 (cyan)
- 6 (blue)
- 7 (purple)
- 8 (magenta)
- 9 (pink)

Sammon Map

"1"

"0"

ISOMAP

"1"

"0"

t-SNE

"0"

"1"

Remember this is __unsupervised__, algorithms do __not__ know the labels.

# Sammon's Map vs. ISOMAP vs. t-SNE

# Sammon's Map vs. ISOMAP vs. t-SNE

Coming Up Next

# T-SNE

# t-Distributed Stochastic Neighbour Embedding

- One key idea in t-SNE:
  - Focus on distance to "neighbours"
    (allow large variance in other distances)



PCA

Tries to preserve larger distances

"crowding" because doesn't focus on small distances

t-SNE

Tries to preserve neighbour distances.

focuses on these

"repulsion" where there are gaps in distances.

# t-Distributed Stochastic Neighbour Embedding

- t-SNE is a special case of MDS (specific $d_1$, $d_2$, and $d_3$ choices):
  - $d_1$: for each $x_i$, compute 'neighbour-ness' of each $x_j$
  - Computation is similar to k-means++, but most weight to close points (Gaussian).
    - Doesn't require explicit graph.

  - $d_2$: for each $z_i$, compute 'neighbour-ness' of each $z_j$.
    - Similar to above, but use student's t (grows really slowly with distance).
    - Avoids 'crowding', because you have a huge range that large distances can fill.

  - $d_3$: Compare $x_i$ and $z_i$ using an entropy-like measure:
    - How much 'randomness' is in probabilities of $x_i$ if you know the $z_i$ (and vice versa)?

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

# t-SNE on Wikipedia Articles

http://jasneetsabharwal.com/assets/files/wiki_tsne_report.pdf

# t-SNE on Product Features

# t-SNE on Leukemia Heterogeneity

# End of Part 4: Latent Factor Models

# End of Part 4: Key Concepts

- We discussed linear latent-factor models:

$$f(W,z) = \sum_{i=1}^{n} \sum_{j=1}^{d} (\langle w_j^j z_i \rangle - x_{ij})^2$$

$$= \sum_{i=1}^{n} \| W^T z_i - x_i \|^2$$

$$= \| ZW - X \|_F^2$$

- Represent 'X' as linear combination of latent factors '$w_c$'.
  - Latent features '$z_i$' give a lower-dimensional version of each '$x_i$'.
  - When k=1, finds direction that minimizes squared orthogonal distance.
- Applications:
  - Outlier detection, dimensionality reduction, data compression, features for linear models, visualization, factor discovery, filling in missing entries.

# End of Part 4: Key Concepts

- We discussed linear latent-factor models:

$$f(W, z) = \sum_{i=1}^{n} \sum_{j=1}^{d} (\langle w_j, z_i \rangle - x_{ij})^2$$

- Principal component analysis (PCA):
  - Often uses orthogonal factors and fits them sequentially (via SVD).
- Non-negative matrix factorization:
  - Uses non-negative factors giving sparsity.
  - Can be minimized with projected gradient.
- Many variations are possible:
  - Different regularizers (sparse coding) or loss functions (robust/binary PCA).
  - Missing values (recommender systems) or change of basis (kernel PCA).

# End of Part 4: Key Concepts

- We discussed multi-dimensional scaling (MDS):
  - Non-parametric method for high-dimensional data visualization.
  - Tries to match distance/similarity in high-/low-dimensions.
    - "Gradient descent on scatterplot points".
- Main challenge in MDS methods is "crowding" effect:
  - Methods focus on large distances and lose local structure.
- Common solutions:
  - Sammon mapping: use weighted cost function.
  - ISOMAP: approximate geodesic distance using via shortest paths in graph.
  - T-SNE: give up on large distances and focus on neighbour distances.

# Summary

- **Different MDS distances/losses/weights** usually gives better results.
- **Manifold learning** focuses on low-dimensional curved structures.
- **ISOMAP** is most common approach:
    - Approximates geodesic distance by shortest path in weighted graph.
- **t-SNE** is promising new data MDS method.

- Next time: deep learning.

# Please Do Course Evaluation!

# Review Questions

- **Q1:** Is MDS sensitive to initialization? Why?

- **Q2:** What is the problem with using linear dimensionality reduction for data on manifold?

- **Q3:** How does ISOMAP compute pair-wise distances among examples?

- **Q4:** What is the key idea behind t-SNE in terms of preserving distances in 2D?

# Does t-SNE always outperform PCA?

- Consider 3D data living on a 2D hyper-plane:



- PCA can perfectly capture the low-dimensional structure.

- T-SNE can capture the local structure, but can "twist" the plane.
  – It doesn't try to get long distances correct.

# Graph Drawing

- A closely-related topic to MDS is graph drawing:
  - Given a graph, how should we display it?
  - Lots of interesting methods: https://en.wikipedia.org/wiki/Graph_drawing

# Bonus Slide: Multivariate Chain Rule

- Recall the <span style="color:blue">univariate chain rule</span>:
$$\frac{d}{dw}\left[ f(g(w)) \right] = f'(g(w))\, g'(w)$$

- The <span style="color:blue">multivariate chain rule</span>:
$$\nabla\left[ f(g(w)) \right] = f'(g(w))\, \nabla g(w)$$
$$\underbrace{\nabla\left[ f(g(w)) \right]}_{d \times 1} = \underbrace{f'(g(w))}_{1 \times 1}\, \underbrace{\nabla g(w)}_{d \times 1}$$

- Example:

$$\nabla\left[ \tfrac{1}{2}(w^\top x_i - y_i)^2 \right]$$

$$= \nabla\left[ f(g(w)) \right]$$

with $g(w) = w^\top x_i - y_i \longrightarrow \nabla g(w) = x_i \longrightarrow$

and $f(r_i) = \tfrac{1}{2} r_i^2 \longrightarrow f'(r_i) = r_i \longrightarrow$

$$\nabla\left[ f(g(w)) \right] = r_i\, x_i$$
$$= (w^\top x_i - y_i) x_i$$

# Bonus Slide: Multivariate Chain Rule for MDS

- General MDS formulation:

$$\underset{Z \in \mathbb{R}^{n \times k}}{\text{argmin}} \sum_{i=1}^{n} \sum_{j=i+1}^{n} g\left( d_1(x_i, x_j), d_2(z_i, z_j) \right)$$

- Using multivariate chain rule we have:

$$\nabla_{z_i} g\left( d_1(x_i, x_j), d_2(z_i, z_j) \right) = g'\left(d_1(x_i, x_j), d_2(z_i, z_j)\right) \nabla_{z_i} d_2(z_i, z_j)$$

- Example:

If $d_1(x_i, x_j) = \|x_i - x_j\|$ and $d_2(z_i, z_j) = \|z_i - z_j\|$ and $g(d_1, d_2) = \frac{1}{2}(d_1 - d_2)^2$

$$\nabla_{z_i} g\left( d_1(x_i, x_j), d_2(z_i, z_j) \right) = -\underbrace{\left( d_1(x_i, x_j) - d_2(z_i, z_j) \right)}_{g'(d_1, d_2)}\left[ -\underbrace{\frac{(z_i - z_j)}{2\|z_i - z_j\|}}_{} \right] \xrightarrow{} \nabla_{z_i} d_2(z_i, z_j)$$

$\hookrightarrow$ Assuming $z_i \neq z_j$

(move distances closer)

(how distance changes in $z$ space)

53

# Latent-Factor Representation of Words

- For natural language, we often represent words by an index.
  - E.g., "cat" is word 124056 among a "bag of words".

- But this may be inefficient:
  - Should "cat" and "kitten" share parameters in some way?

# Latent-Factor Representation of Words

- Latent-factor representation of individual words:
  - Closeness in latent space should indicate similarity.
  - Distances could represent meaning?

- Recent alternative to PCA/NMF is word2vec...

# Using Context

- Consider these phrases:
    - "the <u>cat</u> purred"
    - "the <u>kitten</u> purred"

    - "black <u>cat</u> ran"
    - "black <u>kitten</u> ran"

- Words that occur in the same context likely have similar meanings.

- Word2vec uses this insight to design an MDS distance function.

# Word2Vec

- Two common word2vec approaches:
  1. Try to predict word from surrounding words (continuous bag of words).
  2. Try to predict surrounding words from word (skip-gram).



Figure 1: New model architectures. The CBOW architecture predicts the current word based on the context, and the Skip-gram predicts surrounding words given the current word.

- Train latent-factors to solve one of these supervised learning tasks.

# Word2Vec

- In both cases, each word 'i' is represented by a vector $z_i$.
- In continuous bag of words (CBOW), we optimize the following likelihood:

$$p(x_i \mid x_{surround}) = \prod_{j \in surround} p(x_i \mid x_j) \qquad \text{(independence assumption)}$$

$$= \prod_{j \in surround} \frac{exp(z_i^T z_j)}{\sum_{c=1}^{k} exp(z_c^T z_j)} \qquad \text{(softmax over all words)}$$

- Apply gradient descent to logarithm:
  - Encourages $z_i^T z_j$ to be big for words in same context (making $z_i$ close to $z_j$).
  - Encourages $z_i^T z_j$ to be small for words not appearing in same context (makes $z_i$ and $z_j$ far).
- For CBOW, denominator sums over all words.
- For skip-gram it will be over all possible surrounding words.
  - Common trick to speed things up: sample terms in denominator ("negative sampling").

# Word2Vec Example

- MDS visualization of a set of related words:



- Distances between vectors might represent semantics.

http://sebastianruder.com/secret-word2vec

# Word2Vec

- Subtracting word vectors to find related vectors.

Table 8: *Examples of the word pair relationships, using the best word vectors from Table 4 (Skip-gram model trained on 783M words with 300 dimensionality).*

| Relationship | Example 1 | Example 2 | Example 3 |
|---|---|---|---|
| France - Paris | Italy: Rome | Japan: Tokyo | Florida: Tallahassee |
| big - bigger | small: larger | cold: colder | quick: quicker |
| Miami - Florida | Baltimore: Maryland | Dallas: Texas | Kona: Hawaii |
| Einstein - scientist | Messi: midfielder | Mozart: violinist | Picasso: painter |
| Sarkozy - France | Berlusconi: Italy | Merkel: Germany | Koizumi: Japan |
| copper - Cu | zinc: Zn | gold: Au | uranium: plutonium |
| Berlusconi - Silvio | Sarkozy: Nicolas | Putin: Medvedev | Obama: Barack |
| Microsoft - Windows | Google: Android | IBM: Linux | Apple: iPhone |
| Microsoft - Ballmer | Google: Yahoo | IBM: McNealy | Apple: Jobs |
| Japan - sushi | Germany: bratwurst | France: tapas | USA: pizza |

Table 8 shows words that follow various relationships. We follow the approach described above: the relationship is defined by subtracting two word vectors, and the result is added to another word. Thus for example, *Paris - France + Italy = Rome.* As it can be seen, accuracy is quite good, although

- Word vectors for 157 languages [here](https://arxiv.org/pdf/1301.3781.pdf).

# Multiple Word Prototypes

- What about homonyms and polysemy?
  - The word vectors would need to account for all meanings.

- More recent approaches:
  - Try to cluster the different contexts where words appear.
  - Use different vectors for different contexts.

$$X_{jaguar} = \begin{bmatrix} \cdots \\ \cdots \end{bmatrix} \begin{matrix} z_{j1} \\ z_{j2} \\ z_{j3} \end{matrix}$$

# Multiple Word Prototypes

http://www.socher.org/index.php/Main/ImprovingWordRepresentationsViaGlobalContextAndMultipleWordPrototypes