

Introduction to Information Retrieval
<http://informationretrieval.org>

Flat Clustering: k-means

Slides borrowed from Hinrich Schütze with modifications

Outline

- What is clustering?
- Applications of clustering in information retrieval
- K -means algorithm
- Evaluation of clustering
- How many clusters?

Outline

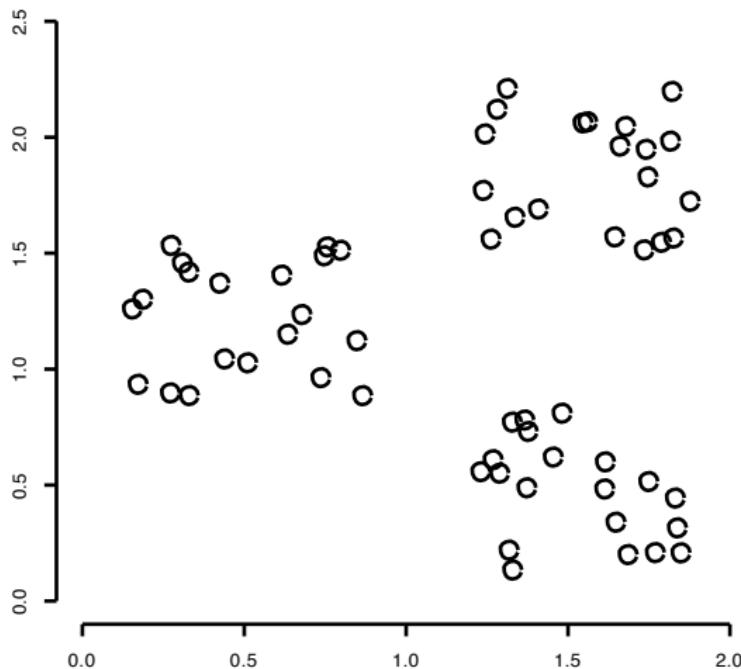
- What is clustering?
-
-
-
-



Clustering: Definition

- (Document) clustering is the process of grouping a set of documents into clusters of similar documents.
- Documents within a cluster should be similar.
- Documents from different clusters should be dissimilar.
- Clustering is the most common form of unsupervised learning.
- Unsupervised = there are no labeled or annotated data.

Data set with clear cluster structure



Propose
algorithm
for finding
the cluster
structure in
this
example

Classification vs. Clustering

- Classification: supervised learning
- Clustering: unsupervised learning
- Classification: Classes are **human-defined** and part of the input to the learning algorithm.
- Clustering: Clusters are **inferred from the data** without human input.
 - However, there are many ways of influencing the outcome of clustering: number of clusters, similarity measure, representation of documents, ...

Outline

- What is clustering?
- Applications of clustering in information retrieval
-
-
-

The cluster hypothesis

Cluster hypothesis. Documents in the same cluster behave similarly with respect to relevance to information needs. All applications of clustering in IR are based (directly or indirectly) on the cluster hypothesis. Van Ribsbergen's original wording (1979): "closely associated documents tend to be relevant to the same requests".

Applications of clustering in IR

application	what is clustered?	benefit
search result clustering	search results	more effective information presentation to user
Scatter-Gather	(subsets of) collection	alternative user interface: “search without typing”
collection clustering	collection	effective information presentation for exploratory browsing
cluster-based retrieval	collection	higher efficiency: faster search

Search result clustering for better navigation

Vivísmo® the Web Advanced Search Help

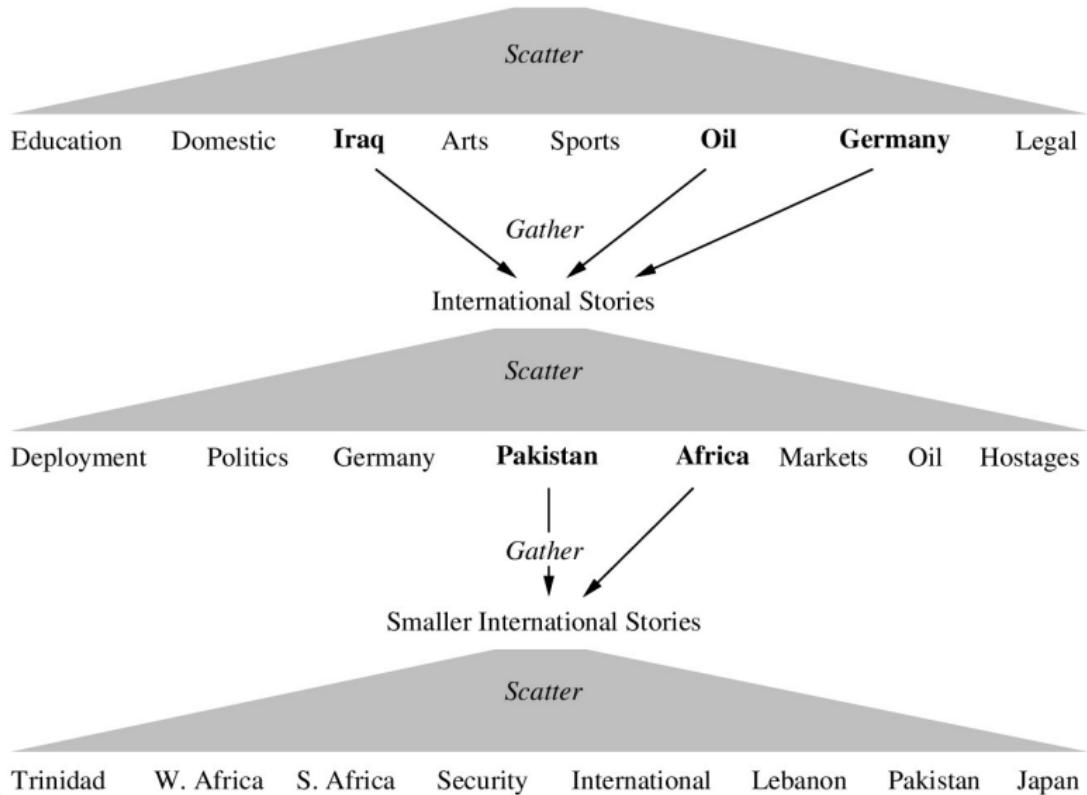
Clustered Results Top 208 results of at least 20,373,974 retrieved for the query **jaguar** ([Details](#))

- ▶ [jaguar \(208\)](#)
 - + ▶ [Cars \(74\)](#)
 - + ▶ [Club \(34\)](#)
 - + ▶ [Cat \(23\)](#)
 - + ▶ [Animal \(13\)](#)
 - + ▶ [Restoration \(10\)](#)
 - + ▶ [Mac OS X \(8\)](#)
 - + ▶ [Jaguar Model \(8\)](#)
 - + ▶ [Request \(5\)](#)
 - + ▶ [Mark Webber \(6\)](#)
 - + ▶ [Maya \(5\)](#)
- ▼ [More](#)

Find in clusters:

1. [Jag-lovers - THE source for all Jaguar information](#) [new window] [frame] [cache] [preview] [clusters]
... Internet! Serving Enthusiasts since 1993 The Jag-lovers Web Currently with 40661 members The Premier **Jaguar** Cars web resource for all enthusiasts Lists and Forums Jag-lovers originally evolved around its ...
www.jag-lovers.org - Open Directory 2, Wisenut 8, Ask Jeeves 8, MSN 9, Looksmart 12, MSN Search 18
2. [Jaguar Cars](#) [new window] [frame] [cache] [preview] [clusters]
[...] redirected to www.jaguar.com
www.jaguarcars.com - Looksmart 1, MSN 2, Lycos 3, Wisenut 6, MSN Search 9, MSN 29
3. <http://www.jaguar.com/> [new window] [frame] [preview] [clusters]
www.jaguar.com - MSN 1, Ask Jeeves 1, MSN Search 3, Lycos 9
4. [Apple - Mac OS X](#) [new window] [frame] [preview] [clusters]
Learn about the new OS X Server, designed for the Internet, digital media and workgroup management
Download a technical factsheet.
www.apple.com/macosx - Wisenut 1, MSN 3, Looksmart 26

Scatter-Gather



Global clustering for navigation: Google News

≡ **Google News** Search for topics, locations & sources

Top stories

- For you
- Favorites
- Saved searches

- U.S.
- World
- Local
- Business
- Technology
- Entertainment
- Sports
- Science
- Health

Language & region
English | United States

Headlines

[More Headlines](#)

Rapper Nipsey Hussle killed in shooting outside his L.A. store
NBCNews.com • 2 hours ago

- **Nipsey Hussle Dead at 33 After Getting Shot in Los Angeles**
TMZ • 1 hour ago
- **Rapper Nipsey Hussle Dead After Being Shot Outside Of His LA Clothing Store**
CBS Los Angeles • 2 hours ago
- **Rapper Nipsey Hussle killed in South L.A. shooting; 2 others wounded**
Los Angeles Times • 15 minutes ago • Local coverage
- **Lakers' LeBron James reacts to Nipsey Hussle dying, shares emotional message**
ClutchPoints • 41 minutes ago

[View full coverage](#)

'It's so easy to Google "Creepy Biden"
POLITICO • 1 hour ago

- **Lucy Flores speaks out on her accusation on Joe Biden**
Cengiz Adabag News • 10 hours ago

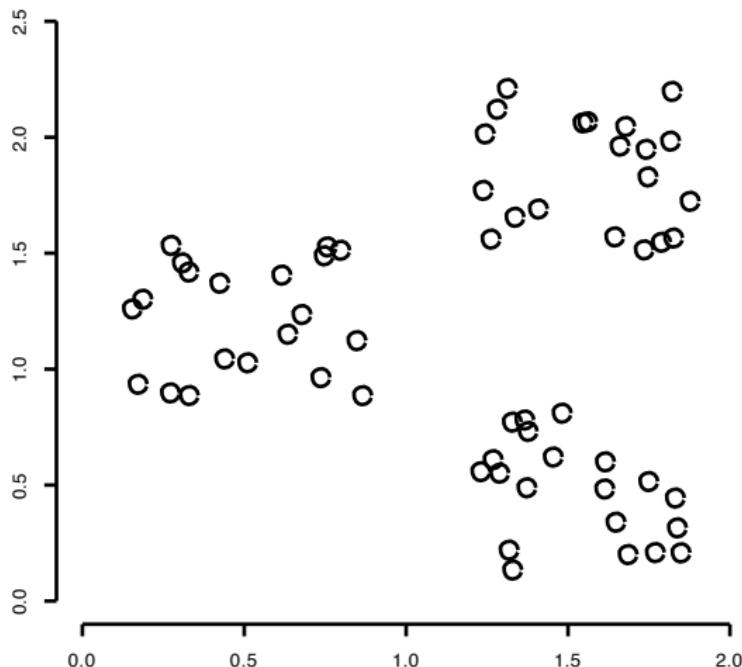
[View full coverage](#)

Trump White House presses threat to close U.S.-Mexico border this
[Image]

Clustering for improving recall

- To improve search recall:
 - Cluster docs in collection a priori
 - When a query matches a doc d , also return other docs in the cluster containing d
- Hope: if we do this: the query “car” will also return docs containing “automobile”
 - Because the clustering algorithm groups together docs containing “car” with those containing “automobile”.
 - Both types of documents contain words like “parts”, “dealer”, “mercedes”, “road trip”.

Data set with clear cluster structure



Propose
algorithm
for finding
the cluster
structure in
this
example

Desiderata for clustering

- General goal: put related docs in the same cluster, put unrelated docs in different clusters.
 - We'll see different ways of formalizing this.
- The number of clusters should be appropriate for the data set we are clustering.
 - Initially, we will assume the number of clusters K is given.
 - Later: Semiautomatic methods for determining K
- Secondary goals in clustering
 - Avoid very small and very large clusters
 - Define clusters that are easy to explain to the user
 - Many others ...

Flat vs. Hierarchical clustering

- Flat algorithms
 - Usually start with a random (partial) partitioning of docs into groups
 - Refine iteratively
 - Main algorithm: *K*-means
- Hierarchical algorithms
 - Create a hierarchy
 - Bottom-up, agglomerative
 - Top-down, divisive



Hard vs. Soft clustering

- Hard clustering: Each document belongs to exactly one cluster.
 - More common and easier to do
- Soft clustering: A document can belong to more than one cluster.
 - Makes more sense for applications like creating browsable hierarchies
 - You may want to put sneakers in two clusters:
 - sports apparel
 - shoes
 - You can only do that with a soft clustering approach.
- We will first cover: flat, hard clustering
- We will then cover: hierarchical, hard clustering
-



Outline

- What is clustering?
- Applications of clustering in information retrieval
- *K*-means algorithm
-
-

Flat algorithms

- Flat algorithms compute a partition of N documents into a set of K clusters.
- Given: a set of documents and the number K
- Find: a partition into K clusters that optimizes the chosen partitioning criterion
- Global optimization: exhaustively enumerate partitions, pick optimal one
 - Not tractable
- Effective heuristic method: K -means algorithm

K-means

- Perhaps the best known clustering algorithm
- Simple, works well in many cases
- Use as default / baseline for clustering documents

Document representations in clustering

- Vector space model
- We measure the distance between vectors by [Euclidean distance](#).

K-means: Basic idea

- Each cluster in K -means is defined by a centroid.
- Objective/ partitioning criterion: minimize the average squared difference from the centroid
- Recall definition of centroid:

$$\underline{\mu}(\omega) = \frac{1}{|\omega|} \sum_{x \in \omega} \underline{x}$$

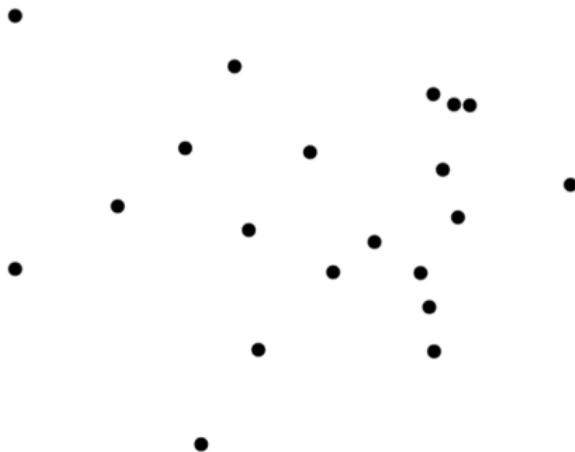
where we use ω to denote a cluster.

- We try to find the minimum average squared difference by iterating two steps:
 - reassignment: assign each vector to its closest centroid
 - recomputation: recompute each centroid as the average of the vectors that were assigned to it in reassignment

K -means pseudocode (μ_k is centroid of ω_k)

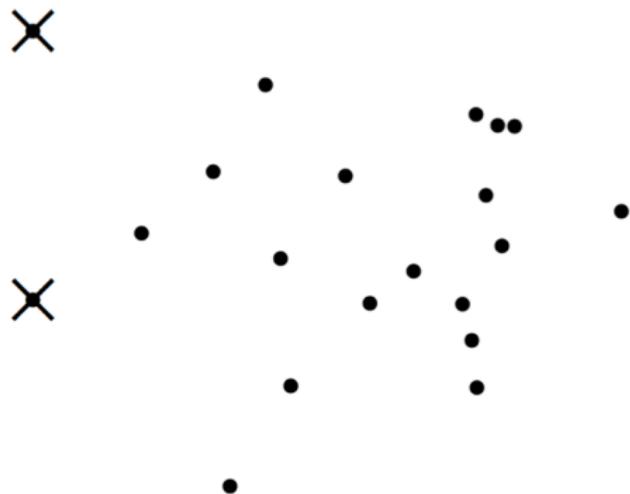
```
K-means( $\{\underline{x}_1, \dots, \underline{x}_N\}, K$ )
1   $(\underline{s}_1, \underline{s}_2, \dots, \underline{s}_K) \leftarrow \text{SelectRandomSeeds}(\{\underline{x}_1, \dots, \underline{x}_N\}, K)$ 
2  for  $k \leftarrow 1$  to  $K$ 
3    do  $\underline{\mu}_k \leftarrow \underline{s}_k$ 
4  while stopping criterion has not been met
5  do for  $k \leftarrow 1$  to  $K$ 
6    do  $\omega_k \leftarrow \{\}$ 
7    for  $n \leftarrow 1$  to  $N$ 
8      do  $j \leftarrow \arg \min_j' |\underline{\mu}_j - \underline{x}_n|$ 
9       $\omega_j \leftarrow \omega_j \cup \{\underline{x}_n\}$  (reassignment of vectors)
10   for  $k \leftarrow 1$  to  $K$ 
11   do  $\underline{\mu}_k \leftarrow \frac{1}{|\omega_k|} \sum_{x \in \omega_k} \underline{x}$  (recomputation of centroids)
12 return  $\{\underline{\mu}_1, \dots, \underline{\mu}_K\}$ 
```

Worked Example : Set of points to be clustered

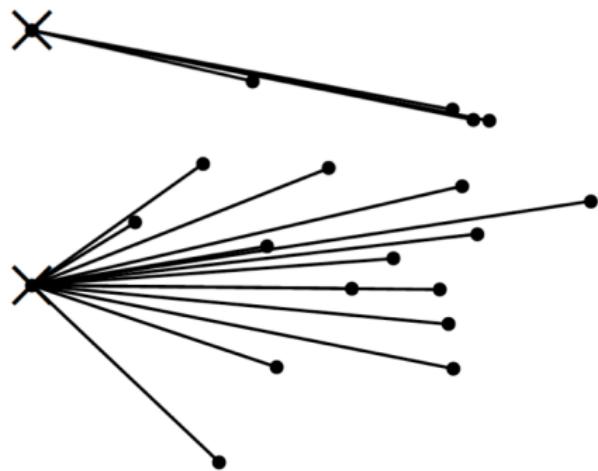


Exercise: (i) Guess what the optimal clustering into two clusters is in this case; (ii) compute the centroids of the clusters

Worked Example: Random selection of initial centroids



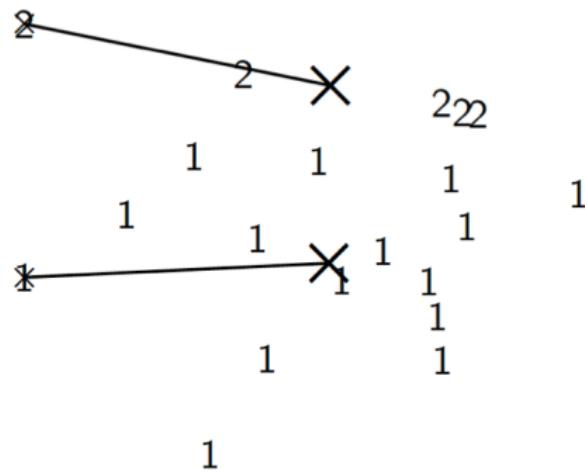
Worked Example: Assign points to closest center



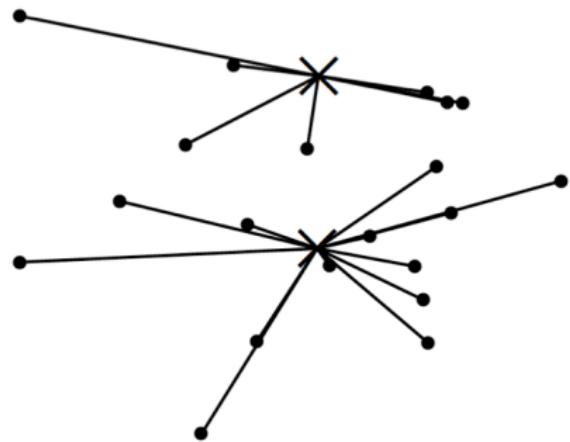
Worked Example: Assignment

$$\begin{array}{r} \times \\ & 2 \\ & 1 & 1 & 1 & 1 & 1 \\ \times & 1 & 1 & 1 & 1 & 1 \\ & & 1 & 1 & 1 & 1 \\ & & & 1 & 1 & 1 \\ & & & & 1 & 1 \\ & & & & & 1 \end{array}$$

Worked Example: Recompute cluster centroids



Worked Example: Assign points to closest centroid

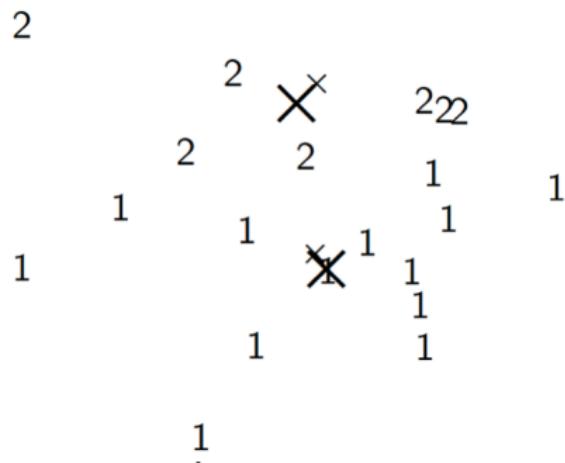


Worked Example: Assignment

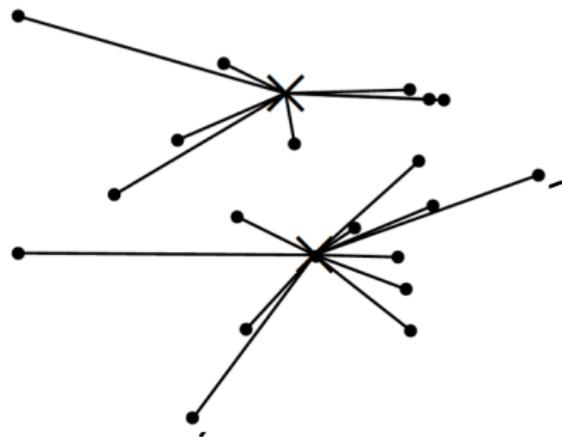
$$\begin{array}{r} 2 \\ \times \quad \quad \quad 22 \\ \hline \begin{array}{r} 2 \\ 2 \\ \times \quad \quad \quad 1 \\ 1 \end{array} \quad \begin{array}{r} 1 \\ 1 \\ 1 \\ 1 \end{array} \end{array}$$

1

Worked Example: Recompute cluster centroids



Worked Example: Assign points to closest centroid

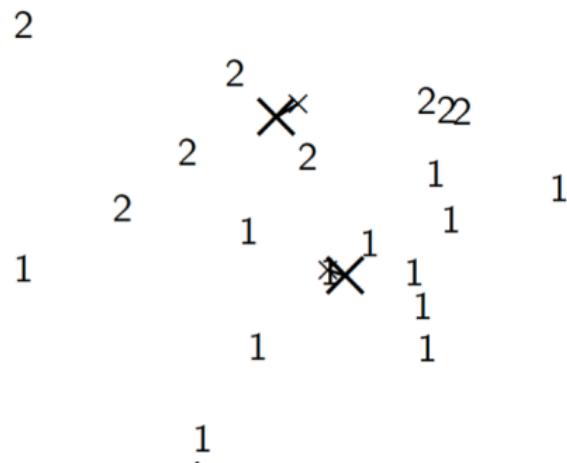


Worked Example: Assignment

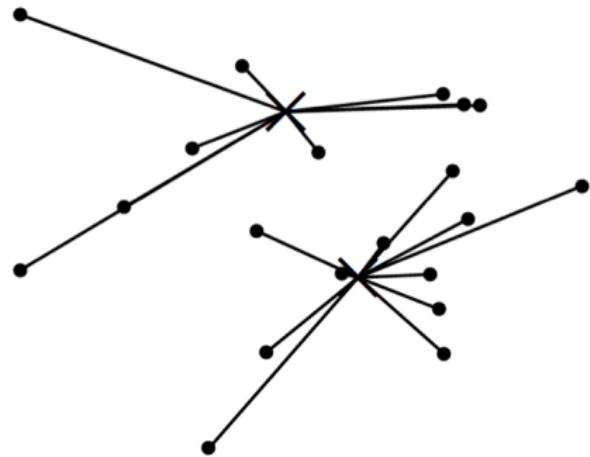
$$\begin{array}{r} 2 \\ \times 22 \\ \hline 2 & 2 \\ & 2 & 1 \\ 1 & & 1 & 1 \\ & 1 & 1 & 1 \\ \hline & & & 1 \end{array}$$

1

Worked Example: Recompute cluster centroids



Worked Example: Assign points to closest centroid

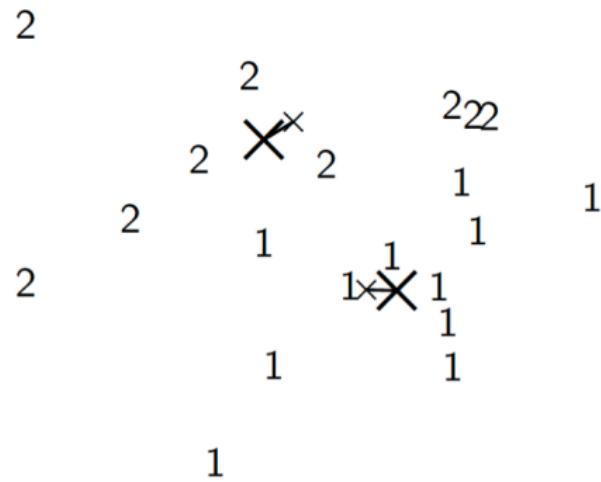


Worked Example: Assignment

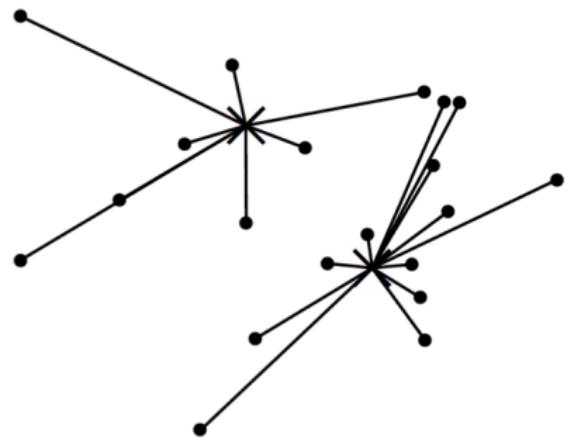
$$\begin{array}{r} 2 \\ \times 2 \\ \hline 2 & 2 \\ 2 & 1 \\ \hline 2 & 1 \\ 1 & 1 \\ \hline 1 & 1 \end{array}$$

1

Worked Example: Recompute cluster centroids



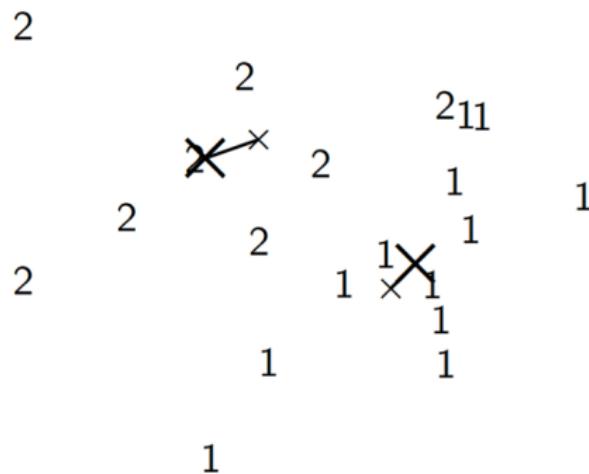
Worked Example: Assign points to closest centroid



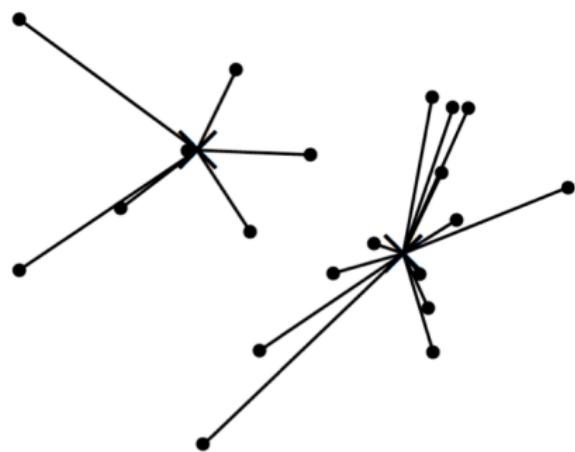
Worked Example: Assignment

$$\begin{array}{r} 2 \\ \times 2 \\ \hline 2 & 2 \\ & 2 \\ \hline 2 & 1 \\ & 1 \\ \hline 1 & 1 \\ & 1 \\ \hline 1 & 1 \\ & 1 \\ \hline 211 & 1 \\ & 1 \\ \hline 211 & 1 \end{array}$$

Worked Example: Recompute cluster centroids



Worked Example: Assign points to closest centroid

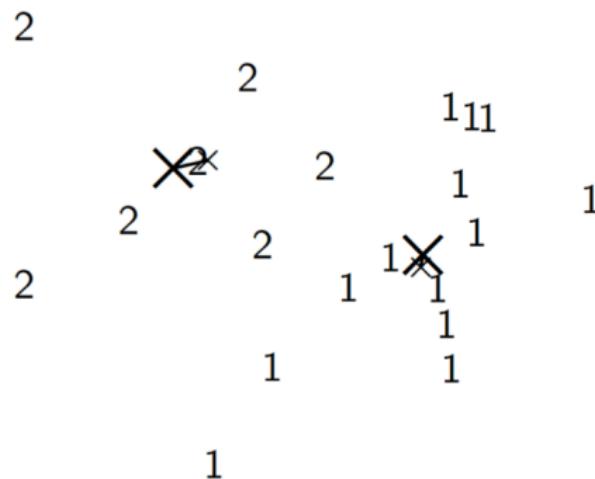


Worked Example: Assignment

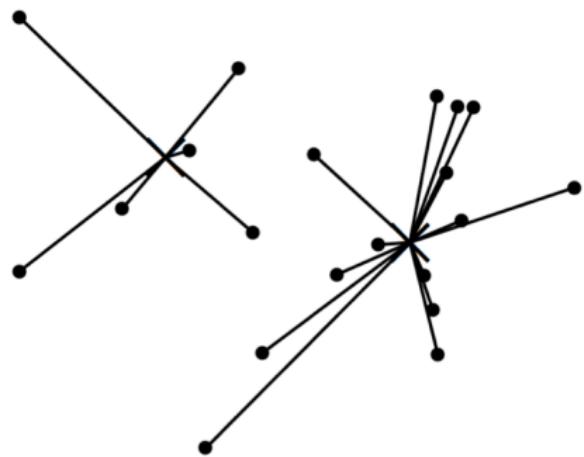
$$\begin{array}{cccccc} & & 2 & & & \\ & & \times & & 1_{11} & \\ & 2 & 2 & 1 & 1 & \\ 2 & & 2 & 1 & 1 & \\ & 2 & 1 & 1 & 1 & \\ & & 1 & 1 & 1 & \\ & & & 1 & & \\ & & & & 1 & \\ & & & & & 1 \end{array}$$

1

Worked Example: Recompute cluster centroids



Worked Example: Assign points to closest centroid



Worked Example: Assignment

2

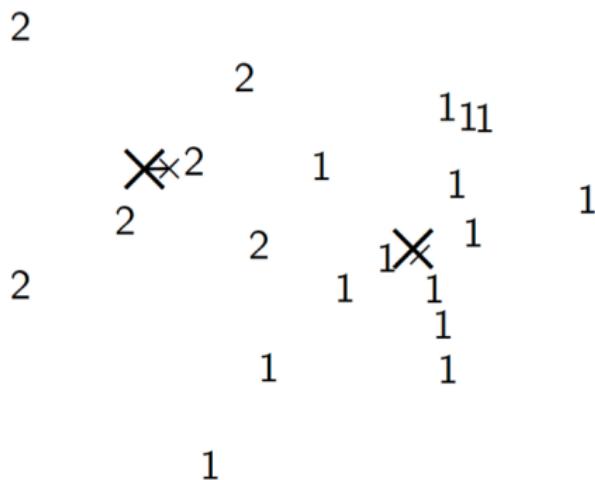
2 2 1 1 1¹¹¹

2 2 1 1 1
~~2~~ 1 1 1 1

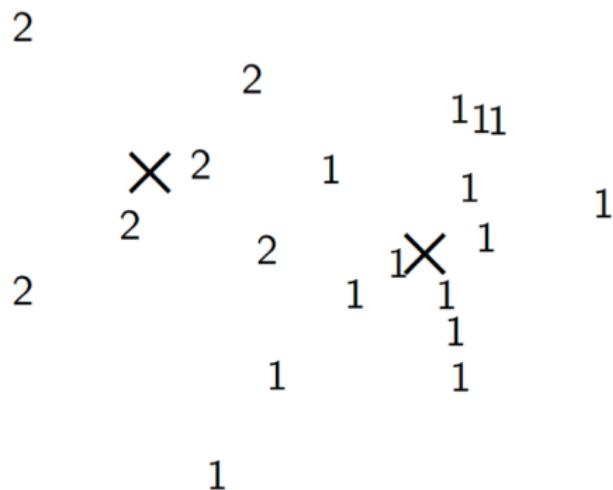
2 1 1 1

1

Worked Example: Recompute cluster centroids



Worked Ex.: Centroids and assignments after convergence



K-means is guaranteed to converge: Proof

- RSS = sum of all squared distances between document vector and closest centroid
- RSS decreases during each reassignment step.
 - because each vector is moved to a closer centroid
- RSS decreases during each recomputation step.
 - see next slide
- There is only a finite number of clusterings.
- Thus: We must reach a fixed point.
- Assumption: Ties are broken consistently.
- Finite set & monotonically decreasing → convergence



Recomputation decreases average distance

$\text{RSS} = \sum_{k=1}^K \text{RSS}_k$ – the residual sum of squares (the “goodness” measure)

$$\text{RSS}_k(\underline{v}) = \sum_{x \in \omega_k} (\underline{v} - \underline{x})^2 = \sum_{x \in \omega_k} \sum_{m=1}^M (v_m - x_m)^2$$

$$\frac{\partial \text{RSS}_k(\underline{v})}{\partial v_m} = \sum_{x \in \omega_k} 2(v_m - x_m) = 0$$

$$v_m = \frac{1}{|\omega_k|} \sum_{x \in \omega_k} x_m$$

The last line is the componentwise definition of the centroid! We minimize RSS_k when the old centroid is replaced with the new centroid. RSS, the sum of the RSS_k , must then also decrease during recomputation.

K -means is guaranteed to converge

- But we don't know how long convergence will take!
- If we don't care about a few docs switching back and forth, then convergence is usually fast (< 10-20 iterations).
- However, complete convergence can take many more iterations.

Optimality of K -means

- Convergence \neq optimality
- Convergence does not mean that we converge to the optimal clustering!
- This is the great weakness of K -means.
- If we start with a bad set of seeds, the resulting clustering can be horrible.

Initialization of K -means

- Random seed selection is just one of many ways K -means can be initialized.
- Random seed selection is not very robust: It's easy to get a suboptimal clustering.
- Better ways of computing initial centroids:
 - Select seeds not randomly, but using some heuristic (e.g., filter out outliers or find a set of seeds that has "good coverage" of the document space)
 - Use hierarchical clustering to find good seeds
 - Select i (e.g., $i = 10$) different random sets of seeds, do a K -means clustering for each, select the clustering with lowest RSS

Time complexity of K -means

- Computing one distance of two vectors is $O(M)$.
- Reassignment step: $O(KNM)$ (we need to compute KN document-centroid distances)
- Recomputation step: $O(NM)$ (we need to add each of the document's $< M$ values to one of the centroids)
- Assume number of iterations bounded by I
- Overall complexity: $O(INKM)$ – linear in all important dimensions
- However: This is not a real worst-case analysis.
- In pathological cases, complexity can be worse than linear.



Outline

- What is clustering?
- Applications of clustering in information retrieval
- *K*-means algorithm
- Evaluation of clustering
-

What is a good clustering?

- Internal criteria
 - Example of an internal criterion: RSS in K -means
- But an internal criterion often does not evaluate the actual utility of a clustering in the application.
- Alternative: External criteria
 - Evaluate with respect to a human-defined classification



External criteria for clustering quality

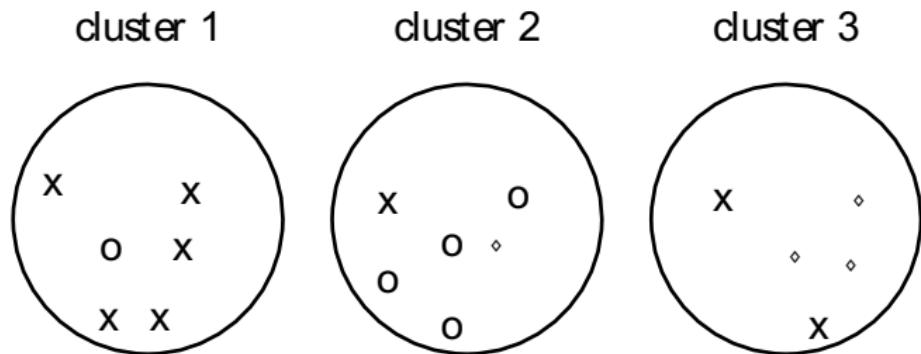
- Based on a gold standard data set, with the class of each instance labeled by human annotators.
- Goal: Clustering should reproduce the classes in the gold standard
- (But we only want to reproduce how documents are divided into groups, not the class labels.)
- First measure for how well we were able to reproduce the classes: **purity**

External criterion: Purity

$$\text{purity}(\Omega, C) = \frac{1}{N} \sum_k \max_j |\omega_k \cap c_j|$$

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_K\}$ is the set of clusters and $C = \{c_1, c_2, \dots, c_J\}$ is the set of classes.
- For each cluster ω_k : find class c_j with most members n_{kj} in ω_k
- Sum all n_{kj} and divide by total number of points

Example for computing purity



To compute

purity: $5 = \max_j |\omega_1 \cap c_j|$ (class x, cluster 1); $4 = \max_j |\omega_2 \cap c_j|$ (class o, cluster 2); and $3 = \max_j |\omega_3 \cap c_j|$ (class \diamond , cluster 3).
Purity is $(1/17) \times (5 + 4 + 3) \approx 0.71$.

Another external criterion: Rand index

- Purity can be increased easily by increasing K – a measure that does not have this problem: Rand index.
- Definition: $RI = \frac{TP + TN}{TP + FP + FN + TN}$
- Based on 2×2 contingency table of all pairs of documents:

	same cluster	different clusters
same class	true positives (TP)	false negatives (FN)
different classes	false positives (FP)	true negatives (TN)
- $TP + FN + FP + TN$ is the total number of pairs.
- $TP + FN + FP + TN = \binom{N}{2}$ for N documents.
- Example: $\binom{17}{2} = 136$ in o/◊/x example
- Each pair is either positive or negative (the clustering puts the two documents in the same or in different clusters) ...
- ... and either “true” (correct) or “false” (incorrect): the clustering decision is correct or incorrect.

Rand Index: Example

As an example, we compute RI for the o/ ◊/ x example. We first compute TP + FP. The three clusters contain 6, 6, and 5 points, respectively, so the total number of “positives” or pairs of documents that are in the same cluster is:

$$TP + FP = \binom{6}{2} + \binom{6}{2} + \binom{5}{2} = 40$$

Of these, the x pairs in cluster 1, the o pairs in cluster 2, the ◊ pairs in cluster 3, and the x pair in cluster 3 are true positives:

$$TP = \binom{5}{2} + \binom{4}{2} + \binom{3}{2} + \binom{2}{2} = 20$$

Thus, $FP = 40 - 20 = 20$. FN and TN are computed similarly. □

Rand measure for the o/ ◊/ x example

	same cluster	different clusters	
same class	TP = 20	FN = 24	RI is then
different classes	FP = 20	TN = 72	

$$(20 + 72) / (20 + 20 + 24 + 72) \approx 0.68.$$

Two other external evaluation measures

- Two other measures
- Normalized mutual information (NMI)
 - How much information does the clustering contain about the classification?
 - Singleton clusters (number of clusters = number of docs) have maximum MI
 - Therefore: normalize by entropy of clusters and classes
- F measure
 - Like Rand, but “precision” and “recall” can be weighted

Evaluation results for the o/ ◊/ x example

	purity	NMI	RI	F_5
lower bound	0.0	0.0	0.0	0.0
maximum	1.0	1.0	1.0	1.0
value for example	0.71	0.36	0.68	0.46

All four measures range from 0 (really bad clustering) to 1 (perfect clustering). □

Outline

- What is clustering?
- Applications of clustering in information retrieval
- *K*-means algorithm
- Evaluation of clustering
- How many clusters?

How many clusters?

- Number of clusters K is given in many applications.
 - E.g., there may be an external constraint on K . Example: In the case of Scatter-Gather, it was hard to show more than 10–20 clusters on a monitor in the 90s.
- What if there is no external constraint? Is there a “right” number of clusters?
- One way to go: define an optimization criterion
 - Given docs, find K for which the optimum is reached.
 - What optimization criterion can we use?
 - We can't use RSS or average squared distance from centroid as criterion: always chooses $K = N$ clusters.

Exercise

- Your job is to develop the clustering algorithms for a competitor to news.google.com
- You want to use K -means clustering.
- How would you determine K ? □

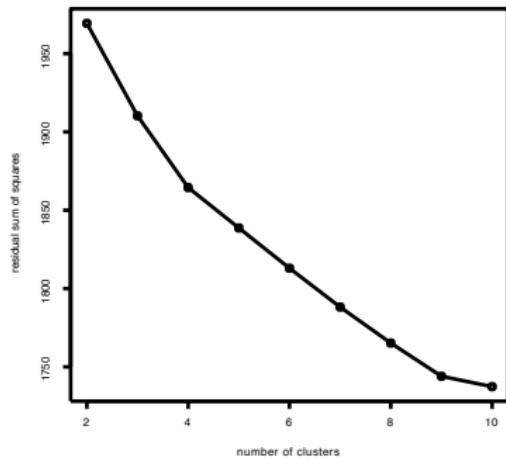
Simple objective function for K : Basic idea

- Start with 1 cluster ($K = 1$)
- Keep adding clusters (= keep increasing K)
- Add a penalty for each new cluster
- Then trade off cluster penalties against average squared distance from centroid
- Choose the value of K with the best tradeoff

Simple objective function for K : Formalization

- Given a clustering, define the cost for a document as (squared) distance to centroid
- Define total **distortion** $\text{RSS}(K)$ as sum of all individual document costs (corresponds to average distance)
- Then: penalize each cluster with a cost λ
- Thus for a clustering with K clusters, total cluster penalty is $K\lambda$
- Define the total cost of a clustering as distortion plus total cluster penalty: $\text{RSS}(K) + K\lambda$
- Select K that minimizes $(\text{RSS}(K) + K\lambda)$
- Still need to determine good value for λ ...

Finding the “knee” in the curve



Pick the number of clusters where
curve “flattens”. Here: 4 or 9.

Resources

- Chapter 16 of IIR
- Resources at <http://cis1.mi.org>
 - Keith van Rijsbergen on the cluster hypothesis (he was one of the originators)
 - Bing/ Y i p p y : search result clustering systems
 - Stirling number: the number of distinct k -clusterings of n items