



Unsupervised Learning

Malte Esders, Christopher Anders, Wojciech Samek
Technische Universität Berlin - Machine Learning Group



This presentation

Presentation repository:

https://github.com/Maltimeore/unsupervised_learning_lecture

Contains code to generate plots

Parts not otherwise marked: BSD licensed





Agenda

Introduction and intuition

- Unsupervised Learning

- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance

- Vector norm notation

- Cosine distance

- Z-scores

Clustering algorithms

- Hierarchical Clustering

- K-Means

- DBSCAN

- More examples

Dimensionality reduction

- Motivation

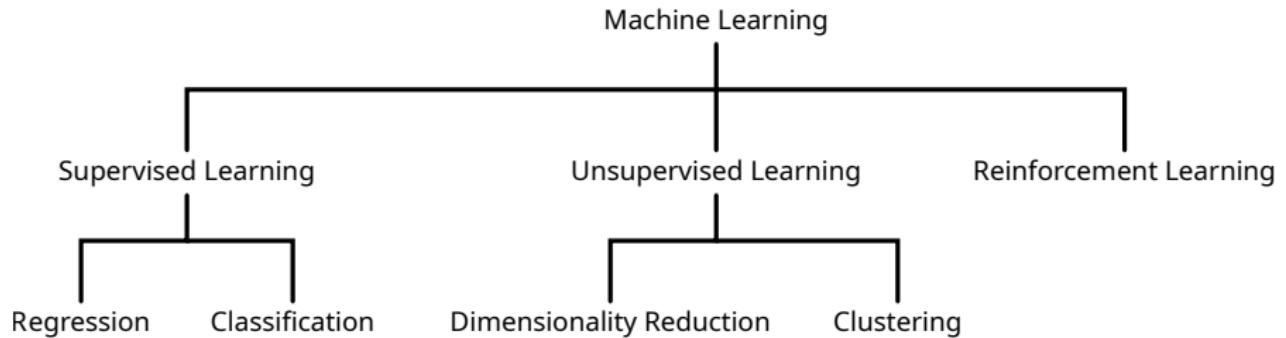
- Principal Component Analysis

- Non-Negative Matrix Factorization

- Nonlinear dimensionality reduction examples



Unsupervised Learning





Unsupervised Learning Motivation

When we're learning to see, nobody's telling us what the right answers are - we just look. Every so often, your mother says 'that's a dog', but that's very little information. You'd be lucky if you got a few bits of information - even one bit per second - that way. The brain's visual system has 10^{14} neural connections. And you only live for 10^9 seconds. So it's no use learning one bit per second. You need more like 10^5 bits per second. And there's only one place you can get that much information: from the input itself.

Geoffrey Hinton, 1996

Genome: $3 * 10^9$ base-pairs (A, T, C, G), with 2 bits per location has $6 * 10^9$ bits.



Unsupervised Learning Motivation

When we're learning to see, nobody's telling us what the right answers are - we just look. Every so often, your mother says 'that's a dog', but that's very little information. You'd be lucky if you got a few bits of information - even one bit per second - that way. The brain's visual system has 10^{14} neural connections. And you only live for 10^9 seconds. So it's no use learning one bit per second. You need more like 10^5 bits per second. And there's only one place you can get that much information: from the input itself.

Geoffrey Hinton, 1996

Genome: $3 * 10^9$ base-pairs (A, T, C, G), with 2 bits per location has $6 * 10^9$ bits.

Yann LeCun's cake analogy (2016):

- Unsupervised learning is the cake itself
- Supervised learning is the icing
- Reinforcement Learning is the cherry on top



Unsupervised Learning Motivation

When we're learning to see, nobody's telling us what the right answers are - we just look. Every so often, your mother says 'that's a dog', but that's very little information. You'd be lucky if you got a few bits of information - even one bit per second - that way. The brain's visual system has 10^{14} neural connections. And you only live for 10^9 seconds. So it's no use learning one bit per second. You need more like 10^5 bits per second. And there's only one place you can get that much information: from the input itself.

Geoffrey Hinton, 1996

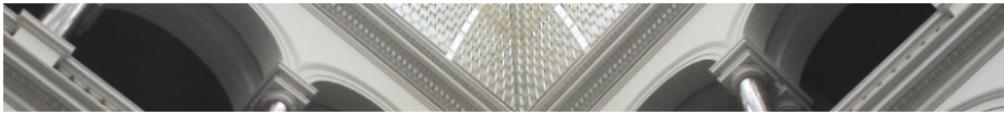
Genome: $3 * 10^9$ base-pairs (A, T, C, G), with 2 bits per location has $6 * 10^9$ bits.

Yann LeCun's cake analogy (2016):

- Unsupervised learning is the cake itself
- Supervised learning is the icing
- Reinforcement Learning is the cherry on top

ChatGPT is unsupervised learning!





Agenda

Introduction and intuition

Unsupervised Learning

Cluster examples

Distance and Similarity measures

Euclidean and Minkowski distance

Vector norm notation

Cosine distance

Z-scores

Clustering algorithms

Hierarchical Clustering

K-Means

DBSCAN

More examples

Dimensionality reduction

Motivation

Principal Component Analysis

Non-Negative Matrix Factorization

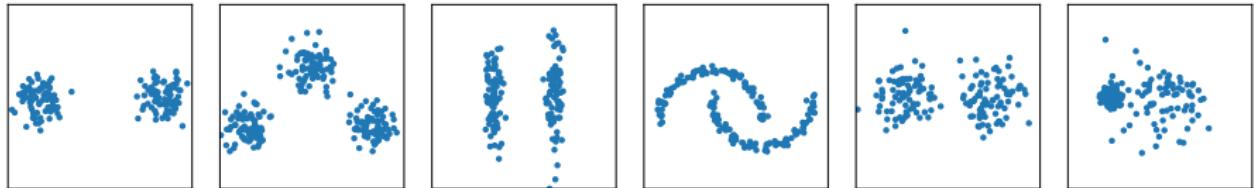
Nonlinear dimensionality reduction examples





Clusters - Intuition

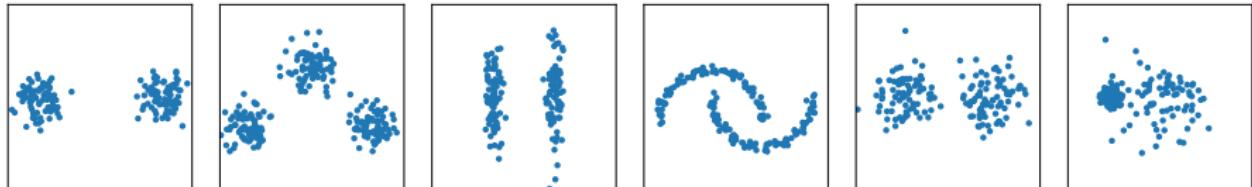
Without labels:



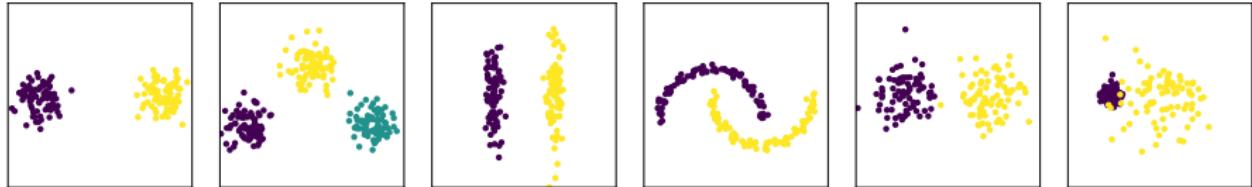


Clusters - Intuition

Without labels:



“True” clusters:



We do not have the “true” clusters with real data!





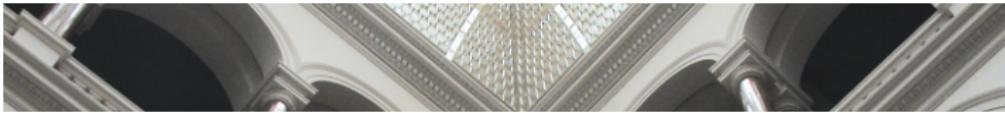
Clusters - Intuition

Clustering is an ill-posed problem: no definition of a cluster.

Some approaches:

- samples within a cluster should be “similar”, samples between clusters “different”
- clusters are regions of densely populated feature space
- samples in the same cluster come from the same underlying process





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples

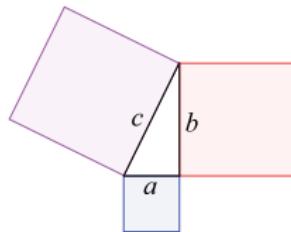




Euclidean distance

Pythagoras formula:

$$c^2 = a^2 + b^2$$
$$c = \sqrt{a^2 + b^2}$$



Wikipedia Pythagorean theorem





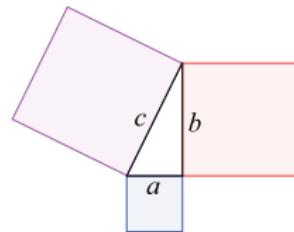
Euclidean distance

Pythagoras formula:

$$c^2 = a^2 + b^2$$
$$c = \sqrt{a^2 + b^2}$$

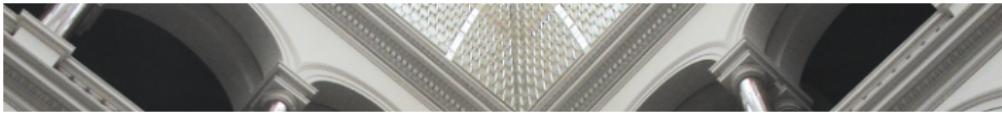
Euclidean distance notation:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\Delta x_1^2 + \Delta x_2^2}$$



Wikipedia Pythagorean theorem





Euclidean distance

Pythagoras formula:

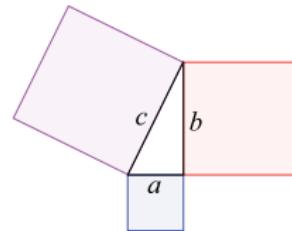
$$c^2 = a^2 + b^2$$
$$c = \sqrt{a^2 + b^2}$$

Euclidean distance notation:

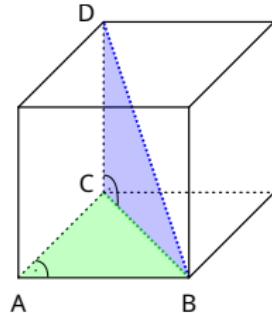
$$d(x, y) = \sqrt{\Delta x_1^2 + \Delta x_2^2}$$

Euclidean distance 3-D:

$$\overline{BC} = \sqrt{\overline{AC}^2 + \overline{AB}^2}$$



Wikipedia Pythagorean theorem





Euclidean distance

Pythagoras formula:

$$c^2 = a^2 + b^2$$
$$c = \sqrt{a^2 + b^2}$$

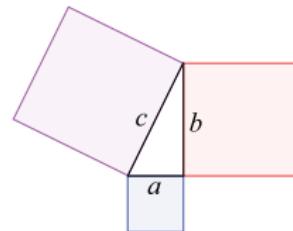
Euclidean distance notation:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\Delta x_1^2 + \Delta x_2^2}$$

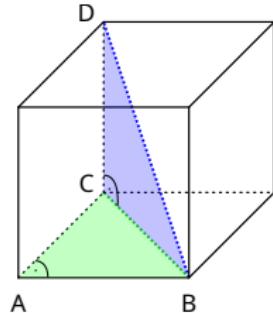
Euclidean distance 3-D:

$$\overline{BC} = \sqrt{\overline{AC}^2 + \overline{AB}^2}$$

$$\overline{BD} = \sqrt{\overline{BC}^2 + \overline{CD}^2}$$



Wikipedia Pythagorean theorem



Euclidean distance

Pythagoras formula:

$$c^2 = a^2 + b^2$$

$$c = \sqrt{a^2 + b^2}$$

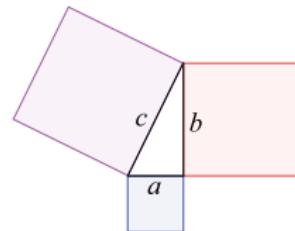
Euclidean distance notation:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\Delta x_1^2 + \Delta x_2^2}$$

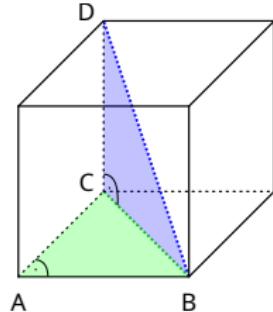
Euclidean distance 3-D:

$$\overline{BC} = \sqrt{\overline{AC}^2 + \overline{AB}^2}$$

$$\overline{BD} = \sqrt{\overline{AC}^2 + \overline{AB}^2 + \overline{CD}^2}$$



Wikipedia Pythagorean theorem





Euclidean distance

Pythagoras formula:

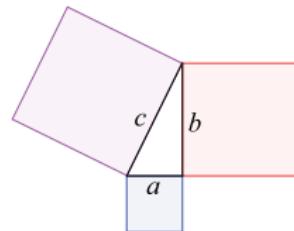
$$c^2 = a^2 + b^2$$
$$c = \sqrt{a^2 + b^2}$$

Euclidean distance notation:

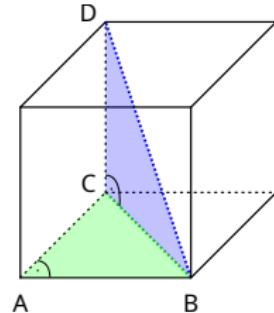
$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\Delta x_1^2 + \Delta x_2^2}$$

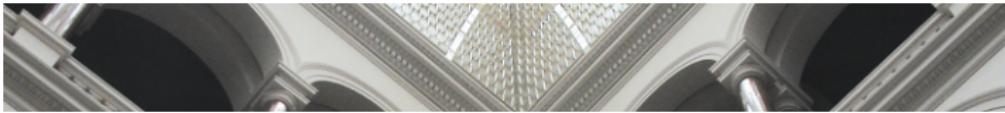
Euclidean distance d-D:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{\sum_{i=1}^d \Delta x_i^2}$$



Wikipedia Pythagorean theorem



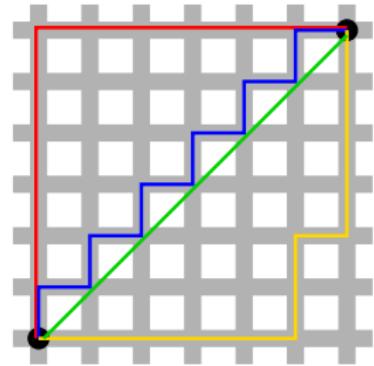


Manhattan distance

“Manhattan” / “Taxicab” distance

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^d |\Delta x_i|$$

Usage: LASSO regularization



Wikipedia Manhattan distance





Minkowski distance

Minkowski distance generalizes Euclidean and Manhattan distance:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |\Delta x_i|^p \right)^{1/p}$$





Minkowski distance

Minkowski distance generalizes Euclidean and Manhattan distance:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |\Delta x_i|^p \right)^{1/p}$$

$p = 1$: Manhattan distance

$p = 2$: Euclidean distance





Minkowski distance

Minkowski distance generalizes Euclidean and Manhattan distance:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |\Delta x_i|^p \right)^{1/p}$$

$p = 1$: Manhattan distance

$p = 2$: Euclidean distance

$p = \infty$:

$$d(\mathbf{x}, \mathbf{y}) = \max_i |\Delta x_i|$$





Minkowski distance

Minkowski distance generalizes Euclidean and Manhattan distance:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |\Delta x_i|^p \right)^{1/p}$$

$p = 1$: Manhattan distance

$p = 2$: Euclidean distance

$p = \infty$:

$$d(\mathbf{x}, \mathbf{y}) = \max_i |\Delta x_i|$$

$p \rightarrow 0$:

$$d(\mathbf{x}, \mathbf{y}) = \min_i |\Delta x_i|$$





Minkowski distance

Minkowski distance generalizes Euclidean and Manhattan distance:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |\Delta x_i|^p \right)^{1/p}$$

$p = 1$: Manhattan distance

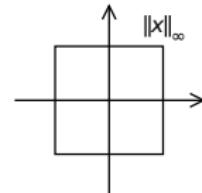
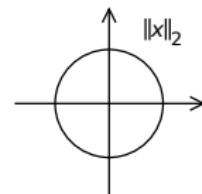
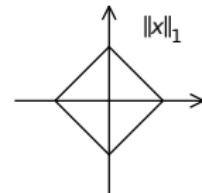
$p = 2$: Euclidean distance

$p = \infty$:

$$d(\mathbf{x}, \mathbf{y}) = \max_i |\Delta x_i|$$

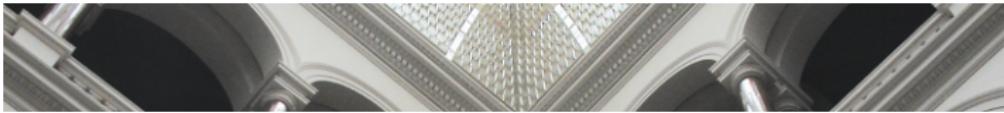
$p \rightarrow 0$:

$$d(\mathbf{x}, \mathbf{y}) = \min_i |\Delta x_i|$$



Wikipedia Norm (mathematics)





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





Vector norm

The distance between vectors is the vector-norm of their difference:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |\Delta x_i|^p \right)^{1/p} = \|\mathbf{x} - \mathbf{y}\|_p$$





Vector norm

The distance between vectors is the vector-norm of their difference:

$$d(\mathbf{x}, \mathbf{y}) = \left(\sum_{i=1}^d |\Delta x_i|^p \right)^{1/p} = \|\mathbf{x} - \mathbf{y}\|_p$$

For $p = 2$, the shorthand is:

$$\|\mathbf{x} - \mathbf{y}\|_2 = \|\mathbf{x} - \mathbf{y}\|$$



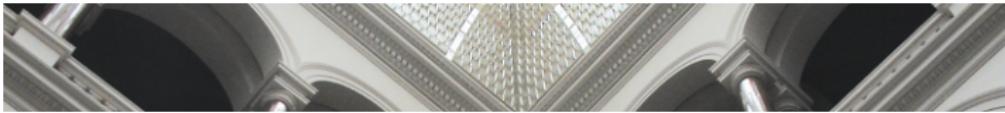


Squared euclidean distance

Squared euclidean distance:

$$||\mathbf{x} - \mathbf{y}||_2^2 = \sum_{i=1}^d |\Delta x_i|^2$$





Squared euclidean distance

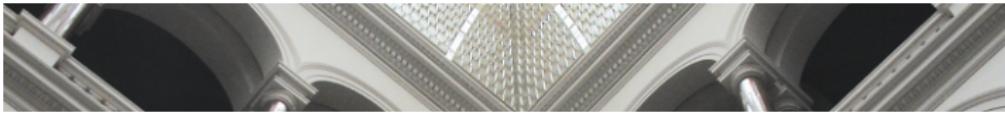
Squared euclidean distance:

$$||\mathbf{x} - \mathbf{y}||_2^2 = \sum_{i=1}^d |\Delta x_i|^2$$

Do not confuse

$$||\mathbf{x} - \mathbf{y}|| = ||\mathbf{x} - \mathbf{y}||_2 \neq ||\mathbf{x} - \mathbf{y}||^2 = ||\mathbf{x} - \mathbf{y}||_2^2$$





Squared euclidean distance

Squared euclidean distance:

$$\|\mathbf{x} - \mathbf{y}\|_2^2 = \sum_{i=1}^d |\Delta x_i|^2$$

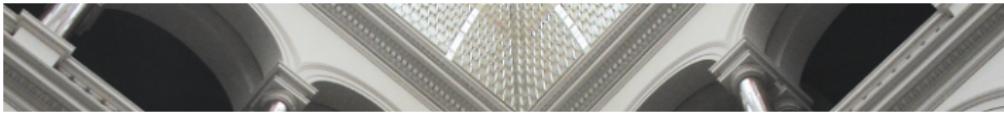
Do not confuse

$$\|\mathbf{x} - \mathbf{y}\| = \|\mathbf{x} - \mathbf{y}\|_2 \neq \|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x} - \mathbf{y}\|_2^2$$

Squared Euclidean distance:

- not a true metric (no triangle inequality)
- overemphasizes big distances (and vice versa)
- makes some optimization problems converge better





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

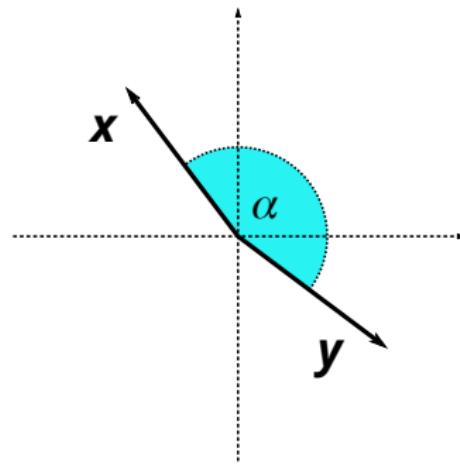
Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples



Cosine distance

Cosine distance only compares angles between vectors.





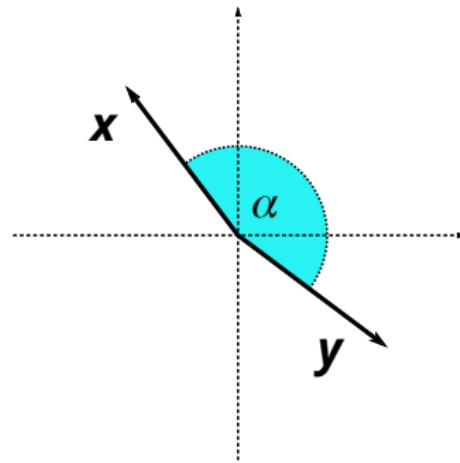
Cosine distance

Cosine distance only compares angles between vectors.

Cosine definition:

$$\mathbf{x}^T \mathbf{y} = |\mathbf{x}| |\mathbf{y}| \cos(\alpha)$$

$$\cos(\alpha) = \frac{\mathbf{x}^T \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$$





Cosine distance

Cosine distance only compares angles between vectors.

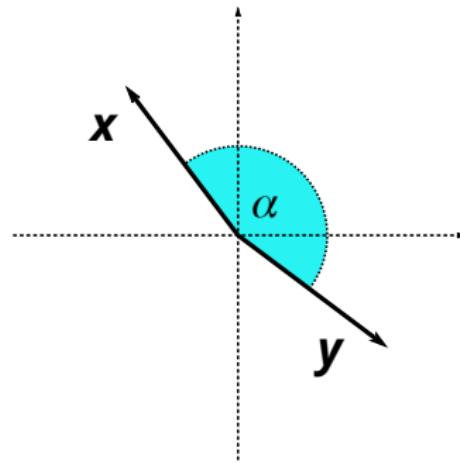
Cosine definition:

$$\mathbf{x}^T \mathbf{y} = |\mathbf{x}| |\mathbf{y}| \cos(\alpha)$$

$$\cos(\alpha) = \frac{\mathbf{x}^T \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$$

Cosine distance:

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{\mathbf{x}^T \mathbf{y}}{|\mathbf{x}| |\mathbf{y}|}$$





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

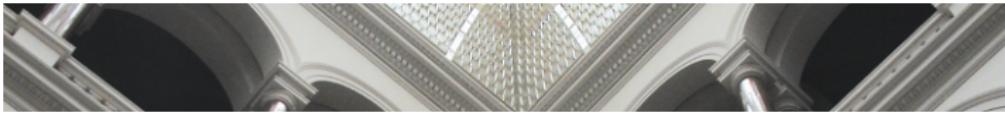
Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





Different scales

Measurements can be at different scales (for instance cm, kg etc.)

We can normalize:

$$z = \frac{x - \mu}{\sigma}$$



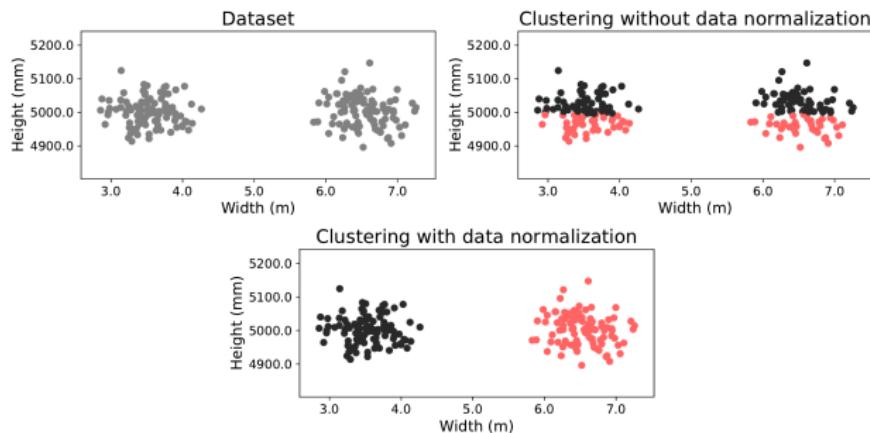


Different scales

Measurements can be at different scales (for instance cm, kg etc.)

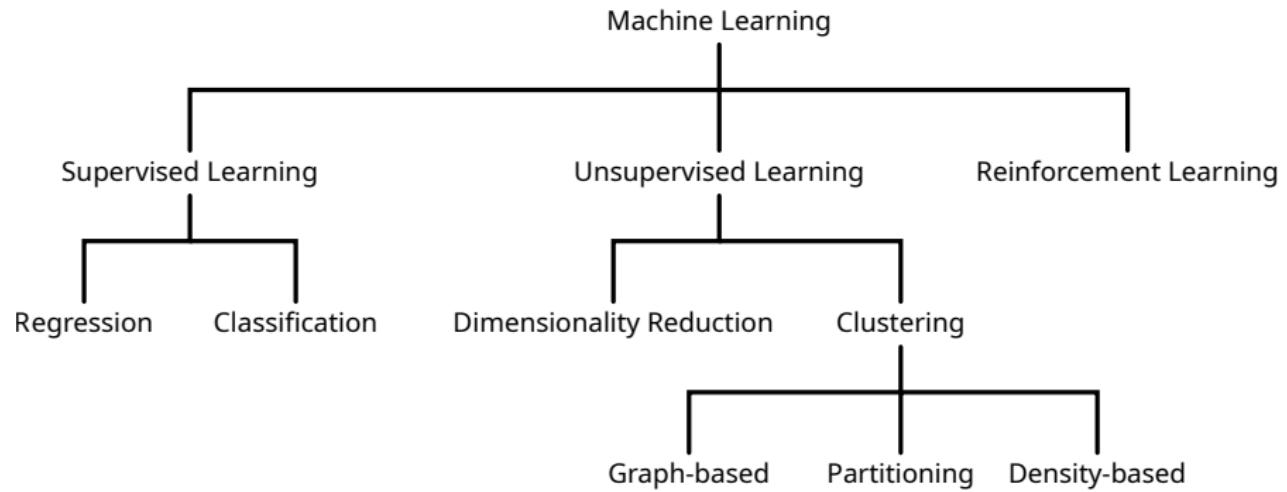
We can normalize:

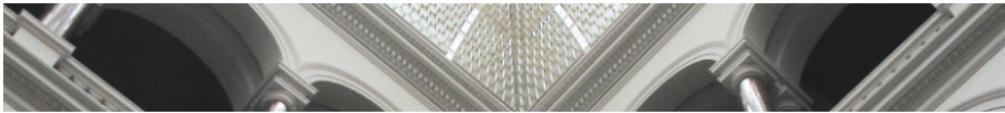
$$z = \frac{x - \mu}{\sigma}$$



Guérin et al. 2017







Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

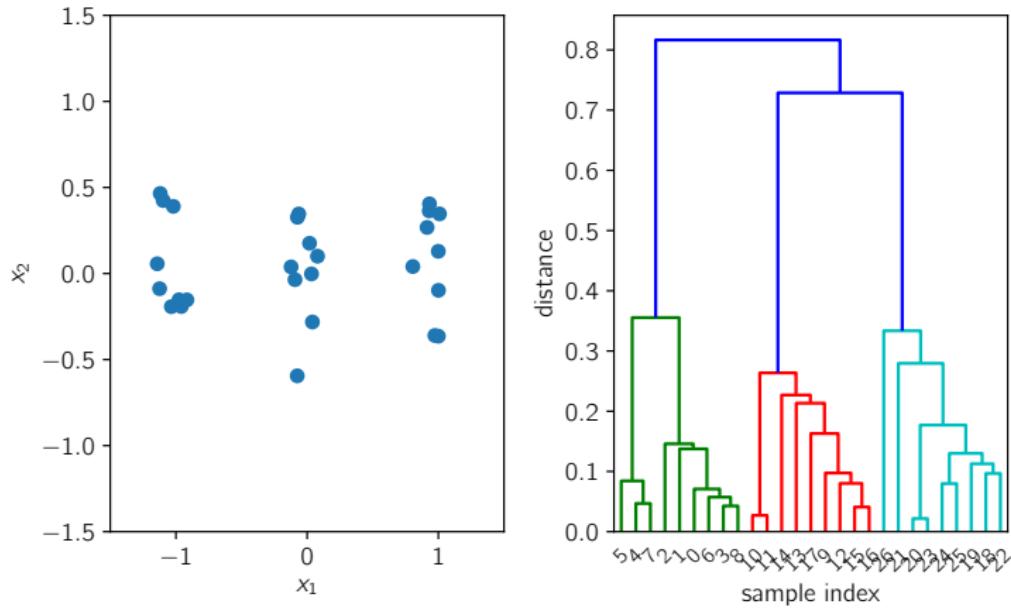
Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





Dendrogram with Euclidean distance





Hierarchical Clustering

- graph-based (no embedding needed!)





Hierarchical Clustering

- graph-based (no embedding needed!)
- builds a hierarchy of clusters, either top down (“divisive”) or bottom up (“agglomerative”)





Hierarchical Clustering

- graph-based (no embedding needed!)
- builds a hierarchy of clusters, either top down (“divisive”) or bottom up (“agglomerative”)
- intuitive





Hierarchical Clustering

- graph-based (no embedding needed!)
- builds a hierarchy of clusters, either top down (“divisive”) or bottom up (“agglomerative”)
- intuitive
- suggests numbers of clusters

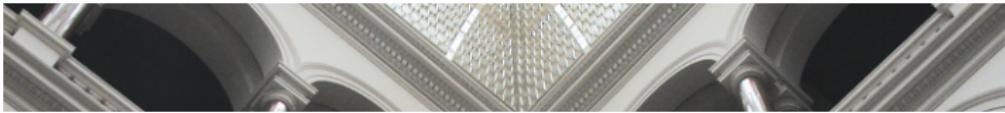




Hierarchical Clustering

- graph-based (no embedding needed!)
- builds a hierarchy of clusters, either top down (“divisive”) or bottom up (“agglomerative”)
- intuitive
- suggests numbers of clusters
- arbitrary cluster shapes





Hierarchical Clustering

- graph-based (no embedding needed!)
- builds a hierarchy of clusters, either top down (“divisive”) or bottom up (“agglomerative”)
- intuitive
- suggests numbers of clusters
- arbitrary cluster shapes
- sensitive to noise





Agglomerative clustering algorithm:

Let \mathcal{D}_i be the sets containing the samples in cluster i .

choose n_{final} clusters

initialize \mathcal{D}_i as singleton clusters

initialize $n_{clusters} = n_{samples}$

while $n_{clusters} > n_{final}$ **do**

| find clusters i and j with smallest distance

| merge \mathcal{D}_i and \mathcal{D}_j

| $n_{clusters} \leftarrow n_{clusters} - 1$

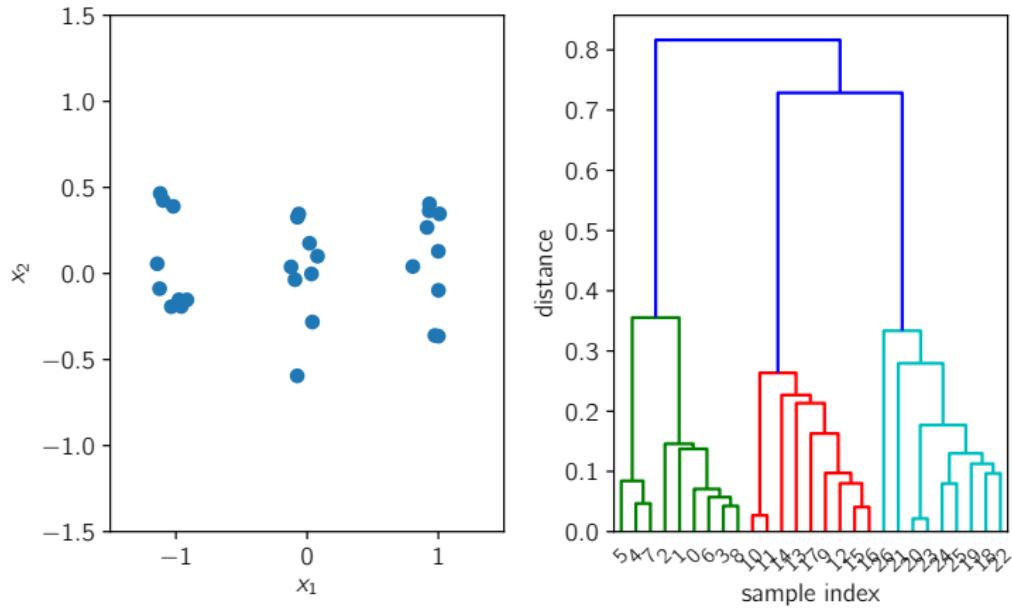
end

return clusters





Dendrogram with Euclidean distance





Cluster distance measures

Let \mathcal{D}_i , \mathcal{D}_j be the sets containing the samples in clusters i and j .

Minimum linkage:

$$d_{min}(\mathcal{D}_i, \mathcal{D}_j) = \min_{\mathbf{x} \in \mathcal{D}_i, \mathbf{x}' \in \mathcal{D}_j} ||\mathbf{x} - \mathbf{x}'||$$





Cluster distance measures

Let \mathcal{D}_i , \mathcal{D}_j be the sets containing the samples in clusters i and j .

Minimum linkage:

$$d_{min}(\mathcal{D}_i, \mathcal{D}_j) = \min_{\mathbf{x} \in \mathcal{D}_i, \mathbf{x}' \in \mathcal{D}_j} ||\mathbf{x} - \mathbf{x}'||$$

Maximum or “complete” linkage:

$$d_{max}(\mathcal{D}_i, \mathcal{D}_j) = \max_{\mathbf{x} \in \mathcal{D}_i, \mathbf{x}' \in \mathcal{D}_j} ||\mathbf{x} - \mathbf{x}'||$$





Cluster distance measures

Let \mathcal{D}_i , \mathcal{D}_j be the sets containing the samples in clusters i and j .

Minimum linkage:

$$d_{min}(\mathcal{D}_i, \mathcal{D}_j) = \min_{\mathbf{x} \in \mathcal{D}_i, \mathbf{x}' \in \mathcal{D}_j} ||\mathbf{x} - \mathbf{x}'||$$

Maximum or “complete” linkage:

$$d_{max}(\mathcal{D}_i, \mathcal{D}_j) = \max_{\mathbf{x} \in \mathcal{D}_i, \mathbf{x}' \in \mathcal{D}_j} ||\mathbf{x} - \mathbf{x}'||$$

Average linkage:

$$d_{avg}(\mathcal{D}_i, \mathcal{D}_j) = \frac{1}{n_i n_j} \sum_{\mathbf{x} \in \mathcal{D}_i} \sum_{\mathbf{x}' \in \mathcal{D}_j} ||\mathbf{x} - \mathbf{x}'||$$

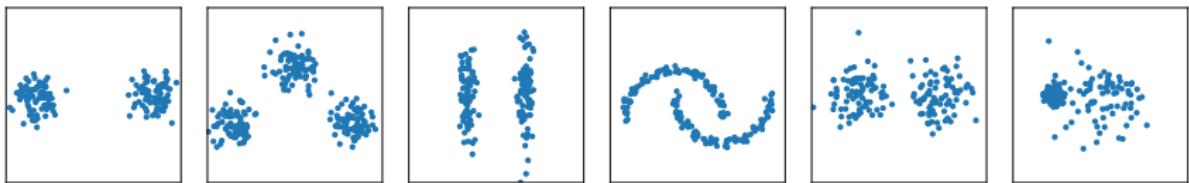
Notation from Duda, Hart, and Stork 2012





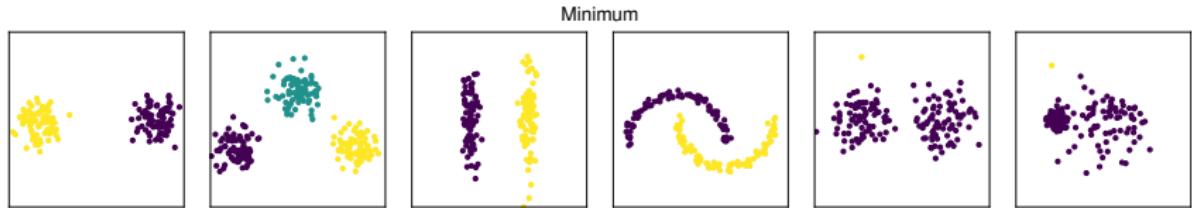
Original data

After clustering, assignments can seem “natural”, good to compare original data:



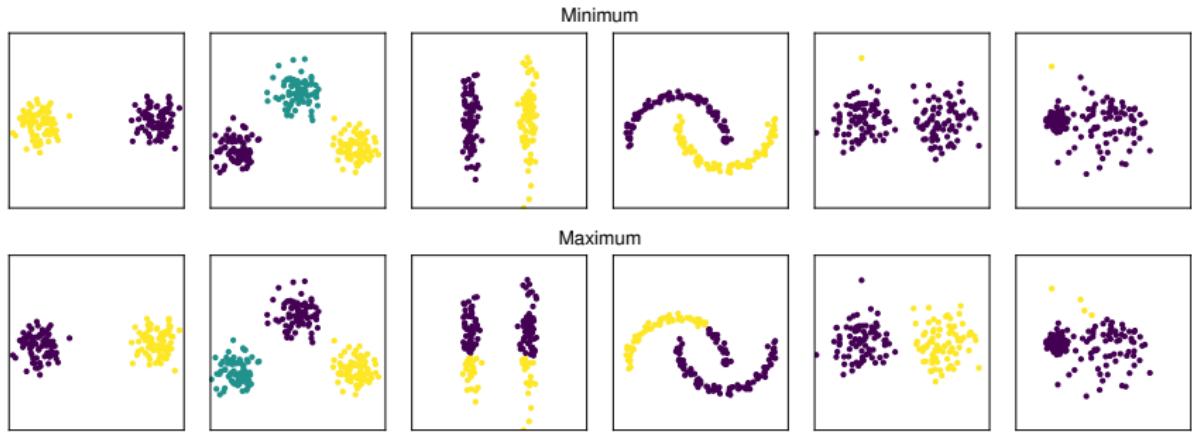


Agglomerative clustering results



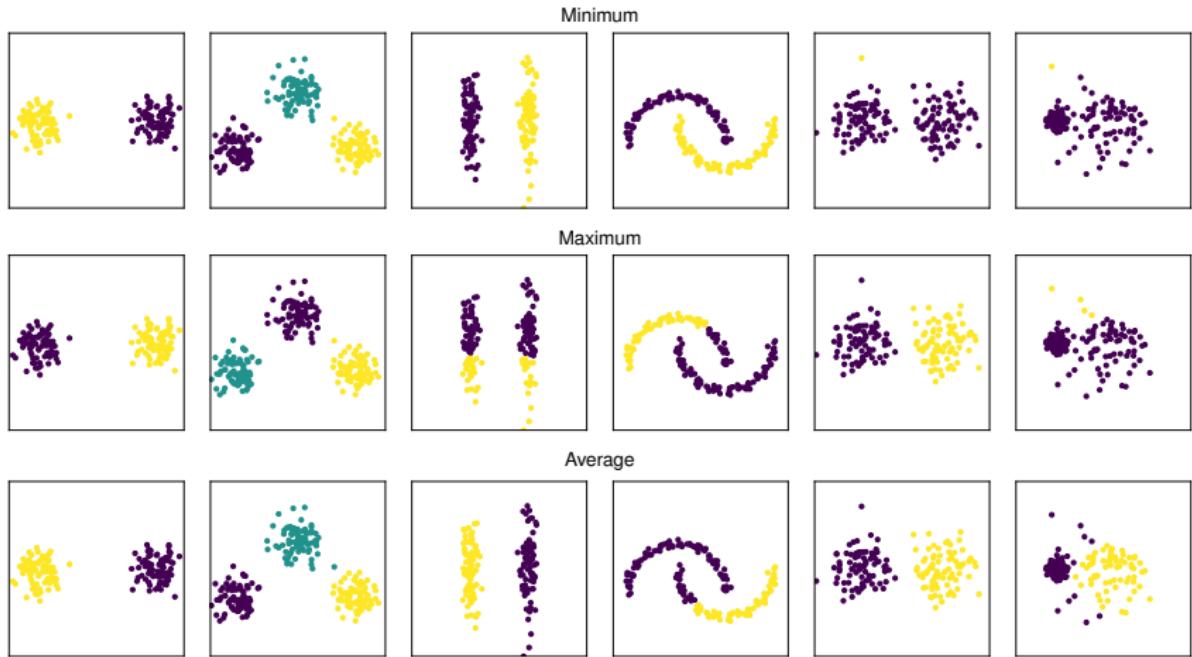


Agglomerative clustering results





Agglomerative clustering results



Tips and tricks

```
class sklearn.cluster.AgglomerativeClustering (n_clusters=2, affinity='euclidean', memory=None,  
connectivity=None, compute_full_tree='auto', linkage='ward', pooling_func=<function mean>)
```

[source]

Agglomerative Clustering

Recursively merges the pair of clusters that minimally increases a given linkage distance.

Read more in the [User Guide](#).

Parameters: `n_clusters` : int, default=2

The number of clusters to find.

`affinity` : string or callable, default: "euclidean"

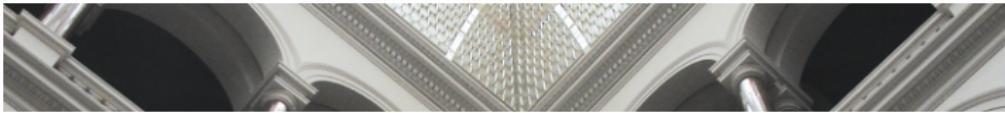
Metric used to compute the linkage. Can be "euclidean", "l1", "l2", "manhattan", "cosine", or 'precomputed'. If linkage is "ward", only "euclidean" is accepted.

Sometimes, documentation can be confusing:

What is the difference between "l2" and "euclidean"?

Or "l1" and "manhattan"?





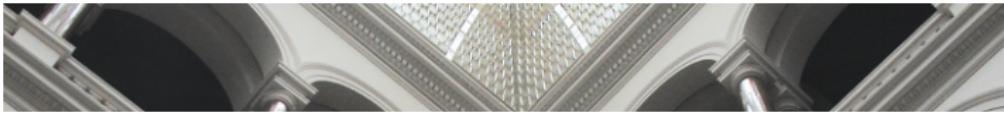
Tips and tricks

Looking into the code:

```
408         elif affinity == 'l2':  
409             # Translate to something understood by scipy  
410             affinity = 'euclidean'  
411         elif affinity in ('l1', 'manhattan'):  
412             affinity = 'cityblock'
```

In the background it does exactly the same.





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

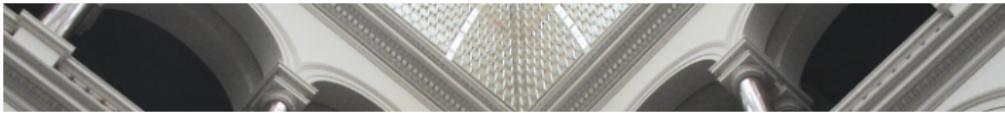
Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





K-Means

Idea: prototype for every cluster

Prototype is mean of all samples in cluster c , \mathcal{D}_c :

$$\mu_c = \frac{1}{n_c} \sum_{i \in \mathcal{D}_c} x_i$$





K-Means

Idea: prototype for every cluster

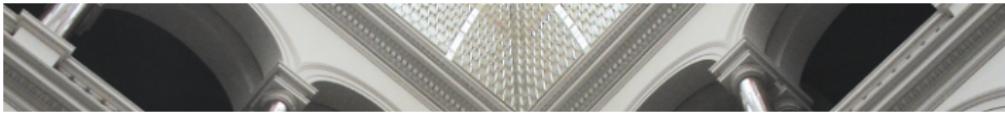
Prototype is mean of all samples in cluster c , \mathcal{D}_c :

$$\mu_c = \frac{1}{n_c} \sum_{i \in \mathcal{D}_c} x_i$$

Minimize distance from prototype to each member in the cluster:

$$\mathcal{L} = \sum_c \sum_{i \in \mathcal{D}_c} d(\mathbf{x}_i, \mu_c)$$





K-Means

Idea: prototype for every cluster

Prototype is mean of all samples in cluster c , \mathcal{D}_c :

$$\mu_c = \frac{1}{n_c} \sum_{i \in \mathcal{D}_c} x_i$$

Minimize distance from prototype to each member in the cluster:

$$\mathcal{L} = \sum_c \sum_{i \in \mathcal{D}_c} d(\mathbf{x}_i, \mu_c)$$

Finding a solution is NP-hard, but in practice straightforward

K-means cannot work with pairwise data such as user ratings





K-Means algorithm

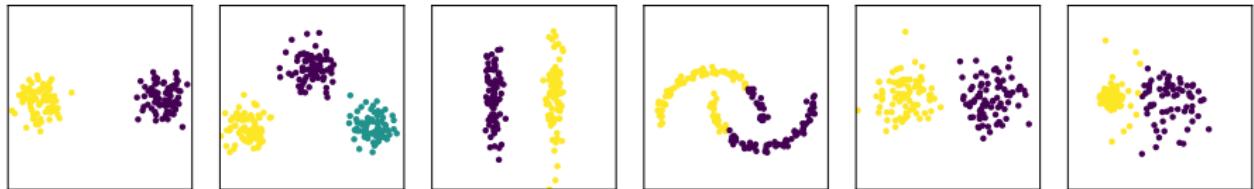
```
initialize  $\mu_c$  to random samples
while clusters change do
    assign  $\mu_c$  to cluster means
    assign samples to clusters
end
return clusters
```





K-Means clusters

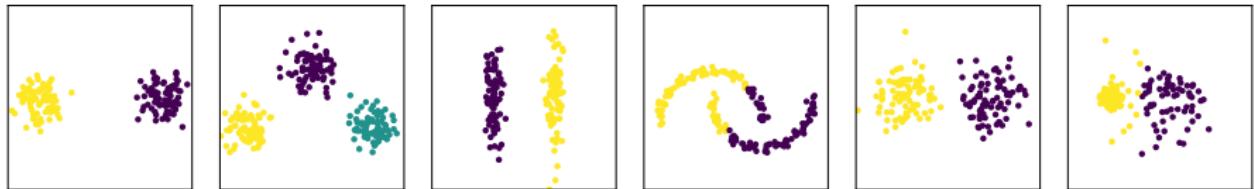
K-Means assignment:



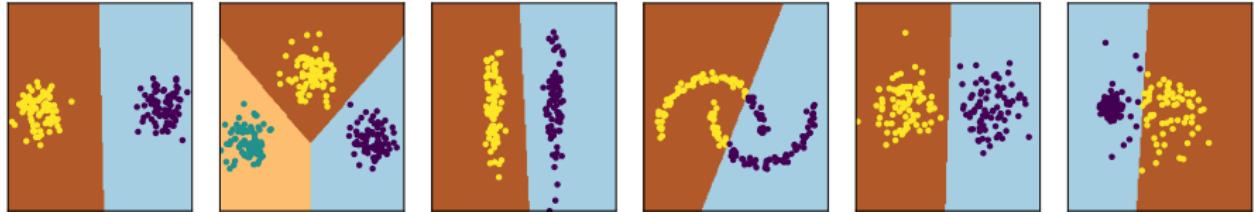


K-Means clusters

K-Means assignment:

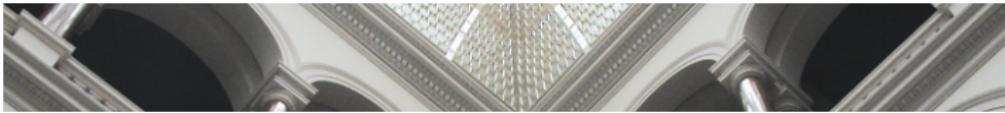


Voronoi tessellation:



Voronoi diagram not possible with hierarchical clustering!





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

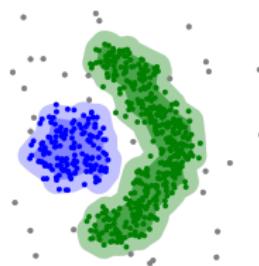
- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





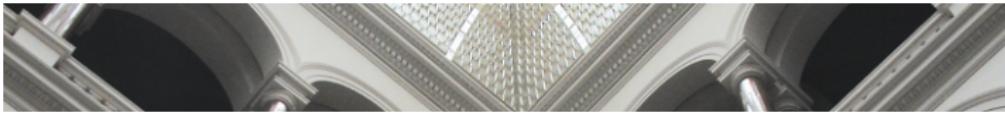
DBSCAN

- clusters are regions of high density separated by regions of low density
- arbitrary cluster shapes



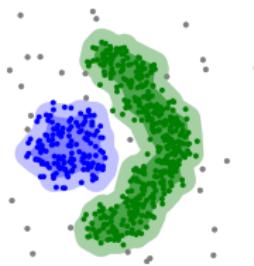
(Wikipedia "DBSCAN")





DBSCAN

- clusters are regions of high density separated by regions of low density
- arbitrary cluster shapes



(Wikipedia "DBSCAN")

- very simple, but effective (2014 SIGKDD test-of-time award (Schubert et al. 2017))
- explicit noise point modeling





DBSCAN

Two parameters:

- ϵ : distance, defines a neighborhood
- MinPts: amount of neighbors to be a core point

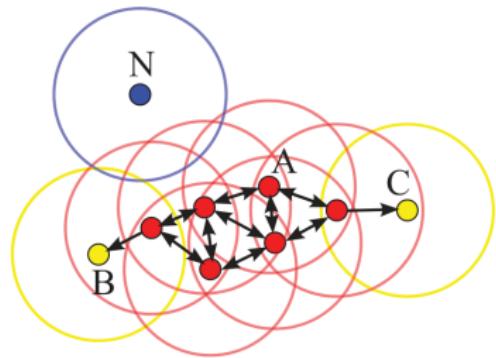




DBSCAN

Two parameters:

- ϵ : distance, defines a neighborhood
- MinPts: amount of neighbors to be a core point



(Wikipedia "DBSCAN")





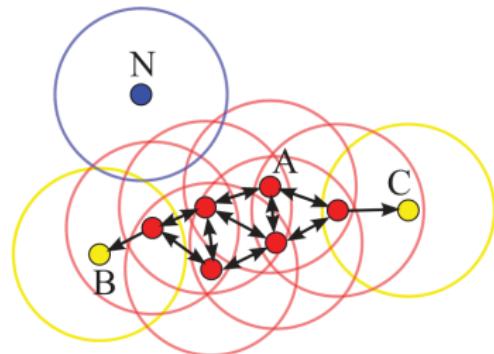
DBSCAN

Two parameters:

- ϵ : distance, defines a neighborhood
- MinPts: amount of neighbors to be a core point

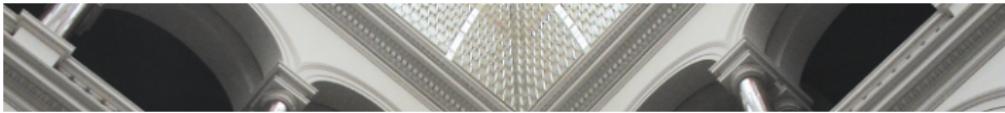
Three types of points:

- core: has MinPts points in its epsilon-neighborhood
- border: within epsilon-neighborhood of core point
- noise: all other points



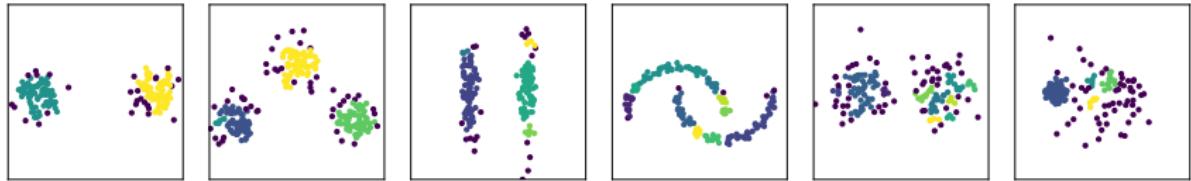
(Wikipedia "DBSCAN")





DBSCAN clustering results

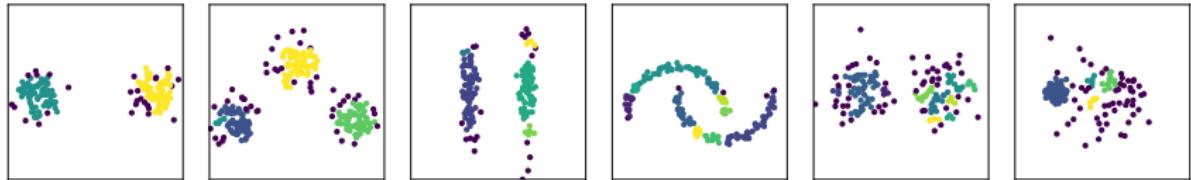
$\epsilon = 0.1, \text{MinPts} = 4$



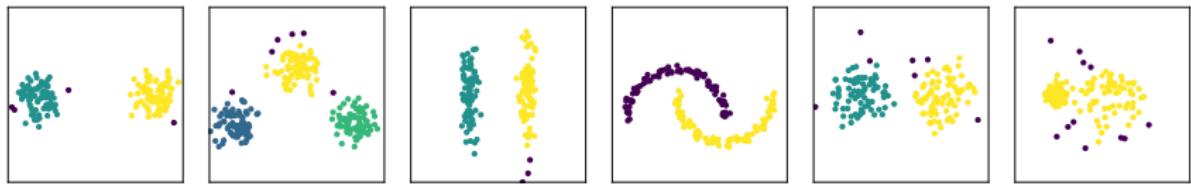


DBSCAN clustering results

$\epsilon = 0.1, \text{MinPts} = 4$



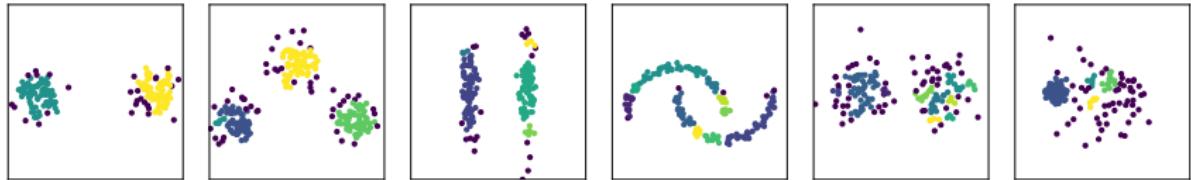
$\epsilon = 0.2, \text{MinPts} = 4$



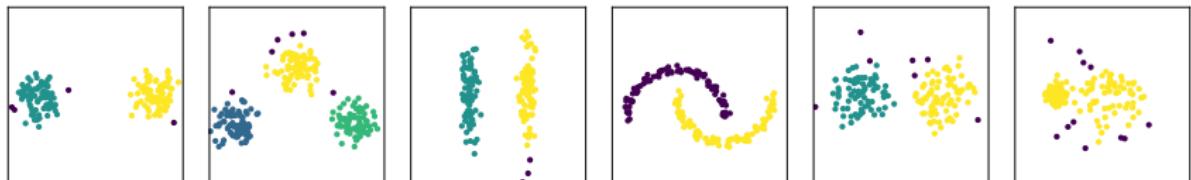


DBSCAN clustering results

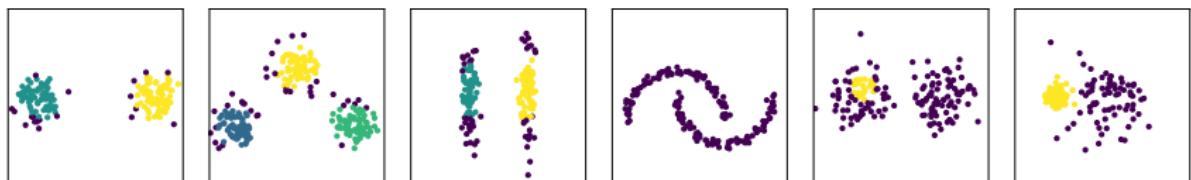
$\epsilon = 0.1, \text{MinPts} = 4$

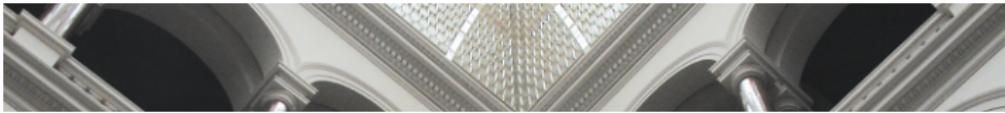


$\epsilon = 0.2, \text{MinPts} = 4$



$\epsilon = 0.2, \text{MinPts} = 20$





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN

More examples

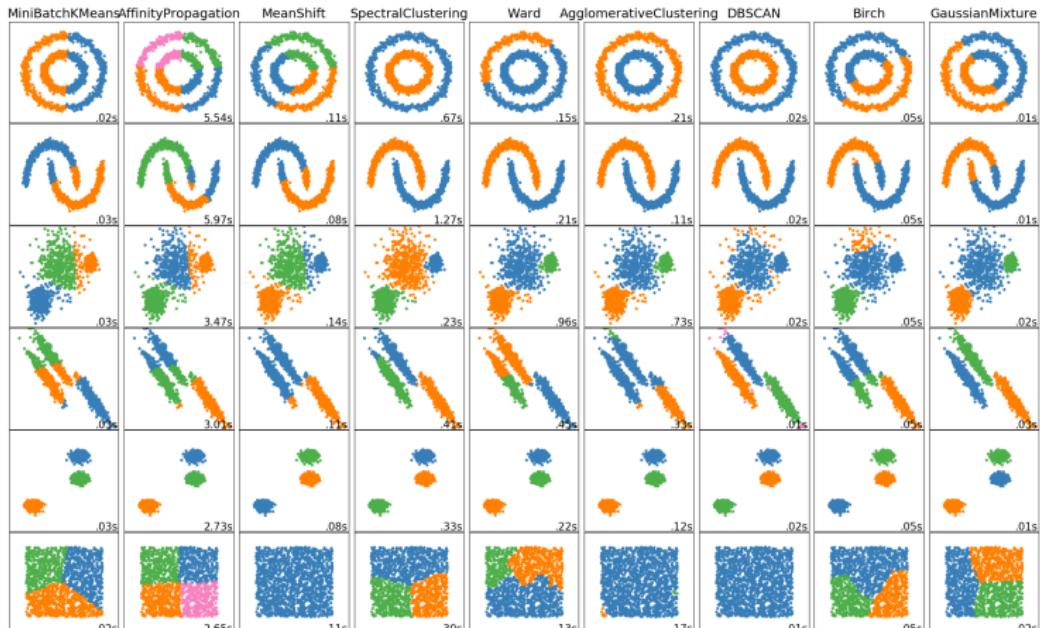
Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples



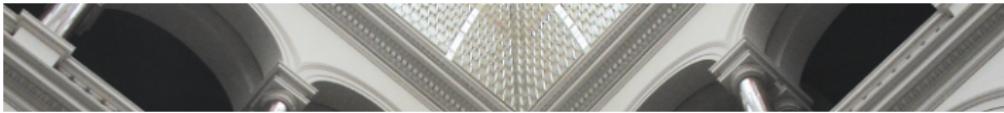


Overview of clustering algorithms



<http://scikit-learn.org/stable/modules/clustering.html> Pedregosa et al. 2011





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

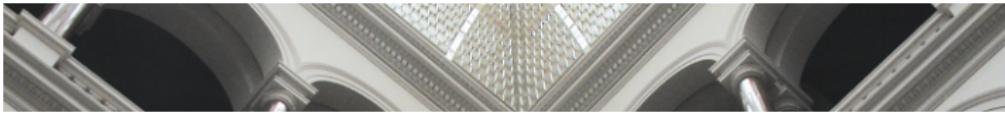
Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





Dimensionality Reduction - Motivation

In many applications, we have

- high-dimensional data
 - reason to believe they lie close to a lower dimensional subspace
- Fewer parameters needed to account for the data properties
hidden causes or latent variables



Dimensionality Reduction - Motivation

In many applications, we have

- high-dimensional data
- reason to believe they lie close to a lower dimensional subspace
- Fewer parameters needed to account for the data properties
hidden causes or latent variables

Examples:

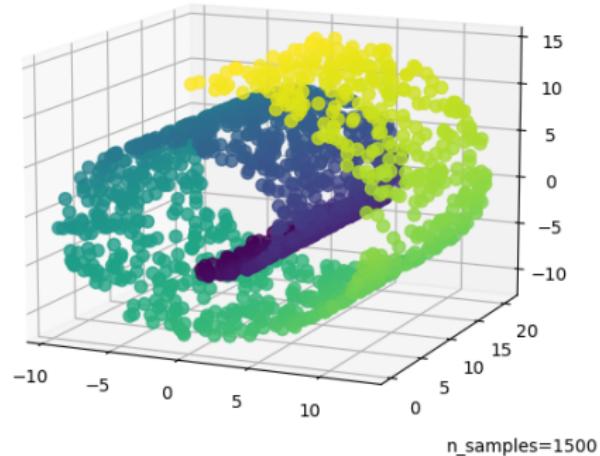
- classify high resolution images
- make a predictive model based on hundreds of customer attributes
- analyse high dimensional neural data
- detect trends in news data





Swiss roll data

Swiss Roll in Ambient Space





Why Dimensionality Reduction

- **Visualization:**

Insights into high-dimensional structures in the data

- **Better Generalization:**

Fewer dimensions → less chances of overfitting

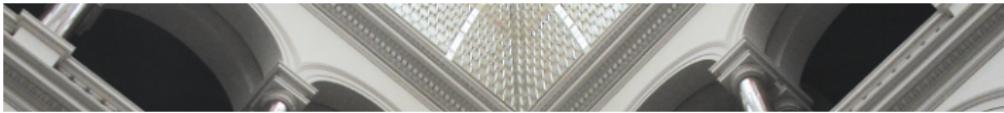
- **Speeding up** learning algorithms:

Most algorithms scale badly with increasing data dimensionality

- **Data compression:**

Less storage requirements





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





The mathematical model for dimensionality reduction

We have

- high-dimensional data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$
 - reason to believe they lie close to a lower dimensional subspace
- $k < D$ parameters needed to account for the data properties
hidden causes or latent variables $H \in \mathbb{R}^{k \times N}$





The mathematical model for dimensionality reduction

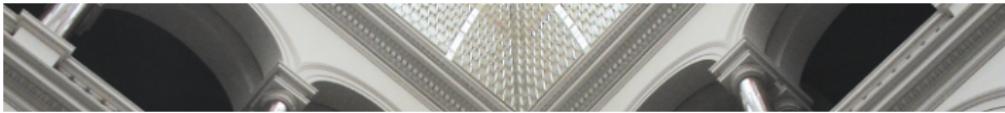
We have

- high-dimensional data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$
 - reason to believe they lie close to a lower dimensional subspace
- $k < D$ parameters needed to account for the data properties
hidden causes or *latent variables* $H \in \mathbb{R}^{k \times N}$

Goal: Find $k < D$ hidden causes $H \in \mathbb{R}^{k \times N}$, that explain the observed data via a mixing $W \in \mathbb{R}^{D \times k}$:

$$X \approx WH$$



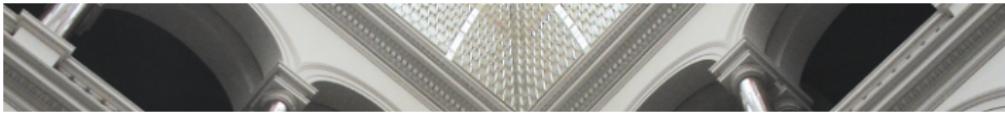


Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal





Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal

$$\text{Var}(\mathbf{w}^\top X) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - E(\mathbf{w}^\top \mathbf{x}))^2$$





Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal

$$\begin{aligned}\text{Var}(\mathbf{w}^\top X) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - E(\mathbf{w}^\top \mathbf{x}))^2 \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - \mathbf{w}^\top E(\mathbf{x}))^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top (\mathbf{x}_n - E(\mathbf{x})))^2\end{aligned}$$





Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal

$$\begin{aligned}\text{Var}(\mathbf{w}^\top X) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - \mathbb{E}(\mathbf{w}^\top \mathbf{x}))^2 \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - \mathbf{w}^\top \mathbb{E}(\mathbf{x}))^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top (\mathbf{x}_n - \mathbb{E}(\mathbf{x})))^2 \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{w}^\top (\mathbf{x}_n - \mathbb{E}(\mathbf{x})) \cdot (\mathbf{x}_n - \mathbb{E}(\mathbf{x}))^\top \mathbf{w}\end{aligned}$$





Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal

$$\begin{aligned}
 \text{Var}(\mathbf{w}^\top X) &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - \mathbb{E}(\mathbf{w}^\top \mathbf{x}))^2 \\
 &= \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - \mathbf{w}^\top \mathbb{E}(\mathbf{x}))^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^\top (\mathbf{x}_n - \mathbb{E}(\mathbf{x})))^2 \\
 &= \frac{1}{N} \sum_{n=1}^N \mathbf{w}^\top (\mathbf{x}_n - \mathbb{E}(\mathbf{x})) \cdot (\mathbf{x}_n - \mathbb{E}(\mathbf{x}))^\top \mathbf{w} \\
 &= \mathbf{w}^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mathbb{E}(\mathbf{x})) \cdot (\mathbf{x}_n - \mathbb{E}(\mathbf{x}))^\top \right)}_{\text{Covariance matrix } S} \mathbf{w}
 \end{aligned}$$





Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal

$$\text{Var}(\mathbf{w}^\top X) = \underbrace{\mathbf{w}^\top \left(\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mathbb{E}(\mathbf{x})) \cdot (\mathbf{x}_n - \mathbb{E}(\mathbf{x}))^\top \right) \mathbf{w}}_{\text{Covariance matrix } S}$$





Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal

$$\begin{aligned}\text{Var}(\mathbf{w}^\top X) &= \mathbf{w}^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mathbf{E}(\mathbf{x})) \cdot (\mathbf{x}_n - \mathbf{E}(\mathbf{x}))^\top \right)}_{\text{Covariance matrix } S} \mathbf{w} \\ &= \mathbf{w}^\top X X^\top \mathbf{w}\end{aligned}$$

where we assume centered data





Principal Component Analysis

We obtained some data $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^\top X$ is maximal

$$\begin{aligned}\text{Var}(\mathbf{w}^\top X) &= \mathbf{w}^\top \underbrace{\left(\frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \mathbf{E}(\mathbf{x})) \cdot (\mathbf{x}_n - \mathbf{E}(\mathbf{x}))^\top \right)}_{\text{Covariance matrix } S} \mathbf{w} \\ &= \mathbf{w}^\top X X^\top \mathbf{w}\end{aligned}$$

where we assume centered data

And: we need to constrain \mathbf{w}





Principal Component Analysis

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^T X$ is maximal

$$\arg \max_{\mathbf{w}} \frac{\mathbf{w}^T S \mathbf{w}}{\mathbf{w}^T \mathbf{w}}$$

This objective function is independent of the scaling of \mathbf{w} .





Principal Component Analysis

PCA finds a direction $\mathbf{w} \in \mathbb{R}^D$ such that the variance of the projected data $\mathbf{w}^T X$ is maximal

$$\arg \max_{\mathbf{w}} \frac{\mathbf{w}^T S \mathbf{w}}{\mathbf{w}^T \mathbf{w}}$$

This objective function is independent of the scaling of \mathbf{w} .

Note the similarity to the objective of Linear Discriminant Analysis!

→ Different covariance matrices, different problem, but: same maths solve it





Principal Component Analysis

$$\arg \max_w \frac{\mathbf{w}^\top S \mathbf{w}}{\mathbf{w}^\top \mathbf{w}}$$

Set the derivative w.r.t \mathbf{w} to zero:

$$\frac{(\mathbf{w}^\top \mathbf{w})2S\mathbf{w} - (\mathbf{w}^\top S\mathbf{w})2\mathbf{w}}{(\mathbf{w}^\top \mathbf{w})^2} = 0$$

$$(\mathbf{w}^\top S\mathbf{w})\mathbf{w} = (\mathbf{w}^\top \mathbf{w})S\mathbf{w}$$

$$\mathbf{w} = S\mathbf{w} \cdot \underbrace{\frac{\mathbf{w}^\top \mathbf{w}}{\mathbf{w}^\top S\mathbf{w}}}_{scalar}$$

$$S\mathbf{w} = \lambda \mathbf{w}$$

This is the standard eigenvalue problem.





Principal Component Analysis

For $S\mathbf{w} = \lambda\mathbf{w}$, we see that the variance in direction \mathbf{w} is given by:

$$\arg \max_{\mathbf{w}} \frac{\mathbf{w}^\top S \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} = \frac{\mathbf{w}^\top \lambda \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} = \lambda$$

The variance of the projected data in an eigendirection \mathbf{w} is given by the corresponding eigenvalue!





Principal Component Analysis

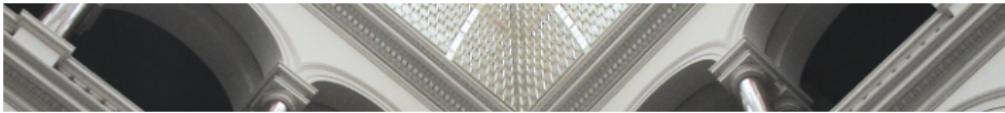
For $\mathbf{S}\mathbf{w} = \lambda\mathbf{w}$, we see that the variance in direction \mathbf{w} is given by:

$$\arg \max_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{S} \mathbf{w}}{\mathbf{w}^T \mathbf{w}} = \frac{\mathbf{w}^T \lambda \mathbf{w}}{\mathbf{w}^T \mathbf{w}} = \lambda$$

The variance of the projected data in an eigendirection \mathbf{w} is given by the corresponding eigenvalue!

The direction of maximal variance in the data is equal to the eigenvector having the largest eigenvalue.





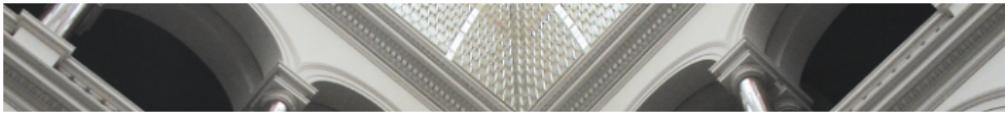
Principal Component Analysis

Require: data $x_1, \dots, x_N \in \mathbb{R}^d$, number of principal components k

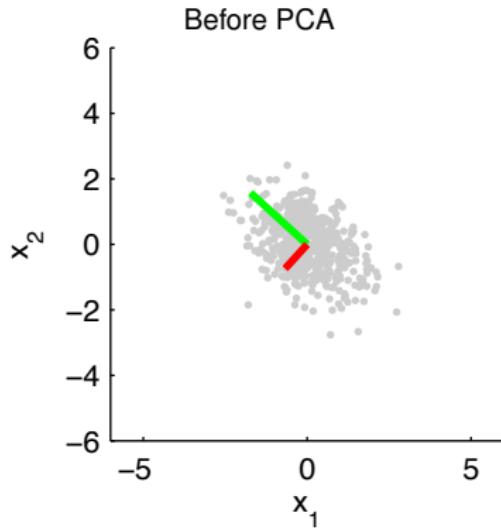
- 1: # Center Data
- 2: $X = X - 1/N \sum_i x_i$
- 3: # Compute Covariance Matrix
- 4: $C = 1/N X X^\top$
- 5: # Compute eigenvectors corresponding to the k largest eigenvalues
- 6: $W = \text{eig}(C)$
- 7: # Project data onto W
- 8: $H = W^\top X$
- 9: **return** W, H

Algorithm 1: Principal Component Analysis



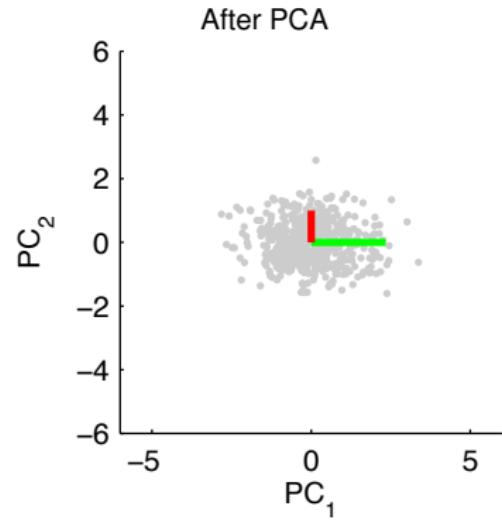
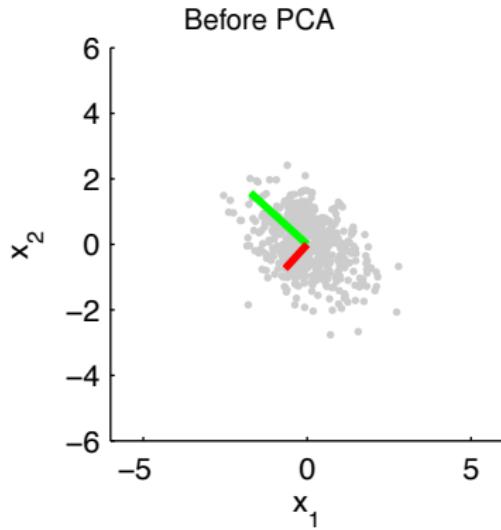


Principal Component Analysis

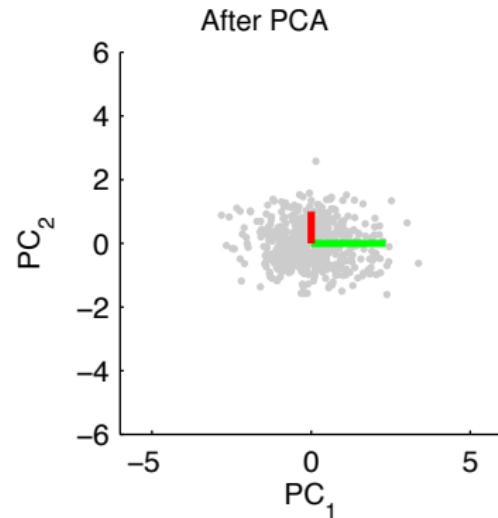
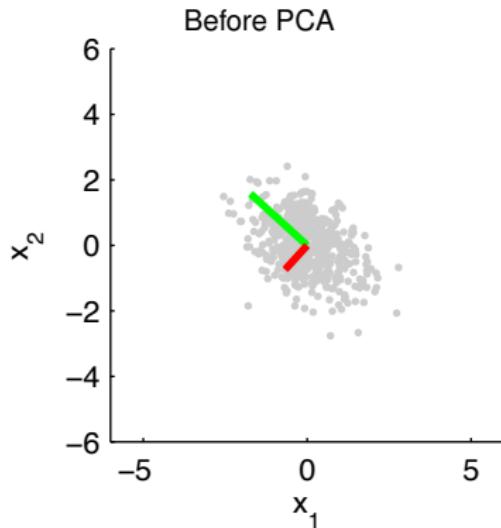




Principal Component Analysis



Principal Component Analysis

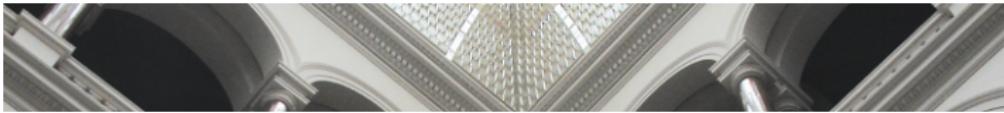


PCA aligns maximum variance directions with standard basis

→ Variance along each dimension is **uncorrelated**

→ Now we can remove each dimension separately





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization
- Nonlinear dimensionality reduction examples





Non-Negative Matrix Factorization

- For some data PCA does not make sense
- Example: Non-negative data
 - Principal directions will have negative entries
 - This can be hard to interpret
- Many data sets are strictly positive
 - Text data
 - Image data
 - Probabilistic data
- Very easy to implement





Non-Negative Matrix Factorization

Given non-negative data $X \in \mathbb{R}_+^{D \times N}$ we want to find $W \in \mathbb{R}_+^{D \times C}$, $H \in \mathbb{R}_+^{C \times N}$ such that

$$\arg \min_{W,H} \|X - WH\|_{\text{Fro}}^2 = \arg \min_{W,H} \sum_{d=1}^D \sum_{n=1}^N (X_{dn} - (WH)_{dn})^2$$





Non-Negative Matrix Factorization

- Local minima





Non-Negative Matrix Factorization

- Local minima
- Gradient descent finds an optimal solution by iterating

$$H \leftarrow H + \eta \left(W^T WH - X^T W \right)$$

$$W \leftarrow W + \eta \left(WHH^T - XH^T \right)$$





Non-Negative Matrix Factorization

- Local minima
- Gradient descent finds an optimal solution by iterating

$$H \leftarrow H + \eta \left(W^T WH - X^T W \right)$$
$$W \leftarrow W + \eta \left(WHH^T - XH^T \right)$$

- By choosing η wisely one can transform the additive updates into multiplicative ones:

$$H = H \odot W^T X \oslash W^T WH$$

$$W = W \odot XH^T \oslash WHH^T$$

where

- is *element-wise* multiplication (in python `**`)
- is *element-wise* division (in python `/`)



NMF Algorithm

Require: data $X = [x_1, \dots, x_N] \in \mathbb{R}_+^{D \times N}$, number of factors k

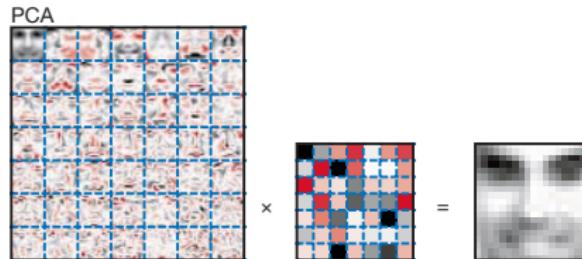
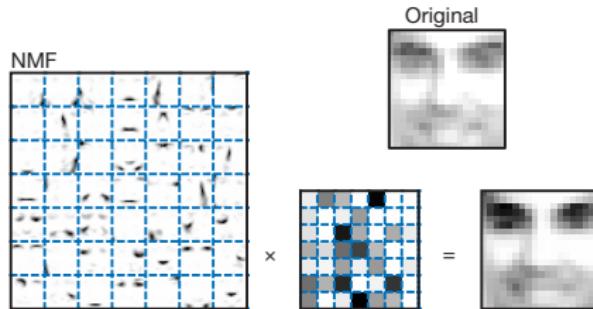
```
1: # Initialize  $W \in \mathbb{R}_+^{D \times k}$ ,  $H \in \mathbb{R}_+^{k \times N}$  randomly
2: # Add a small constant  $\epsilon = 10^{-19}$  to  $X$  to avoid zero-divisions
3: for it  $\leq$  Iterations do
4:    $H = H \odot W^\top X \oslash W^\top WH$ 
5:    $W = W \odot XH^\top \oslash WHH^\top$ 
6: end for
7: return  $W, H$ 
```

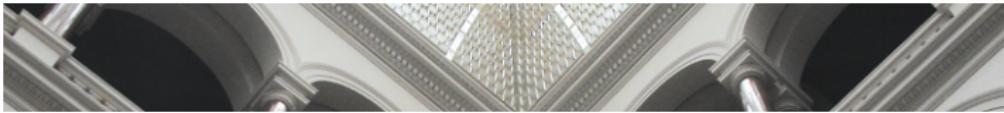
Algorithm 2: Non-negative Matrix Factorization





NMF Applications - Face Parts





Agenda

Introduction and intuition

- Unsupervised Learning
- Cluster examples

Distance and Similarity measures

- Euclidean and Minkowski distance
- Vector norm notation
- Cosine distance
- Z-scores

Clustering algorithms

- Hierarchical Clustering
- K-Means
- DBSCAN
- More examples

Dimensionality reduction

- Motivation
- Principal Component Analysis
- Non-Negative Matrix Factorization

Nonlinear dimensionality reduction examples



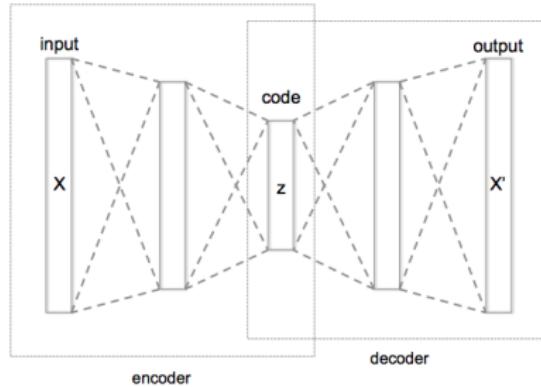


Autoencoder

Two parts:

Encoder: A neural network that forces the input into a low-dimensional "code"

Decoder: A neural network that tries to reconstruct the input from the "code"



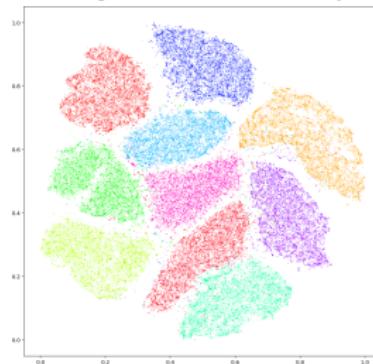
Wikipedia autoencoder





t-SNE (t-distributed Stochastic Neighbor Embedding)

Unlike Autoencoder, non-parametric, i.e. can not easily embed new points.
Local similarities in the data preserved, higher distances not preserved.



Wikipedia t-SNE. A visualization of MNIST in two t-SNE dimensions.





References

- Duda, Richard O, Peter E Hart, and David G Stork (2012). **Pattern classification.** John Wiley & Sons.
- Guérin, Joris et al. (Mar. 2017). "Clustering for Different Scales of Measurement - the Gap-Ratio Weighted K-means Algorithm". [en. In: arXiv:1703.07625 \[cs, stat\]](#). arXiv: 1703.07625. URL: <http://arxiv.org/abs/1703.07625> (visited on 08/29/2018).
- Pedregosa, Fabian et al. (2011). "Scikit-learn: Machine learning in Python". In: **Journal of machine learning research** 12.Oct, pp. 2825–2830.
- Schubert, Erich et al. (2017). "DBSCAN revisited, revisited: why and how you should (still) use DBSCAN". In: **ACM Transactions on Database Systems (TODS)** 42.3, p. 19.

