

UNIVERSITY OF BERGEN

INFO 284 – Machine Learning

Dimension Reduction

Bjørnar Tessem

UNIVERSITY OF BERGEN



Learning in AI

Any **component** of an intelligent system can be improved by learning from data. Learning depends on four factors:

- which component is to be improved
- what prior knowledge the agent has
- **what representation is used for the data and the component**
- **what feedback is available to learn from**

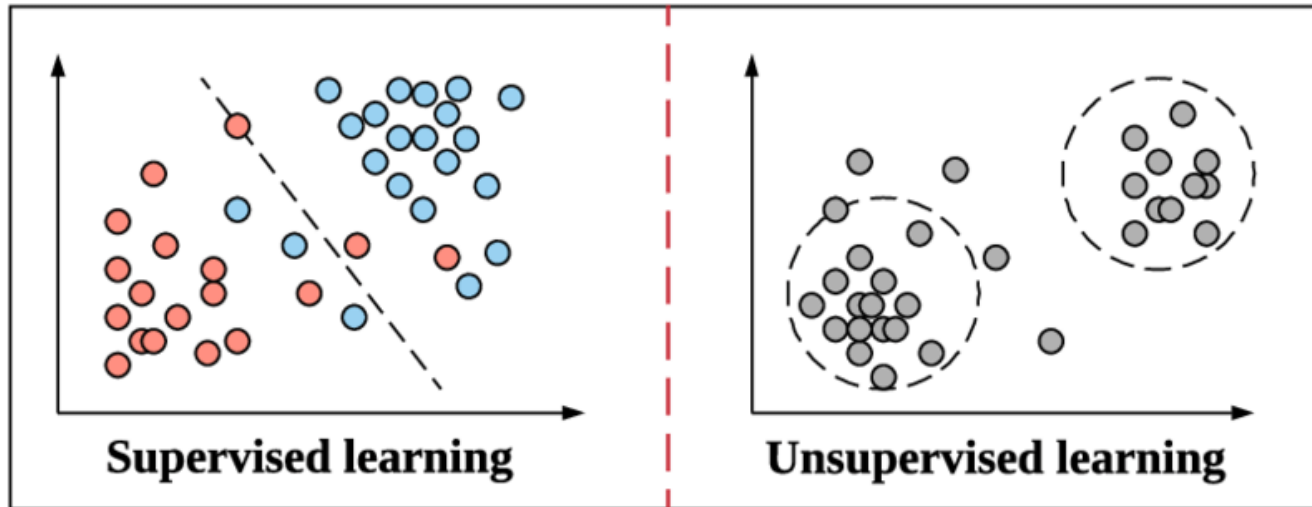


Available feedback

- Supervised learning - correct/wrong
- Unsupervised learning - no feedback
- Self-supervised learning
 - Generate labels from original data
 - Predict some parts of the data from other parts
- Reinforcement learning - reward/punishment



Unsupervised learning



Types of unsupervised learning

- Clustering
 - Find inherent groupings in the data
 - grouping customers by purchasing behaviour
 - news by topic
- Unsupervised transformations
 - Create new representations of the data
 - easier to visualise
 - easier to learn from
 - Dimensionality reduction

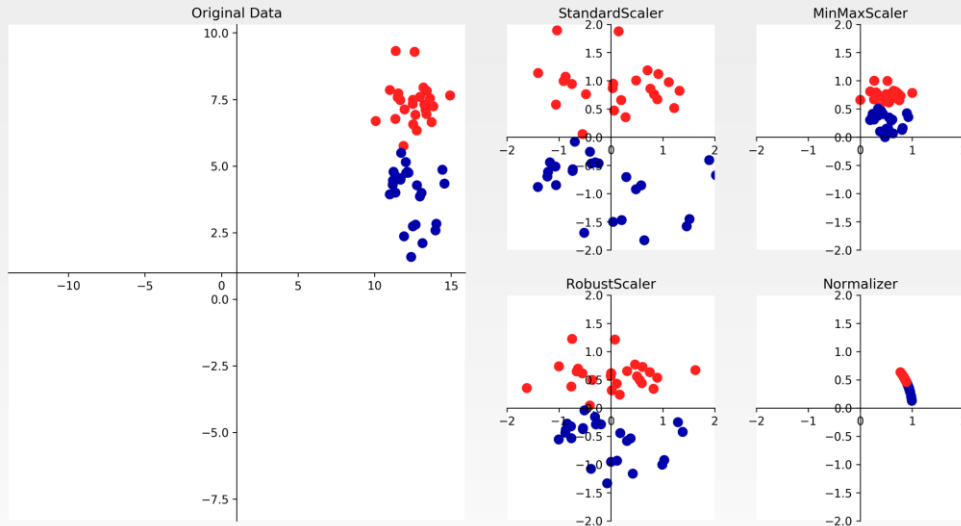


Dimensionality reduction

- Scaling
 - Ensuring all features are on same/similar scale
- Principal Component Analysis
 - Extract numerical features to represent data
- Non-negative Matrix Factorisation
 - Extract **positive valued** numerical features
- t-distributed Stochastic Neighbor Embedding
 - Two- or three-dimensional representations for visualisation



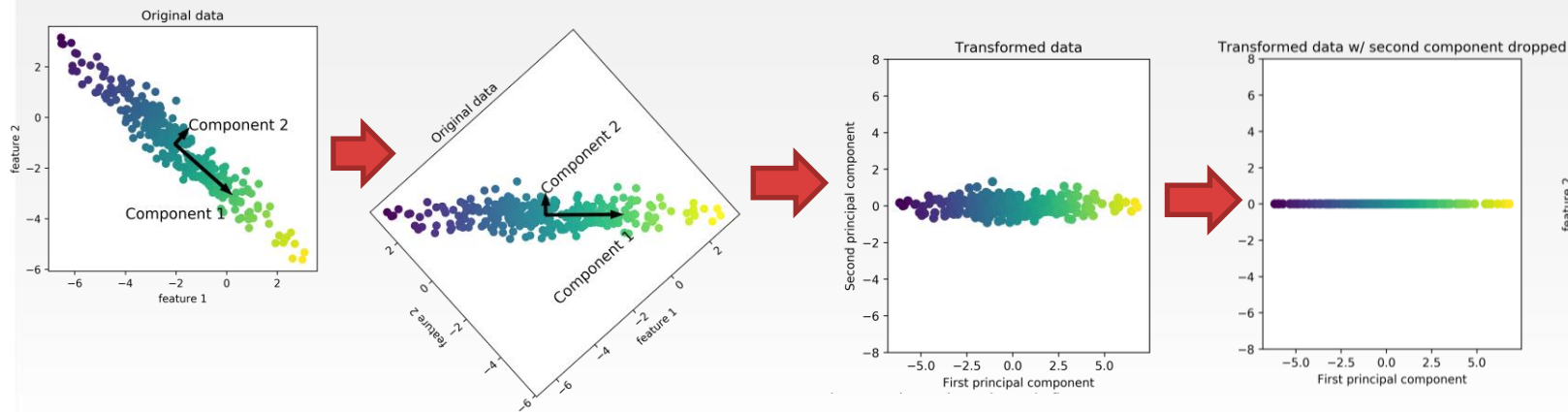
Scaling



- Standard scaler
 - transform to standard Gaussian
- MinMax scaler
 - transform to numbers in $[0,1]$
- Robust scaler
 - transform based on quartiles
- Normalizer
 - transform so that data vector has length (norm) 1

Principal Component Analysis

- Compute alternative feature set
- Rotate dataset so that the new features are statistically uncorrelated



Mathematics

- How much to rotate in each dimension?
- Data: $\begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}_n$ n datapoints, m features
 $\quad \quad \quad m$
- Covariance matrix: $X^T X \rightarrow \begin{bmatrix} c \\ \vdots \\ c \end{bmatrix}_m$
 $\quad \quad \quad m$
- Eigenvectors: Find w-vectors such that $\lambda w = Cw$
 - λ -s are called eigenvalues
 - w-vectors are called eigenvectors



Use

- Place all eigenvectors in own matrix $W = [w_1, w_2, \dots, w_m]$
 - rank by λ (eigenvalue)
- Transform all datapoints x to
$$t = xW$$
 - To reduce dimension use only the n first eigenvectors of W
 - The chosen eigenvectors are called principal components
 - t becomes an n -dimensional vector called scores
 - W columns are called loadings



Visualisation

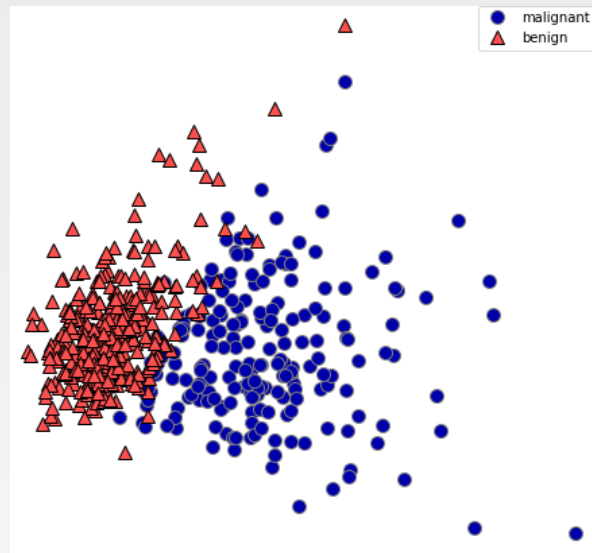
- Pair plots

```
from sklearn.decomposition import PCA
# keep the first two principal components of the data
pca = PCA(n_components=2)
# fit PCA model to breast cancer data
pca.fit(X_scaled)

# transform data onto the first two principal components
X_pca = pca.transform(X_scaled)
```

Original shape: (569, 30)

Reduced shape: (569, 2)



- The two axes in the plot are often not easy to interpret
 - uncorrelated combinations of the original features



PCA for feature extraction

- Find a representation of data that is better suited to analysis than the raw representation.



PCA in scikit-learn for feature extraction

- Problem: Does a previously unseen face belong to a known person from the database?
- Approach 1: build a classifier where each person is a separate class
 - Problem:
 - many different people
 - not that many pictures of each => hard to train a classifier
- Approach 2: Classifier that **can** be trained is kNN.
 - Computing distances between pixels is not very informative
 - Using PCA of images can help

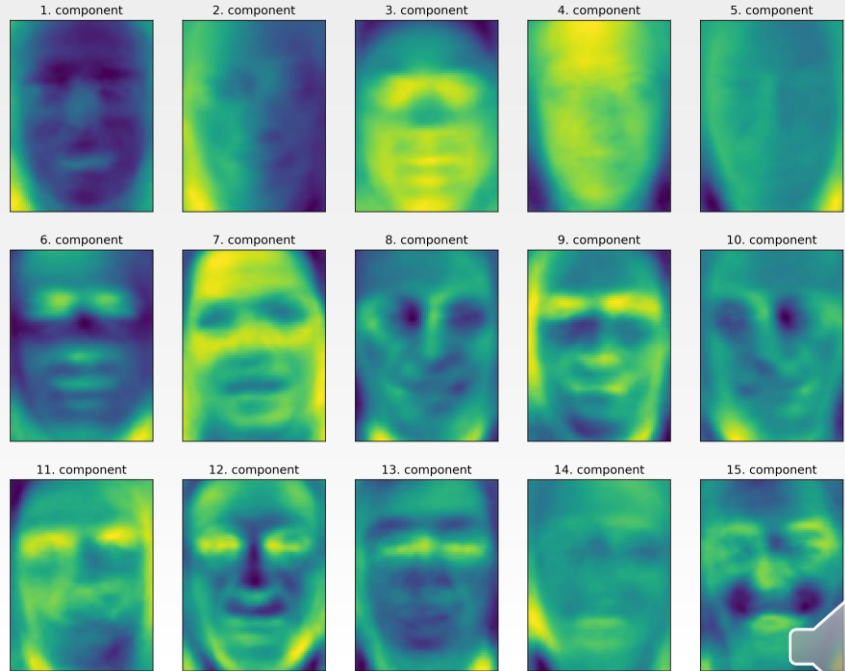


PCA for feature extraction

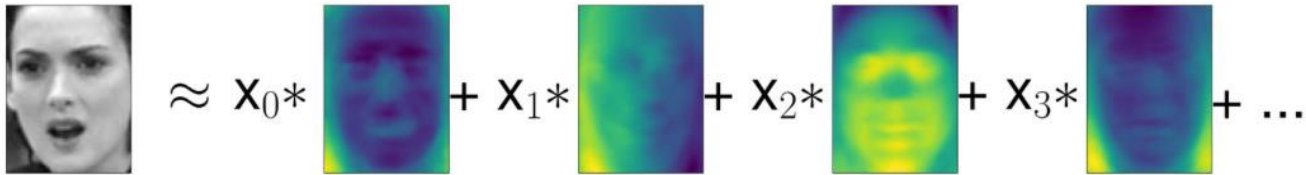
```
pca = PCA(n_components=100, whiten=True, random_state=0).fit(X_train)
X_train_pca = pca.transform(X_train)
X_test_pca = pca.transform(X_test)
```

```
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train_pca, y_train)
```

Scaling principal
components



Schematic view of PCA



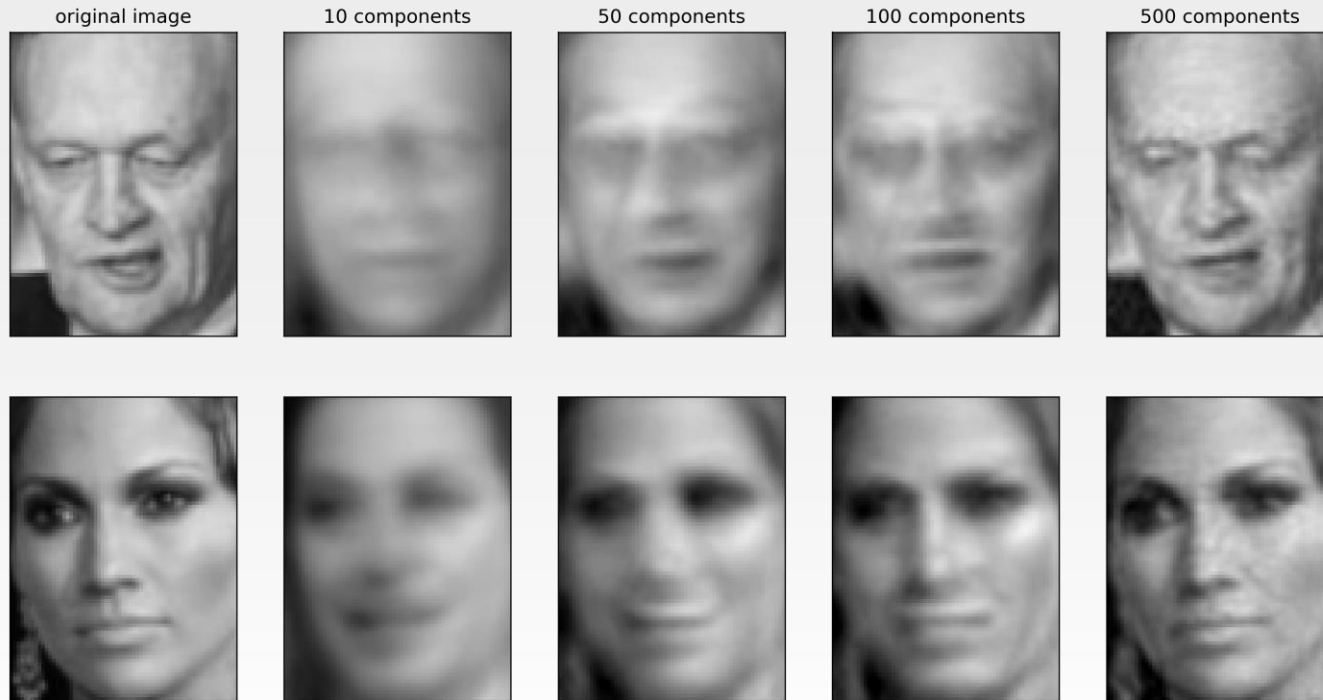
The diagram illustrates the decomposition of a grayscale image of a woman's face into a weighted sum of principal components. On the left is the original grayscale image. To its right is an approximation symbol followed by a series of terms: X_0 multiplied by a color heatmap, plus X_1 multiplied by a color heatmap, plus X_2 multiplied by a color heatmap, plus X_3 multiplied by a color heatmap, followed by an ellipsis. Each heatmap represents a principal component, showing different features of the face in a color-coded manner (blue for low intensity, yellow for high intensity).

$$\approx X_0 * \text{Component}_0 + X_1 * \text{Component}_1 + X_2 * \text{Component}_2 + X_3 * \text{Component}_3 + \dots$$

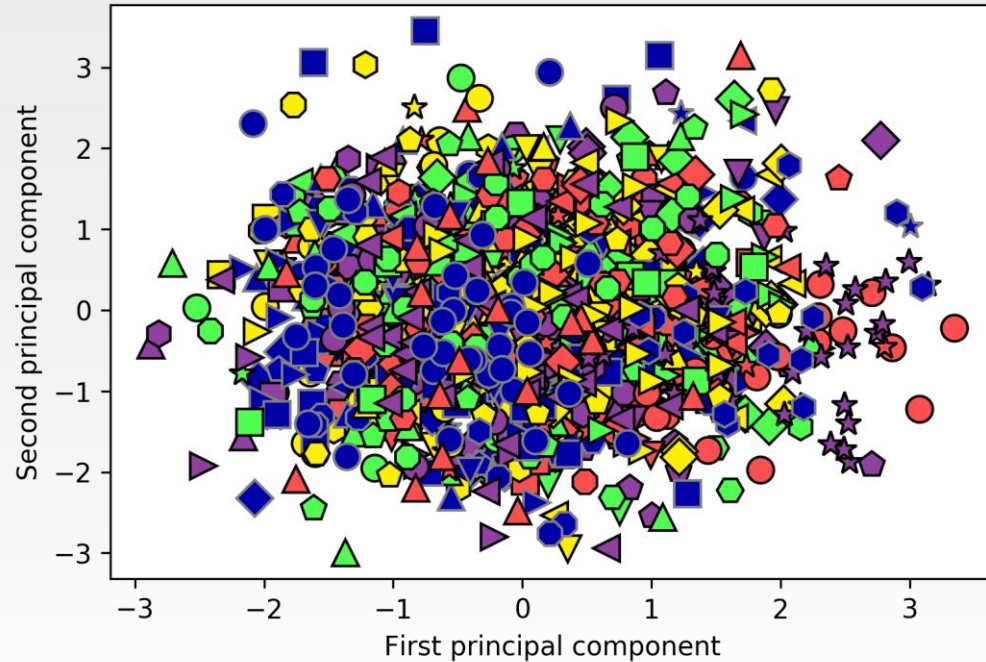
Figure 3-10. Schematic view of PCA as decomposing an image into a weighted sum of components



Reconstructing images



Plot of first two principle components



Non-negative Matrix Factorization

- Establish new set of features
- Similar to Principal Component Analysis
 - PCA: orthogonal components
- NMF: positive components and coefficients
 - Not orthogonal
 - $t = x \cdot W$
 - x and W only positive entries



Non-negative Matrix Factorisation

$$\begin{bmatrix} & W \\ \begin{bmatrix} \square & \square \\ \square & \square \\ \square & \square \\ \square & \square \end{bmatrix} & \times & \begin{bmatrix} & H \\ \begin{bmatrix} \square & \square & \square & \square & \square & \square \end{bmatrix} \end{bmatrix} \approx \begin{bmatrix} & V \\ \begin{bmatrix} \square & \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square & \square \\ \square & \square & \square & \square & \square & \square \end{bmatrix} \end{bmatrix}$$

The matrix \mathbf{V} (the original dataset) is represented by the two smaller matrices \mathbf{W} and \mathbf{H} . $\mathbf{W} \times \mathbf{H}$ approximately reconstruct \mathbf{V} .

We want to transfer \mathbf{V} to \mathbf{H} to reduce dimensionality.

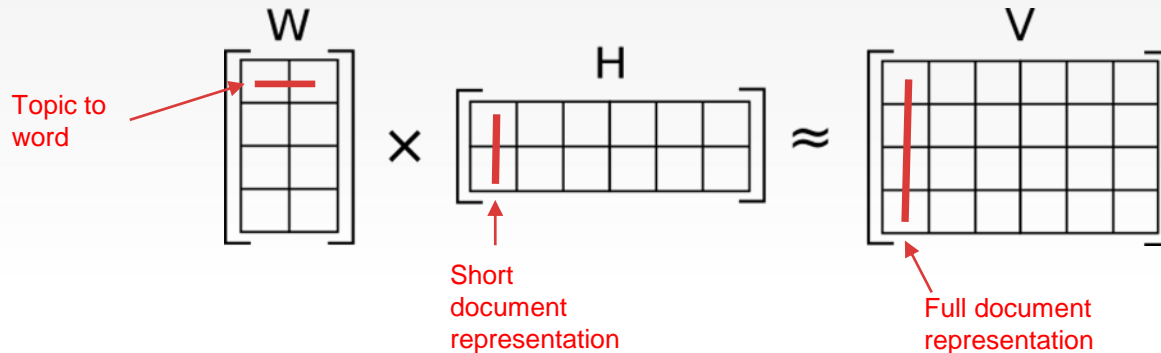
Each row in \mathbf{H} represent the scores on a particular «hidden» feature. Each column is a datapoint.

Each column in \mathbf{W} tells how we transform a «hidden» feature to the observed data found in \mathbf{V}

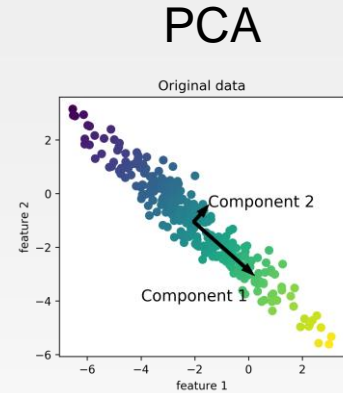
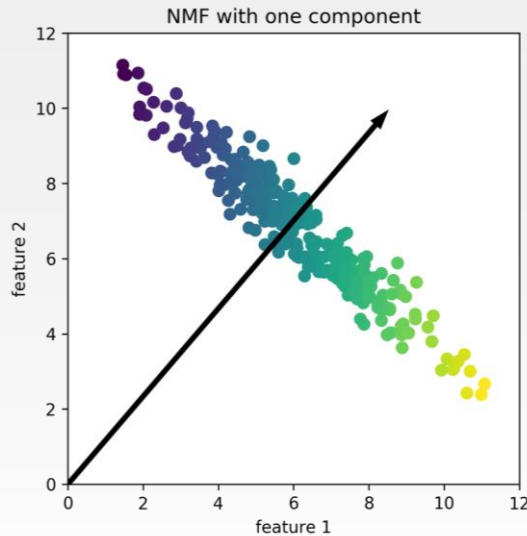
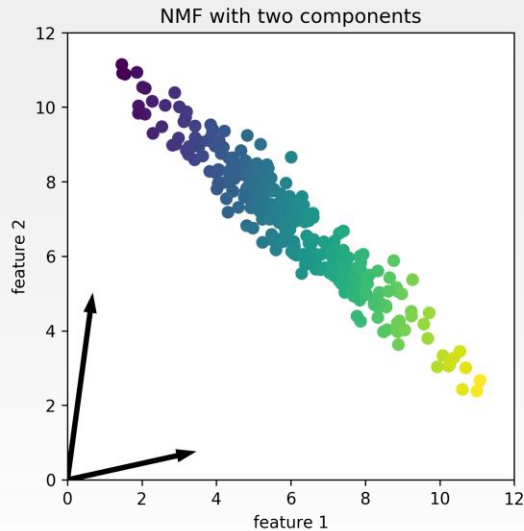


Topic modeling

- Popular approach in document analysis
- Uses NMF technique
- Represent documents as one-hot-coded document representations (or TF-IDF)
- Convert to a combination of a few topics.



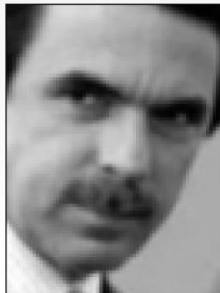
Non-negative Matrix Factorisation



Faces in MNF components



High scorers on component 3



High scorers on component 7



t--distributed Stochastic Neighbor Embedding

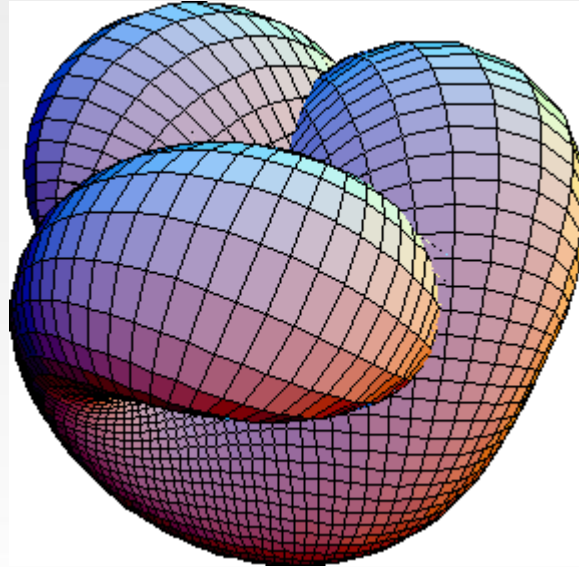
- t-SNE
 - tool for visualisation of many-dimensional data sets
- PCA maintains variance in data
 - Pairwise distances in data are preserved
- t-SNE
 - Preserve only short pairwise distances
 - Close data points in full feature space features will be close in 2(3) dimensions
 - «Trained» by gradient descent
 - New data → new training



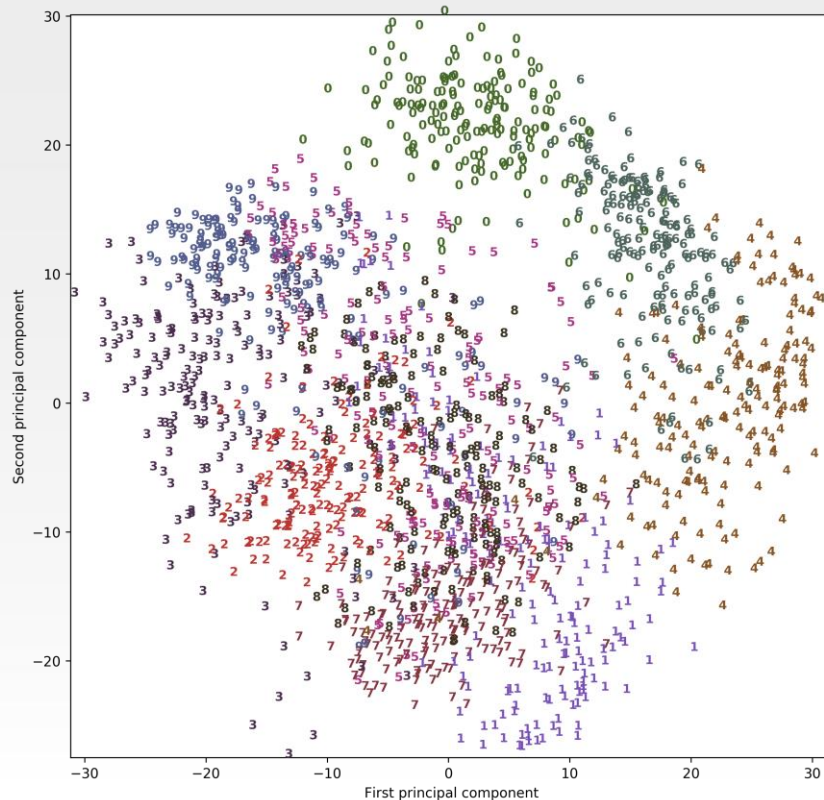
Manifolds

Data is often placed within particular continuous subsets (surfaces) of the full \mathbb{R}^n feature space.

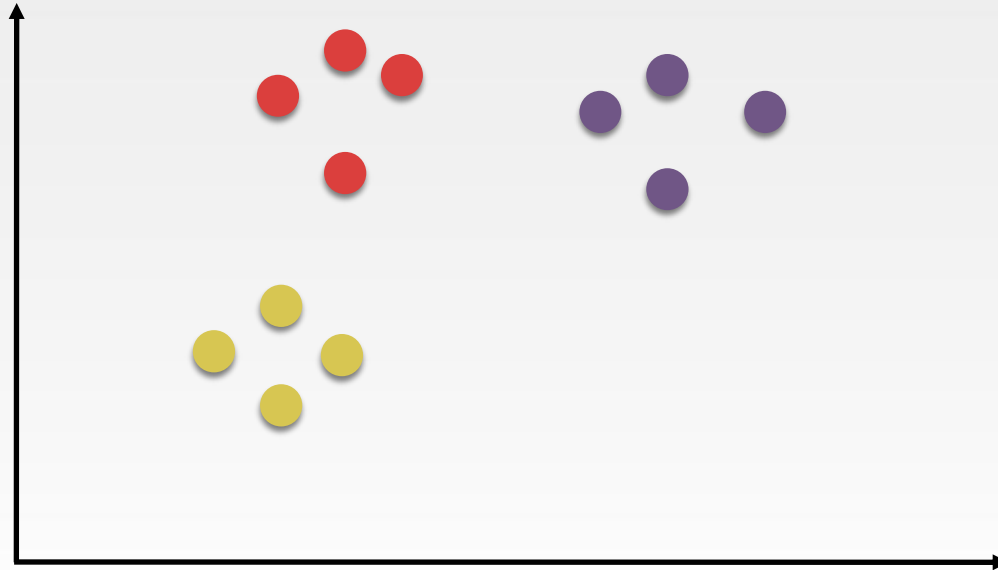
Visualise by reducing this data MANIFOLD to two-dimensions



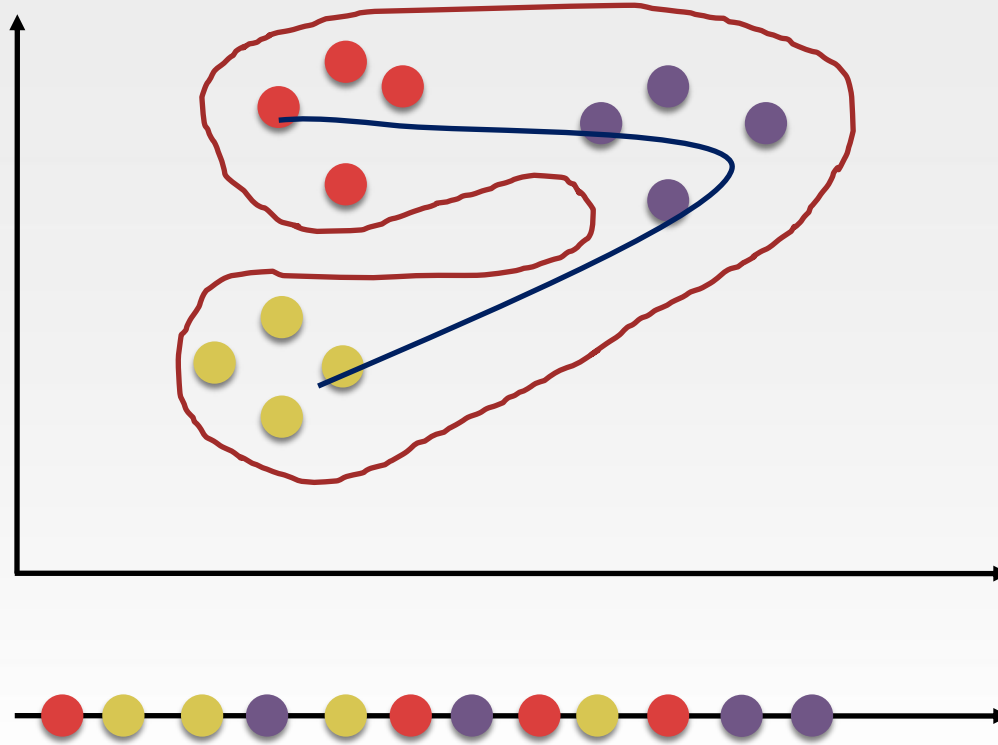
Plot of only the first two principle components MNIST



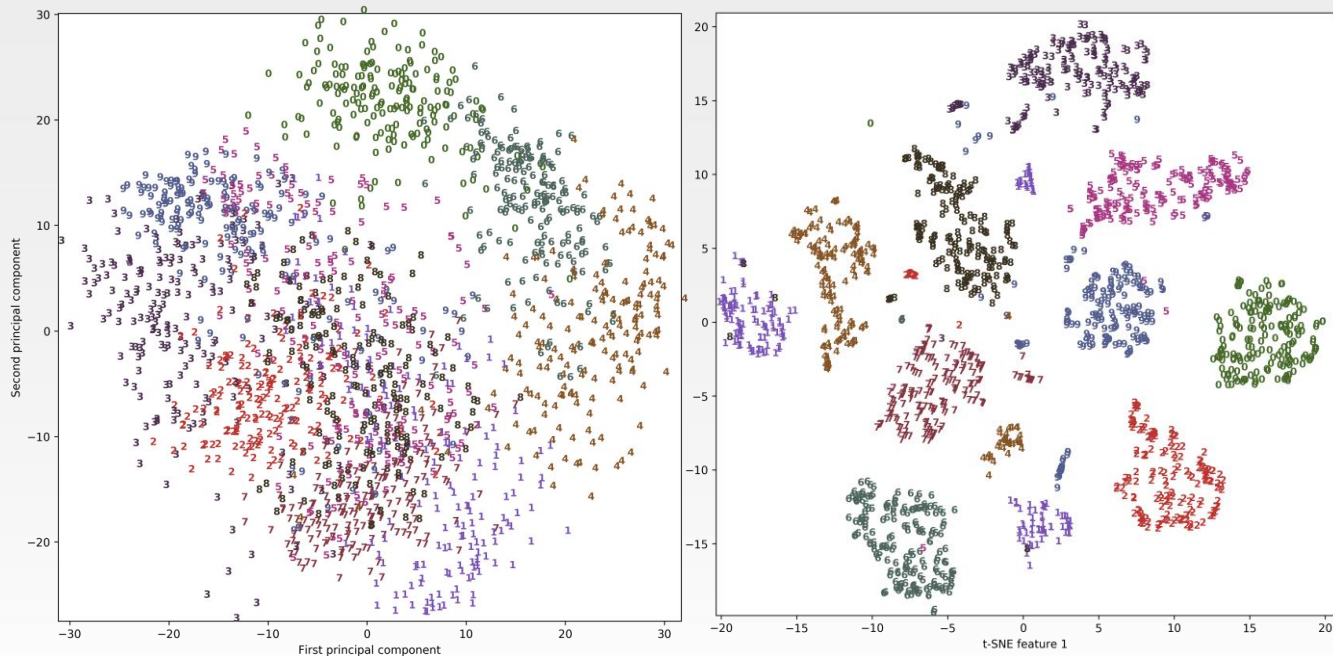
Basic ideas of t-SNE



Basic ideas of t-SNE



Plot MNIST using t-SNE



<https://www.youtube.com/watch?v=RJVL80Gg3IA>





uib.no

