

Dimensionality Reduction: PCA and low-D embeddings

Applied Machine Learning Derek Hoiem

McInnes et al. (UMAP, 2000): Visualization of 30,000,000 integers as represented by binary vectors of prime divisibility, colored by integer value of the point

So far... KNN and retrieval

K-nearest neighbor classification and regression

Measuring and understanding error

- Fast retrieval and clustering
 - A few carry over slides coming up...

Advantages of Hierarchical K-Means

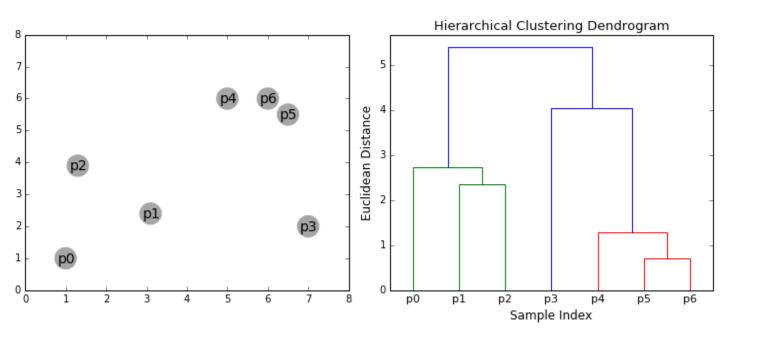
- Fast cluster training
 - With a branching factor of 10, can cluster into 1M clusters by clustering into 10 clusters ~111,111 times, each time using e.g. 10K data points
 - Vs. e.g. clustering 1B data points into 1M clusters
 - Kmeans is O(K*N*D) per iteration so this is a 900,000x speedup!
- Fast lookup
 - Find cluster number in O(log(K)*D) vs. O(K*D)
 - 16,667x speedup in the example above

Are there any disadvantages of hierarchical Kmeans?

Yes, the assignment might not be quite as good, but often usually isn't a huge deal since K means is used to approximate data points with centroid anyway

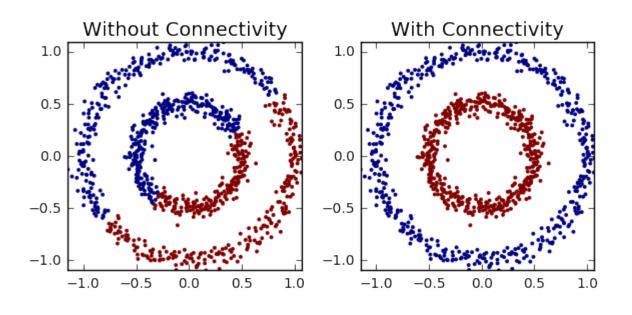
Agglomerative clustering

- Iteratively merge the two most similar points or clusters
 - Can use various distance measures
 - Can use different "linkages", e.g. distance of nearest points in two clusters or the cluster averages
 - Ideally the minimum distance between clusters should increase after each merge (e.g. if using the distance between cluster centers)
 - Number of clusters can be set based on when the cost to merge increases suddenly



Agglomerative clustering

 With good choices of linkage, agglomerative clustering can reflect the data connectivity structure ("manifold")



Clustering based on distance of 5 nearest neighbors between clusters

Applications of clustering

- K-means
 - Quantization (codebooks for image generation)
 - Search
 - Data visualization (show the average image of clusters of images)
- Hierarchical K-means
 - Fast search (document / image search)
- Agglomerative clustering
 - Finding structures in the data (image segmentation, grouping camera locations together)

https://tinyurl.com/AML441-L4

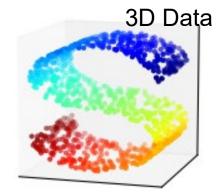


This class – dimensionality reduction

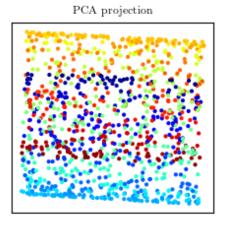
- Goal: We want to represent high dimensional data with fewer dimensions, e.g. for:
 - Compression: reduced storage, faster retrieval
 - Visualization: plot in two dimensions

- A good dimensionality reduction can be defined in different ways
 - Be able to reproduce the original data
 - Preserve discriminative features
 - Preserve the neighborhood structure

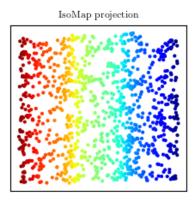
This class – PCA and manifolds



- Linear projection
 - PCA: Principal Components Analysis
 - Reduce dimension while preserving variance of data



- Embedding/manifold learning
 - MDS (multidimensional scaling), IsoMap, t-SNE, UMAP
 - Preserve point distances and/or local structure



- Vectors and matrices
- Translation
- Projection
- Scaling
- Rotation
- Rank
- Eigenvectors/eigenvalues
- SVD

Vector can represent a data point x or a **projection** w onto a **coordinate**

- $w^T x$ projects data point x onto the axis defined by w
- E.g. suppose d=2
 - $\mathbf{w} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ selects the first value of \mathbf{x}
 - $w = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ adds the two values of x together

Matrix can represent a set of data points or a set of projection vectors

Translation is a transformation that adds a constant value to each coordinate

• E.g. centering is $x_c = x - \mu_x$, where μ_x is the mean of x

Scaling is a transformation that multiplies each coordinate by a constant value

• E.g.
$$W = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}$$
 will double the value of the first coordinate of x when applied as Wx (blanks are assumed to be zero)

Rotation is a set of projections that preserves the distances between points and distances to the origin

- Each row and column represents a basis vector
- The basic vectors must have a unit norm and be orthogonal to each other: $\mathbf{r}_i^T \mathbf{r}_i = 1$, $\mathbf{r}_i^T \mathbf{r}_i = 0$, $\forall (i, j \neq i)$

Rank of matrix M is the number of linearly independent vectors in the rows or columns of M

• Rank(M) is at most the smaller of the number of rows and columns in M

•
$$Rank \begin{pmatrix} \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \end{pmatrix} = 2$$

• $Rank \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} = 1$ because one of the rows (or columns) can be composed of a weighted sum of the others

Eigenvector v and corresponding **eigenvalue** λ of matrix M are defined by having the special property $Mv = \lambda v$

 Eigenvectors characterize matrices, and appear as part of a solution to many linear algebra problems

SVD (singular value decomposition) is a factorization of a matrix \boldsymbol{A} into $\boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^T$, where:

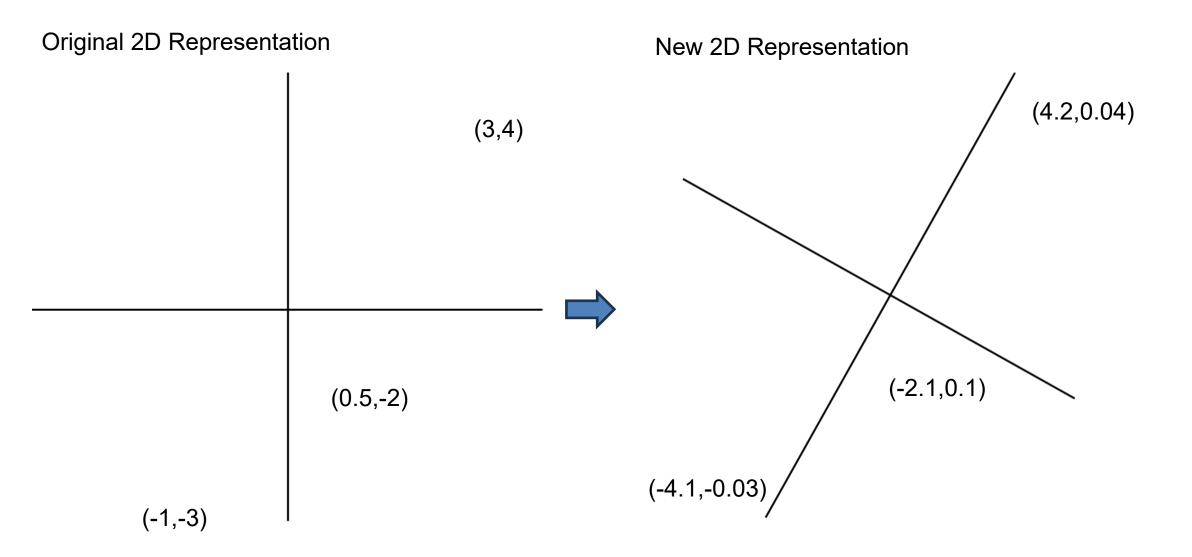
- Columns of U are eigenvectors of AA^T
- Columns of V are eigenvectors of A^TA
- Σ is a diagonal (scaling) matrix of singular values, which are square roots of the eigenvalues of both AA^T and A^TA

PCA

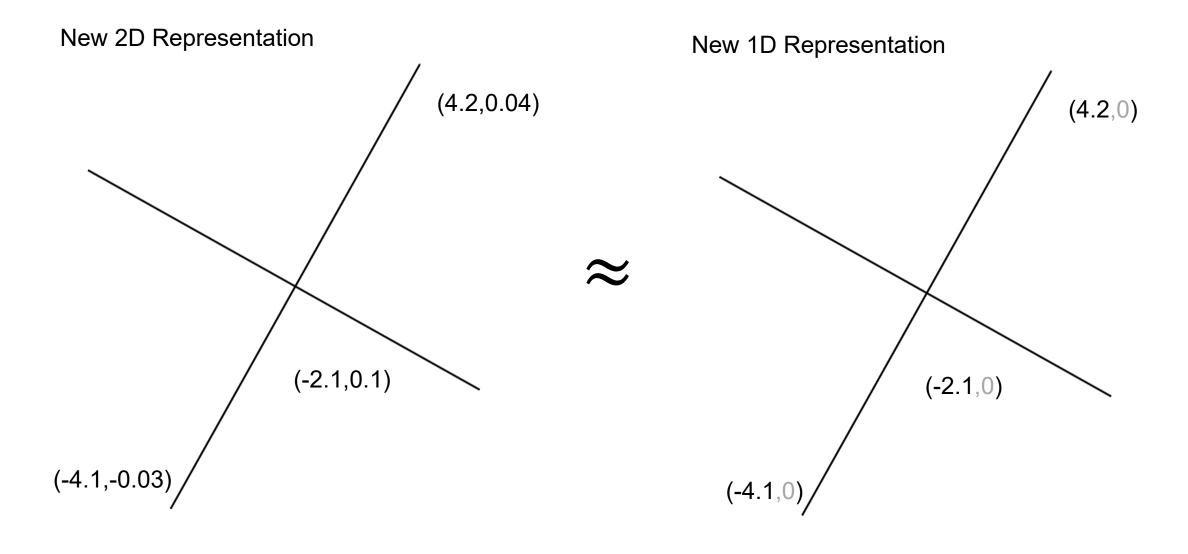
- General dimensionality reduction technique
 - Finds major orthogonal directions of variation

- Preserves most of variance with a much more compact representation
 - Lower storage requirements
 - Faster matching/retrieval
 - Easier to work in low dimensions, e.g. for probability estimation

Center and rotate the axes

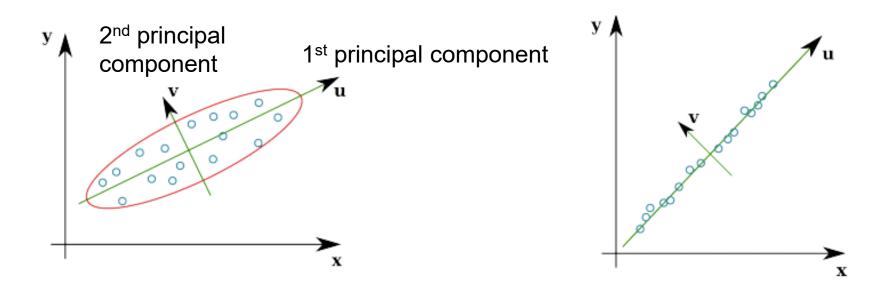


Now data is well-approximated by dropping the coordinates with the least variance



Principal Component Analysis

- Given a point set $\{\vec{\mathbf{p}}_j\}_{j=1...P}$, in an M-dim space, PCA finds a basis such that
 - The most variation is in the first basis vector
 - The second most, in the second vector that is orthogonal to the first vector
 - The third...



Principal Component Analysis (PCA)

- Given: N data points x₁, ..., x_N in R^d
- We want to find a new set of features that are linear combinations of original ones:

$$u(\mathbf{x}_i) = \mathbf{u}^T(\mathbf{x}_i - \mathbf{\mu})$$

 (μ) : mean of data points)

 Choose unit vector u in R^d that captures the most data variance

Principal Component Analysis

 $\mathbf{u}^{\mathrm{T}} \Sigma \mathbf{u}$

Direction that maximizes the variance of the projected data:

$$\begin{aligned} & \underbrace{\frac{1}{N} \sum_{i=1}^{N} \mathbf{u}^{\! \mathrm{T}} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{u}^{\! \mathrm{T}} (\mathbf{x}_i - \boldsymbol{\mu}))^{\! \mathrm{T}}}_{\text{Subject to } ||\mathbf{u}|| = 1} \\ & = \mathbf{u}^{\! \mathrm{T}} \left[\sum_{i=1}^{N} (\mathbf{x}_i - \boldsymbol{\mu}) (\mathbf{x}_i - \boldsymbol{\mu})^{\! \mathrm{T}} \right] \mathbf{u}} \end{aligned}$$

The direction that maximizes the variance is the eigenvector associated with the largest eigenvalue of Σ (can be derived using Raleigh's quotient or Lagrange multiplier)

PCA in Python

```
print(X.shape)
x_{mu} = np.mean(X, axis=0)
X_cov = np.matmul((X-x_mu).transpose(), X-x_mu)/X.shape[0]
[lam, v] = np.linalg.eig(X_cov)
print(lam[:3])
v = v.transpose()
from sklearn.decomposition import PCA
pca_transform = PCA().fit(X)
print(np.max(np.abs(v[0]-pca_transform.components_[0])))
print(pca_transform.explained_variance_[:3])
```

Compute PCA components as eigenvectors of covariance matrix

Compute PCA using the PCA function

```
(50000, 784)

[5.0695131 3.7301149 3.25871794]

3.9811903773667723e-16

[5.06961449 3.7301895 3.25878311]
```

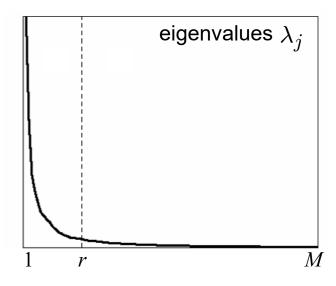
Printout shows that eigenvalues are the same as the explained variance, and the first component is identical (up to numerical precision)

Principal Component Analysis

First r < M principal component vectors provide an approximate basis that minimizes the mean-squared-error (MSE) of reconstructing the original points

Choosing subspace dimension r:

- look at decay of the eigenvalues as a function of r
- Larger r means lower expected error in the subspace data approximation



Example on aligned faces

 $\mathbf{x}_1,...,\mathbf{x}_N$ are pixel values of each face

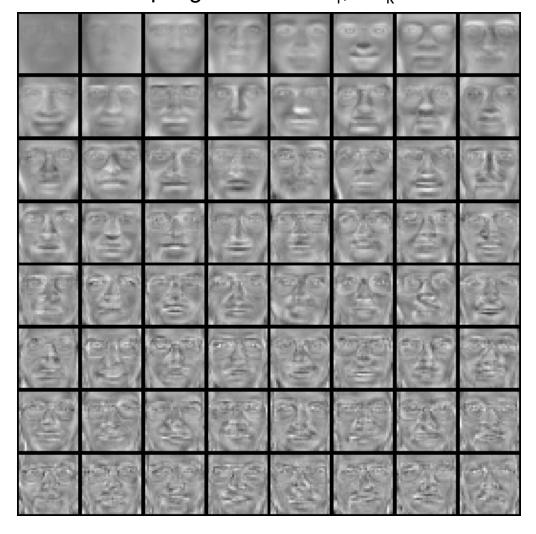


PCA of aligned face images (called "eigenfaces")

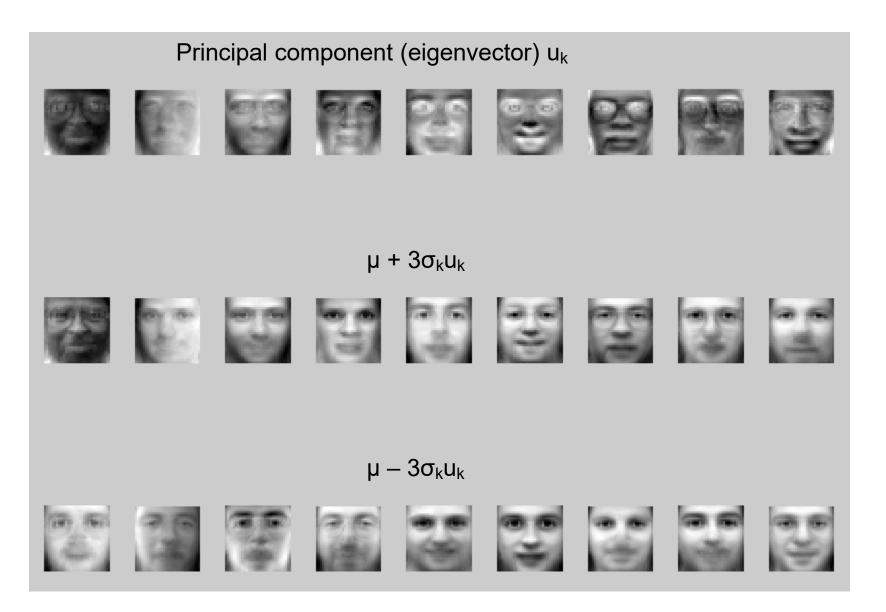
Top eigenvectors: $u_1, \dots u_k$

Mean: µ





Visualization of eigenfaces (appearance variation)



Representation and reconstruction

Face x in "face space" coordinates:



$$\mathbf{x} \to [\mathbf{u}_1^{\mathrm{T}}(\mathbf{x} - \mu), \dots, \mathbf{u}_k^{\mathrm{T}}(\mathbf{x} - \mu)]$$

$$= w_1, \dots, w_k$$

Representation and reconstruction

• Face x in "face space" coordinates:

$$\mathbf{x} \to [\mathbf{u}_1^{\mathrm{T}}(\mathbf{x} - \mu), \dots, \mathbf{u}_k^{\mathrm{T}}(\mathbf{x} - \mu)]$$

$$= w_1, \dots, w_k$$

• Reconstruction:

Reconstruction



After computing eigenfaces using 400 face images from ORL face database

Note

Preserving variance (minimizing MSE) does not necessarily lead to perceptually good reconstruction.

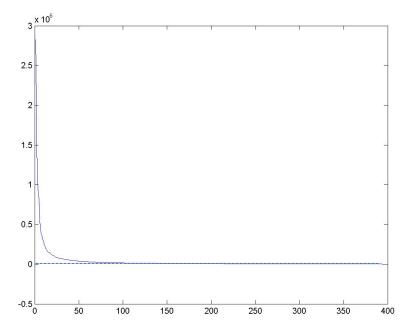






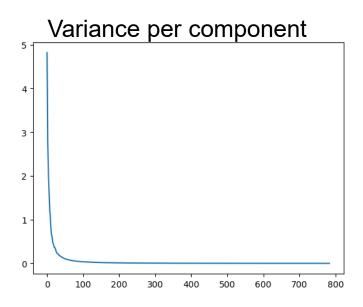




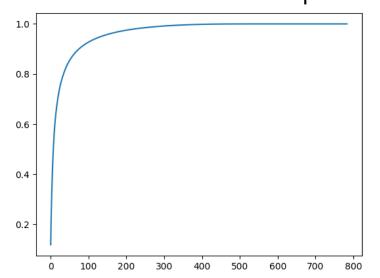


Plot of eigenvalues for each eigenvector of the covariance matrix, equivalent to the variance contained along each principal component

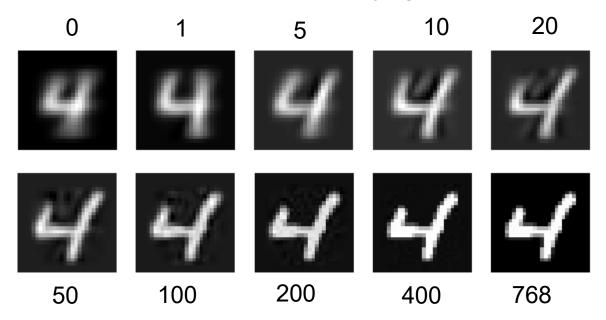
Another example, representing MNIST 4's



Cumulative % variance explained



Reconstructions with varying # PCs

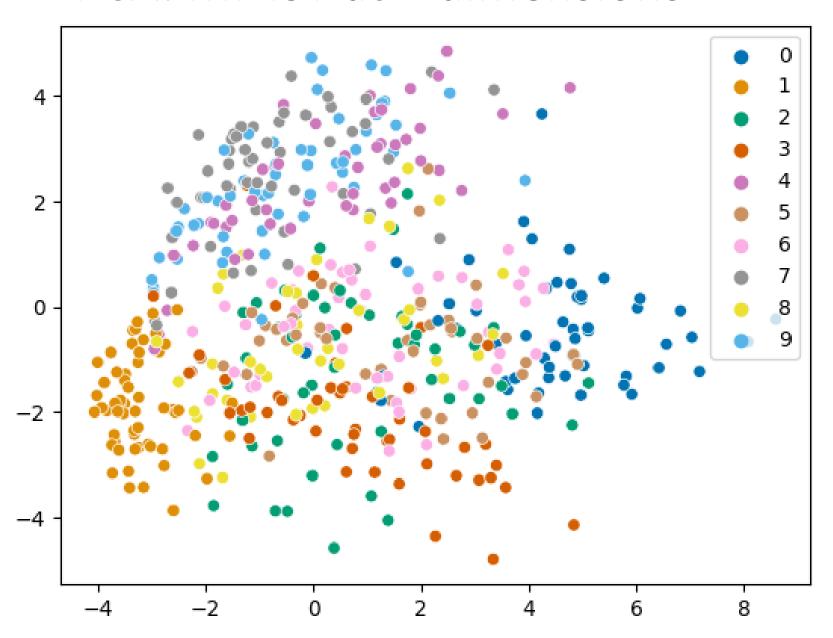


```
from sklearn.decomposition import PCA
x4 = x_train[y_train==4]
pca = PCA().fit(x4)
pca_x4 = pca.transform(x4)
plt.plot(pca.explained_variance_)
plt.show()
plt.plot(np.cumsum(pca.explained_variance_ratio_))
plt.show()
ims = np.zeros((10, 28*28))
nc = np.int32([0, 1, 5, 10, 20, 50, 100, 200, 400, 768])
for i in range(len(nc)):
    c = nc[i]
    ims[i] = pca.mean_ + np.matmul(pca_x4[0][:c],pca.components_[:c])
display_mnist(ims, 1, 10)
```

https://tinyurl.com/AML441-L5



PCA: MNIST at 2 dimensions



Note: I'm only plotting a subset

Non-Linear Scaling and Manifold Estimation

We may care less about being able to reconstruct each data point than representing the relationships between data points

- MDS: Preserve Euclidean pairwise distances
- Non-metric MDS: Preserve distance orderings
- ISOMAP: Define distances in terms of "geodesic" (graph-based) similarity

MultiDimensional Scaling (MDS)

- For all data points, solve for new coordinate positions that preserve some input set of pairwise distances
- Classic case (equations on right) uses Euclidean distance and has closed form solution
- More generally, distance can be defined arbitrarily (e.g. from user surveys)
- Major drawback is that the solution is most influenced by points that are far from each other

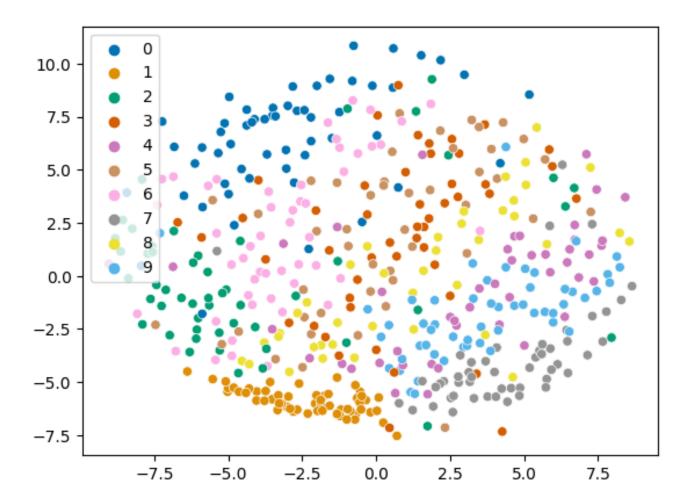
Solve for y that minimizes $\sum_{ij} \left(D_{ij}^{(2)}(\mathbf{x}) - D_{ij}^{(2)}(\mathbf{y}) \right)^2$

where $D_{ij}^{(2)}(\mathbf{x}) = (\mathbf{x}_i - \mathbf{x}_j)^T (\mathbf{x}_i - \mathbf{x}_j)$

x are old coordinates, y are new coordinates

MDS on MNIST

- 30 PCA components
- MDS to 2 dimensions



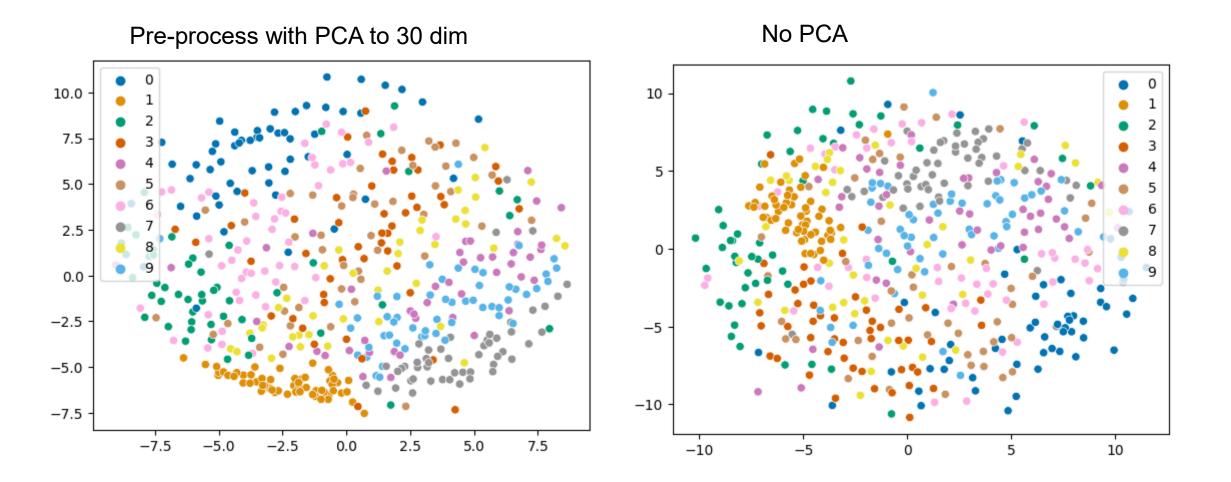
```
from sklearn.manifold import MDS

pca = PCA()
pca.fit(x_train)
x_pca = pca.transform(x_train)

ind = train_indices['s']
x_mds = MDS(n_components=2, normalized_stress='auto').fit_transform(x_pca[ind, :30])
sns.scatterplot(x=x_mds[ind,0],y=x_mds[ind,1], hue=y_train[ind], palette="colorblind")
```

Note: For MDS and others, manifolds are fit on 500 pts for speed

MDS on MNIST



Non-metric MDS

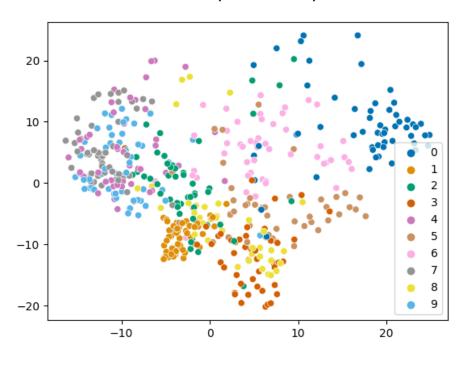
- Optimize position of data points so that Euclidean distance preserves the ordering of input pairwise distances
- Requires only an order of dissimilarities

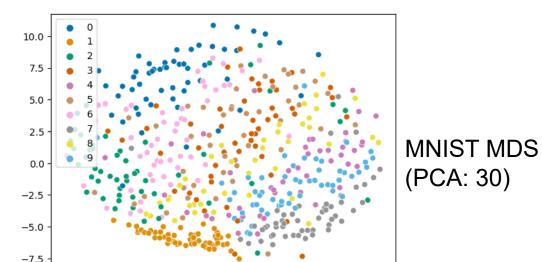
Slow because this is a complicated optimization

ISOMAP

- Same as MDS but define distance in a graph
 - Compute adjacency graph (e.g. 5 nearest neighbors)
 - Distance is shortest path in graph between two points

MNIST ISOMAP (PCA: 30)





t-SNE

Map to 2 or 3 dimensions while preserving similarities of nearby points

- 1. Assign probability p_{ij} that pairs of points are similar, e.g. with Gaussian weighted distance
- 2. Define similarity q_{ij} in new coordinates
- 3. Minimize KL divergence between p and q (i.e. they should have similar distributions) using gradient descent

For $i \neq j$, define

$$p_{j|i} = rac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma_i^2)}{\sum_{k
eq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2/2\sigma_i^2)} \ p_{ij} = rac{p_{j|i} + p_{i|j}}{2N}$$

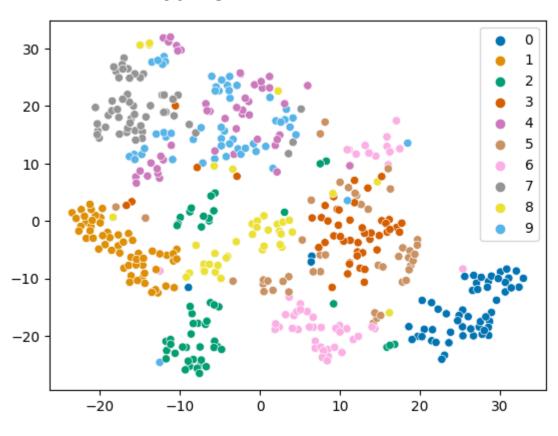
$$q_{ij} = rac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_k \sum_{l
eq k} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$

$$ext{KL}\left(P \parallel Q
ight) = \sum_{i
eq j} p_{ij} \log rac{p_{ij}}{q_{ij}}$$

t-SNE

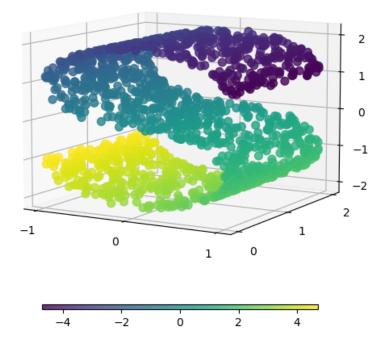
 In high dimensions, each point tends to be similarly distant to many points, so we often use PCA before applying t-SNE

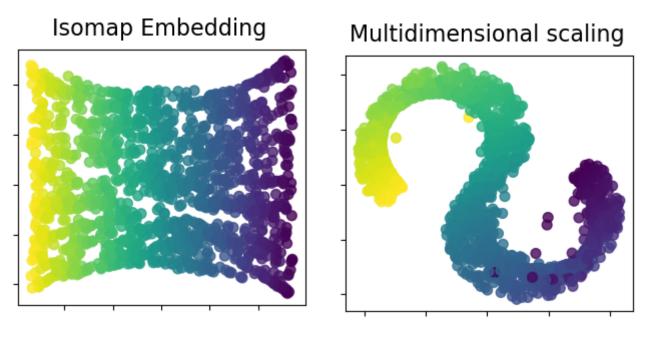
t-SNE on 30 PCA dimensions of MNIST



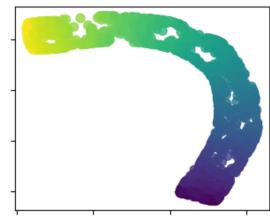
Comparison

Original S-curve samples

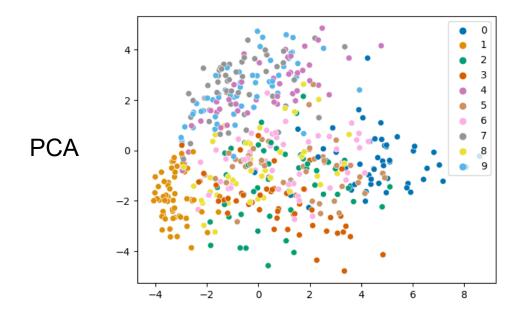


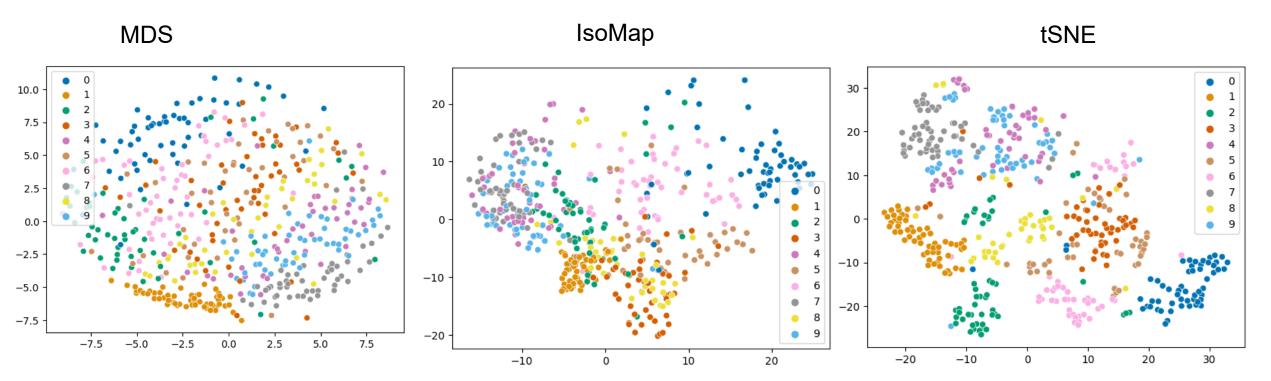


T-distributed Stochastic Neighbor Embedding



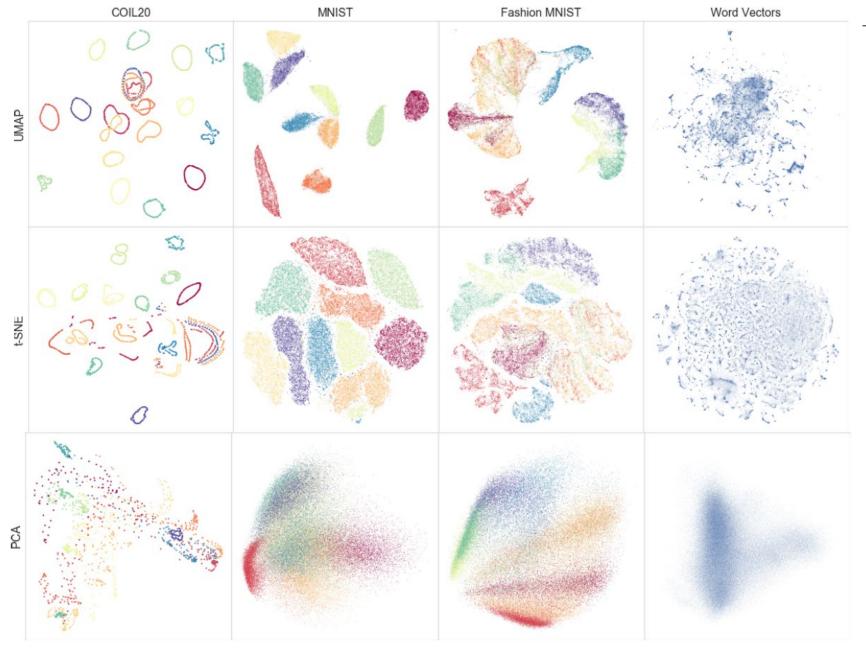
Comparison on MNIST





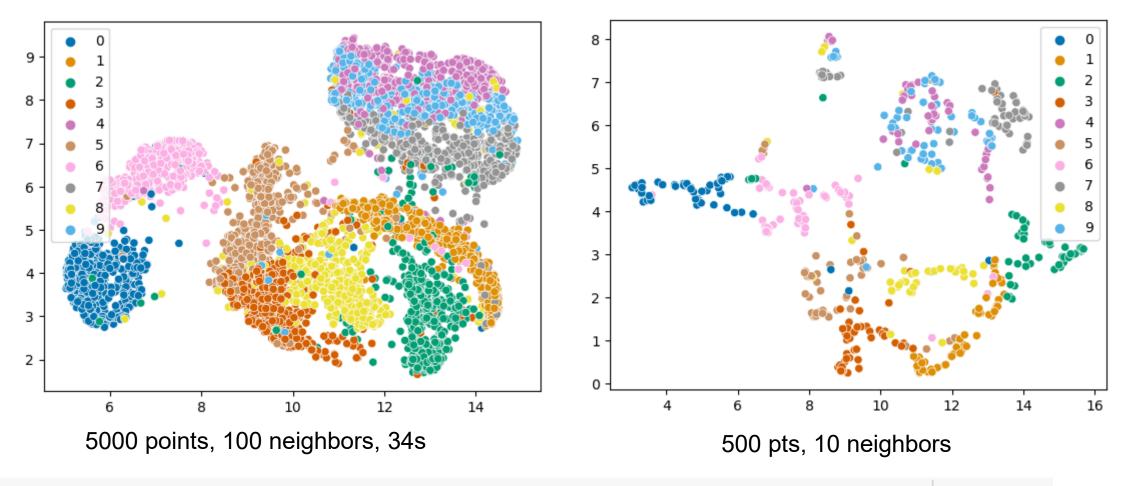
UMAP (McInnes, 2020)

- Assumes data is uniformly distributed on an underlying manifold that is locally connected. Goal is to preserve that local structure.
- Algorithm: relatively complicated, incorporates many ideas from other methods, has strong mathematical foundations
- Hyperparameters:
 - number of neighbors to consider
 - dimension of target embedding
 - desired separation between close points
 - number of training epochs

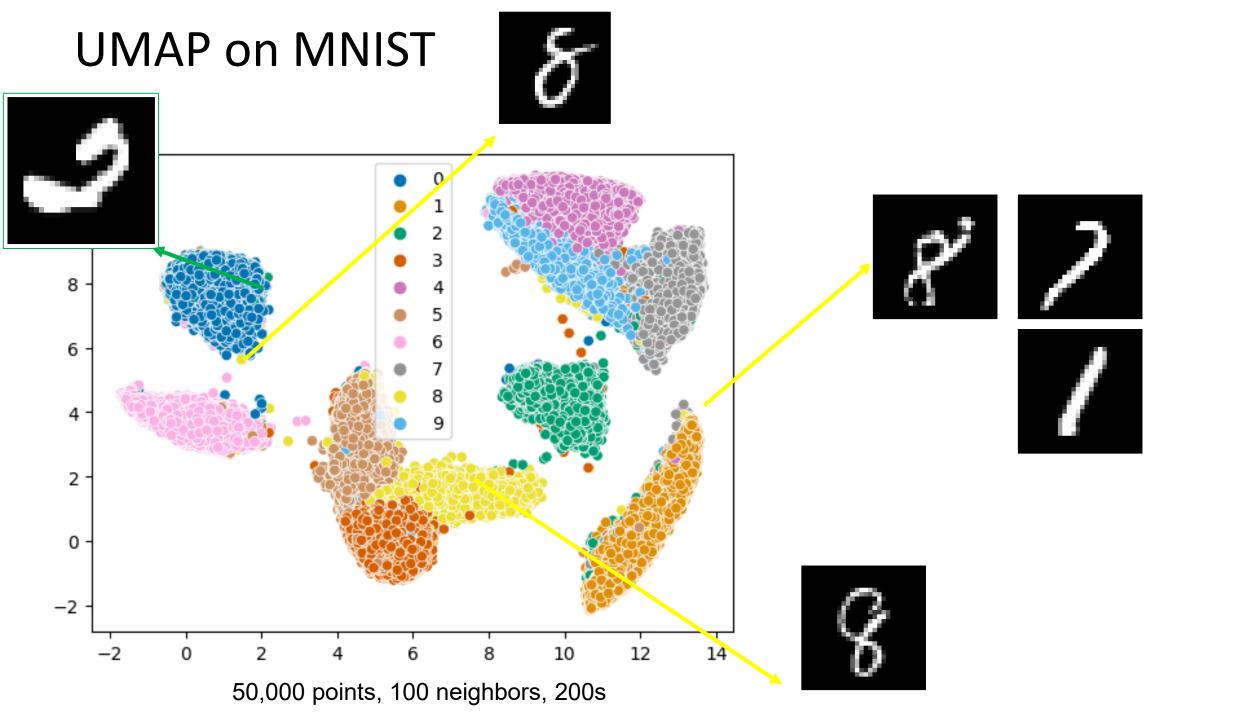


	UMAP	t-SNE	Isomap
Pen Digits	9s	17s	2s
(1797x64)			
COIL20	12s	22s	58s
(1440x16384)			
COIL100	85s	810s	3210s
(7200x49152)			32103
scRNA	28s	258s	0220
(21086x1000)	208	2308	923s
		=	
Shuttle	94s	714s	_
(58000x9)			
MNIST	87s	1450s	_
(70000x784)			
F-MNIST	65s	934s	_
(70000x784)			
Flow	102s	1135s	_
(100000x17)			
Google News	361s	16906s	
(200000x300)	3018	107008	_
(

UMAP on MNIST



```
#!pip install umap-learn
from umap import UMAP
ind = train_indices['s']
x_umap = UMAP(n_components=2, n_neighbors=100).fit_transform(x[ind, :30])
sns.scatterplot(x=x_umap[ind,0],y=x_umap[ind,1], hue=y_train[ind], palette="colorblind")
```

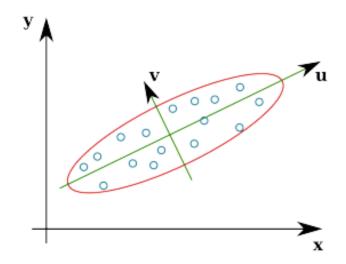


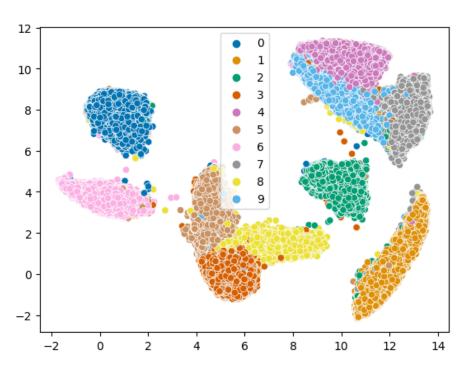
https://tinyurl.com/AML441-L5



Things to remember

- PCA reduces dimensions by linear projection
 - Preserves variance to reproduce data as well as possible, according to mean squared error
 - May not preserve local connectivity structure or discriminative information
- Other methods try to preserve relationships between points
 - MDS: preserve pairwise distances
 - IsoMap: MDS but using a graph-based distance
 - t-SNE: preserve a probabilistic distribution of neighbors for each point (also focusing on closest points)
 - UMAP: incorporates k-nn structure, spectral embedding, and more to achieve good embeddings relatively quickly





Next class: Linear Regression and Regularization