

Applied Data Analysis (CS401)



Lecture 10 Unsupervised learning 2018/11/22



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Robert West



Announcements

- Course [evaluations](#) are being collected (participate by Sunday)
- Poster session: Jan 21, 22, or 23 -- vote [here](#)!
- HW 4 was due yesternight
- HW 3 grades are out
- Project milestone 2 due on Sun, Nov 25, 23:59
- Tomorrow's lab session:
 - Project office hours (sign up [here](#))
 - Guest talk: Ryan Faulkner
(Google **DeepMind**, London)



Give us feedback on this lecture here:

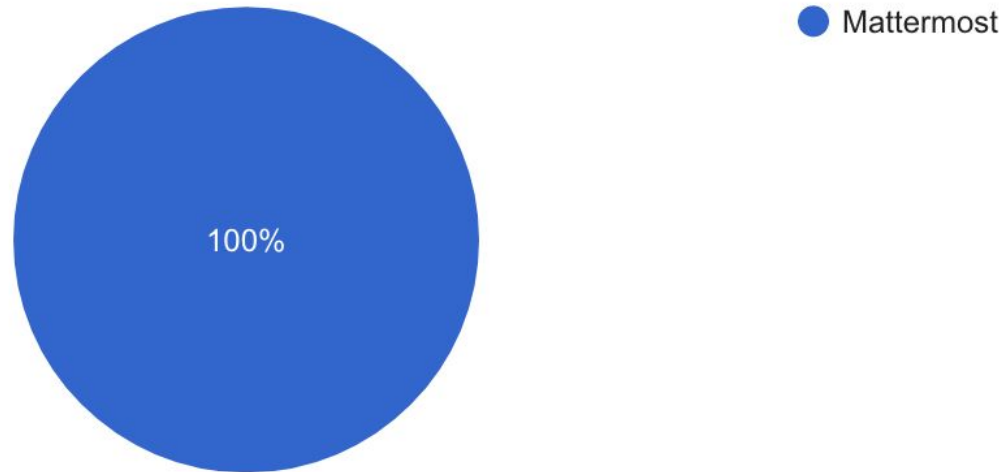
<https://go.epfl.ch/ada2018-lec10-feedback>

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more (fewer) details?
- ...

Last week's feedback form [\[link\]](#)

What's your favorite messaging system?

2 responses

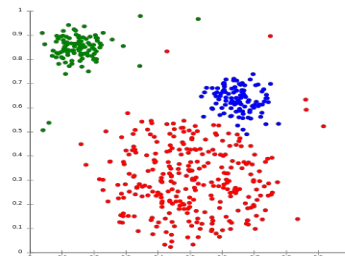


Machine Learning

- **Supervised:** We are given input/output samples (x, y) which we relate with a function $y = f(x)$. We would like to “learn” f , and evaluate it on new data. Types:
 - **Classification:** y is discrete (class labels).
 - **Regression:** y is continuous, e.g. linear regression.
- **Unsupervised:** Given only samples x of the data, we compute a function f such that $y = f(x)$ is a “simpler” representation.
 - Discrete y : **clustering**
 - Continuous y : **dimensionality reduction** (e.g., matrix factorization, unsupervised neural networks)

The clustering problem

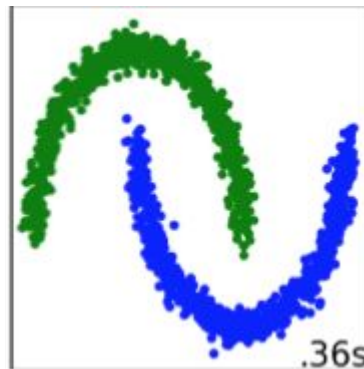
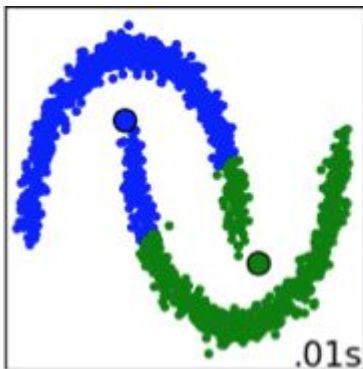
- Given a **set of points**, with a notion of **distance** between points, **group the points** into some number of ***clusters***, such that
 - members of a cluster are close (i.e., similar) to each other
 - members of different clusters are dissimilar
- **Usually:**
 - Points are in a high-dimensional space
 - Similarity is defined using a distance measure
 - Euclidean, cosine, Jaccard, edit distance, ...



Characteristics of clustering methods

Quantitative: scalability (many samples), dimensionality (many features)

Qualitative: types of features (numerical, categorical, etc.), type of shapes (spheres, hyperplanes, manifolds, etc.)

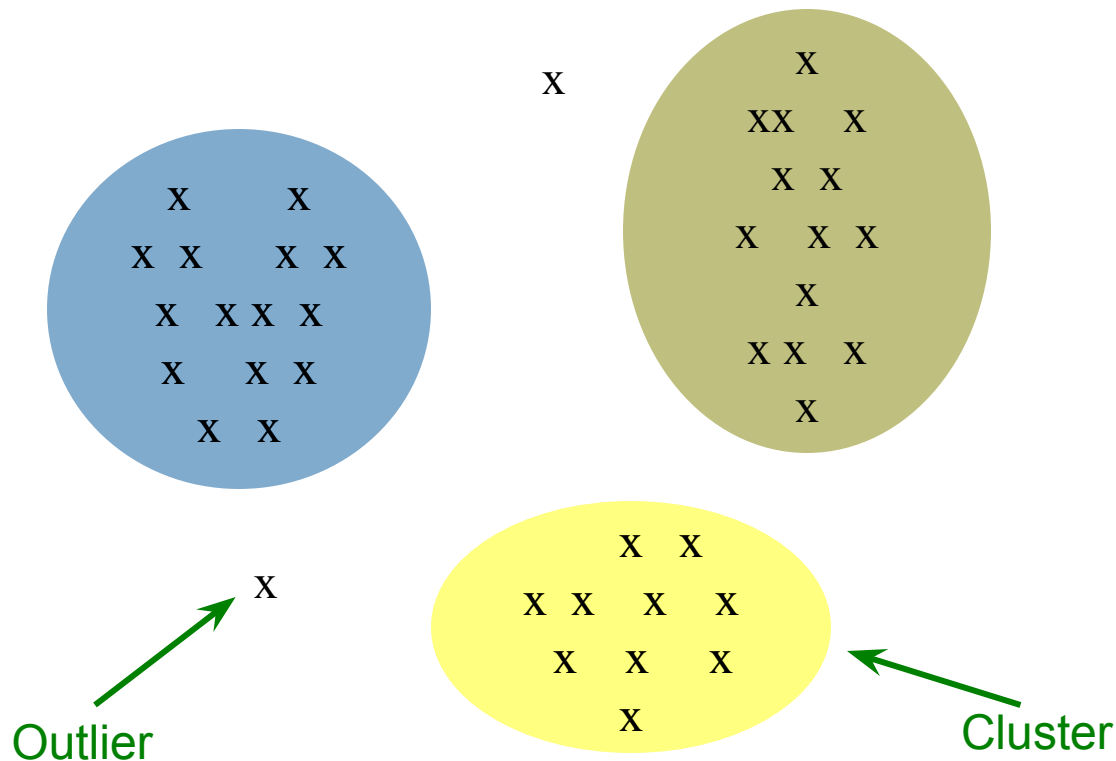


Characteristics of clustering methods

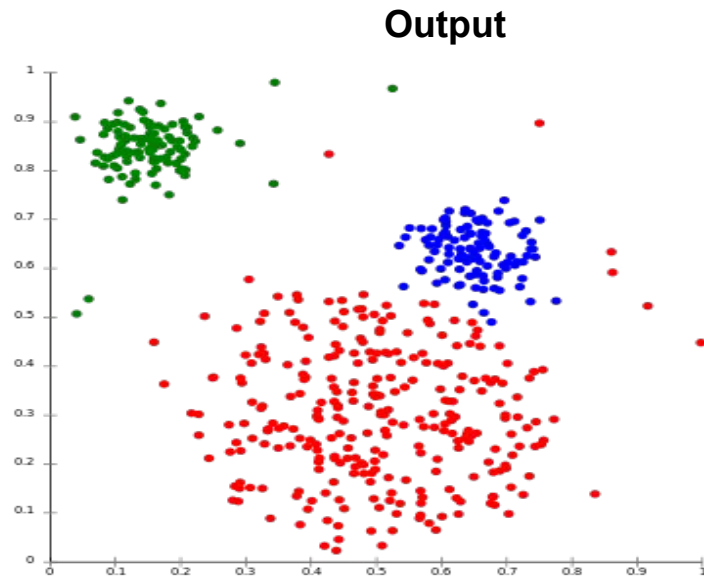
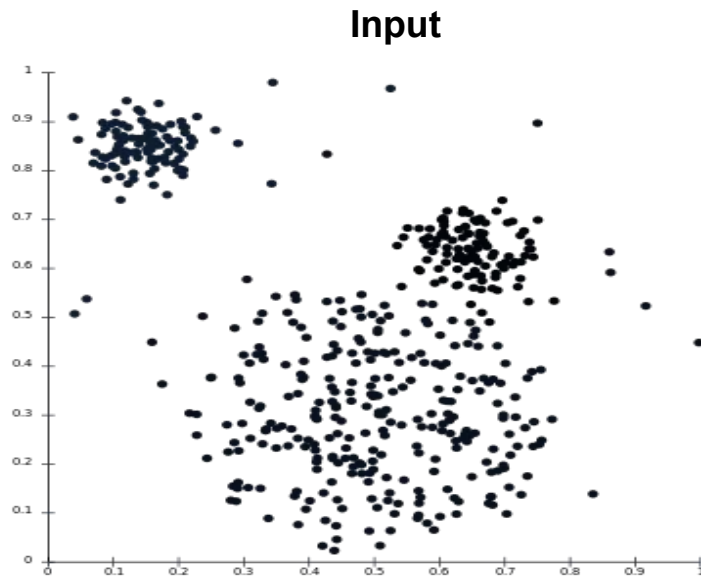
Robustness: sensitivity to noise and outliers, sensitivity to the processing order

User interaction: incorporation of user constraints (e.g., number of clusters, max size of clusters), interpretability and usability

Example: clusters & outliers

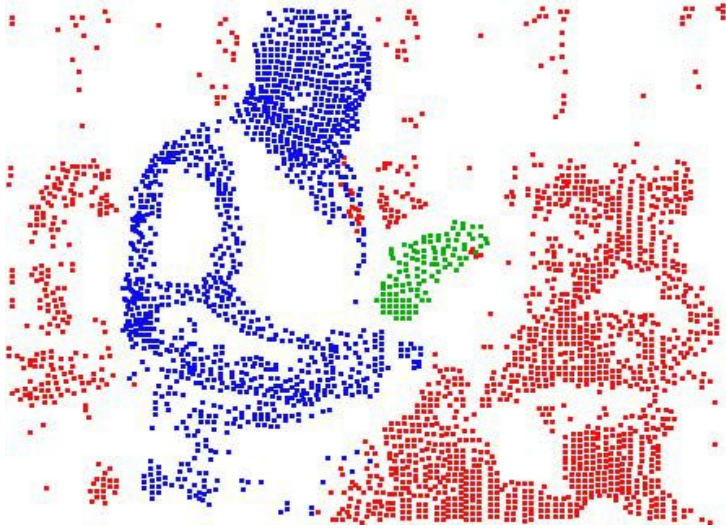


A typical clustering example



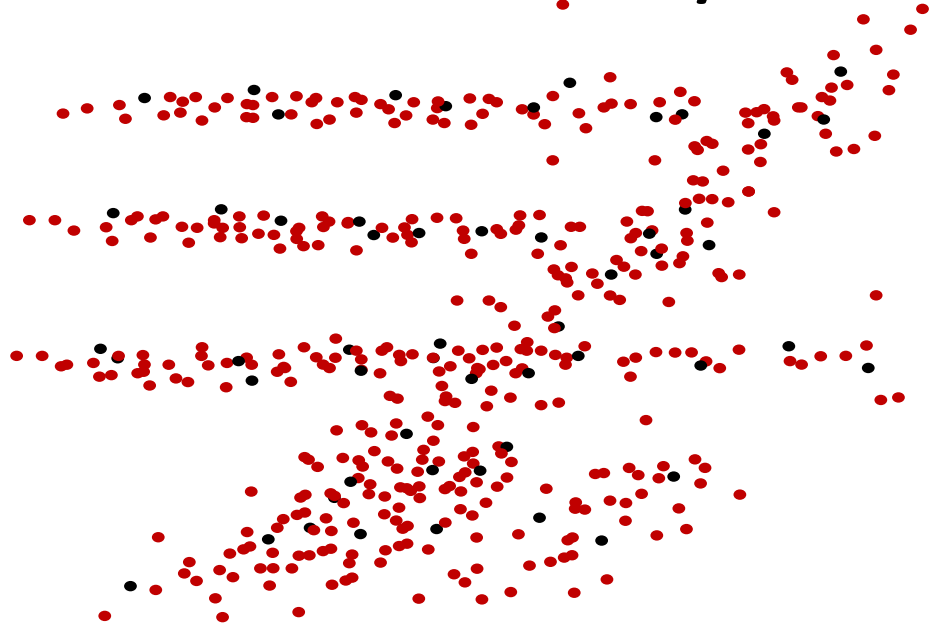
Note: Above is 2D; real scenarios often much more high-dimensional, e.g., 10,000-dimensional for 100x100 images.

Clustering for segmentation



Note: We break an image into regions of points with similar features (Brox and Malik, ECCV 2010).

Condensation/compression



Here we don't require that clusters extract meaningful structure, but that they give a coarse-grained version of the data.

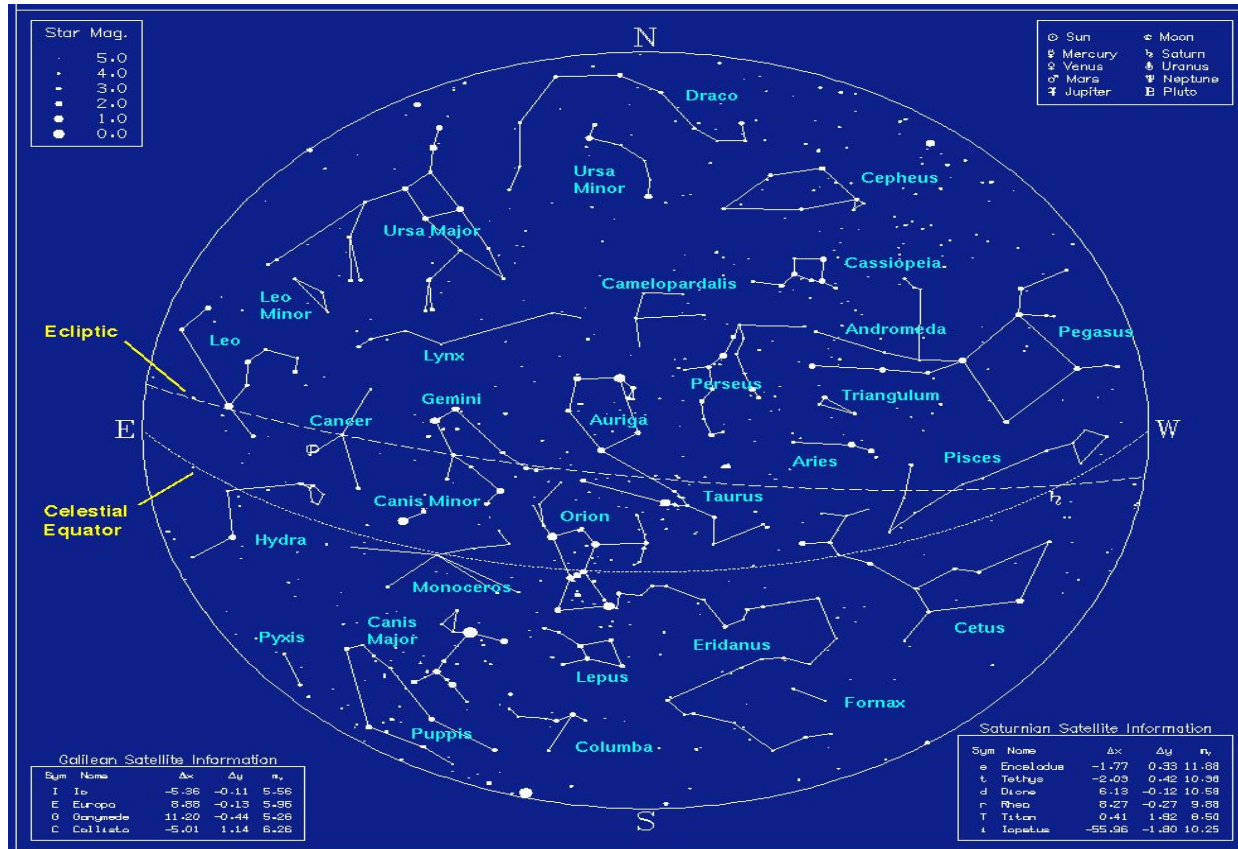
Beware of “cluster bias”!

- Human beings conceptualize the world through categories represented as *exemplars* (Rosch 73, Estes 94).



- We tend to see cluster structure whether it is there or not.
- Works well for dogs, but...

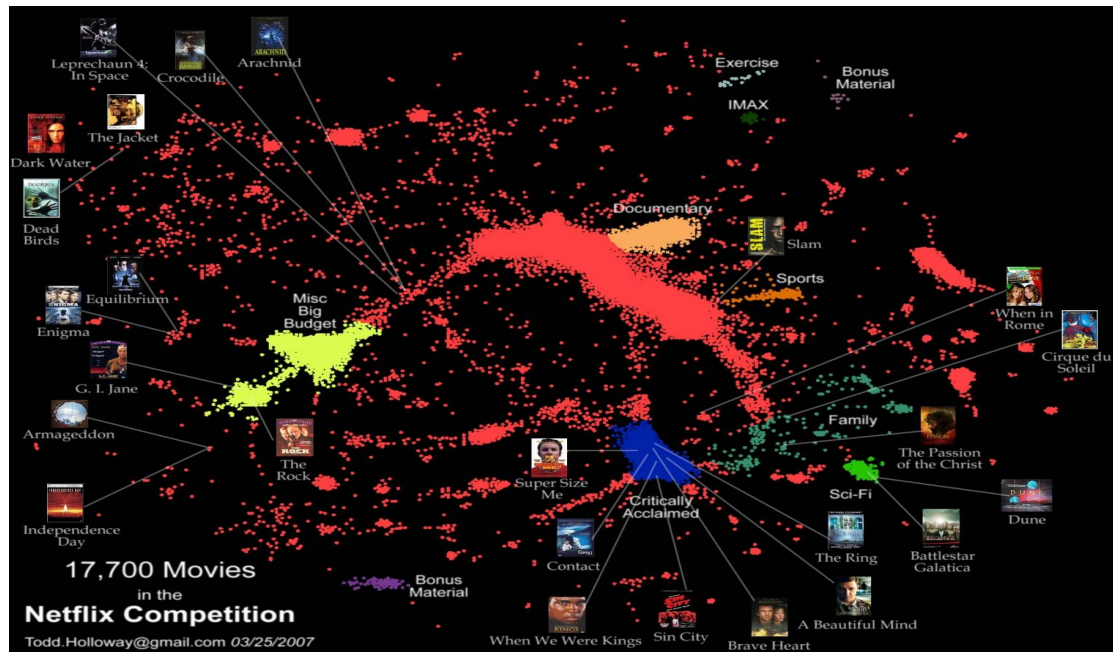
Cluster bias



“Cluster bias”

- **Clustering is used more than it should be**, because people assume an underlying domain has discrete classes in it.
- This is especially true for characteristics of people, e.g., Myers-Briggs personality types like “ENTP”.
- In reality the underlying data is usually **continuous**.
- Continuous models (matrix factorization, “soft” clustering, kNN) tend to do better (cf. next slide)

Netflix



- More of a continuum than discrete clusters
- Other methods (e.g., matrix factorization, kNN) may do better than discrete cluster models

Terminology

- **Hierarchical clustering:** clusters form a tree-shaped hierarchy. Can be computed bottom-up or top-down.
- **Flat clustering:** no inter-cluster structure.
- **Hard clustering:** items assigned to a unique cluster.
- **Soft clustering:** cluster membership is a probability distribution over all clusters

Clustering is a hard problem!



Why is it hard?

- Clustering in two dimensions looks easy
 - Clustering small amounts of data looks easy
 - And in most cases, looks are not deceiving
-
- Many applications involve not 2, but 10 or 10,000 dimensions
 - **High-dimensional spaces look different:** Almost all pairs of points are at about the same distance
("curse of dimensionality", cf. lecture 8)

Clustering problem: galaxies

- A catalog of 2 billion “sky objects” represents objects by their radiation in 7 dimensions (frequency bands)
- Problem: Cluster into similar objects, e.g., galaxies, nearby stars, quasars, etc.
- Sloan Digital Sky Survey [\[link\]](#)



Clustering problem: music CDs

- **Intuitively: Music divides into categories, and customers prefer a few categories**
 - But what are categories really?
- Represent a CD by a set of customers who bought it
- Similar CDs have similar sets of customers, and vice-versa

Clustering problem: music CDs

Space of all CDs:

- Think of a space with one dimension for each customer
 - Values in a dimension may be 0 or 1 only
 - A CD is a point in this space (x_1, x_2, \dots, x_k) ,
where $x_i = 1$ iff the i^{th} customer bought the CD
- For Amazon, the dimension is tens of millions
- **Task:** Find clusters of similar CDs

Clustering problem: documents

Finding topics:

- Represent a document by a vector (x_1, x_2, \dots, x_k) , where $x_i = 1$ iff the i^{th} word (in some order) appears in the document
- **Documents with similar sets of words may be about the same topic**

Cosine, Jaccard, and Euclidean

- In both examples (CDs, documents) we have a choice when we thinking of data points as sets of features (users, words):
 - **Sets as vectors:**
 - Measure similarity by **Euclidean distance**
 - Measure similarity by the **cosine distance**
 - **Sets as sets:** Measure similarity by the [Jaccard distance](#)

Overview: Methods of clustering

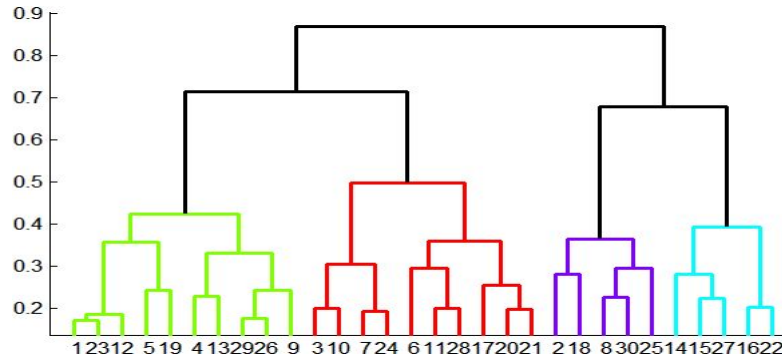
■ Hierarchical:

▪ **Agglomerative** (bottom up):

- Initially, each point is a cluster
- Repeatedly combine the two “nearest” clusters into one

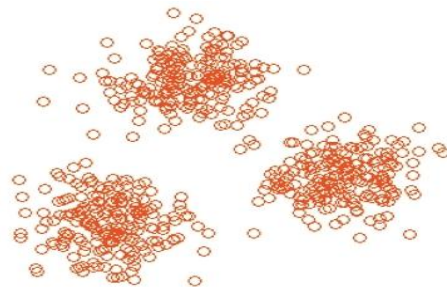
▪ **Divisive** (top down):

- Start with one cluster and recursively split it



■ Point assignment:

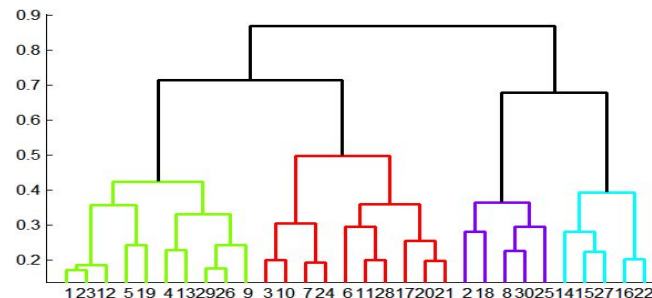
- Maintain a set of clusters
- Points belong to “nearest” cluster



Agglomerative hierarchical clustering

■ Key operation:

Repeatedly combine two nearest clusters



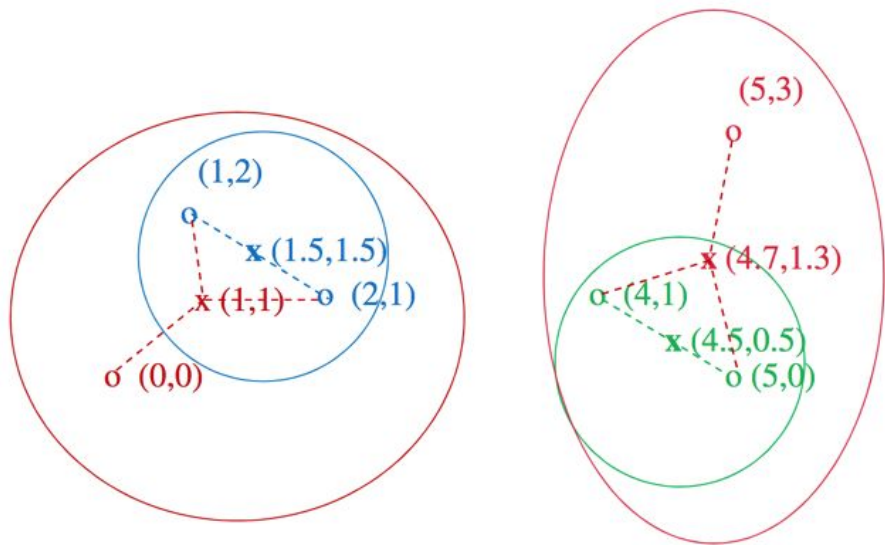
■ Three important questions:

- **1)** How do you represent a cluster of more than one point?
- **2)** How do you determine the “nearness” of clusters?
- **3)** When to stop combining clusters?

Agglomerative hierarchical clustering

- **Key operation: Repeatedly combine two nearest clusters**
- **(1) How to represent a cluster of many points?**
 - **Euclidean case:** each cluster has a **centroid** = average of its points
 - What about non-Euclidean case?
- **(2) How to determine “nearness” of clusters?**
 - Measure cluster distances by distances of centroids
 - What about non-Euclidean case?

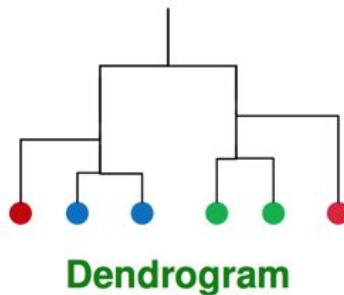
Example: Hierarchical clustering



Data:

σ ... data point

\bar{x} ... centroid



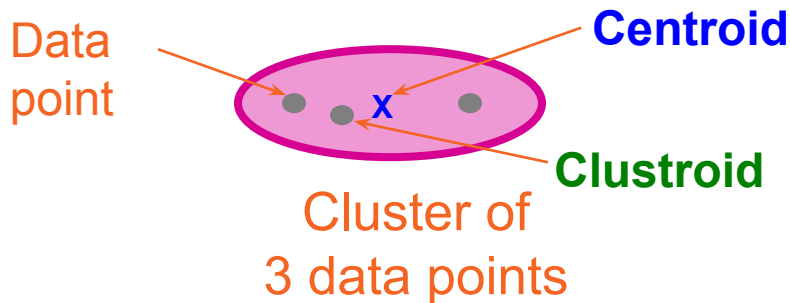
Non-Euclidean case: clustroids

- (1) How to represent a cluster of many points?

clustroid = point “closest” to other points

- Possible meanings of “closest”:

- Smallest maximum distance to other points
- Smallest average distance to other points (a.k.a. **medoid**)
- Smallest sum of squares of distances to other points
 - For distance metric d clustroid c of cluster C is:



Centroid is the avg. of all (data)points in the cluster. This means centroid is an “artificial” point.

Clustroid is an **existing** (data)point that is “closest” to all other points in the cluster.

Defining “nearness” of clusters

■ (2) How do you determine the “nearness” of clusters?

▪ Approach 1:

Intercluster distance = minimum of the distances between any two points, one from each cluster; or average of distances; or distance between centroids/clustroids; etc.

▪ Approach 2:

Pick a notion of “**cohesion**” (“tightness”) of clusters

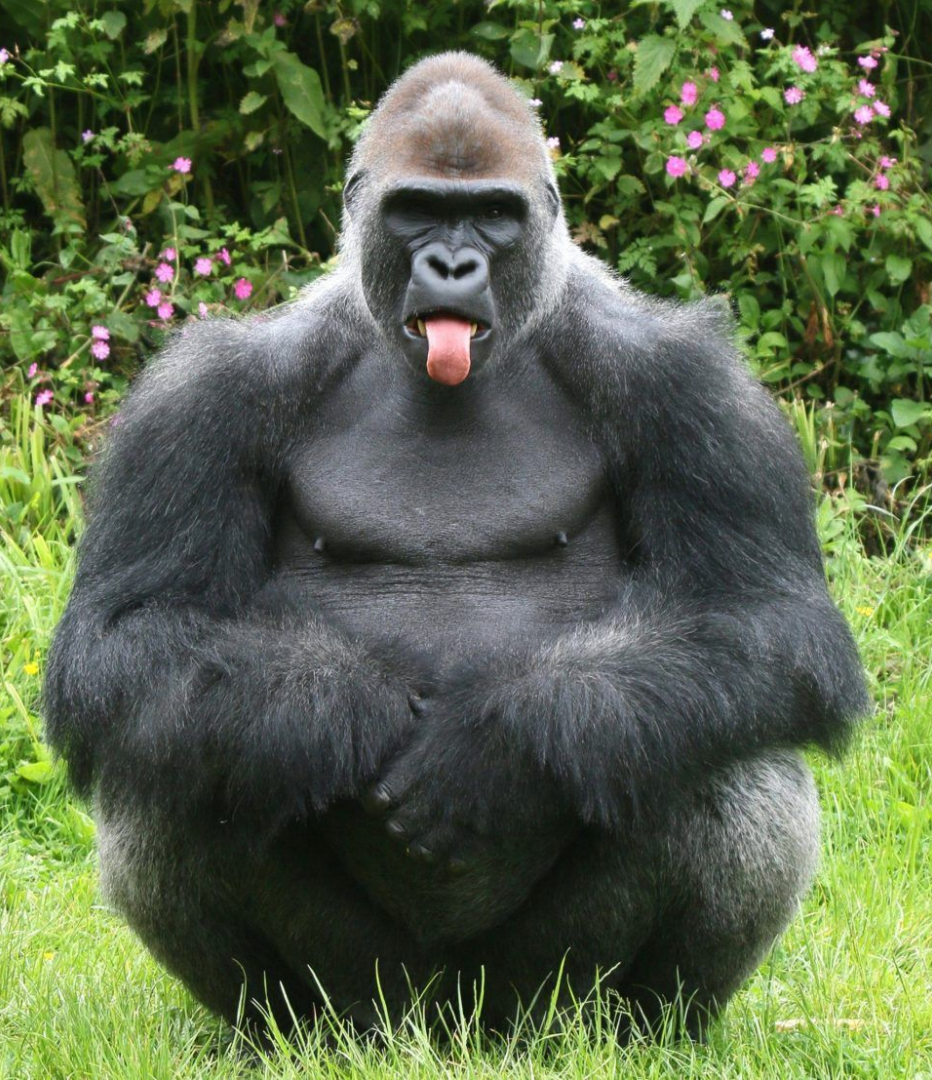
- Nearness of clusters = cohesion of their *union*

Cohesion

- **Approach 2.1:** Use the **diameter** of the merged cluster = maximum distance between points in the cluster
- **Approach 2.2:** Use the **average distance** between points in the cluster

Implementation

- **Naïve implementation of hierarchical clustering:**
 - At each step, compute pairwise distances between all pairs of clusters, then merge
 - $O(N^3)$
- Careful implementation using priority queue can reduce time to $O(N^2 \log N)$
 - **Still too expensive for really big datasets that do not fit in memory**



K-means

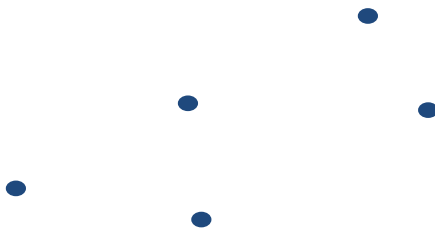
The gorilla among the point-assignment
clustering algorithms

K-means clustering

- A simple greedy algorithm (usually called Lloyd's algorithm)
- Locally minimizes the distance (usually Euclidean) from data points to their respective centroids:
 - **Find the closest cluster centroid** for each item, and assign it to that cluster.
 - **Recompute the cluster centroid** (the mean of items in the cluster) for each cluster.

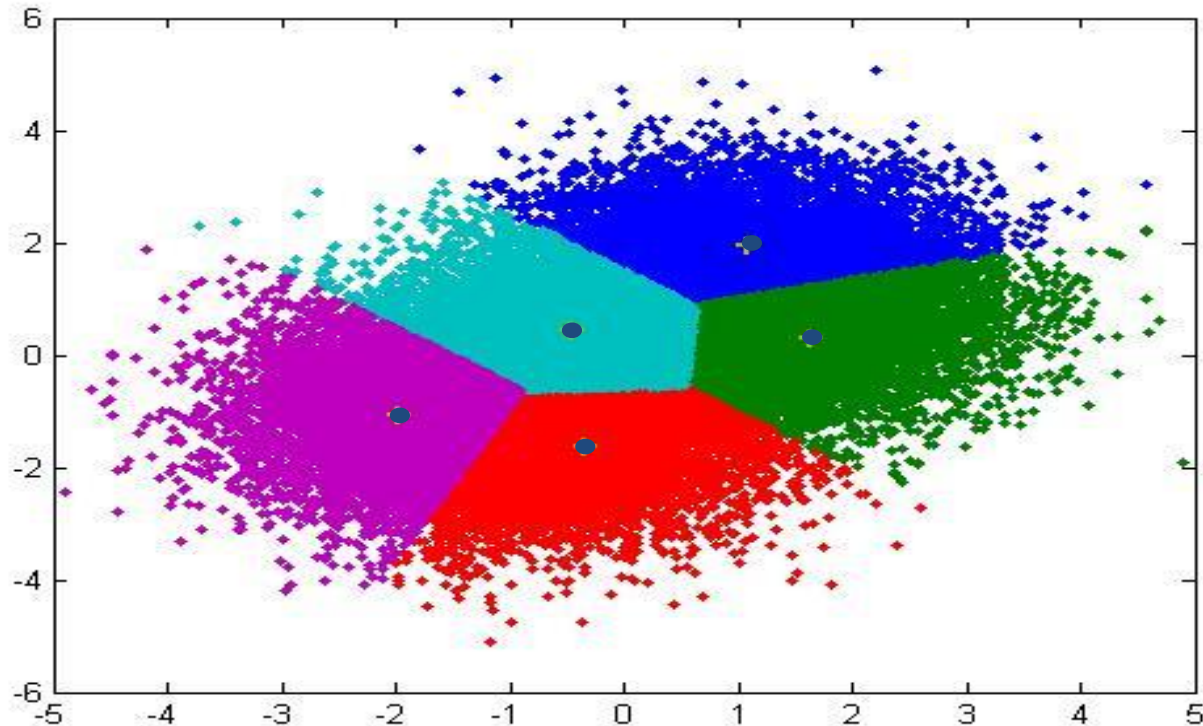
K-means clustering

Cluster centers – can pick by sampling the input data.



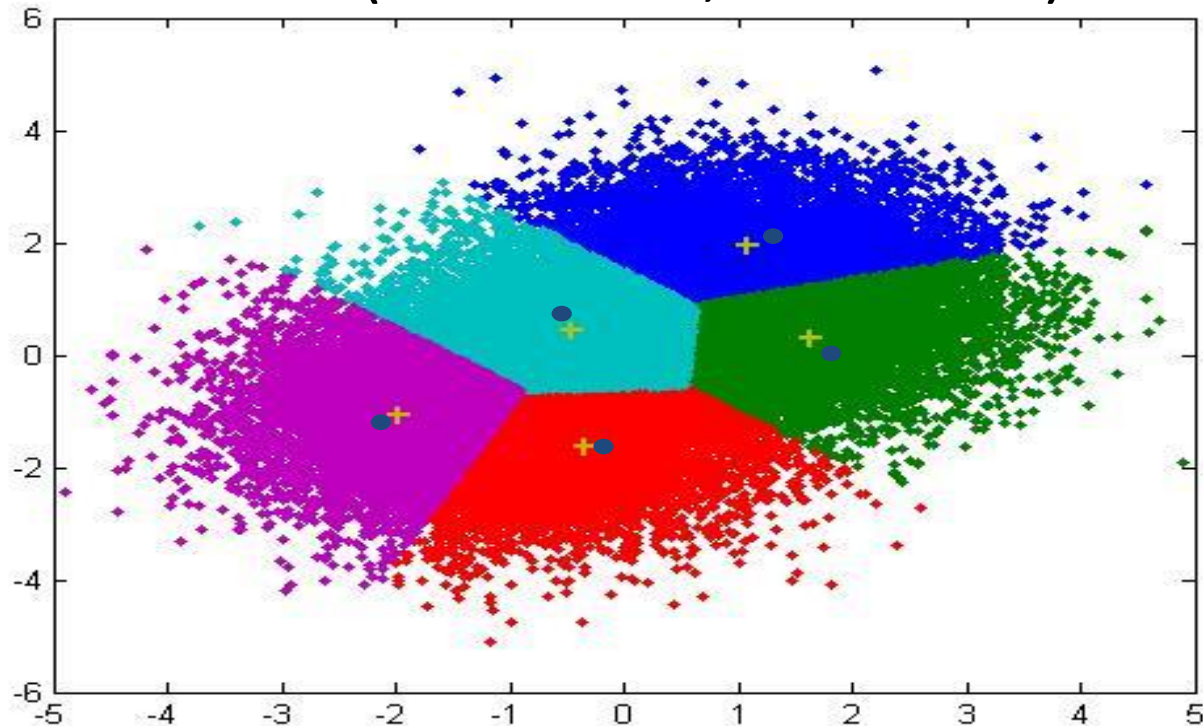
K-means clustering

Assign points to closest centroid



K-means clustering

Recompute centers (old = cross, new = dot)



K-means clustering

Iterate:

- For fixed number of iterations
- Until no change in assignments
- Until small change in cluster “tightness” (sum of distances from points to centroids)

K-means initialization

We need to pick some points for the first round of the algorithm:

- **Random sample:** Pick a random subset of k points from the dataset.
- **K-Means++:** Iteratively construct a random sample with good spacing across the dataset.

Note: Finding an optimal k-means clustering is NP-hard. The above help avoid bad configurations.

K-means++ [\[link\]](#)

Start: Choose first cluster center at random from the data points

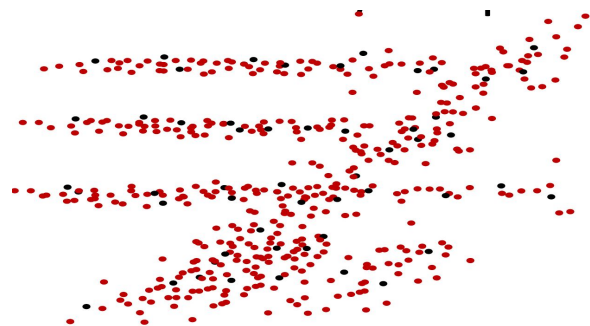
Iterate:

- For every remaining data point x , compute the distance $D(x)$ from x to the closest cluster center.
- Choose a remaining point x randomly with probability proportional to $D(x)^2$, and make it a new cluster center.

Intuitively, this finds a sample of widely-spaced points avoiding “collapsing” of the clustering into a few internal centers.

K-means properties

- It's a greedy algorithm with random setup – **solution isn't optimal** and varies significantly with different initial points.
- Very simple convergence proofs.
- **Performance is $O(nk)$ per iteration**, not bad and can be heuristically improved.
n = total features in the dataset, k = number clusters
- Many variants, e.g.
 - Fixed-size clusters
 - Soft clustering
- Works well for data condensation/compression.



K-means drawbacks

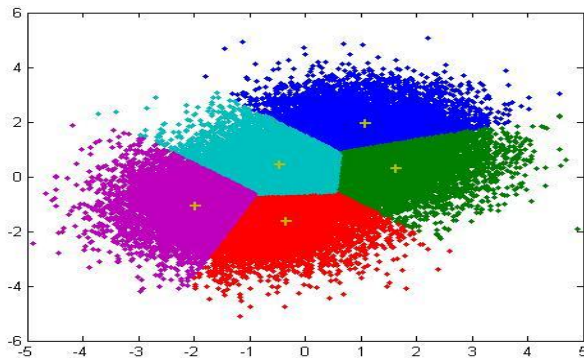
Often terminates at a **local optimum** (mitigated by smart initialization, e.g., k-means++)

Need a notion of **mean**

Need to specify **k** (number of clusters) in advance

Doesn't handle **noisy data and outliers** well

Clusters **only** have **convex shapes**



How to choose k?

Execute for $k = 1, 2, 3 \dots$

$b(i)$: distance to
closest
other centroid

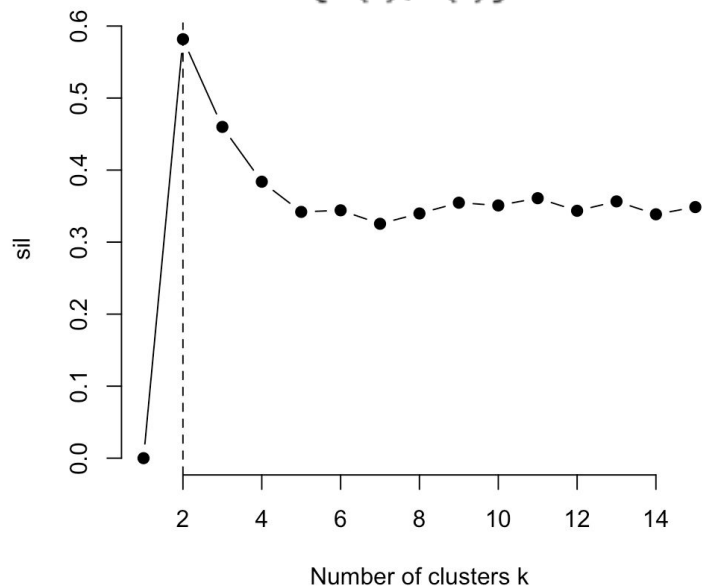
$a(i)$: distance to
own centroid

For each data point i , compute “silhouette” $s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$


S = average of $s(i)$ over all i

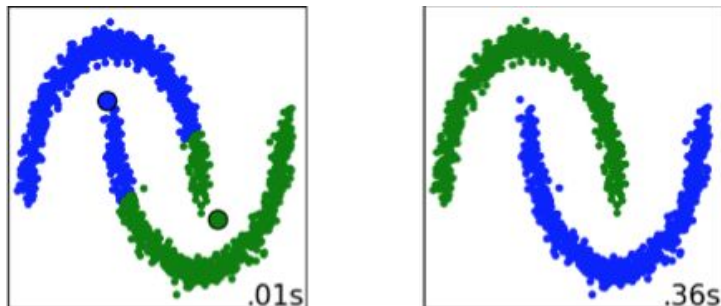
Plot S against k

Pick k for which S is greatest

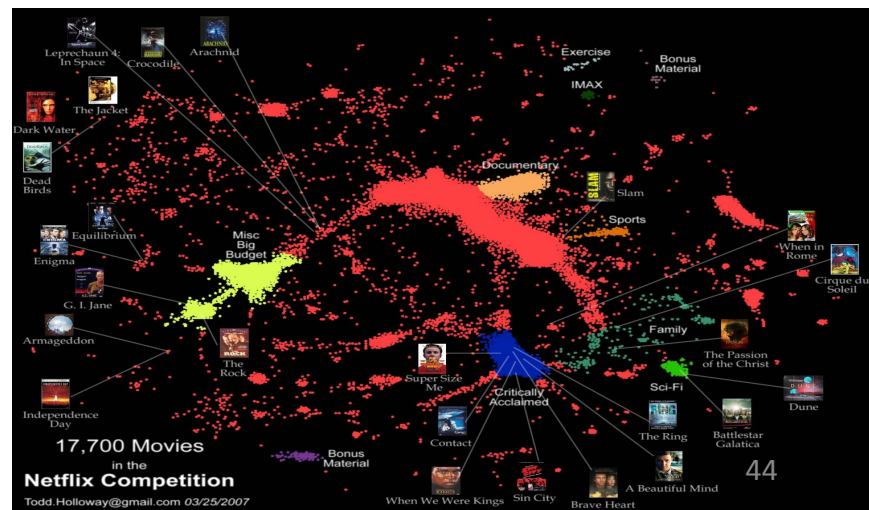


DBSCAN

- “**Density-based spatial clustering of applications with noise**”
 - Motivation: Centroid-based clustering methods like k-means favor clusters that are spherical, and have great difficulty with anything else
- 

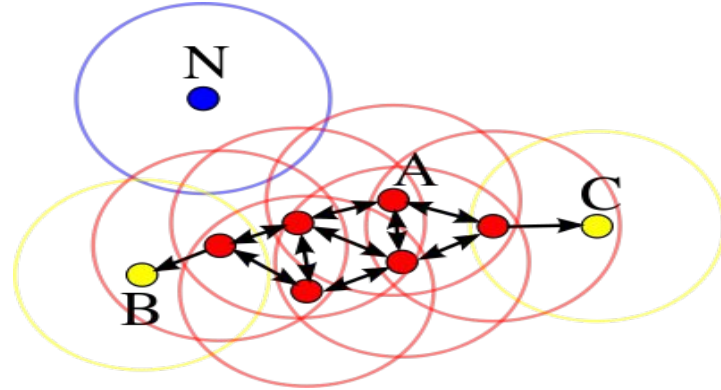


- But with real data we have:



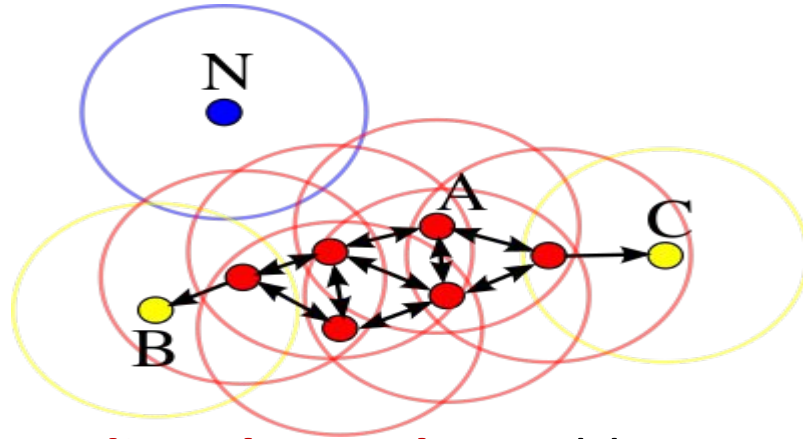
DBSCAN

- DBSCAN performs density-based clustering, and follows the shape of dense neighborhoods of points.



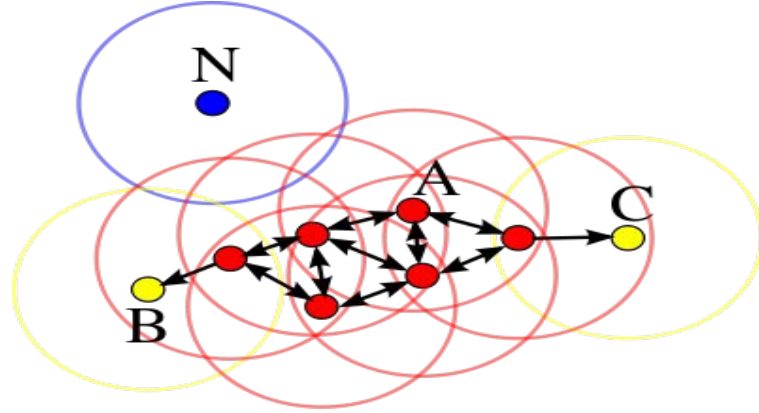
- Core points** have at least minPts neighbors in a sphere of diameter ϵ around them.
- The red points here are core points with at least $\text{minPts} = 3$ neighbors in an ϵ -sphere around them.

DBSCAN



- **Core points** can **directly reach** neighbors in their ϵ -sphere
- From non-core points, no other points can be reached (by def.)
- Point q is **density-reachable** from p if there is a series of points $p = p_1, p_2, \dots, p_n = q$ such that p_{i+1} is directly reachable from p_i
- All points not density-reachable from any other points are **outliers**

DBSCAN clusters



- Points p, q are **density-connected** if there is a point o such that both p and q are density-reachable from o .
- A **cluster** is a set of points which are **mutually density-connected**.
- That is, if a point is density-reachable from a cluster point, it is part of the cluster as well.
- In the above figure, red points in A are mutually density-reachable; B and C are density-connected; N is an outlier.

DBSCAN algorithm

```
DBSCAN(DB, dist, eps, minPts) {
```

```
    C = 0
```

```
    for each point P in database DB {
```

```
        if label(P) ≠ undefined then continue
```

```
        Neighbors N = RangeQuery(DB, dist, P, eps)
```

```
        if |N| < minPts then {
```

```
            label(P) = Noise
```

```
            continue
```

```
        }
```

```
        C = C + 1
```

```
        label(P) = C
```

```
        Seed set S = N \ {P}
```

```
        for each point Q in S {
```

```
            if label(Q) = Noise then label(Q) = C
```

```
            if label(Q) ≠ undefined then continue
```

```
            label(Q) = C
```

```
            Neighbors N = RangeQuery(DB, dist, Q, eps)
```

```
            if |N| ≥ minPts then {
```

```
                S = S ∪ N
```

```
            }
```

```
        }
```

```
    }
```

```
}
```

```
/* Cluster counter */
```

```
/* Previously processed in inner loop */
```

```
/* Find neighbors */
```

```
/* Density check */
```

```
/* Label as Noise */
```

```
/* next cluster label */
```

```
/* Label initial point */
```

```
/* Neighbors to expand */
```

```
/* Process every seed point */
```

```
/* Change Noise to border point */
```

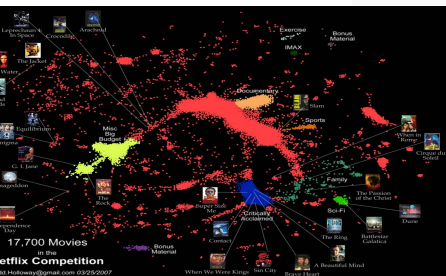
```
/* Previously processed */
```

```
/* Label neighbor */
```

```
/* Find neighbors */
```

```
/* Density check */
```

```
/* Add new neighbors to seed set */
```



DBSCAN performance

- DBSCAN uses all-pairs point distances, but using an efficient indexing structure, each RangeQuery takes $O(\log n)$ time
- The algorithm overall can be made to run in **$O(n \log n)$**
- Fast neighbor search becomes progressively harder (higher constants) in higher dimensions

Give us feedback on this lecture here:

<https://go.epfl.ch/ada2018-lec10-feedback>

- What did you (not) like about this lecture?
- What was (not) well explained?
- On what would you like more (fewer) details?
- ...