

CPSC 340: Machine Learning and Data Mining

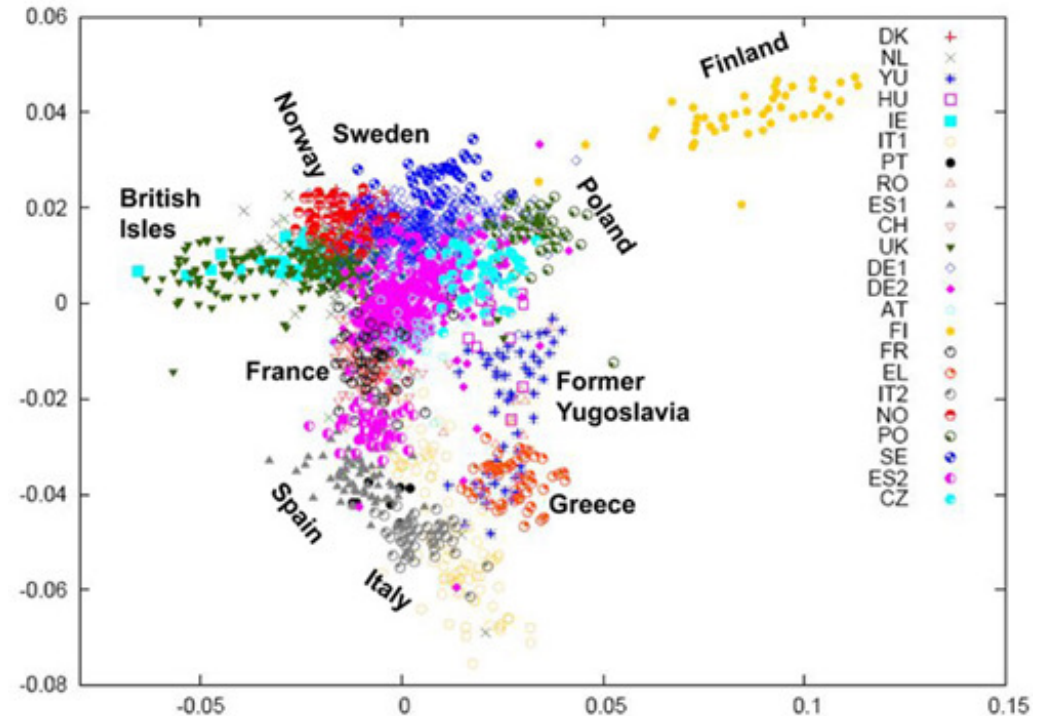
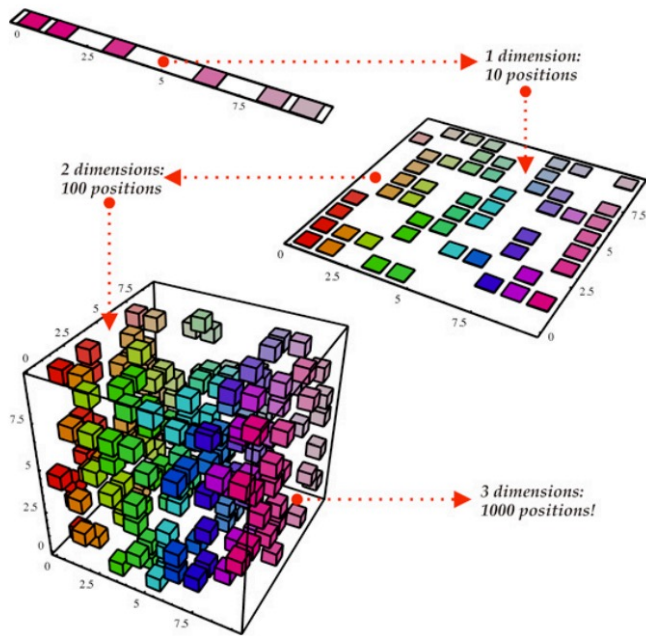
Multi-Dimensional Scaling

Last Time: Variants of PCA

- Solve the PCA objective function by alternative minimization and gradient descent.
- Variants of PCA: robust PCA, binary PCA, regularized PCA.
- Non-negative matrix factorization, topic modeling.
- We discussed **recommender systems**:
 - Predicting what ratings users have for different products.
 - **content-based filtering (supervised)**: Extract features of users and products, and use these to predict rating.
 - **collaborative filtering (unsupervised)**: Methods that **only looks at ratings**, not features of products.

Visualization High-Dimensional Data

- PCA for visualizing high-dimensional data:
 - Use PCA ‘W’ matrix to linearly transform data to get the z_i values.
 - And then we plot the z_i values as locations in a scatterplot.



Visualization High-Dimensional Data

- PCA for **visualizing high-dimensional** data:
 - Use PCA ‘W’ matrix to **linearly transform data** to get the z_i values.
 - And then we plot the z_i values as locations in a scatterplot.
- An common alternative is **multi-dimensional scaling (MDS)**:
 - **Directly optimize the pixel locations of the z_i values.**
 - “Gradient descent on the points in a scatterplot”.
 - Needs a “cost” function saying how “good” the z_i locations are.

- Traditional **MDS cost function**:

$$f(z) = \sum_{i=1}^n \sum_{j=i+1}^n$$

sum over
pairs of
examples

$$\left(\|z_i - z_j\| - \|x_i - x_j\| \right)^2$$

distance in
scatterplot

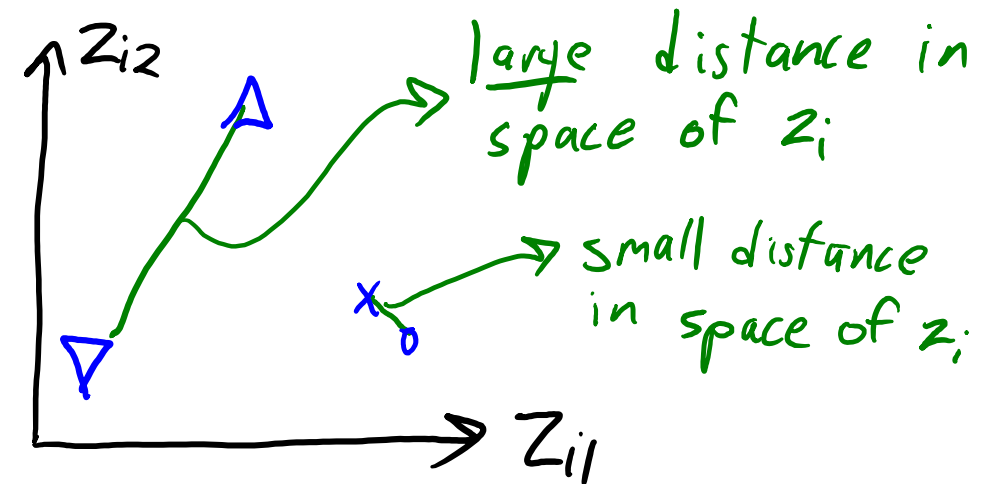
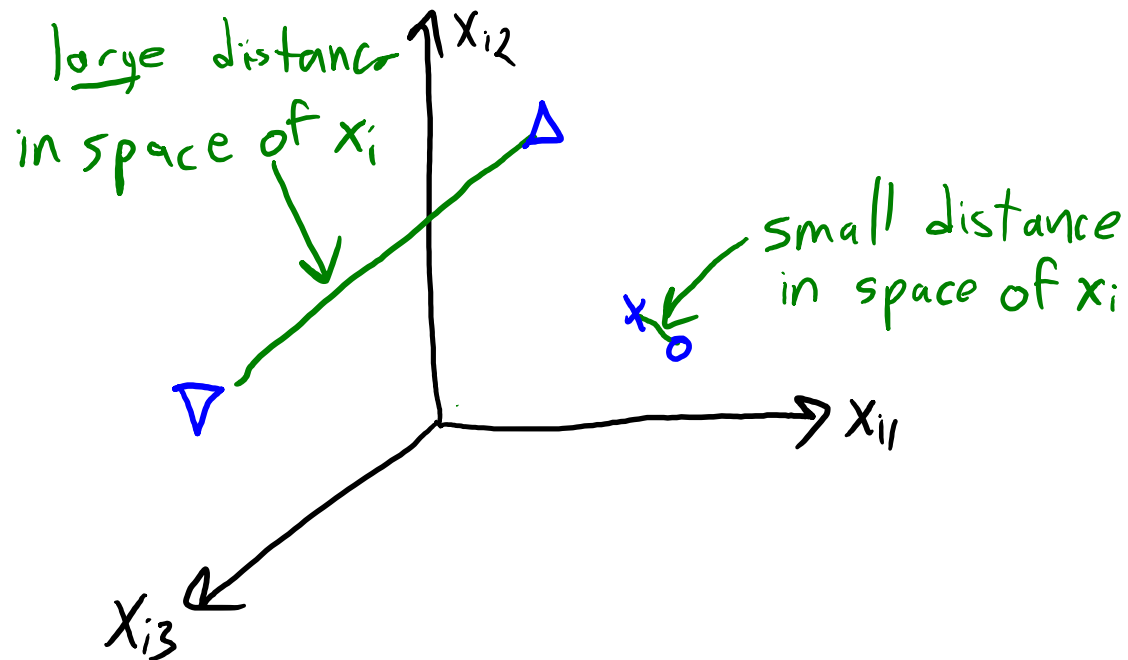
Distance between points
in original 'd' dimensions

Try to make scatterplot
distances match high-dimensional
distance

Multi-Dimensional Scaling

- Multi-dimensional scaling (MDS):
 - Directly optimize the final locations of the z_i values.

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



Multi-Dimensional Scaling

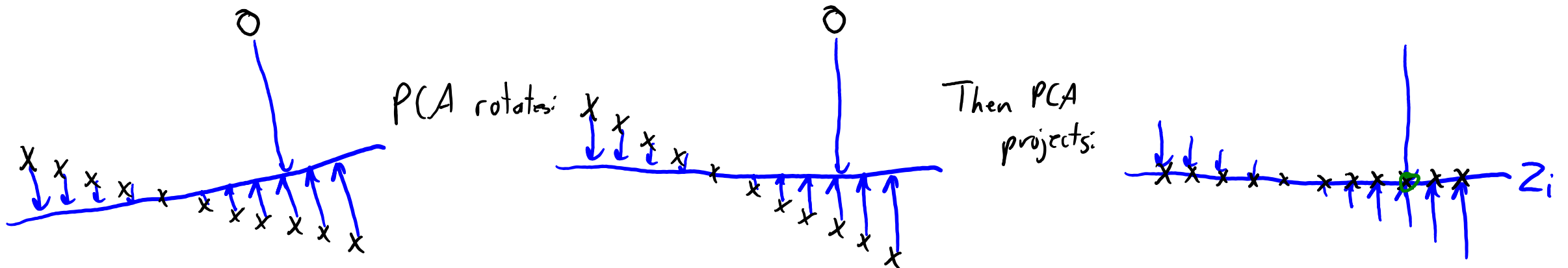
- Multi-dimensional scaling (MDS):

- Directly optimize the final locations of the z_i values.

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

- Non-parametric dimensionality reduction and visualization:

- No 'W': just trying to make z_i preserve high-dimensional distances between x_i .



Multi-Dimensional Scaling

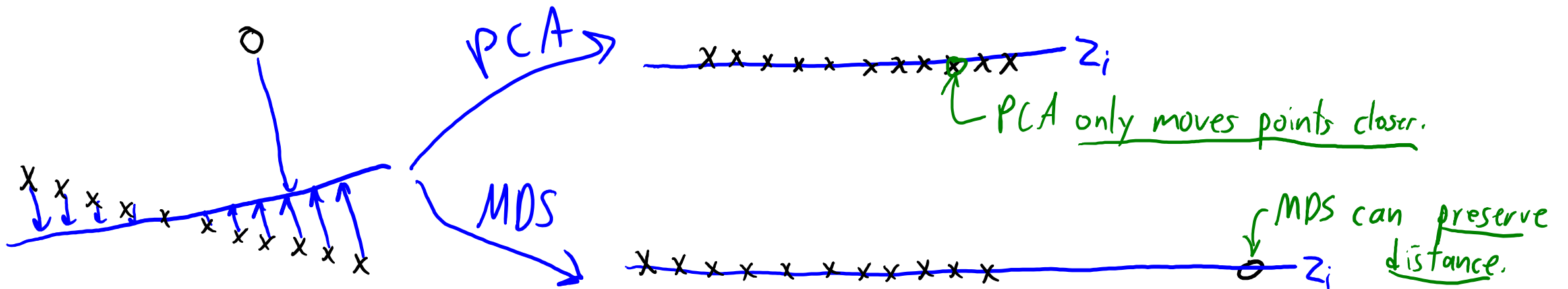
- Multi-dimensional scaling (MDS):

- Directly optimize the final locations of the z_i values.

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

- Non-parametric dimensionality reduction and visualization:

- No 'W': just trying to make z_i preserve high-dimensional distances between x_i .



Multi-Dimensional Scaling

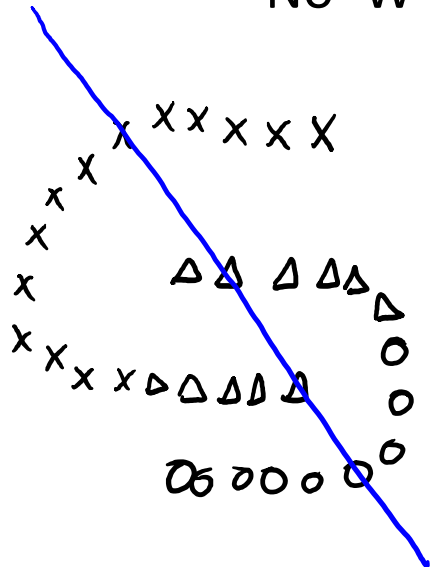
- Multi-dimensional scaling (MDS):

- Directly optimize the final locations of the z_i values.

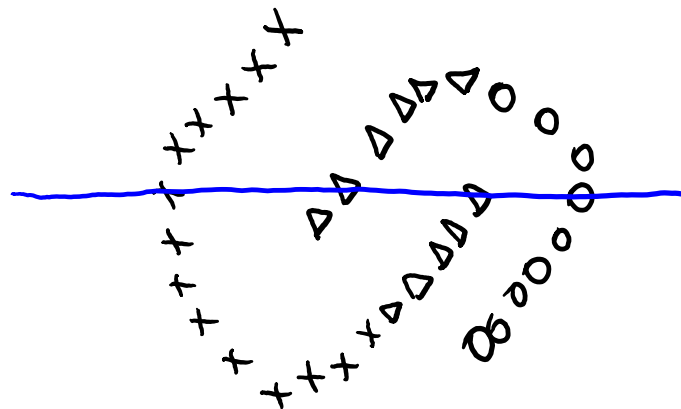
$$f(z) = \sum_{i=1}^n \sum_{j=i+1}^n (\|z_i - z_j\| - \|x_i - x_j\|)^2$$

- Non-parametric dimensionality reduction and visualization:

- No 'W': just trying to make z_i preserve high-dimensional distances between x_i .



PCA rotation:



PCA projection:



Multi-Dimensional Scaling

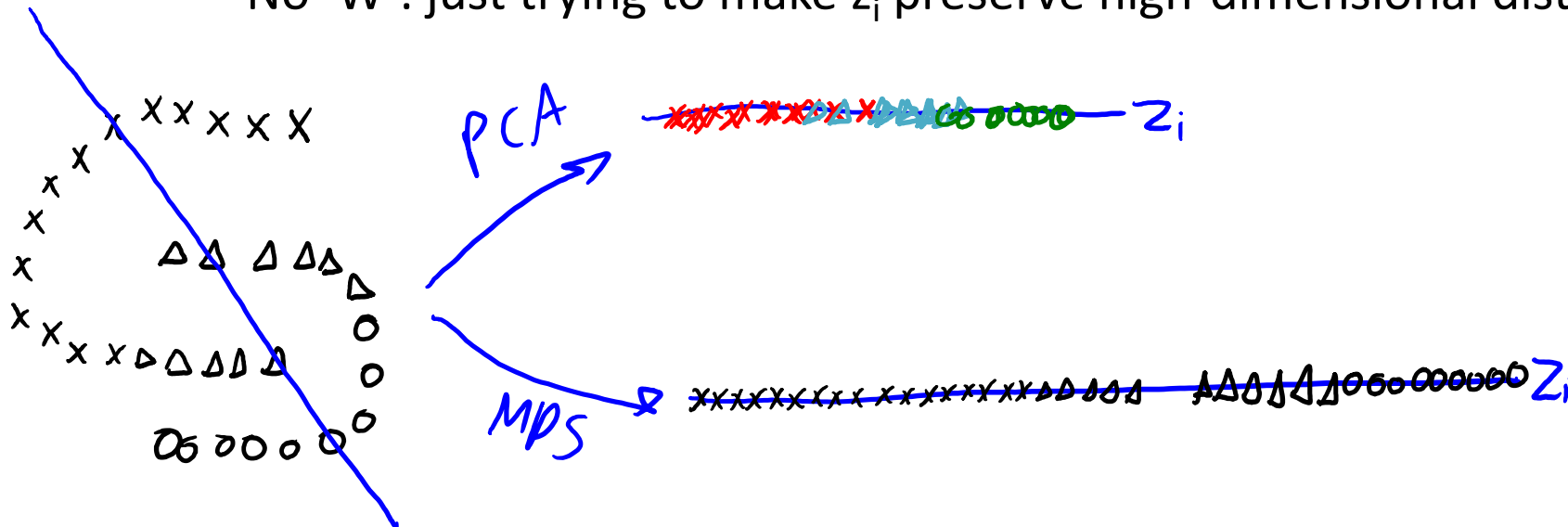
- Multi-dimensional scaling (MDS):

- Directly optimize the final locations of the z_i values.

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

- 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):

- Directly optimize the final locations of the z_i values.

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

- 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.

Different MDS Cost Functions

- The default MDS objective function using the Euclidean distance:

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

- We could consider **different distances/similarities**:

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

- Where the functions are **not necessarily the same**:

- d_1 is the high-dimensional distance we want to match.
- d_2 is the low-dimensional distance we can control.
- d_3 controls how we compare high-/low-dimensional distances.

PCA is a Special MDS

- Let d_1 and d_2 be the dot product, and d_3 the square distance

Note, $\mathbf{K} = \mathbf{X}\mathbf{X}^T$ is called the Gram matrix, which measures the dot-product between data points (centered)

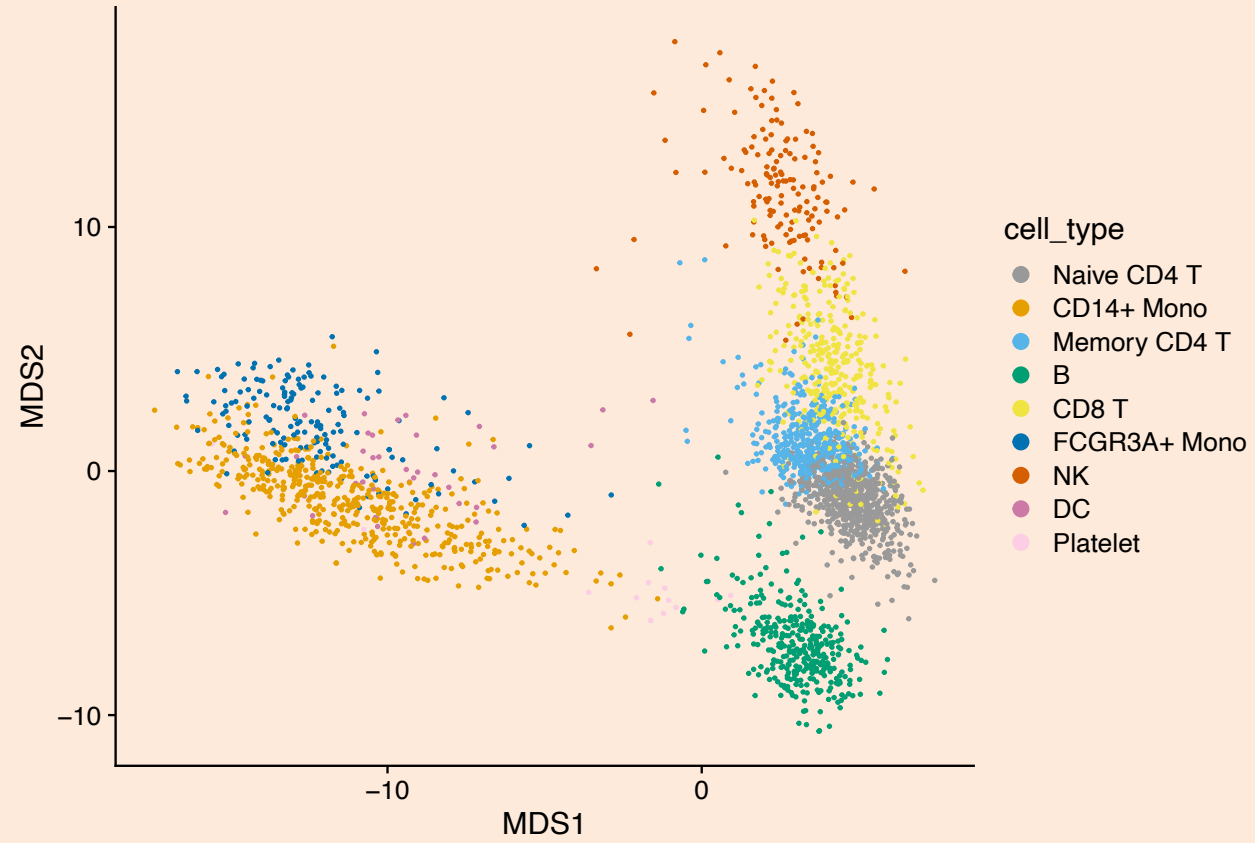
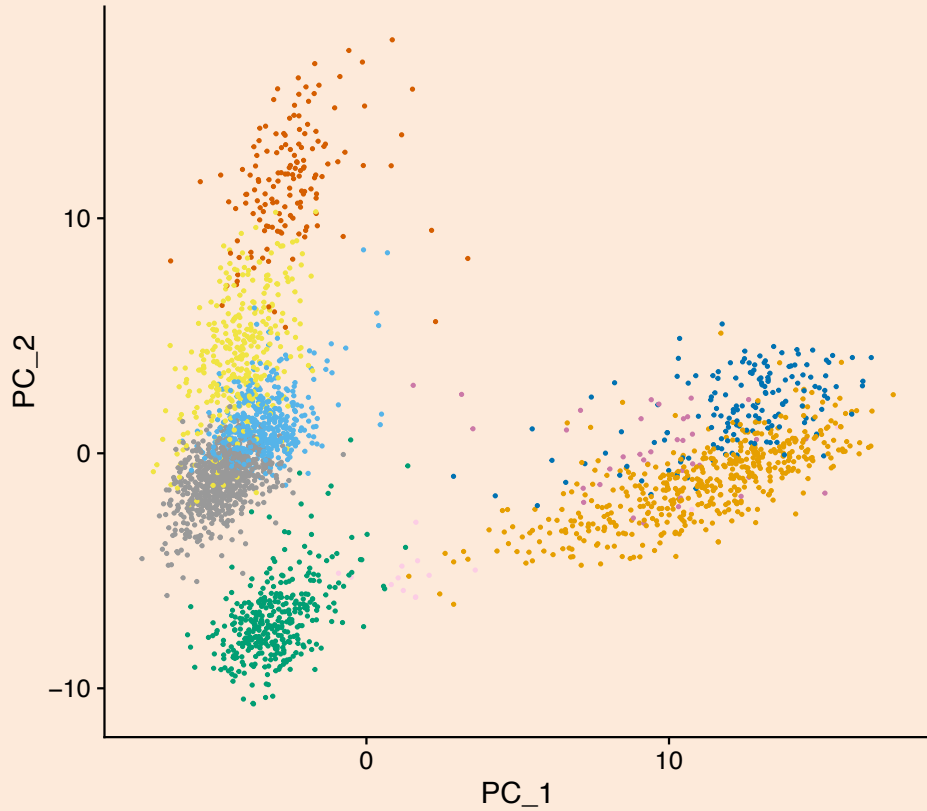
We may be interested in minimizing the distortions in distances after projecting \mathbf{x} to \mathbf{z} , e.g., $\sum_{i,j} (\mathbf{x}_i^T \mathbf{x}_j - \mathbf{z}_i^T \mathbf{z}_j)^2 = \|\mathbf{K} - \mathbf{Z}\mathbf{Z}^T\|_F^2$

Using SVD, we can factorize, $\mathbf{X}\mathbf{X}^T = \mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{V}^T\mathbf{V}\mathbf{\Lambda}^{1/2}\mathbf{U}^T = \mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{\Lambda}^{1/2}\mathbf{U}^T$

The best rank- K approximation to \mathbf{K} is $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^T = \mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{\Lambda}^{1/2}\mathbf{U}^T$

Thus, we can get $\mathbf{Z} = \mathbf{U}\mathbf{\Lambda}^{1/2}$ (PCA latent representation and the classic multi-dimensional scaling)

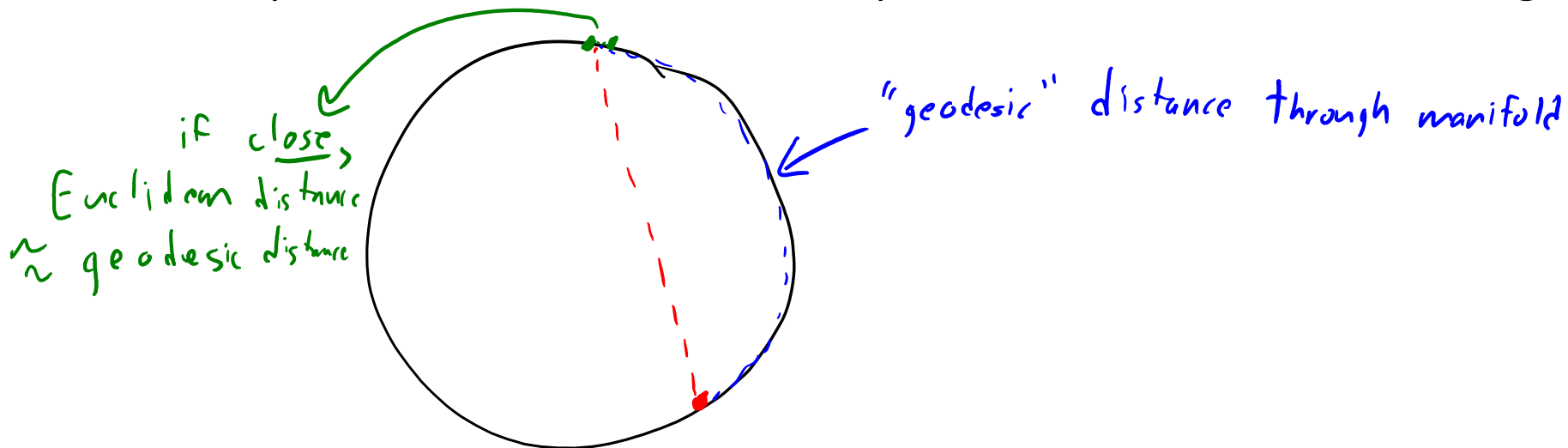
PCA is a Special MDS



Next Topic: *t*-SNE

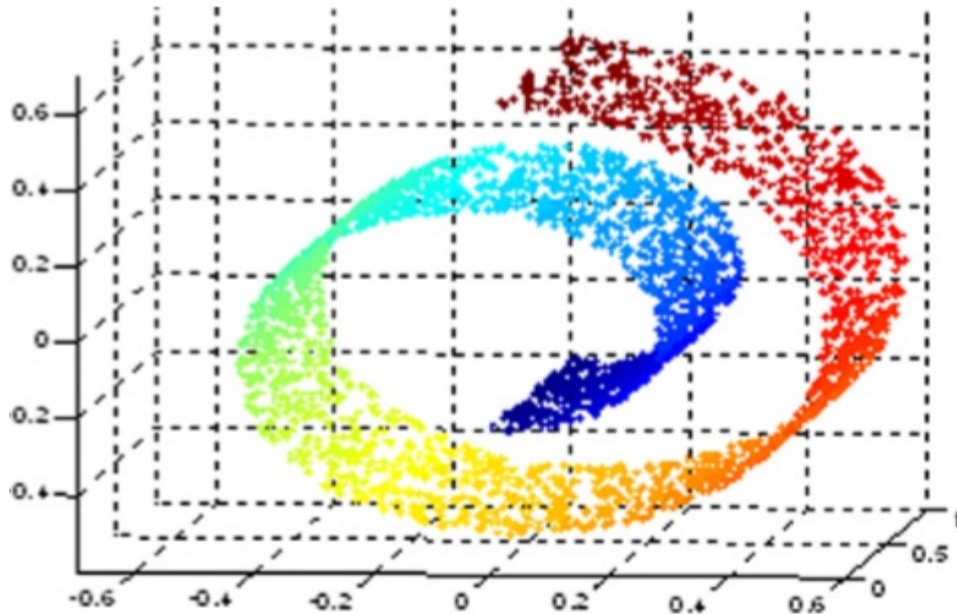
Data on Manifolds

- Consider data that lives on a **low-dimensional “manifold”**.
 - Where **Euclidean distances make sense “locally”**.
 - But **Euclidean distances may not make sense “globally”**.
 - Wikipedia example: Surface of the Earth is “locally” flat.
 - Euclidean distance accurately measures distance “along the surface” locally.
 - For far points Euclidean distance is a poor measure of distance “along the surface”.

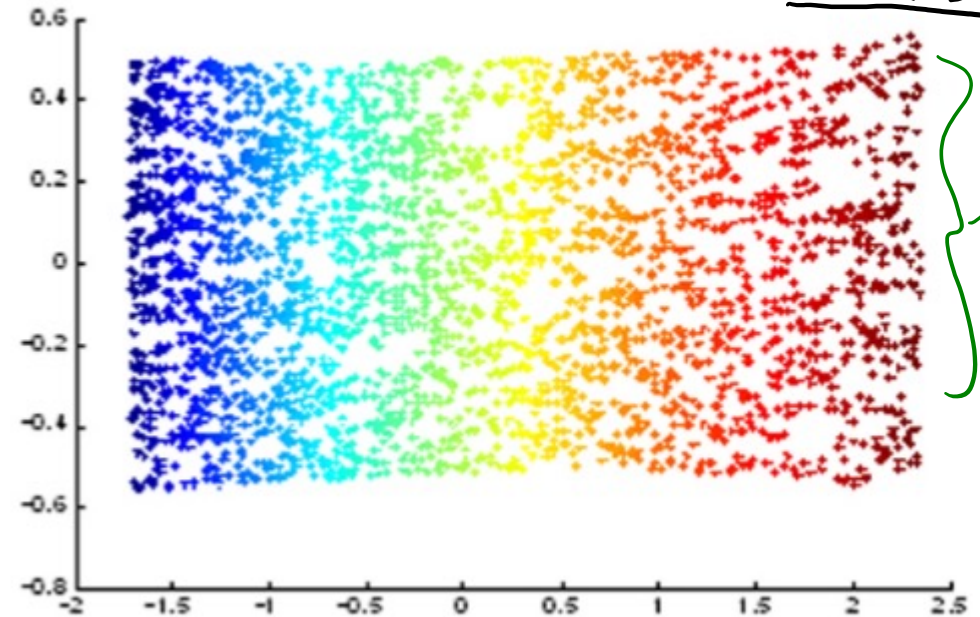


Data on Manifolds

- Consider data that lives on a **low-dimensional “manifold”**.
 - Where **Euclidean distances make sense “locally”**.
 - But **Euclidean distances may not make sense “globally”**.
- Example is the ‘Swiss roll’:



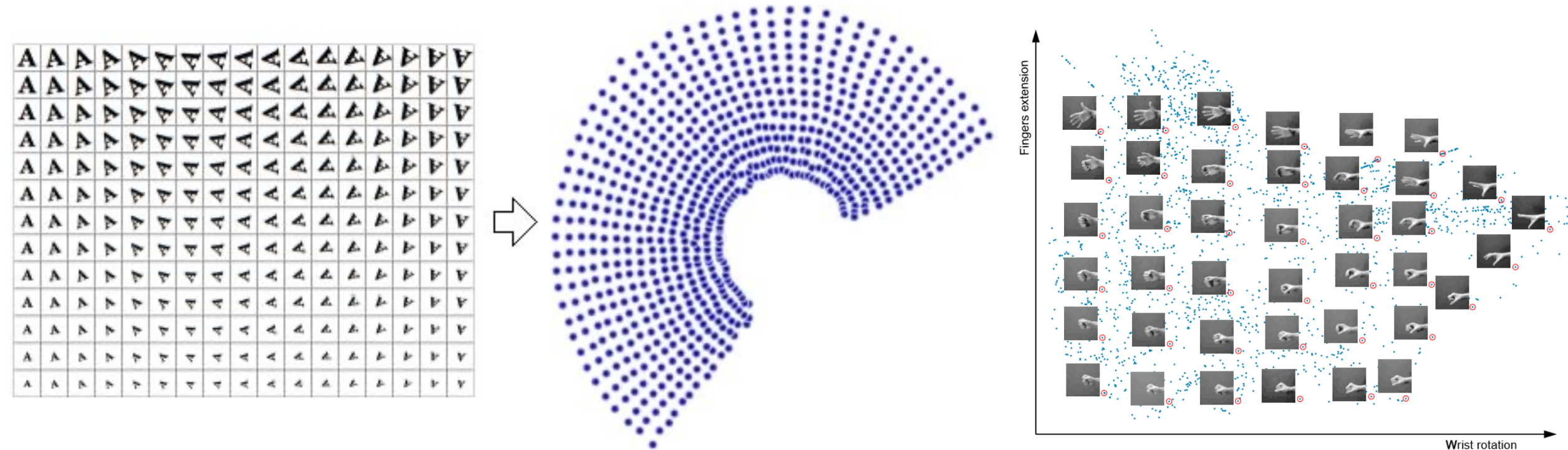
We want an MDS method that visualizes manifolds



This visualization "unrolls" the Swiss roll.

Example: Manifolds in Image Space

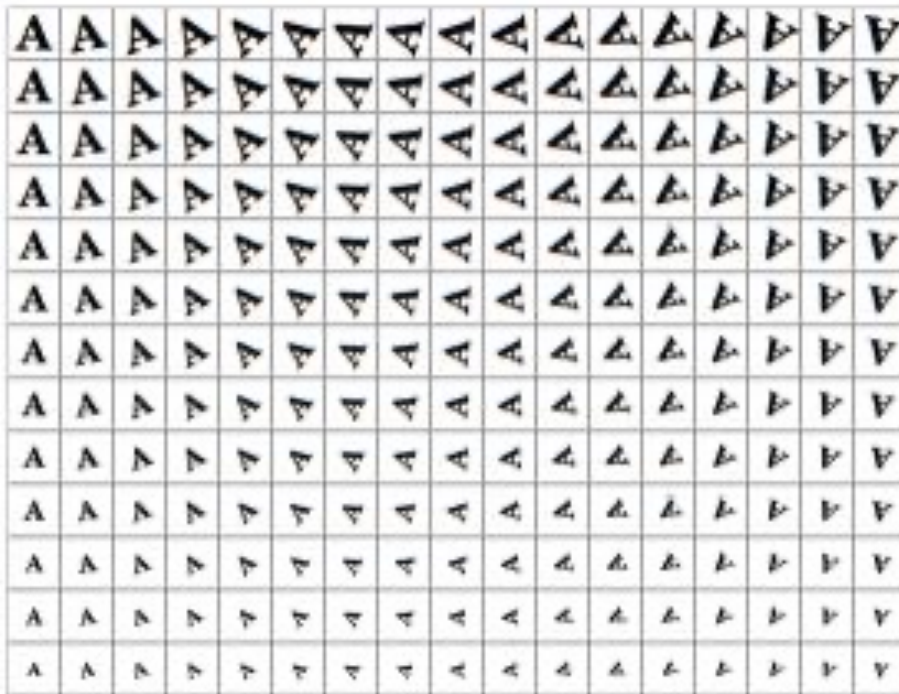
- Slowly-varying image transformations exist on a manifold:



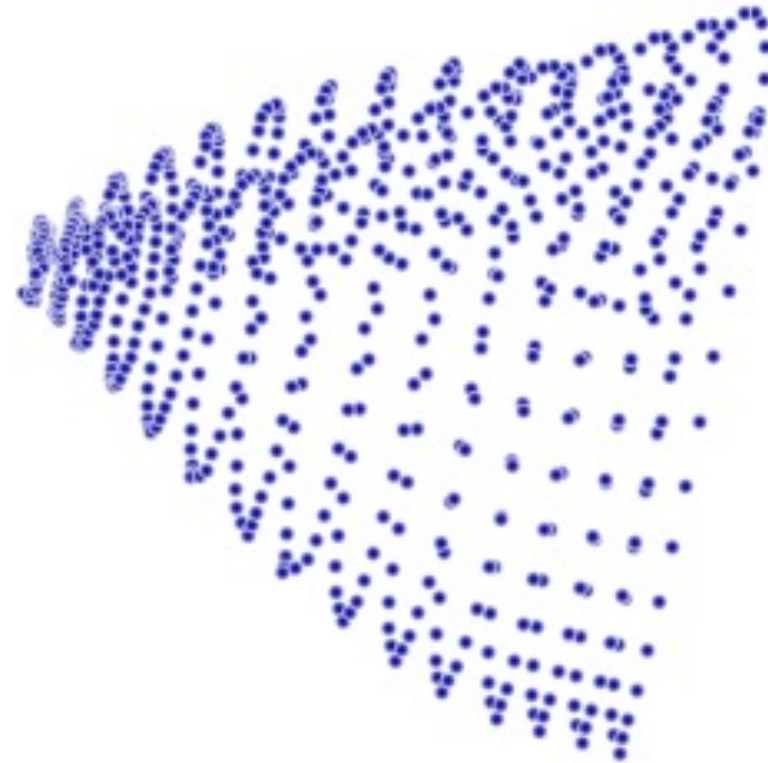
- “Neighbouring” images are close in Euclidean distance.
 - But distances between very-different images are not reliable.

Learning Manifolds

- With usual distances, **PCA/MDS do not discover non-linear manifolds.**



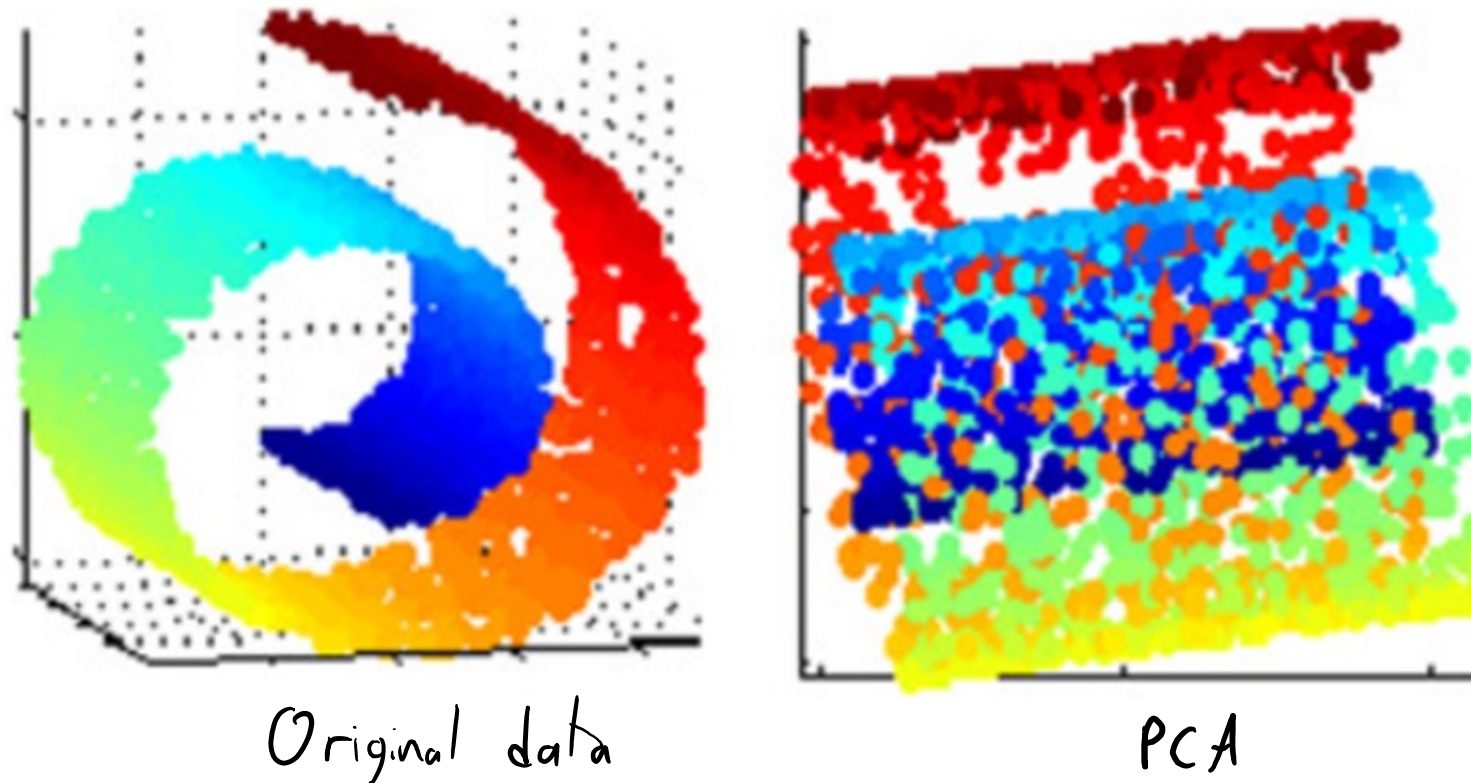
Original data



PCA

Learning Manifolds

- With usual distances, **PCA/MDS do not discover non-linear manifolds.**



- We could use **change of basis** or **kernels**: but **still need to pick basis.**

Sammon's Mapping

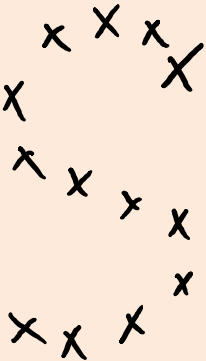
- Challenge for most MDS models: they **focus on large distances**.
 - Leads to “crowding” effect like with PCA.
- Early attempt to address this is **Sammon's mapping**:
 - **Weighted MDS** so large/small distances are more comparable.

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

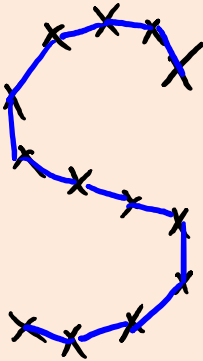
- Denominator **reduces focus on large distances**.

ISOMAP

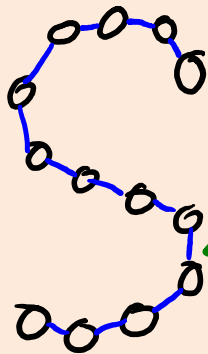
- **ISOMAP** is latent-factor model for visualizing data on manifolds:



find "neighbours"
of each point



Represent points
and neighbours
as a weighted
graph.



"weight" on each
edge is distance
between points

Approximate geodesic distance
by shortest path through
graph.

ISOMAP z_i values in 1D or 2D

Run MDS
with these
approximate geodesic distances.

$$D = \begin{bmatrix} 0 & 1 & 2 & 3 & \dots \\ 1 & 0 & 1 & 2 & \dots \\ 2 & 1 & 0 & 1 & \dots \\ 3 & 2 & 1 & 0 & \dots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix}$$

Sammon's Map vs. ISOMAP vs. PCA (MNIST)



Sammon Map

ISOMAP

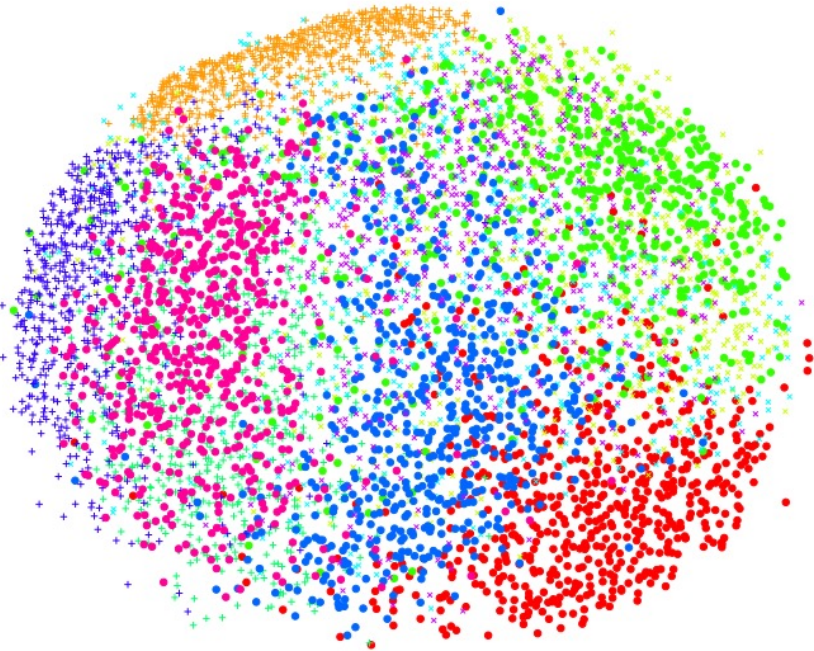
PCA



Sammon's Map vs. ISOMAP vs. t -SNE (MNIST)



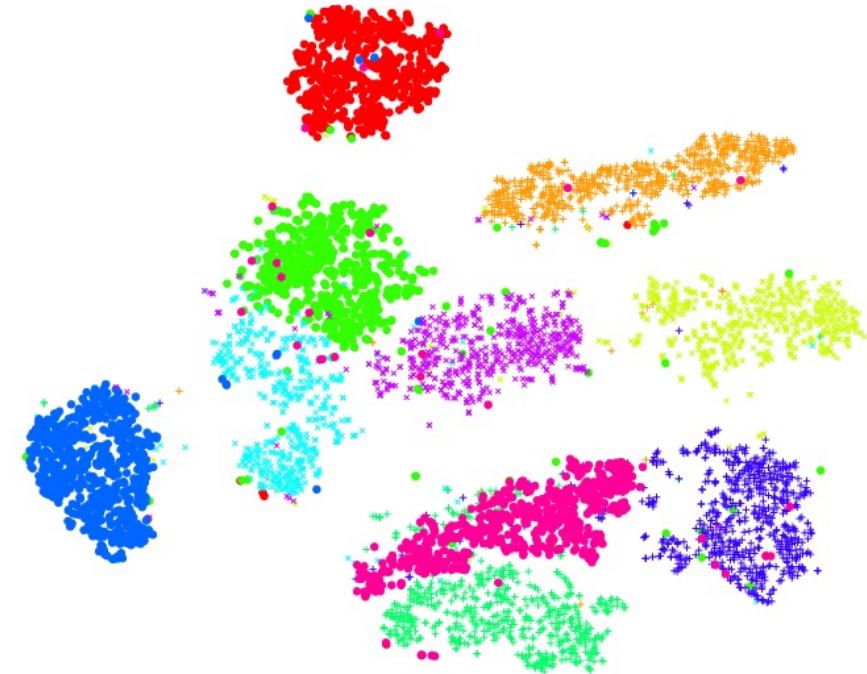
Sammon Map



ISOMAP

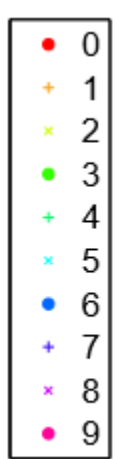


t -SNE

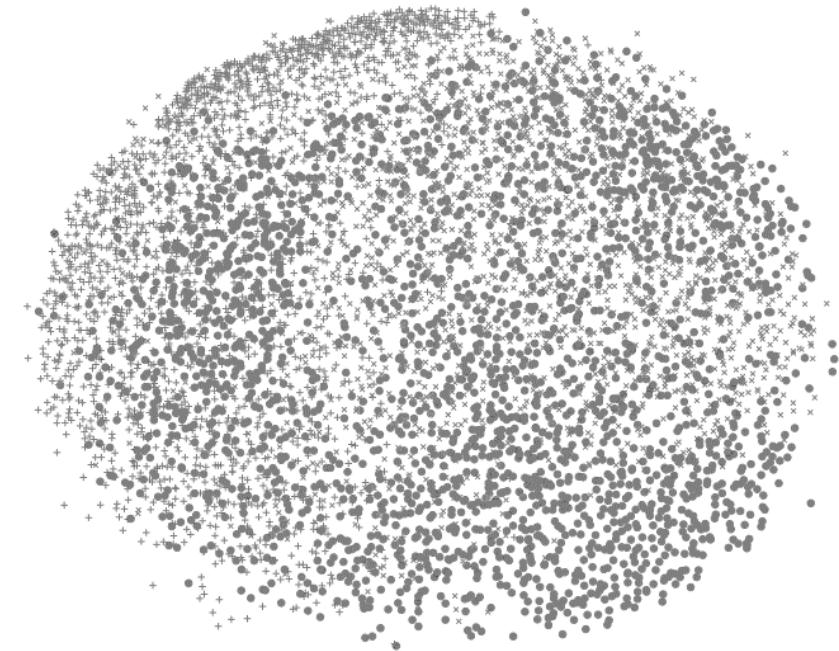


- A 'modern' way to visualize manifolds and clusters is t -SNE.

Sammon's Map vs. ISOMAP vs. t-SNE (MNIST)



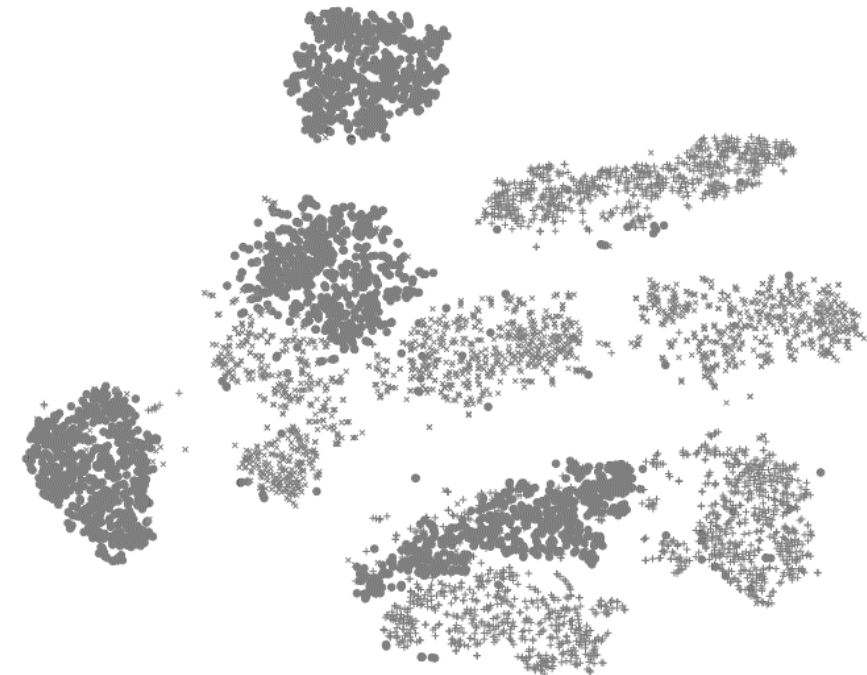
Sammon Map



ISOMAP



t-SNE

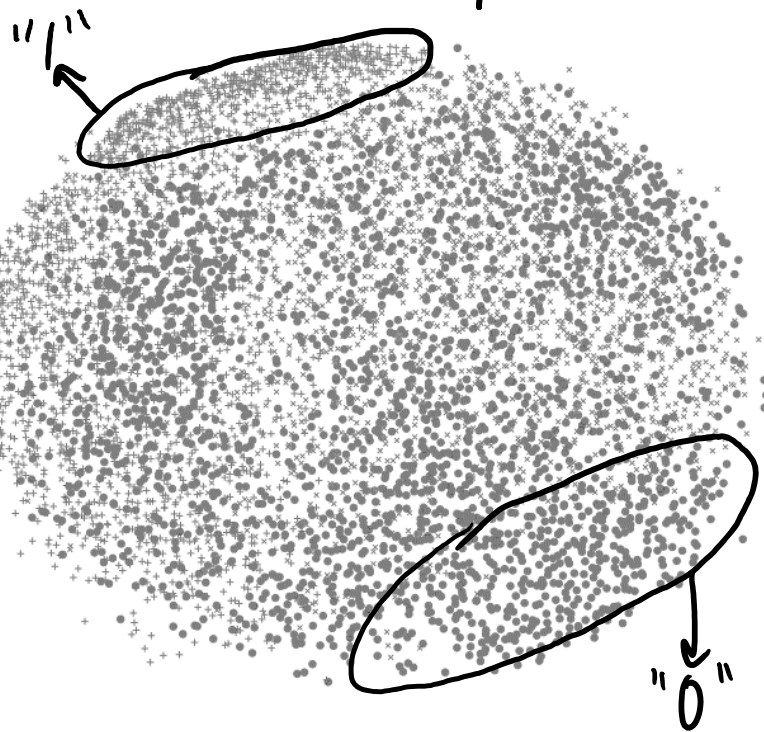


Remember this is unsupervised, algorithms do not know the labels.

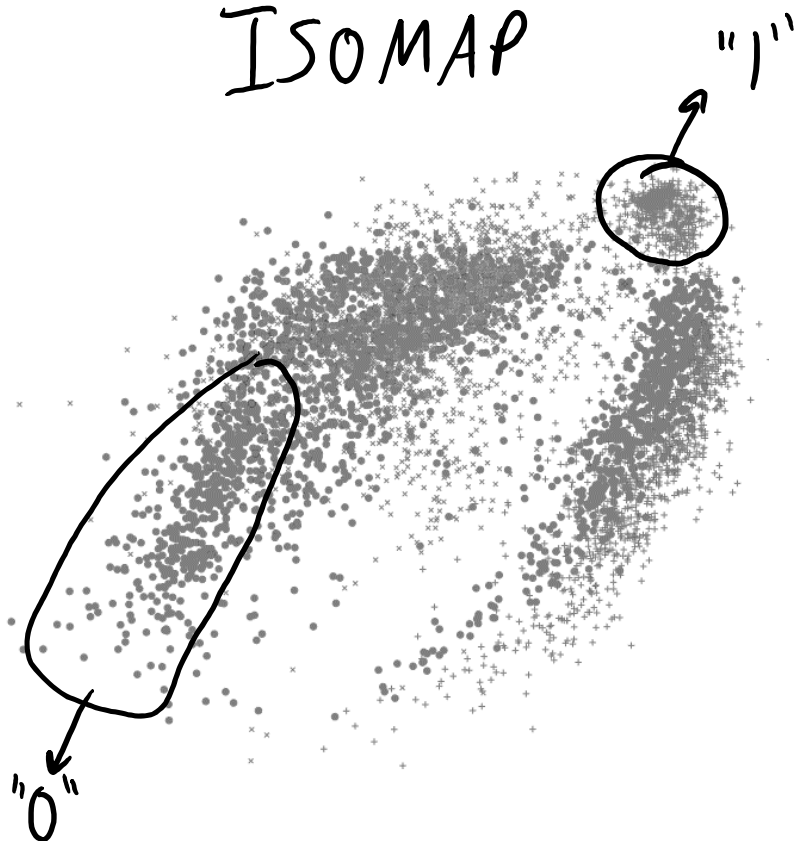
Sammon's Map vs. ISOMAP vs. t-SNE (MNIST)

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

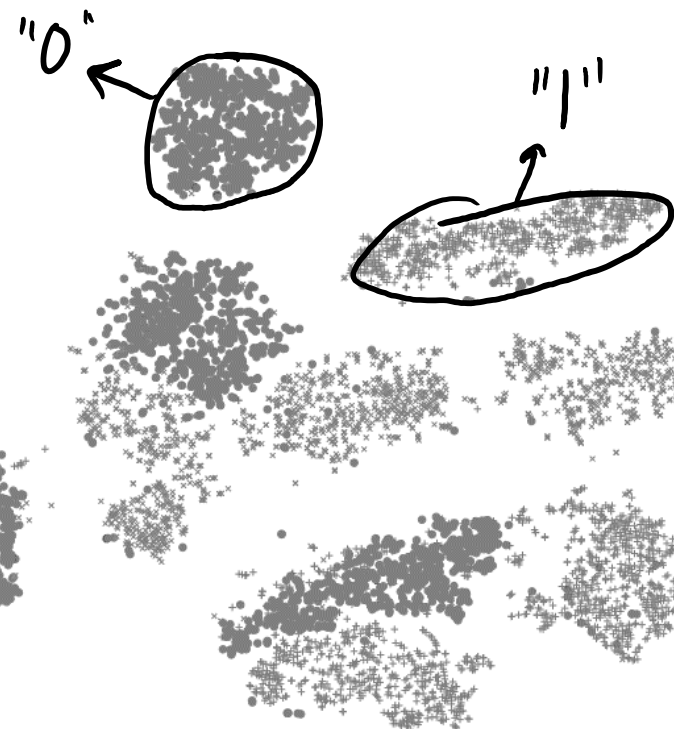
Sammon Map



ISOMAP



t-SNE

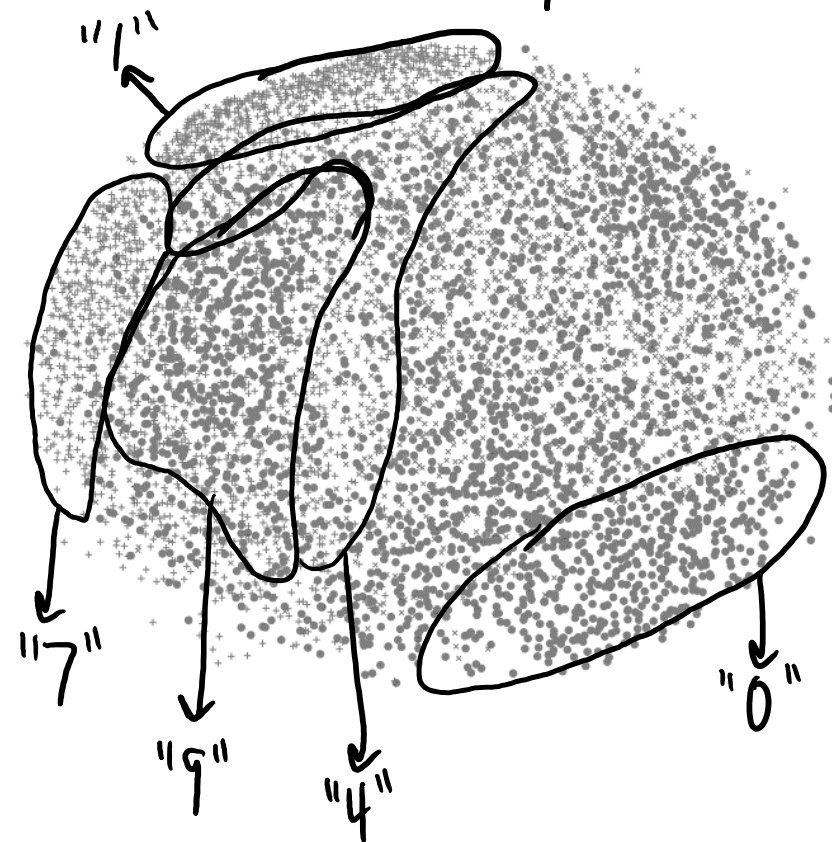


Remember this is unsupervised, algorithms do not know the labels.

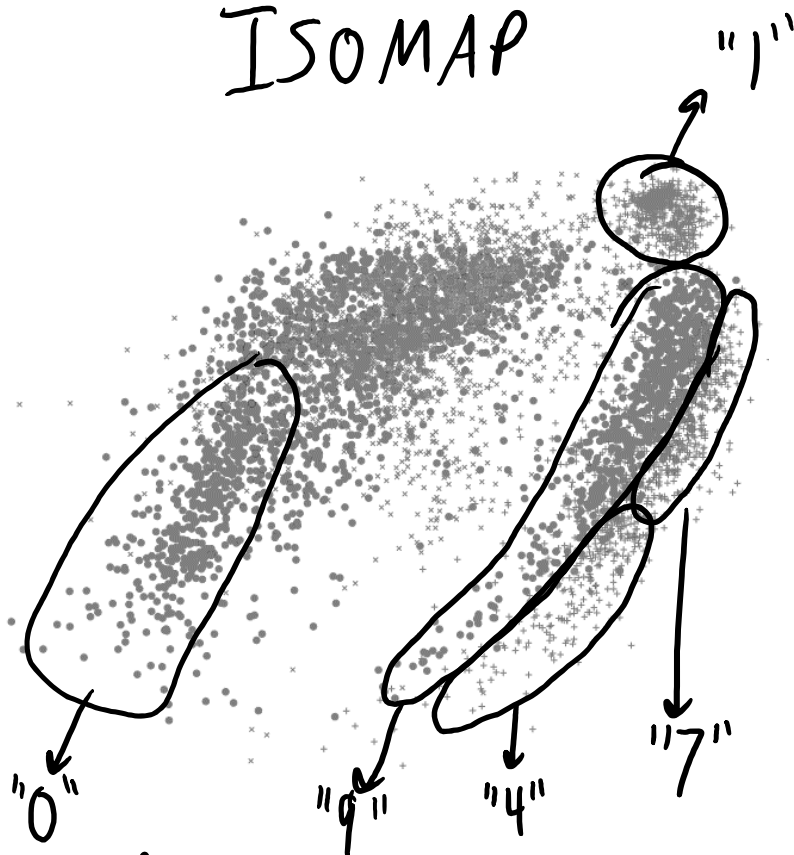
Sammon's Map vs. ISOMAP vs. t-SNE (MNIST)

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

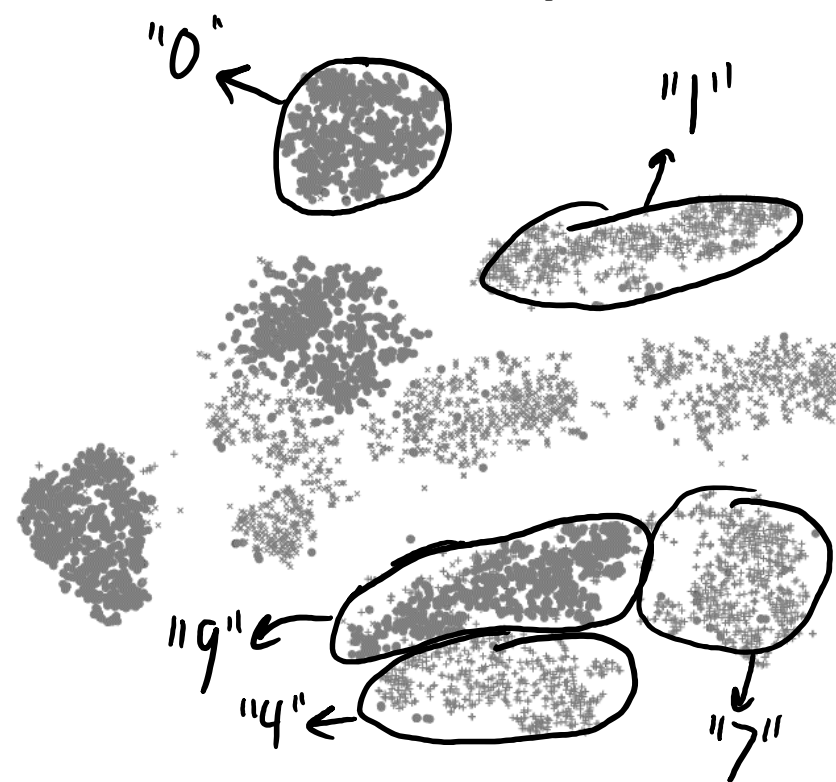
Sammon Map



ISOMAP



t-SNE

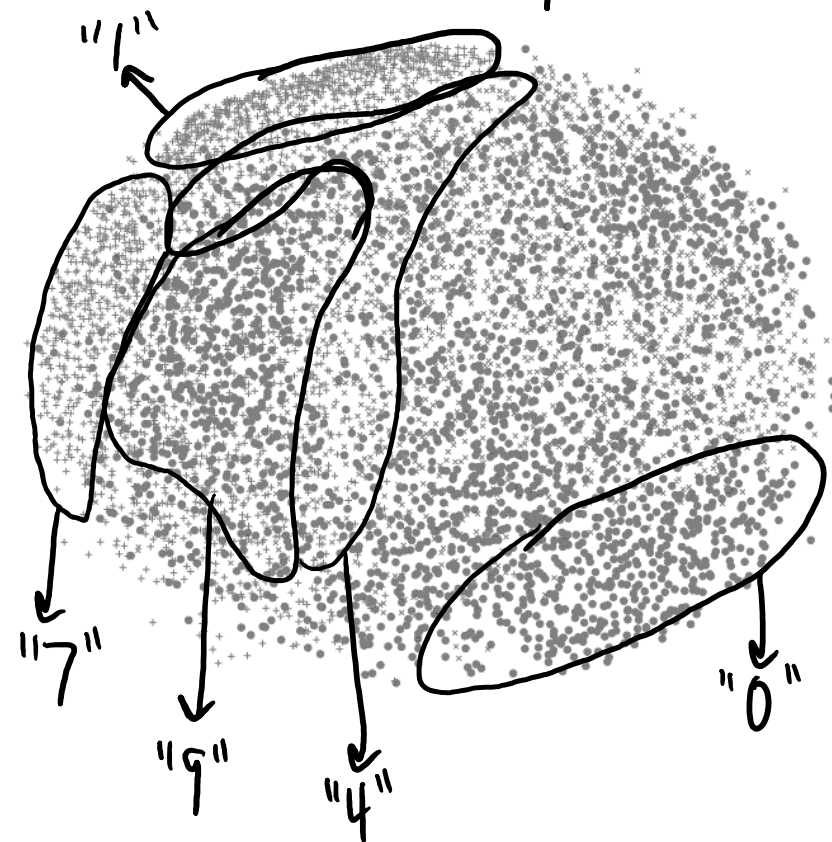


Remember this is unsupervised, algorithms do not know the labels.

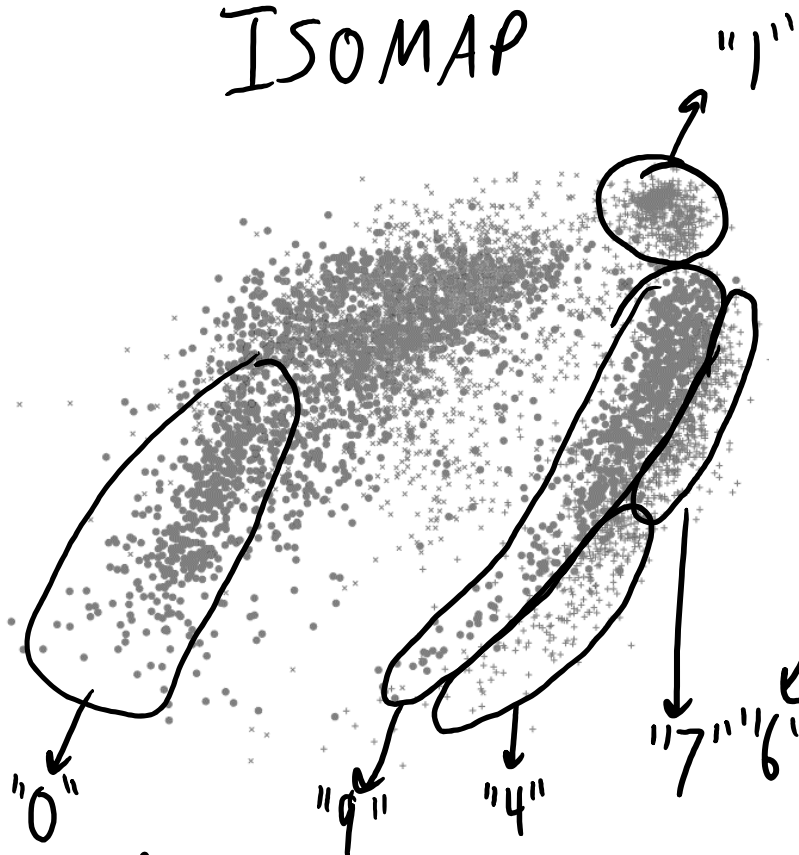
Sammon's Map vs. ISOMAP vs. t-SNE (MNIST)

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

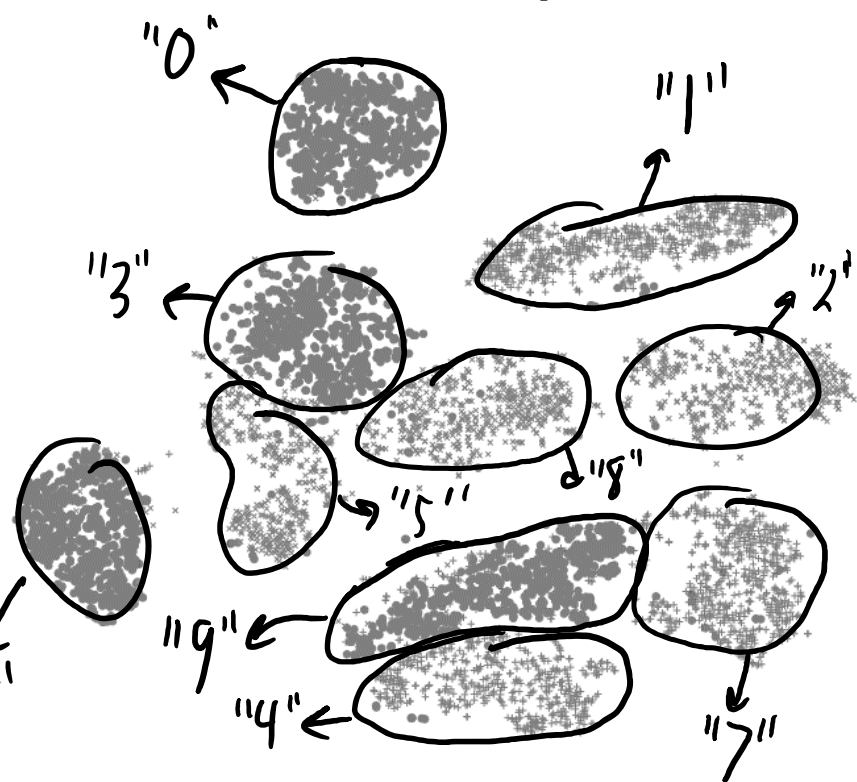
Sammon Map



ISOMAP



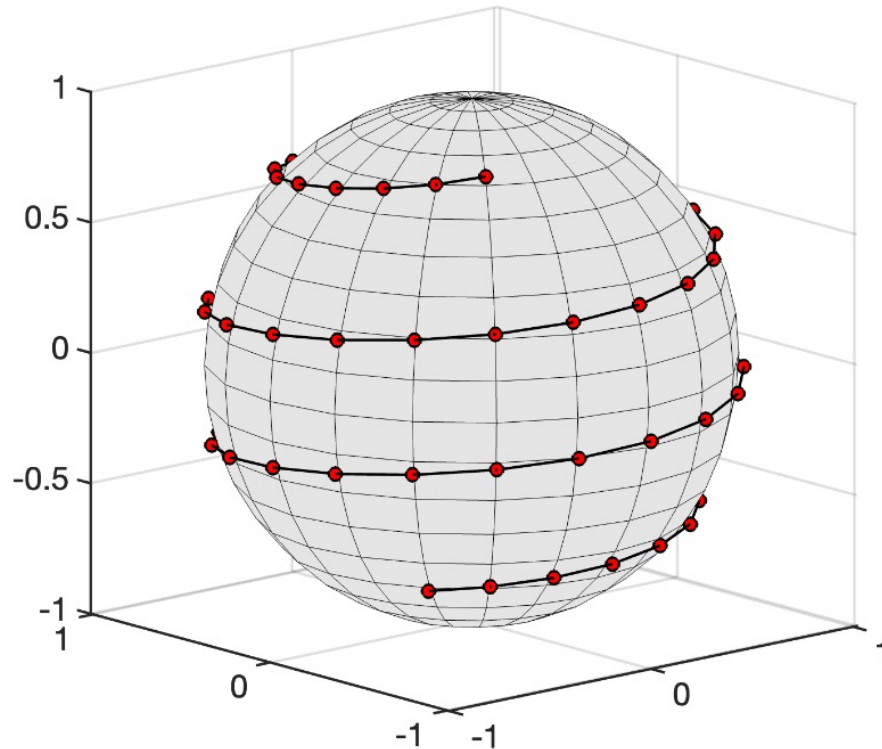
t-SNE



Remember this is unsupervised, algorithms do not know the labels.

t -Distributed Stochastic Neighbour Embedding

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



t -Distributed Stochastic Neighbour Embedding

Visualizing Data using t-SNE

Laurens van der Maaten

TiCC

Tilburg University

P.O. Box 90153, 5000 LE Tilburg, The Netherlands

LVDMAATEN@GMAIL.COM

Geoffrey Hinton

Department of Computer Science

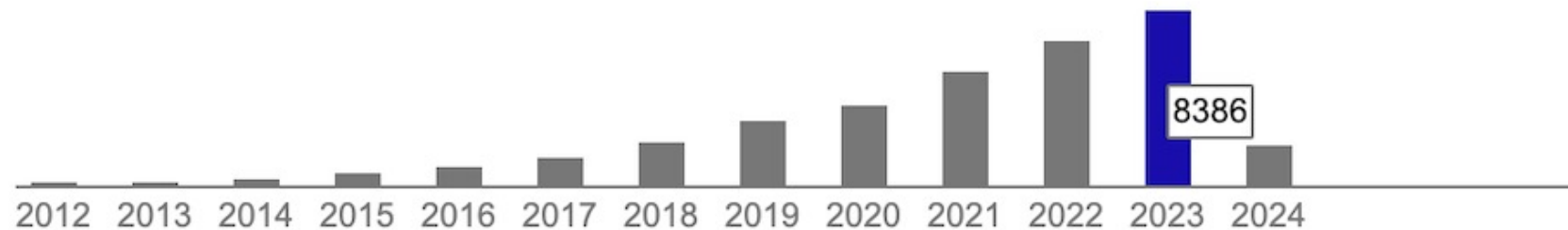
University of Toronto

6 King's College Road, M5S 3G4 Toronto, ON, Canada

HINTON@CS.TORONTO.EDU

Editor: Yoshua Bengio

Cited by 35319



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

The Loss Function

- The Kullback–Leibler divergence between the affinity (similarity) matrix \mathbf{P} from the high-dimensional data and the affinity matrix from the low-dimensional data \mathbf{Q}

$$\text{KL} (P \parallel Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

- The loss is optimized via gradient descent.
- Keep nearby data points in the high-dimensional space nearby in the low-dimensional space, **while push all data points in the low-dimensional space apart from each other.**

The High-dimensional Affinity Matrix

- The high-dimensional affinity matrix

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)}$$

- Typically symmetrize and normalize to be a probability mass function

$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2N}$$

- The data point-dependent parameter σ_i is adaptively calculated to achieve the desired *perplexity* (30 by default)

$$2^{-\sum_j p_{j|i} \log_2 p_{j|i}}$$

The Low-dimensional Affinity Matrix

- The low-dimensional affinity matrix q_{ij} :

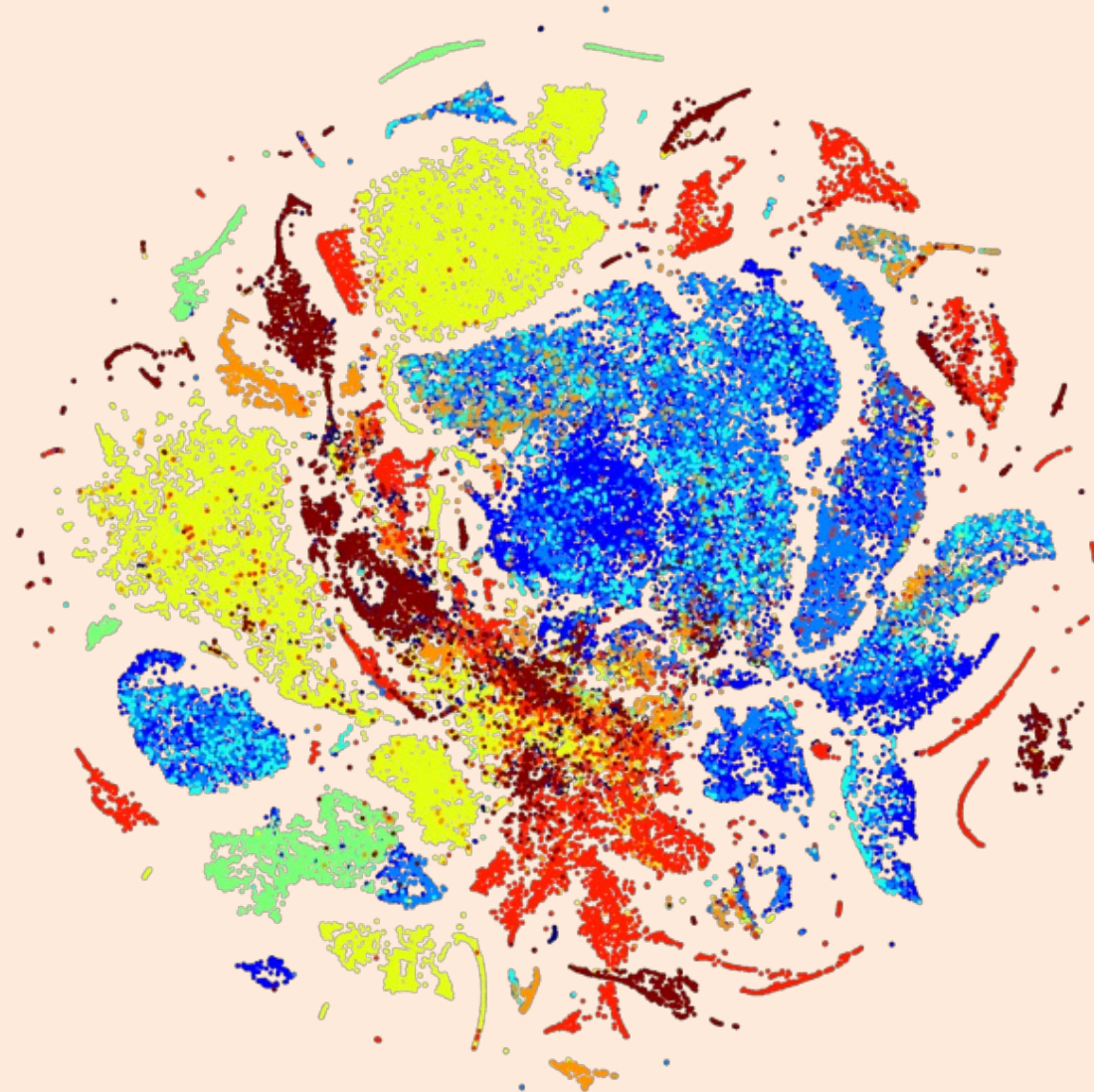
$$\frac{(1 + \|\mathbf{z}_i - \mathbf{z}_j\|^2 / \nu)^{-\frac{\nu+1}{2}}}{\sum_{k, k \neq i} (1 + \|\mathbf{z}_i - \mathbf{z}_k\|^2 / \nu)^{-\frac{\nu+1}{2}}}$$

- Here we typically use the Student's t distribution with $\nu=1$ (the Cauchy distribution to measure the similarity between points).

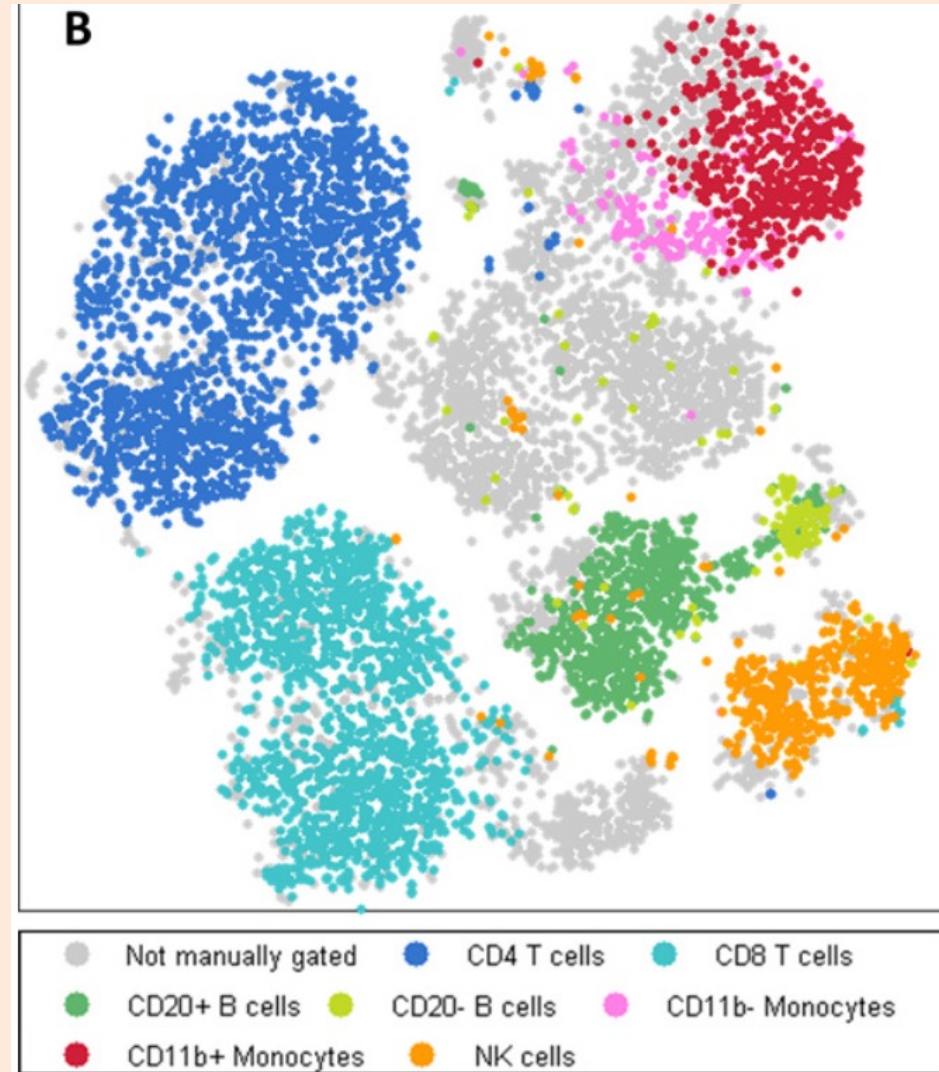
Other Details

- T-SNE is sensitive to initialization, typically we initial Z by PCA.
- To compute the **high-dimensional affinity** matrix can be slow $O(n^2d)$, we use approximate k-NN search to only compute the affinities between a point and its k-NNs ($k = 3 * \textit{perplexity}$).
- We set both p_{ij} and $q_{ij} = 0$ (only pairwise similarities are of interest).
- Optimization trick (early exaggeration – multiply the attractive force by 12 for the first 250 iterations).
- Speedup calculating the repulsive force (the FFT-SNE algorithm).
- A more recent nonlinear dimension reduction tool: UMAP ([Uniform manifold approximation and projection for dimension reduction](#), published in 2018, >10k citations).

t -SNE on Product Features



t -SNE on Leukemia Heterogeneity

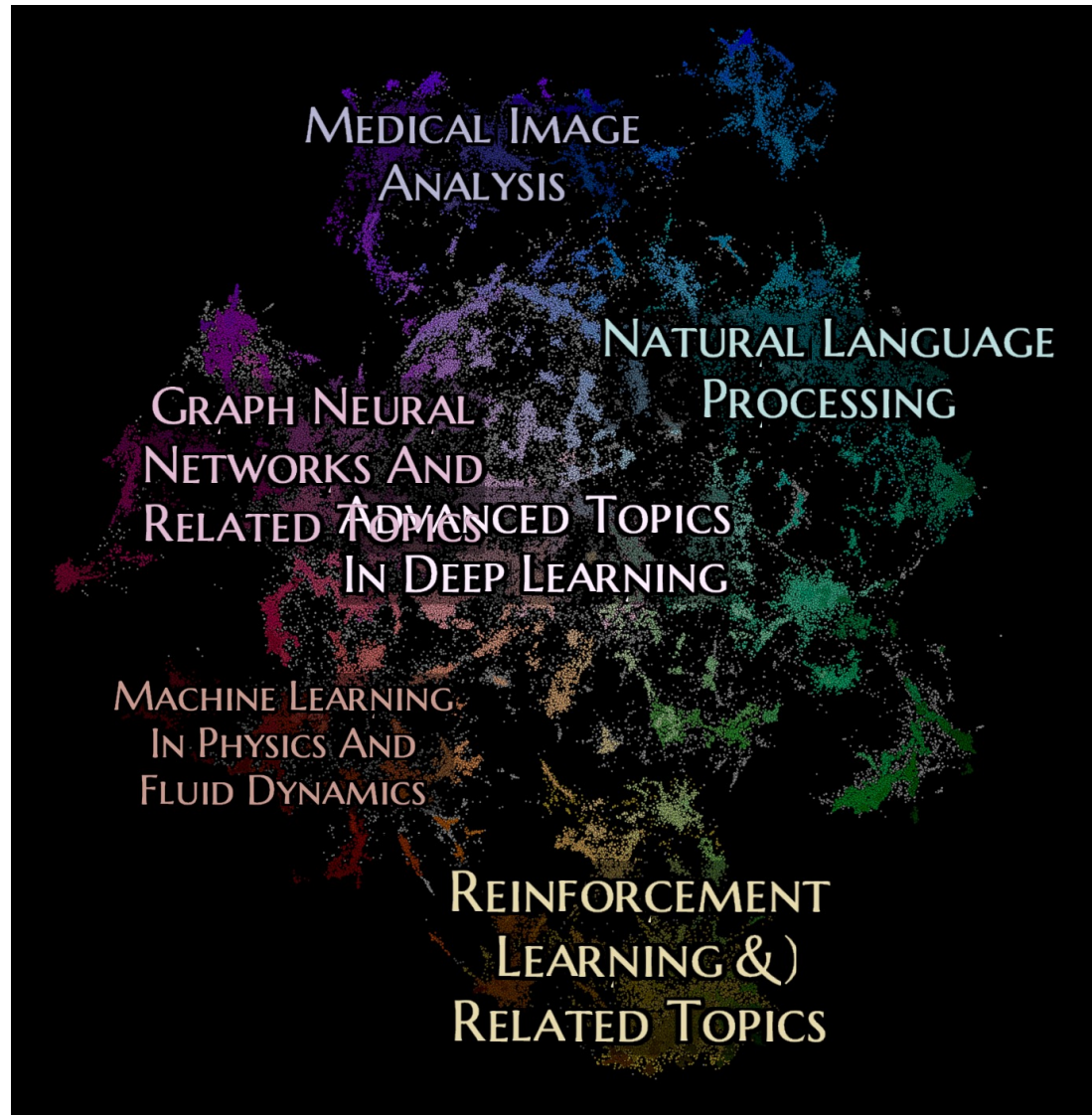


UMAP on Mouse Brain Data

- 4 million single-cell transcriptomes from adult mouse brain labeled by source brain region.



ArXiv Machine Learning Landscape



https://lmcinnes.github.io/datamapplot_examples/ArXiv_data_map_example.htm

Next Topic: Word2Vec

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” features be **related** in some way?
- We want a **latent-factor representation** of individual words:
 - Closeness in latent space should indicate similarity.
 - Distances could represent meaning?
- Recent alternative to PCA is **word2vec**...

Using Context

- Consider these phrases:
 - “the cat purred”
 - “the kitten purred”

 - “black cat ran”
 - “black kitten ran”
- Words that occur in the same context likely have similar meanings.
- **Word2vec** uses this insight to design an **MDS distance function**.

Word2Vec (Continuous Bag of Words)

- A common **word2vec** approaches (called **continuous bag of words**):
 - Each **word 'i'** is represented by a vector of real numbers z_i .
 - Training data: sentence fragments with **“hidden” middle word**:
 - “We introduce basic principles and techniques in”
 - “the fields of data mining and machine”
 - “tools behind the emerging field of data”
 - “techniques are now running behind the scenes”
 - “discover patterns and make predictions in various”
 - “the core data mining and machine learning”
 - “with motivating applications from a variety of”
 - Train so that z_i of **“hidden” words** are similar to z_i of surrounding words.

Word2Vec (Continuous Bag of Words)

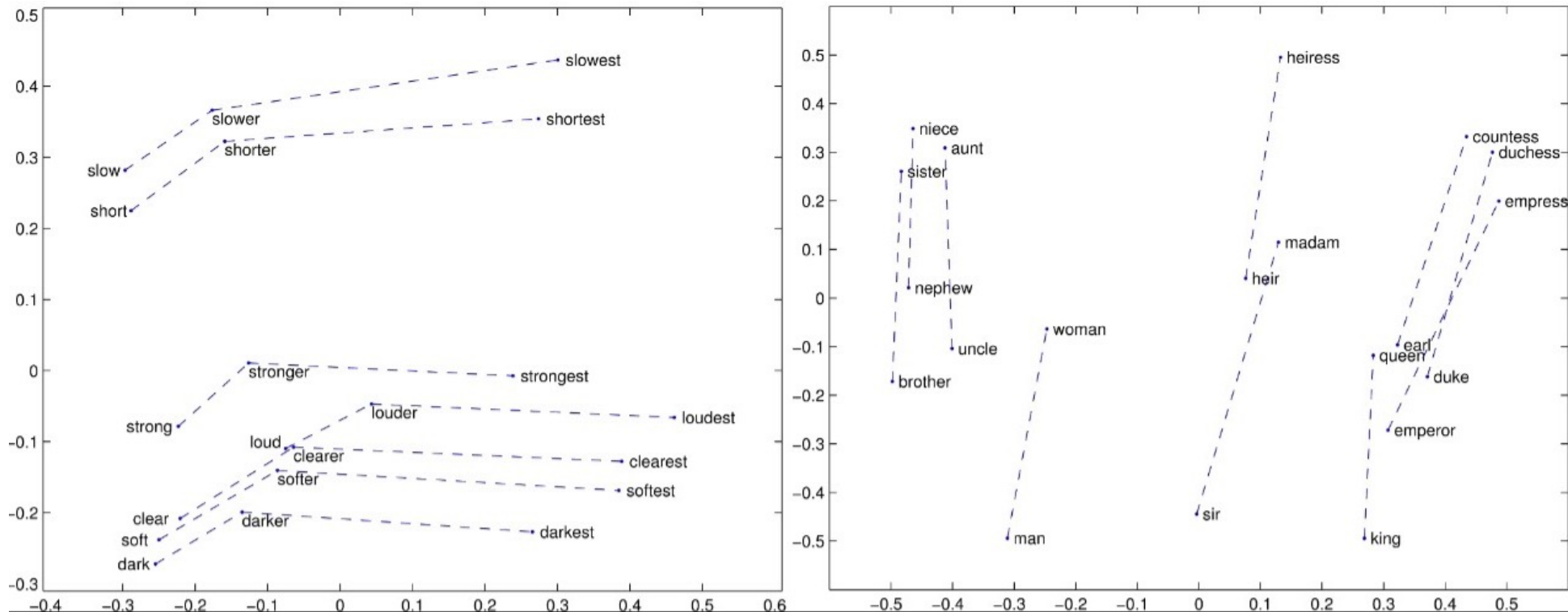
- Continuous bag of words model probability of middle word 'i' as:

$$\prod_{j \in \text{surrounding words}} \frac{\exp(z_i^T z_j)}{\sum_{c=1}^{\# \text{ words}} \exp(z_c^T z_j)}$$

- We use gradient descent on negative logarithm of these probabilities:
 - Makes $z_i^T z_j$ big for words appearing in same context (making z_i close to z_j).
 - Makes $z_i^T z_j$ small for words not appearing together (makes z_i and z_j far).
- Once trained, you use these z_i as features for language tasks.
 - Tends to work much better than bag of words.
 - Allows you to get useful features of words from unlabeled text data.

Word2Vec Example

- MDS visualization of a set of related words:



- Distances between vectors might represent semantics.

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](#).

Summary

- **Multi-dimensional scaling** is a non-parametric latent-factor model.
- **Different MDS distances/losses/weights** usually gives better results.
- **Manifold**: space where local Euclidean distance is accurate.
 - Structured data like images often form manifolds in space.
- ***t*-SNE** is an MDS method focusing on matching small distances.
- **Word2vec**:
 - Latent-factor (continuous) representation of words.
 - Based on predicting word from its context (or context from word).
- Next time: Neural Networks.

Word2Vec (Skip-Gram)

- A common **word2vec** approaches (**skip gram**):
 - Each **word 'i'** is represented by a vector of real numbers z_i .
 - Training data: sentence fragments with “hidden” **surrounding** word:
 - ~~“We introduce basic principles and techniques in”~~
 - ~~“the fields of data mining and machine”~~
 - ~~“tools behind the emerging field of data”~~
 - ~~“techniques are now running behind the scenes”~~
 - ~~“discover patterns and make predictions in various”~~
 - ~~“the core data mining and machine learning”~~
 - ~~“with motivating applications from a variety of”~~
 - Train so that z_i of “hidden” words are similar to z_i of surrounding words.
 - Uses same probability as continuous bag of words.
 - But **denominator sums over all possible surrounding** words (often just sample terms for speed).

Stochastic Gradient for SVDfeature

- Common approach to fitting SVDfeature is **stochastic gradient**.
- Previously you saw stochastic gradient for supervised learning:
 - Choose a random example 'i'
 - Update parameters 'w' using gradient of example 'i'
- **Stochastic gradient for SVDfeature** (formulas as bonus):
 - Choose a random user 'i' and a random product 'j'
 - Update β , β_i , β_j , w , z_i , and w^j based on their gradient for this user-product.

Updated every time



SVDfeature with SGD: the gory details

Objective: $\frac{1}{2} \sum_{(i,j) \in R} (\hat{y}_{ij} - y_{ij})^2$ with $\hat{y}_{ij} = \beta + \beta_i + \beta_j + w^T x_{ij} + (w^j)^T z_i$

Update based on random (i,j) :

$$\beta = \beta - \alpha r_{ij}$$

$$\beta_i = \beta_i - \alpha r_{ij}$$

$$\beta_j = \beta_j - \alpha r_{ij}$$

Updates are the same,

but ' β ' is always update while β_i and β_j are only updated for the specific user + product

$$w = w - \alpha r_{ij} x_{ij} \leftarrow \text{Updated every time.}$$

$$z_i = z_i - \alpha r_{ij} w^j$$

$$w^j = w^j - \alpha r_{ij} z_i$$

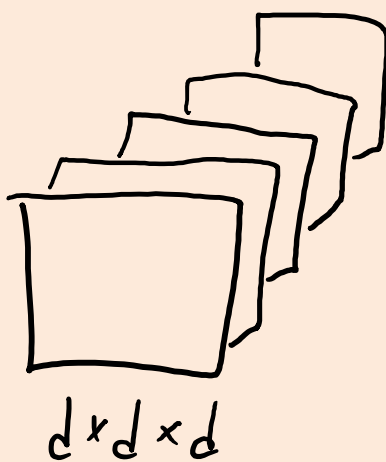
Updated for specific user and product.

(Adding regularization adds an extra term)

Tensor Factorization

- Tensors are higher-order generalizations of matrices:

Scalar $\alpha = []_{1 \times 1}$ Vector $\alpha = []_{d \times 1}$ Matrix $A = []_{d \times d}$ Tensor $A = []_{d \times d \times d}$



- Generalization of matrix factorization is **tensor factorization**:

$$y_{ijm} \approx \sum_{c=1}^k w_{jc} z_{ic} v_{mc}$$

- Useful if there are other relevant variables:
 - Instead of ratings based on {user,movie}, ratings based {user,movie,group}.
 - Useful if you have groups of users, or if ratings change over time.

Field-Aware Matrix Factorization

- **Field-aware factorization machines (FFMs):**
 - Matrix factorization with multiple z_i or w_c for each example or part.
 - You choose which z_i or w_c to use based on the value of feature.
- Example from “click through rate” prediction:
 - E.g., predict whether “male” clicks on “nike” advertising on “espn” page.
 - A previous matrix factorization method for the 3 factors used:

$$w_{espn} w_{nike} + w_{espn} w_{male} + w_{nike} w_{male}$$
$$w_{espn}^A w_{nike}^P + w_{espn}^G w_{male}^P + w_{nike}^G w_{male}^A$$

- FFMs could use:
 - w_{espn}^A is the factor we use when multiplying by an advertiser’s latent factor.
 - w_{espn}^G is the factor we use when multiplying by a group’s latent factor.
- This approach has won some Kaggle competitions ([link](#)), and has shown to work well in production systems too ([link](#)).

Warm-Starting

- We've used data $\{X,y\}$ to fit a model.
- We now have **new training data** and **want to fit new and old data**.
- Do we need to re-fit from scratch?
- This is the **warm starting** problem.
 - It's easier to warm start some models than others.

Easy Case: K-Nearest Neighbours and Counting

- K-nearest neighbours:

- KNN just stores the training data, so just **store the new data**.

- Counting-based models:

- Models that base predictions on frequencies of events.

- E.g., naïve Bayes.

- Just **update the counts**:
$$p(\text{"vicodin"} | \text{"spam"}) = \frac{\text{count of } \{\text{vicodin, spam}\} \text{ in } \underline{\text{new and old data}}}{\text{count of "spam" in } \underline{\text{new and old data}}}$$

- Decision trees with fixed rules: just update counts at the leaves.

Medium Case: L2-Regularized Least Squares

- L2-regularized least squares is obtained from linear algebra:

$$w = (X^T X + \lambda I)^{-1} (X^T y)$$

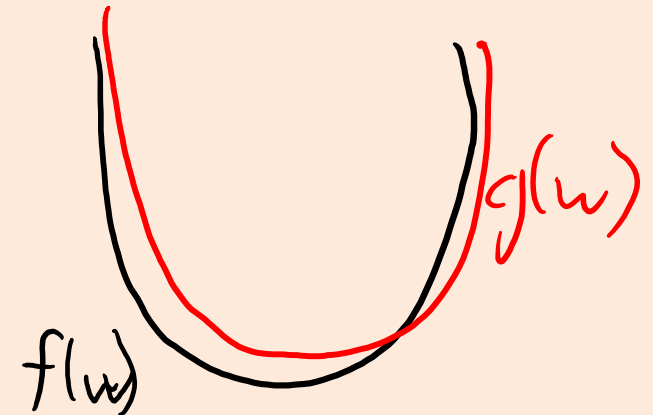
- Cost is $O(nd^2 + d^3)$ for ‘n’ training examples and ‘d’ features.
- Given one new point, we need to compute:
 - $X^T y$ with one row added, which costs $O(d)$.
 - Old $X^T X$ plus $x_i x_i^T$, which costs $O(d^2)$.
 - Solution of linear system, which costs $O(d^3)$.
 - So cost of adding ‘t’ new data point is $O(td^3)$.
- With “matrix factorization updates”, can reduce this to $O(td^2)$.
 - Cheaper than computing from scratch, particularly for large d.

Medium Case: Logistic Regression

- We fit **logistic regression** by **gradient descent** on a convex function.
- With new data, convex function $f(w)$ changes to new function $g(w)$.

$$f(w) = \sum_{i=1}^n f_i(w) \qquad g(w) = \sum_{i=1}^{n+1} f_i(w)$$

- If we don't have much more data, 'f' and 'g' will be "close".
 - Start gradient descent on 'g' with minimizer of 'f'.
 - You can show that it **requires fewer iterations**.



Hard Cases: Non-Convex/Greedy Models

- For **decision trees**:
 - “Warm start”: continue splitting nodes that haven’t already been split.
 - “Cold start”: re-fit everything.
- Unlike previous cases, this **won’t in general give same result as re-fitting**:
 - New data points might lead to **different splits** higher up in the tree.
- Intermediate: usually do warm start but occasionally do a cold start.
- Similar heuristics/conclusions for other non-convex/greedy models:
 - **K-means clustering**.
 - **Matrix factorization** (though you can continue PCA algorithms).

Different MDS Cost Functions

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

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

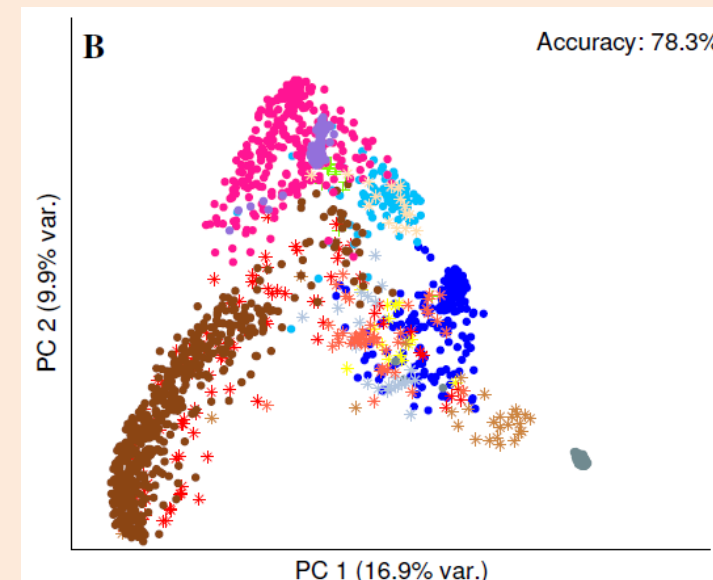
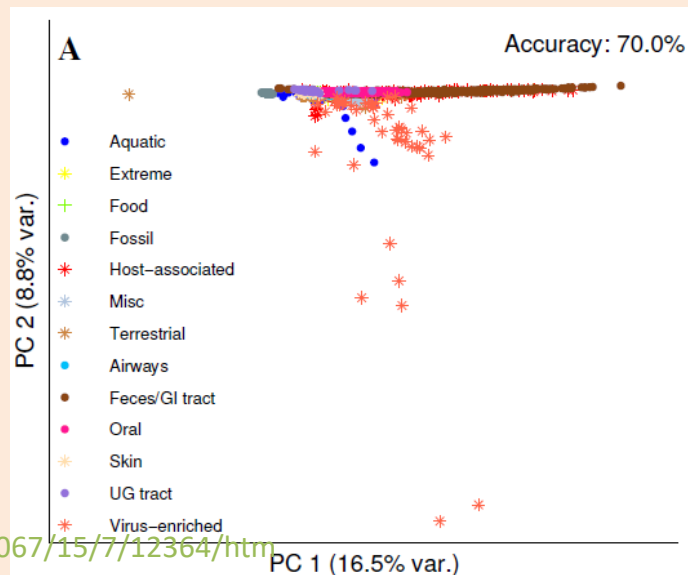
- A possibility is **“classic” MDS** with $d_1(x_i, x_j) = x_i^T x_j$ and $d_2(z_i, z_j) = z_i^T z_j$.
 - We obtain **PCA in this special case** (centered x_i , d_3 as the squared L2-norm).
 - Not a great choice because it's **a linear model**.

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 $d_2(z_i, z_j) = ||z_i - z_j||$.
 - The z_i approximate the high-dimensional L_1 -norm distances.



Sammon's Mapping

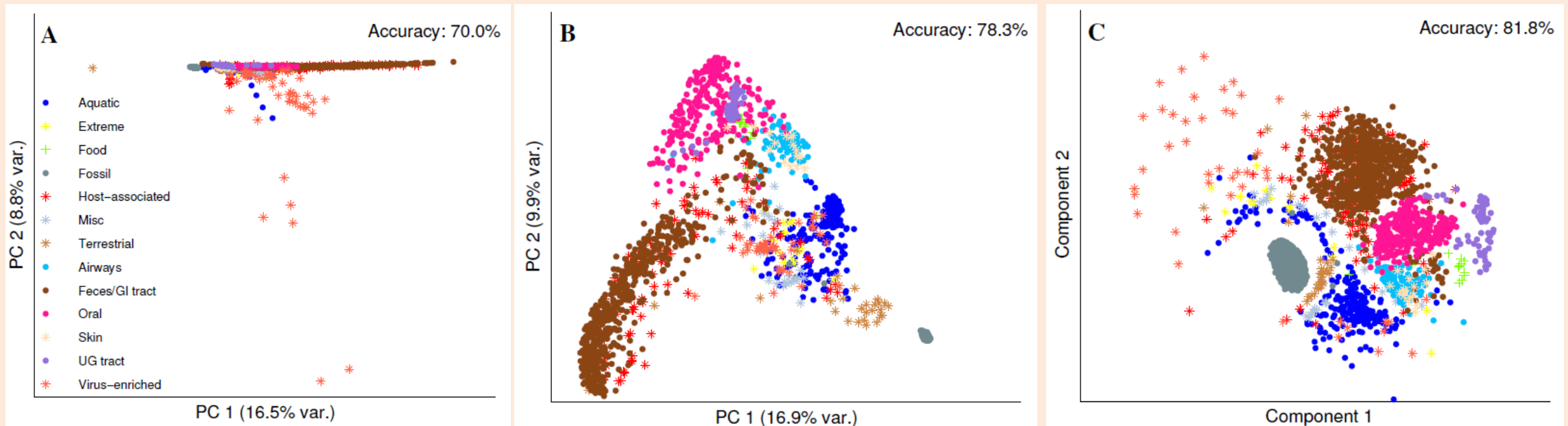
- Challenge for most MDS models: they **focus on large distances**.
 - Leads to “crowding” effect like with PCA.
- Early attempt to address this is **Sammon's mapping**:
 - **Weighted MDS** so large/small distances are more comparable.

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

- Denominator **reduces focus on large distances**.

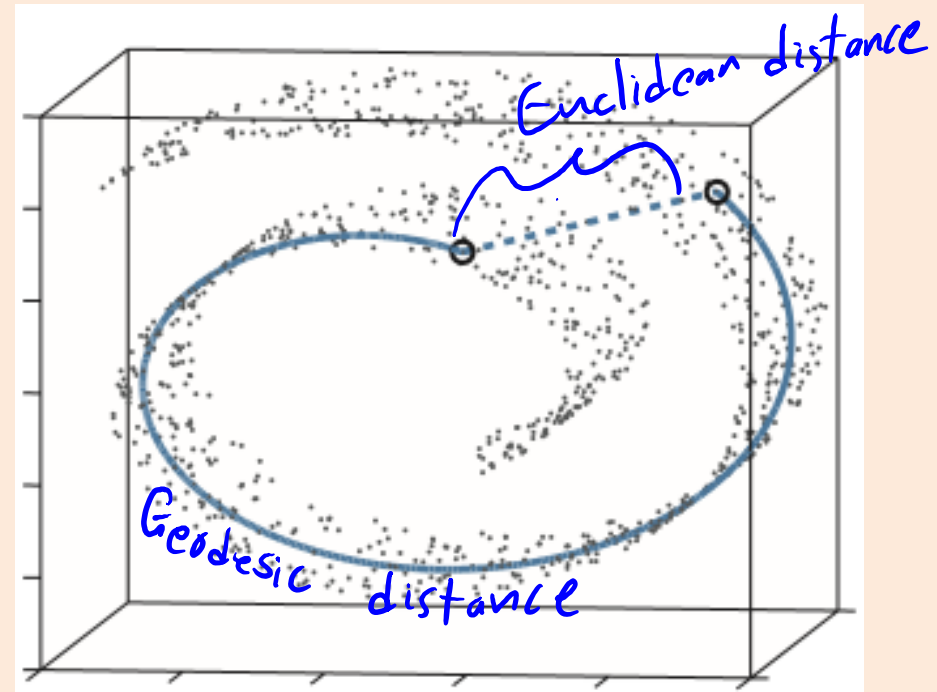
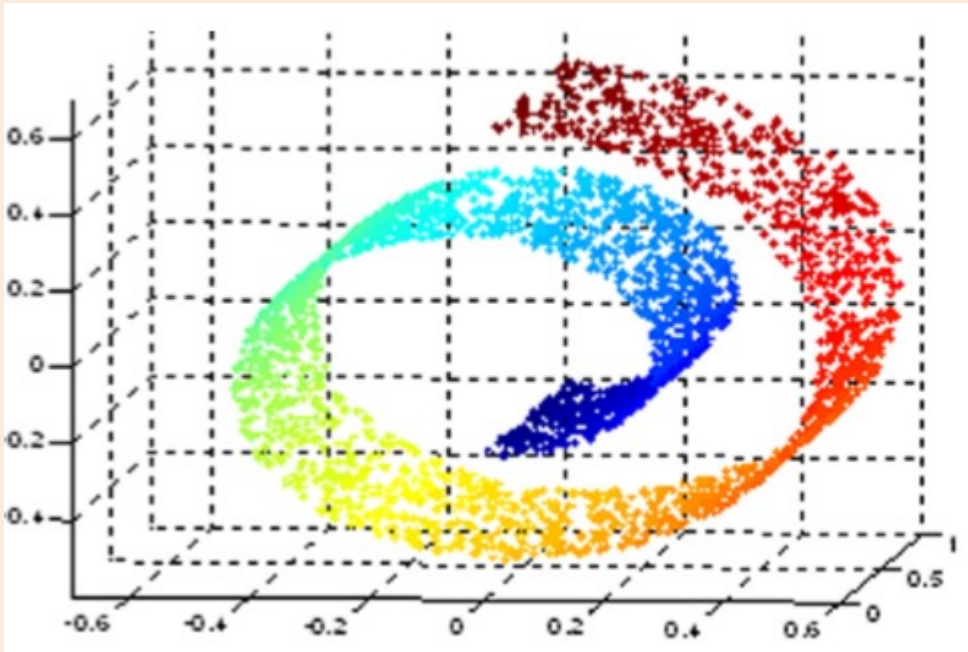
Sammon's Mapping

- Challenge for most MDS models: they **focus on large distances**.
 - 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.



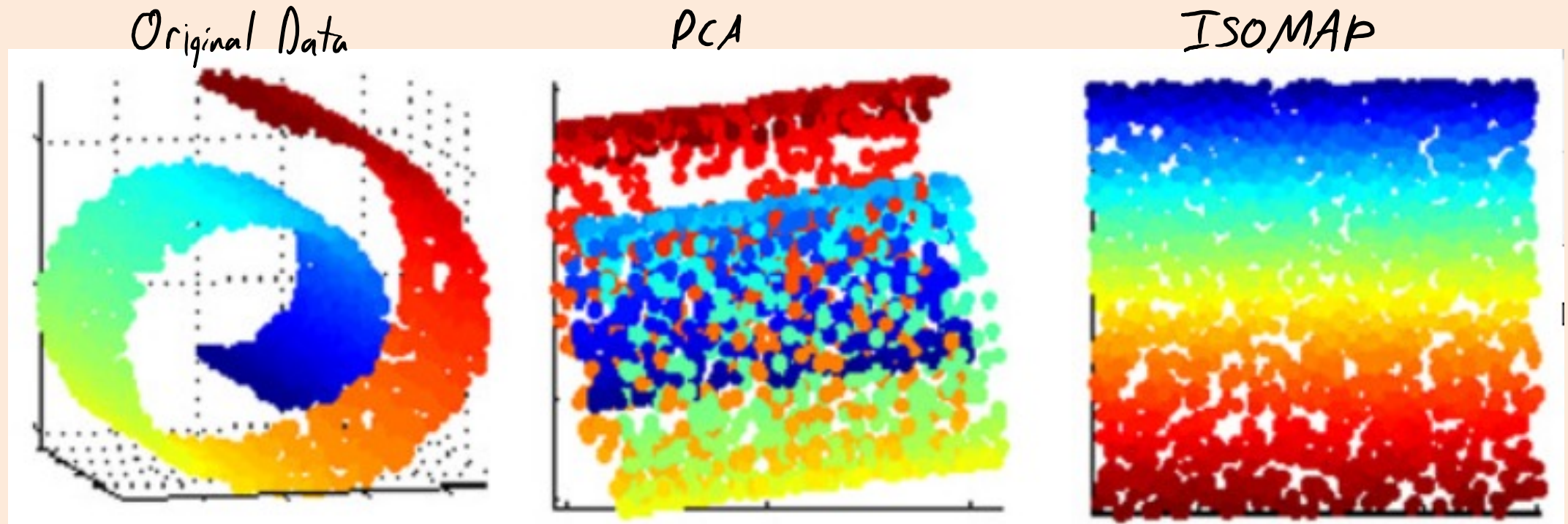
Geodesic Distance on Manifolds

- 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 **distance *through* the manifold**.



ISOMAP

- ISOMAP can “unwrap” the roll:
 - Shortest paths are approximations to geodesic distances.



- Sensitive to having the right graph:
 - Points off of manifold and gaps in manifold cause problems.

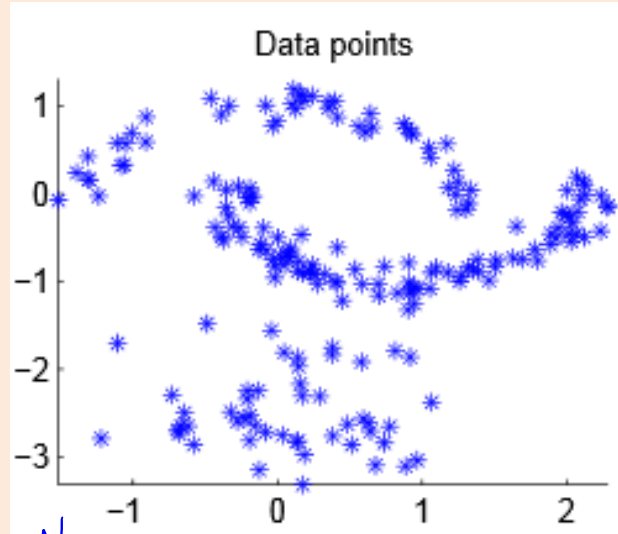
Constructing Neighbour Graphs

- 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**:
 - Approach 1 (“**epsilon graph**”): connect x_i to all x_j within some threshold ϵ .
 - Like we did with density-based clustering.
 - Approach 2 (“**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 2 (“**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 .

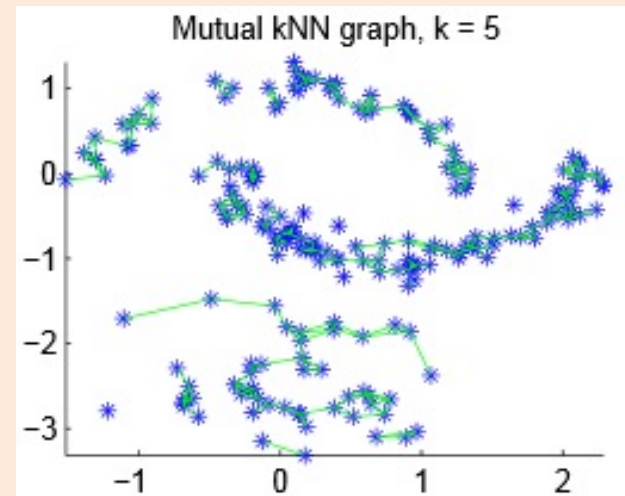
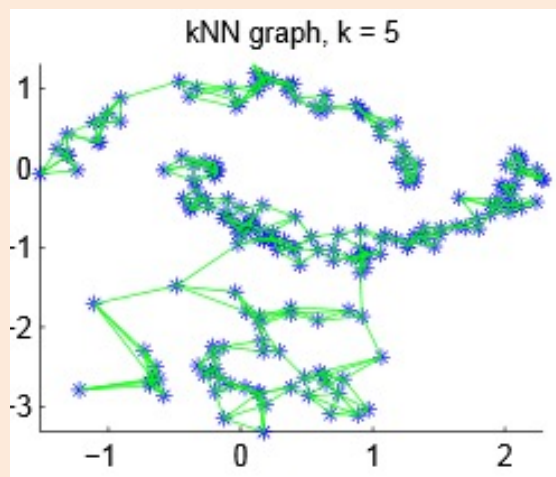
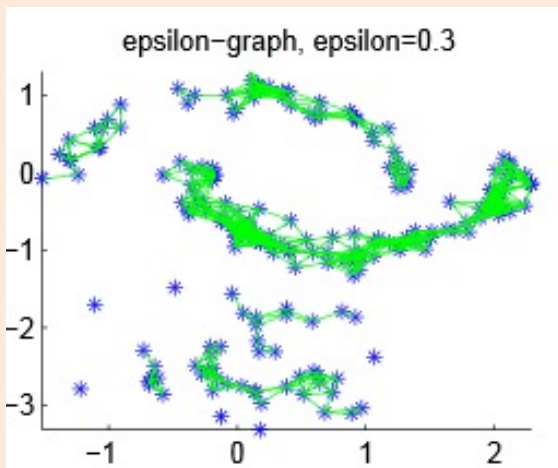
Converting from Features to Graph

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

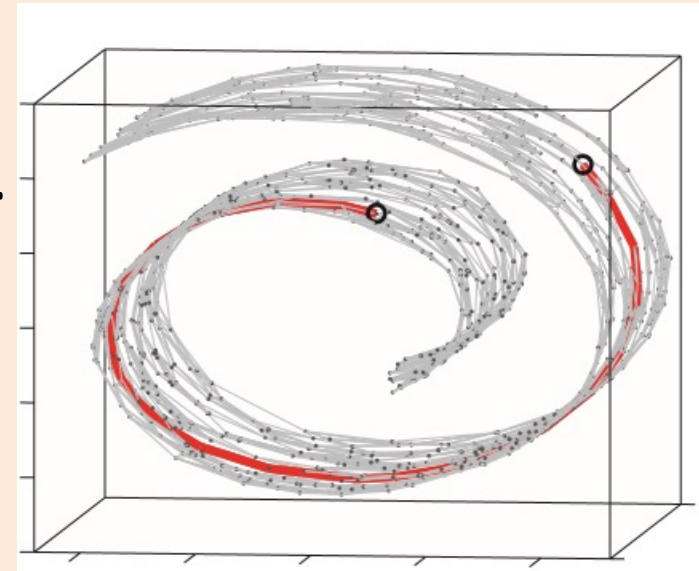


add edge if
 i and j
are kNNs
of each other.



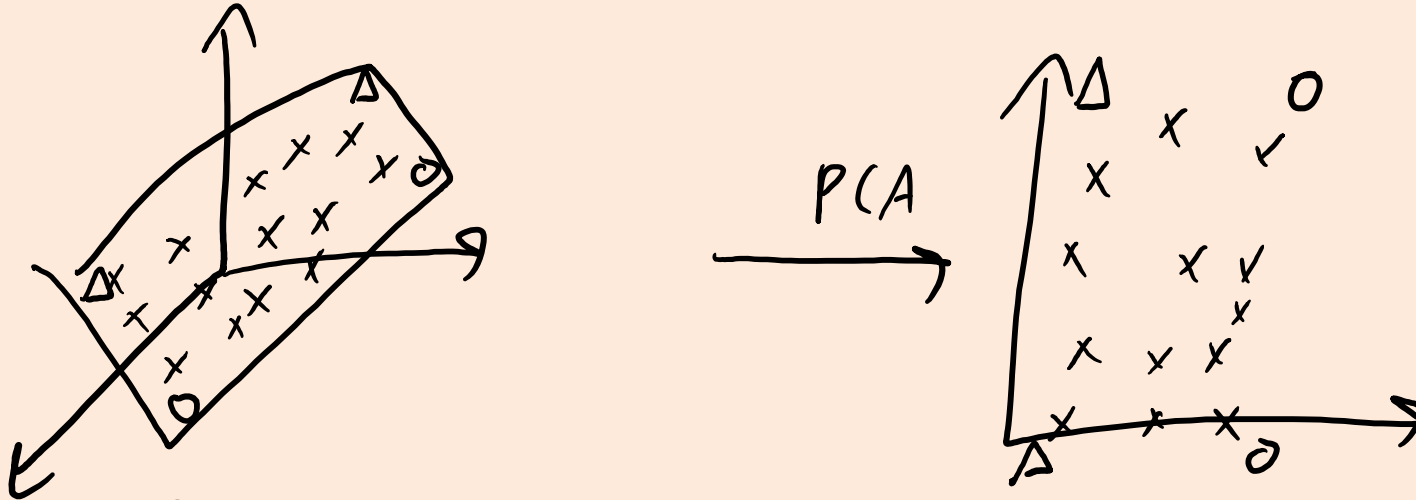
ISOMAP

- **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.

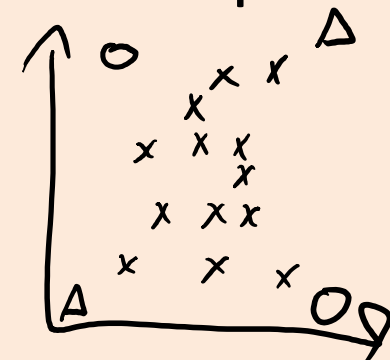


Does t-SNE always outperform PCA?

- 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.



t-SNE

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 **univariate chain rule**:

$$\frac{d}{dw} [f(g(w))] = f'(g(w)) g'(w)$$

- The **multivariate chain rule**:

$$\underbrace{\nabla [f(g(w))]}_{1 \times 1} = \underbrace{f'(g(w))}_{1 \times 1} \underbrace{\nabla g(w)}_{d \times 1}$$

- Example:

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

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

$$\text{with } g(w) = w^T x_i - y_i$$

$$\text{and } f(r_i) = \frac{1}{2} r_i^2$$

$$\nabla g(w) = x_i$$

$$f'(r_i) = r_i$$

$$\nabla [f(g(w))] = r_i x_i$$

$$= (w^T x_i - y_i) x_i$$

Bonus Slide: Multivariate Chain Rule for MDS

- General **MDS** formulation:

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

- Using **multivariate chain rule** we have:

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

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

$$\nabla_{z_i} g(d_1(x_i, x_j), d_2(z_i, z_j)) = \underbrace{-(d_1(x_i, x_j) - d_2(z_i, z_j))}_{g'(d_1, d_2)} \left[\underbrace{-\frac{(z_i - z_j)}{2\|z_i - z_j\|}}_{\text{(how distance changes in } z \text{ space)}} \right] \rightarrow \nabla_{z_i} d_2(z_i, z_j)$$

↳ Assuming $z_i \neq z_j$

(move distances closer)

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**.

$$X_{j \text{ dagger}} \approx \begin{bmatrix} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{bmatrix} \begin{matrix} z_{j1} \\ z_{j2} \\ z_{j3} \end{matrix}$$

