

Introduction to Information Retrieval

Hinrich Schütze and Christina Lioma
Chapter 16: Flat Clustering
(Slides modified by Scott Sanner)

1

Overview

- ① Recap
- ② Clustering: Introduction
- ③ Clustering in IR
- ④ *K*-means
- ⑤ Evaluation
- ⑥ How many clusters?

2

Outline

- ① Recap
- ② Clustering: Introduction
- ③ Clustering in IR
- ④ *K*-means
- ⑤ Evaluation
- ⑥ How many clusters?

3

Recap

- Up until now...
 - We have built document classifiers given labeled data
 - Classification known as supervised learning
- Today...
 - What if we don't have labeled data?
 - We don't know the class labels?
 - Can we still assign reasonable labels?
 - We could try to cluster documents (cluster=class)
 - What would make a human readable label?
 - Clustering known as unsupervised learning

4

Outline

- ① Recap
- ② Clustering: Introduction
- ③ Clustering in IR
- ④ *K*-means
- ⑤ Evaluation
- ⑥ How many clusters?

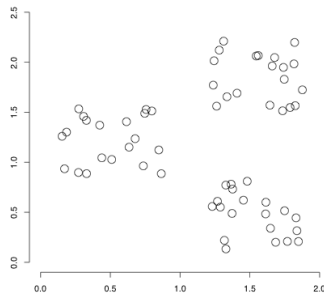
5

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.

6

Data set with clear cluster structure



Propose algorithm for finding the cluster structure in this example

7

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, . . .

8

Outline

- 1 Recap
- 2 Clustering: Introduction
- 3 Clustering in IR
- 4 K-means
- 5 Evaluation
- 6 How many clusters?

9

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 Rijsbergen's original wording: "closely associated documents tend to be relevant to the same requests".

10

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

deduplication
source clustering
diversification

11

Search result clustering for better navigation

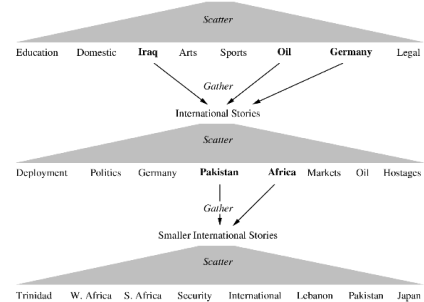
12

Clustering for better navigation

- Ambiguity like Jaguar a bad example
- Better case of ambiguity: London Heathrow
 - What is my intent?
 - Checkin
 - Ground transport
 - Food availability
 - Security

13

Scatter-Gather



14

Global navigation: Yahoo

15

Global navigation: MESH (upper level)

MeSH Tree Structures - 2008

- [Return to Entry Page](#)
1. Anatomy (A)
 2. Organisms (B)
 3. Diseases (C)
 - 3.1 Infectious Diseases and Microbes (C.001) +
 - 3.2 Parasitic Diseases (C.002) +
 - 3.3 Virology (C.003) +
 - 3.4 Immunology (C.004) +
 - 3.5 Allergies (C.005) +
 - 3.6 Digestive System Diseases (C.006) +
 - 3.7 Respiratory System Diseases (C.007) +
 - 3.8 Cardiovascular Diseases (C.008) +
 - 3.9 Nervous System Diseases (C.009) +
 - 3.10 Endocrine System Diseases (C.010) +
 - 3.11 Male Reproductive Diseases (C.011) +
 - 3.12 Female Reproductive Diseases and Pregnancy Complications (C.012) +
 - 3.13 Musculoskeletal Diseases (C.013) +
 - 3.14 Skin and Connective Tissue Diseases (C.014) +
 - 3.15 Sensory and Motor System Diseases (C.015) +
 - 3.16 Mental and Behavioral Disorders (C.016) +
 - 3.17 Endocrine System Diseases (C.017) +
 - 3.18 Reproductive System Diseases (C.018) +
 - 3.19 Disorders of Environmental Origin (C.019) +
 - 3.20 Nutritional Disorders (C.020) +
 - 3.21 Pathological Conditions, Signs and Symptoms (C.021) +
 4. Chemicals and Drugs (D)
 5. Analytical, Diagnostic and Therapeutic Techniques and Equipment (E)
 6. Psychology and Psychiatry (F)
 7. Biological Sciences (G)
 8. Natural Sciences (H)
 9. Anthropology, Education, Sociology and Social Phenomena (I)
 10. Technology, Industry, Agriculture (J)
 11. Humanities (K)

16

Global navigation: MESH (lower level)

- [Neoplasia \[C04\]](#)
- Cysts [C04.182] +
 - Hematomas [C04.445] +
 - Neoplasia by Histologic Type [C04.557]
 - Benign Neoplasms [C04.557.001] +
 - Benign Neoplasms by Site [C04.557.002] +
 - Benign Neoplasms by System [C04.557.003] +
 - Benign Neoplasms by Tissue [C04.557.004] +
 - Benign Neoplasms by Organ [C04.557.005] +
 - Benign Neoplasms by Cell [C04.557.006] +
 - Benign Neoplasms by Tissue [C04.557.007] +
 - Benign Neoplasms by Organ [C04.557.008] +
 - Benign Neoplasms by Cell [C04.557.009] +
 - Benign Neoplasms by Tissue [C04.557.010] +
 - Benign Neoplasms by Organ [C04.557.011] +
 - Benign Neoplasms by Cell [C04.557.012] +
 - Benign Neoplasms by Tissue [C04.557.013] +
 - Benign Neoplasms by Organ [C04.557.014] +
 - Benign Neoplasms by Cell [C04.557.015] +
 - Benign Neoplasms by Tissue [C04.557.016] +
 - Benign Neoplasms by Organ [C04.557.017] +
 - Benign Neoplasms by Cell [C04.557.018] +
 - Benign Neoplasms by Tissue [C04.557.019] +
 - Benign Neoplasms by Organ [C04.557.020] +
 - Benign Neoplasms by Cell [C04.557.021] +
 - Benign Neoplasms by Tissue [C04.557.022] +
 - Benign Neoplasms by Organ [C04.557.023] +
 - Benign Neoplasms by Cell [C04.557.024] +
 - Benign Neoplasms by Tissue [C04.557.025] +
 - Benign Neoplasms by Organ [C04.557.026] +
 - Benign Neoplasms by Cell [C04.557.027] +
 - Benign Neoplasms by Tissue [C04.557.028] +
 - Benign Neoplasms by Organ [C04.557.029] +
 - Benign Neoplasms by Cell [C04.557.030] +
 - Benign Neoplasms by Tissue [C04.557.031] +
 - Benign Neoplasms by Organ [C04.557.032] +
 - Benign Neoplasms by Cell [C04.557.033] +
 - Benign Neoplasms by Tissue [C04.557.034] +
 - Benign Neoplasms by Organ [C04.557.035] +
 - Benign Neoplasms by Cell [C04.557.036] +
 - Benign Neoplasms by Tissue [C04.557.037] +
 - Benign Neoplasms by Organ [C04.557.038] +
 - Benign Neoplasms by Cell [C04.557.039] +
 - Benign Neoplasms by Tissue [C04.557.040] +
 - Benign Neoplasms by Organ [C04.557.041] +
 - Benign Neoplasms by Cell [C04.557.042] +
 - Benign Neoplasms by Tissue [C04.557.043] +
 - Benign Neoplasms by Organ [C04.557.044] +
 - Benign Neoplasms by Cell [C04.557.045] +
 - Benign Neoplasms by Tissue [C04.557.046] +
 - Benign Neoplasms by Organ [C04.557.047] +
 - Benign Neoplasms by Cell [C04.557.048] +
 - Benign Neoplasms by Tissue [C04.557.049] +
 - Benign Neoplasms by Organ [C04.557.050] +
 - Benign Neoplasms by Cell [C04.557.051] +
 - Benign Neoplasms by Tissue [C04.557.052] +
 - Benign Neoplasms by Organ [C04.557.053] +
 - Benign Neoplasms by Cell [C04.557.054] +
 - Benign Neoplasms by Tissue [C04.557.055] +
 - Benign Neoplasms by Organ [C04.557.056] +
 - Benign Neoplasms by Cell [C04.557.057] +
 - Benign Neoplasms by Tissue [C04.557.058] +
 - Benign Neoplasms by Organ [C04.557.059] +
 - Benign Neoplasms by Cell [C04.557.060] +
 - Benign Neoplasms by Tissue [C04.557.061] +
 - Benign Neoplasms by Organ [C04.557.062] +
 - Benign Neoplasms by Cell [C04.557.063] +
 - Benign Neoplasms by Tissue [C04.557.064] +
 - Benign Neoplasms by Organ [C04.557.065] +
 - Benign Neoplasms by Cell [C04.557.066] +
 - Benign Neoplasms by Tissue [C04.557.067] +
 - Benign Neoplasms by Organ [C04.557.068] +
 - Benign Neoplasms by Cell [C04.557.069] +
 - Benign Neoplasms by Tissue [C04.557.070] +
 - Benign Neoplasms by Organ [C04.557.071] +
 - Benign Neoplasms by Cell [C04.557.072] +
 - Benign Neoplasms by Tissue [C04.557.073] +
 - Benign Neoplasms by Organ [C04.557.074] +
 - Benign Neoplasms by Cell [C04.557.075] +
 - Benign Neoplasms by Tissue [C04.557.076] +
 - Benign Neoplasms by Organ [C04.557.077] +
 - Benign Neoplasms by Cell [C04.557.078] +
 - Benign Neoplasms by Tissue [C04.557.079] +
 - Benign Neoplasms by Organ [C04.557.080] +
 - Benign Neoplasms by Cell [C04.557.081] +
 - Benign Neoplasms by Tissue [C04.557.082] +
 - Benign Neoplasms by Organ [C04.557.083] +
 - Benign Neoplasms by Cell [C04.557.084] +
 - Benign Neoplasms by Tissue [C04.557.085] +
 - Benign Neoplasms by Organ [C04.557.086] +
 - Benign Neoplasms by Cell [C04.557.087] +
 - Benign Neoplasms by Tissue [C04.557.088] +
 - Benign Neoplasms by Organ [C04.557.089] +
 - Benign Neoplasms by Cell [C04.557.090] +
 - Benign Neoplasms by Tissue [C04.557.091] +
 - Benign Neoplasms by Organ [C04.557.092] +
 - Benign Neoplasms by Cell [C04.557.093] +
 - Benign Neoplasms by Tissue [C04.557.094] +
 - Benign Neoplasms by Organ [C04.557.095] +
 - Benign Neoplasms by Cell [C04.557.096] +
 - Benign Neoplasms by Tissue [C04.557.097] +
 - Benign Neoplasms by Organ [C04.557.098] +
 - Benign Neoplasms by Cell [C04.557.099] +
 - Benign Neoplasms by Tissue [C04.557.100] +
 - Malignant Neoplasms [C04.557.101] +
 - Malignant Neoplasms by Site [C04.557.102] +
 - Malignant Neoplasms by System [C04.557.103] +
 - Malignant Neoplasms by Tissue [C04.557.104] +
 - Malignant Neoplasms by Organ [C04.557.105] +
 - Malignant Neoplasms by Cell [C04.557.106] +
 - Malignant Neoplasms by Tissue [C04.557.107] +
 - Malignant Neoplasms by Organ [C04.557.108] +
 - Malignant Neoplasms by Cell [C04.557.109] +
 - Malignant Neoplasms by Tissue [C04.557.110] +
 - Malignant Neoplasms by Organ [C04.557.111] +
 - Malignant Neoplasms by Cell [C04.557.112] +
 - Malignant Neoplasms by Tissue [C04.557.113] +
 - Malignant Neoplasms by Organ [C04.557.114] +
 - Malignant Neoplasms by Cell [C04.557.115] +
 - Malignant Neoplasms by Tissue [C04.557.116] +
 - Malignant Neoplasms by Organ [C04.557.117] +
 - Malignant Neoplasms by Cell [C04.557.118] +
 - Malignant Neoplasms by Tissue [C04.557.119] +
 - Malignant Neoplasms by Organ [C04.557.120] +
 - Malignant Neoplasms by Cell [C04.557.121] +
 - Malignant Neoplasms by Tissue [C04.557.122] +
 - Malignant Neoplasms by Organ [C04.557.123] +
 - Malignant Neoplasms by Cell [C04.557.124] +
 - Malignant Neoplasms by Tissue [C04.557.125] +
 - Malignant Neoplasms by Organ [C04.557.126] +
 - Malignant Neoplasms by Cell [C04.557.127] +
 - Malignant Neoplasms by Tissue [C04.557.128] +
 - Malignant Neoplasms by Organ [C04.557.129] +
 - Malignant Neoplasms by Cell [C04.557.130] +
 - Malignant Neoplasms by Tissue [C04.557.131] +
 - Malignant Neoplasms by Organ [C04.557.132] +
 - Malignant Neoplasms by Cell [C04.557.133] +
 - Malignant Neoplasms by Tissue [C04.557.134] +
 - Malignant Neoplasms by Organ [C04.557.135] +
 - Malignant Neoplasms by Cell [C04.557.136] +
 - Malignant Neoplasms by Tissue [C04.557.137] +
 - Malignant Neoplasms by Organ [C04.557.138] +
 - Malignant Neoplasms by Cell [C04.557.139] +
 - Malignant Neoplasms by Tissue [C04.557.140] +
 - Malignant Neoplasms by Organ [C04.557.141] +
 - Malignant Neoplasms by Cell [C04.557.142] +
 - Malignant Neoplasms by Tissue [C04.557.143] +
 - Malignant Neoplasms by Organ [C04.557.144] +
 - Malignant Neoplasms by Cell [C04.557.145] +
 - Malignant Neoplasms by Tissue [C04.557.146] +
 - Malignant Neoplasms by Organ [C04.557.147] +
 - Malignant Neoplasms by Cell [C04.557.148] +
 - Malignant Neoplasms by Tissue [C04.557.149] +
 - Malignant Neoplasms by Organ [C04.557.150] +
 - Malignant Neoplasms by Cell [C04.557.151] +
 - Malignant Neoplasms by Tissue [C04.557.152] +
 - Malignant Neoplasms by Organ [C04.557.153] +
 - Malignant Neoplasms by Cell [C04.557.154] +
 - Malignant Neoplasms by Tissue [C04.557.155] +
 - Malignant Neoplasms by Organ [C04.557.156] +
 - Malignant Neoplasms by Cell [C04.557.157] +
 - Malignant Neoplasms by Tissue [C04.557.158] +
 - Malignant Neoplasms by Organ [C04.557.159] +
 - Malignant Neoplasms by Cell [C04.557.160] +
 - Malignant Neoplasms by Tissue [C04.557.161] +
 - Malignant Neoplasms by Organ [C04.557.162] +
 - Malignant Neoplasms by Cell [C04.557.163] +
 - Malignant Neoplasms by Tissue [C04.557.164] +
 - Malignant Neoplasms by Organ [C04.557.165] +
 - Malignant Neoplasms by Cell [C04.557.166] +
 - Malignant Neoplasms by Tissue [C04.557.167] +
 - Malignant Neoplasms by Organ [C04.557.168] +
 - Malignant Neoplasms by Cell [C04.557.169] +
 - Malignant Neoplasms by Tissue [C04.557.170] +
 - Malignant Neoplasms by Organ [C04.557.171] +
 - Malignant Neoplasms by Cell [C04.557.172] +
 - Malignant Neoplasms by Tissue [C04.557.173] +
 - Malignant Neoplasms by Organ [C04.557.174] +
 - Malignant Neoplasms by Cell [C04.557.175] +
 - Malignant Neoplasms by Tissue [C04.557.176] +
 - Malignant Neoplasms by Organ [C04.557.177] +
 - Malignant Neoplasms by Cell [C04.557.178] +
 - Malignant Neoplasms by Tissue [C04.557.179] +
 - Malignant Neoplasms by Organ [C04.557.180] +
 - Malignant Neoplasms by Cell [C04.557.181] +
 - Malignant Neoplasms by Tissue [C04.557.182] +
 - Malignant Neoplasms by Organ [C04.557.183] +
 - Malignant Neoplasms by Cell [C04.557.184] +
 - Malignant Neoplasms by Tissue [C04.557.185] +
 - Malignant Neoplasms by Organ [C04.557.186] +
 - Malignant Neoplasms by Cell [C04.557.187] +
 - Malignant Neoplasms by Tissue [C04.557.188] +
 - Malignant Neoplasms by Organ [C04.557.189] +
 - Malignant Neoplasms by Cell [C04.557.190] +
 - Malignant Neoplasms by Tissue [C04.557.191] +
 - Malignant Neoplasms by Organ [C04.557.192] +
 - Malignant Neoplasms by Cell [C04.557.193] +
 - Malignant Neoplasms by Tissue [C04.557.194] +
 - Malignant Neoplasms by Organ [C04.557.195] +
 - Malignant Neoplasms by Cell [C04.557.196] +
 - Malignant Neoplasms by Tissue [C04.557.197] +
 - Malignant Neoplasms by Organ [C04.557.198] +
 - Malignant Neoplasms by Cell [C04.557.199] +
 - Malignant Neoplasms by Tissue [C04.557.200] +

17

Navigational hierarchies: Manual vs. automatic creation

- Note: Yahoo/MESH are **not** examples of clustering.
- But they are well known examples for using a global hierarchy for navigation.
- Some examples for global navigation/exploration based on clustering:
 - Cartia
 - Themespaces
 - Google News

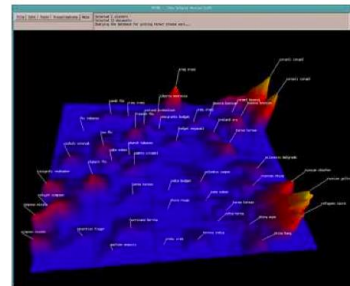
18

Global navigation combined with visualization (1)



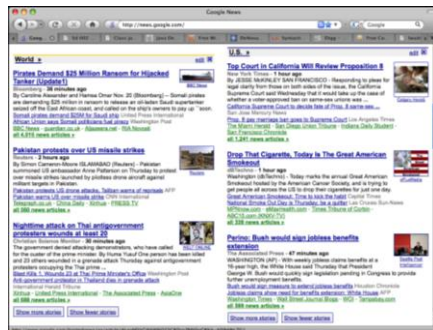
19

Global navigation combined with visualization (2)

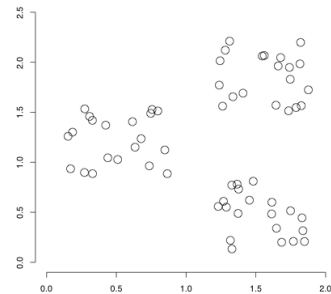


20

Google News: automatic clustering gives an effective news presentation metaphor



Data set with clear cluster structure



Propose algorithm
for finding the
cluster structure
in this example

22

Desiderata for clustering

- General goal: put related docs in the same cluster, put unrelated docs in different clusters.
 - How do we formalize 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 . . .

23

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

24

Hard vs. Soft clustering

- **Hard clustering:** Each document in **exactly one** cluster.
 - More common and easier to do
- **Soft clustering:** A document can be in **more than one** cluster.
 - Makes more sense for 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 do **flat, hard clustering only** in this class.
- See IIR 16.5, IIR 17, IIR 18 for soft clustering and hierarchical clustering

25

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**

26

Outline

- 1 Recap
- 2 Clustering: Introduction
- 3 Clustering in IR
- 4 **K-means**
- 5 Evaluation
- 6 How many clusters?

27

K-means

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

28

Document representations in clustering

- Vector space model
- K-means can use any distance metric
 - Intuition often given via Euclidean distance (i.e., this lecture)
 - But for documents you want to use (1 – cosine similarity)
 - Why?

29

K-means

- 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:

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

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

30

K-means algorithm

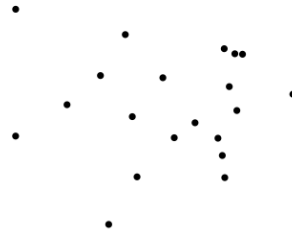
```

K-MEANS( $\{\vec{x}_1, \dots, \vec{x}_N\}, K$ )
1   $(\vec{s}_1, \vec{s}_2, \dots, \vec{s}_K) \leftarrow \text{SELECTRANDOMSEEDS}(\{\vec{x}_1, \dots, \vec{x}_N\}, K)$ 
2  for  $k \leftarrow 1$  to  $K$ 
3  do  $\vec{\mu}_k \leftarrow \vec{s}_k$ 
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  $j \leftarrow \arg \min_j |\vec{\mu}_j - \vec{x}_n|$ 
9   $\omega_j \leftarrow \omega_j \cup \{\vec{x}_n\}$  (reassignment of vectors)
10 for  $k \leftarrow 1$  to  $K$ 
11 do  $\vec{\mu}_k \leftarrow \frac{1}{|\omega_k|} \sum_{\vec{x} \in \omega_k} \vec{x}$  (recomputation of centroids)
12 return  $\{\vec{\mu}_1, \dots, \vec{\mu}_K\}$ 

```

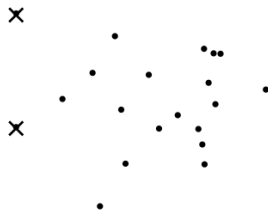
31

Worked Example: Set of to be clustered



32

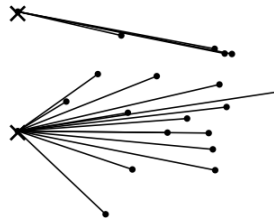
Worked Example: Random selection of initial centroids



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

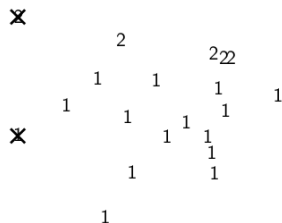
33

Worked Example: Assign points to closest center



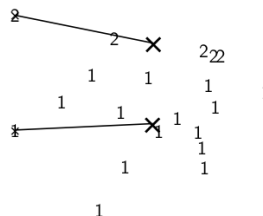
34

Worked Example: Assignment



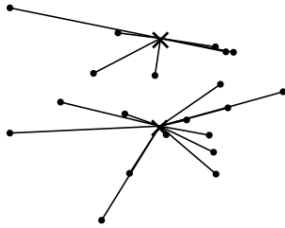
35

Worked Example: Recompute cluster centroids



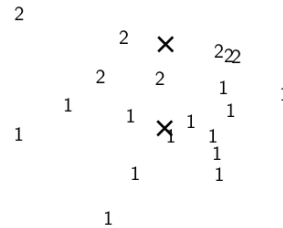
36

Worked Example: Assign points to closest centroid



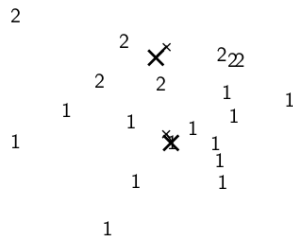
37

Worked Example: Assignment



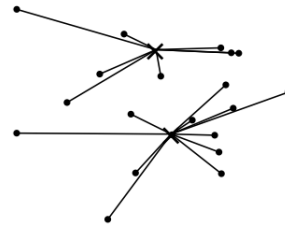
38

Worked Example: Recompute cluster centroids



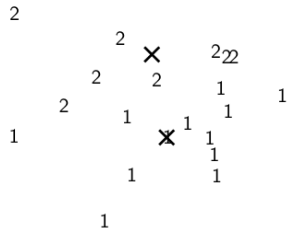
39

Worked Example: Assign points to closest centroid



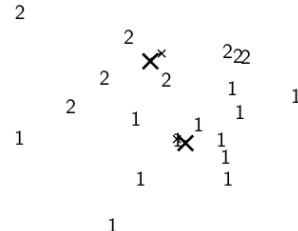
40

Worked Example: Assignment



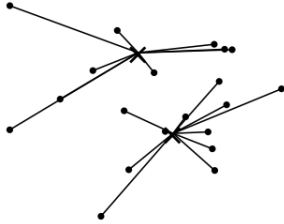
41

Worked Example: Recompute cluster centroids



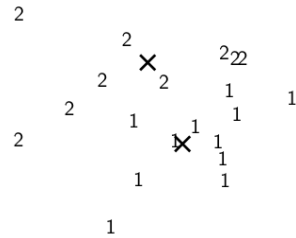
42

Worked Example: Assign points to closest centroid



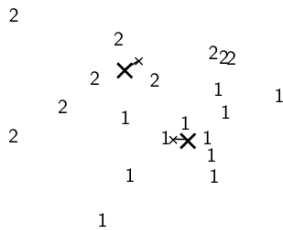
43

Worked Example: Assignment



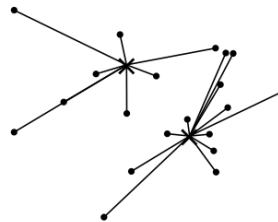
44

Worked Example: Recompute cluster centroids



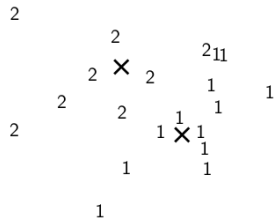
45

Worked Example: Assign points to closest centroid



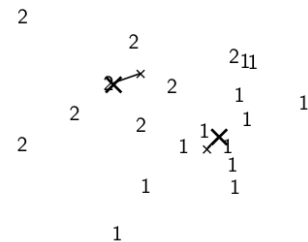
46

Worked Example: Assignment



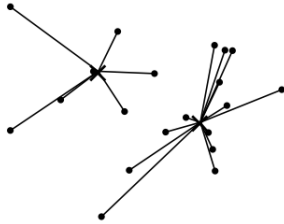
47

Worked Example: Recompute cluster centroids



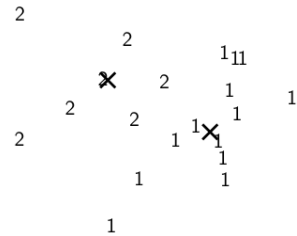
48

Worked Example: Assign points to closest centroid



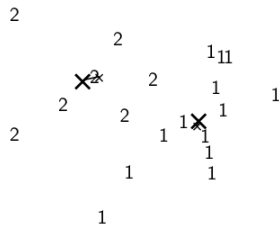
49

Worked Example: Assignment



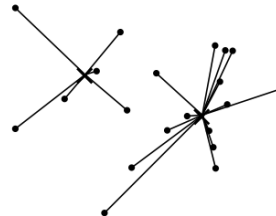
50

Worked Example: Recompute cluster centroids



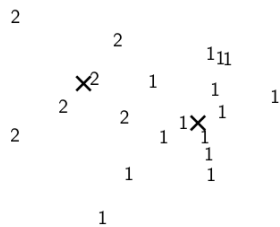
51

Worked Example: Assign points to closest centroid



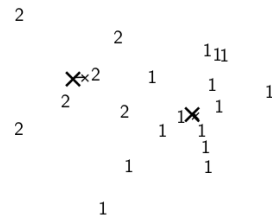
52

Worked Example: Assignment



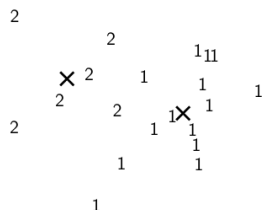
53

Worked Example: Recompute cluster centroids



54

Worked Ex.: Centroids and assignments after convergence



55

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 are only a finite number of clusterings.
- Thus: We must reach a fixed point.
- Assumption: Ties are broken consistently.

56

Recomputation decreases average distance

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

$$RSS_k(\vec{v}) = \sum_{\vec{x} \in \omega_k} \|\vec{v} - \vec{x}\|^2 = \sum_{\vec{x} \in \omega_k} \sum_{m=1}^M (v_m - x_m)^2$$

$$\frac{\partial RSS_k(\vec{v})}{\partial v_m} = \sum_{\vec{x} \in \omega_k} 2(v_m - x_m) = 0$$

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

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.

Mean also reduces RSS for cosine distance – spherical K-means! ⁵⁷

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.

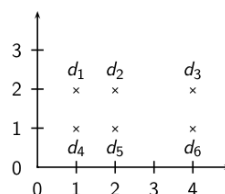
58

Optimality of K-means

- 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.

59

Convergence Exercise: Suboptimal clustering



- What is the optimal clustering for $K = 2$?
- Do we converge on this clustering for arbitrary seeds d_i, d_j ?

60

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

61

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: $O(IKNM)$ – linear in all important dimensions
- However: This is not a real worst-case analysis.
- In pathological cases, complexity can be worse than linear.

62

Outline

- 1 Recap
- 2 Clustering: Introduction
- 3 Clustering in IR
- 4 K-means
- 5 Evaluation
- 6 How many clusters?

63

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

64

External criteria for clustering quality

- Based on a gold standard data set, e.g., the Reuters collection we also used for the evaluation of classification
- 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**

65

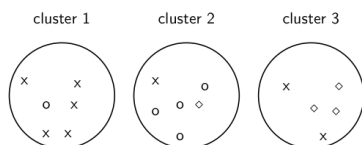
External criterion: Purity

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

- $\Omega = \{\omega_1, \omega_2, \dots, \omega_k\}$ 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

66

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 o, cluster 3). Purity is $(1/17) \times (5 + 4 + 3) \approx 0.71$.

67

Rand index

Definition: $RI = \frac{TP+TN}{TP+FP+FN+TN}$

Based on 2x2 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.
- There are $\binom{N}{2}$ pairs for N documents.
- Example: $\binom{17}{2} = 136$ in o/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.

68

Rand Index: Example

As an example, we compute RI for the o/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 o 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.

69

Rand measure for the o/o/x example

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

RI is then

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

70

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

71

Evaluation results for the o/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).

72

Internal criteria: PMI

- For **text documents**, coherency of top cluster unigrams can be a measure of coherency
- How to measure coherency without labels?
 - Compute average pointwise mutual informatin (PMI) of top-10 unigrams in each cluster (averaged over clusters)

$$\text{pmi}(x; y) \equiv \log \frac{p(x, y)}{p(x)p(y)} = \log \frac{p(x|y)}{p(x)} = \log \frac{p(y|x)}{p(y)}.$$

- Probabilities $p(\cdot)$ computed from background sources (e.g., Wikipedia)

Outline

- 1 Recap
- 2 Clustering: Introduction
- 3 Clustering in IR
- 4 K-means
- 5 Evaluation
- 6 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 ?

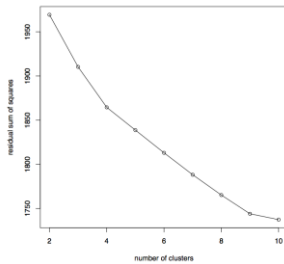
Simple objective function for K (1)

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

Simple objective function for K (2)

- 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 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 $\lambda \dots$

Finding the “knee” in the curve



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

79

Take-away today

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

80

Resources

- Chapter 16 of IIR
- Resources at <http://ifnlp.org/ir>
 - K-means example
 - Keith van Rijsbergen on the cluster hypothesis (he was one of the originators)
 - Bing/Carrot2/Clusty: search result clustering

81