

Clustering and Similarity: Retrieving Documents



Emily Fox & Carlos Guestrin

Machine Learning Specialization

University of Washington

Retrieving documents of interest

Document retrieval

- Currently reading article you like



Document retrieval

- Currently reading article you like
- **Goal:** Want to find similar article



Document retrieval



Challenges

- How do we measure similarity?
- How do we search over articles?

是否有效率高的找法?
(暴力解需看過所有文章)



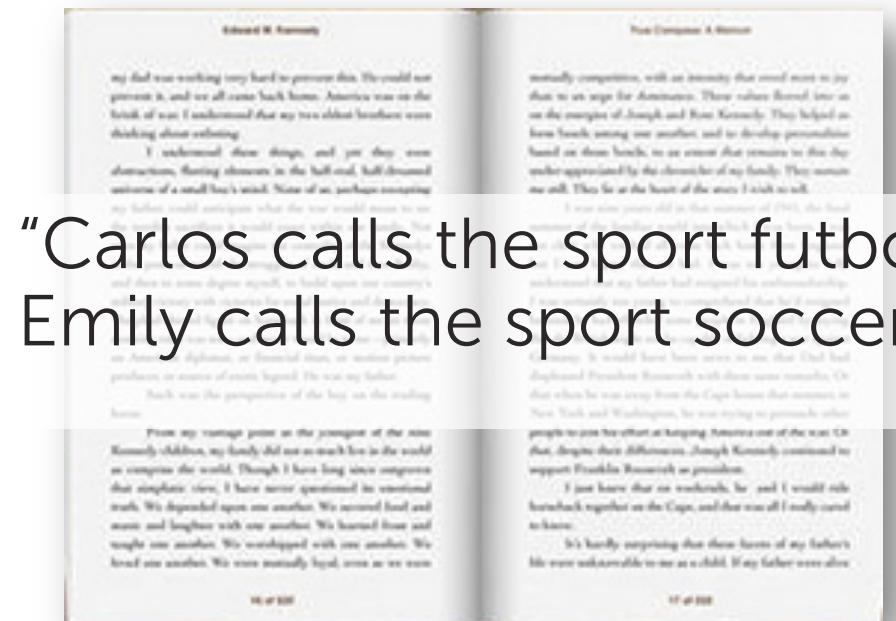
Word count representation for measuring similarity

Word count document representation

- Bag of words model
 - Ignore order of words
 - Count # of instances of each word in vocabulary



Carlos the tree calls sport cat futbol dog soccer Emily



Measuring similarity

相似度：向量內積



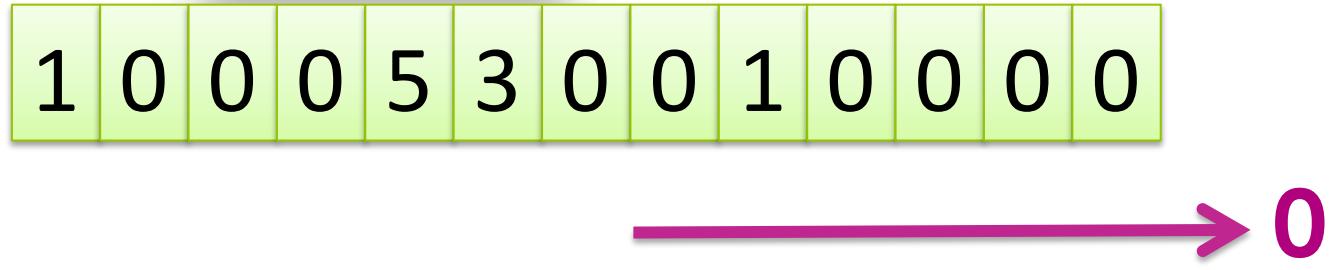
$$1 * 3$$



$$\begin{aligned} &+ \\ &5 * 2 \\ &= 13 \end{aligned}$$



Measuring similarity



Issues with word counts – Doc length

用 word count 會受文章
長度影響



Solution = normalize



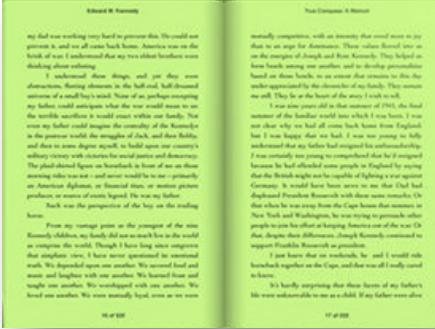
1	0	0	0	5	3	0	0	1	0	0	0	0
---	---	---	---	---	---	---	---	---	---	---	---	---

$$\sqrt{1^2 + 5^2 + 3^2 + 1^2}$$

1					5	3		1				
/	0	0	0	/	/	0	0	/	0	0	0	0
6				6	6		6					

Prioritizing important words with tf-idf

Issues with word counts – Rare words



Common words in doc: “the”, “player”, “field”, “goal”

Dominate rare words like: “futbol”, “Messi”

Document frequency

- What characterizes a **rare word**?
 - Appears **infrequently** in the corpus
- Emphasize words appearing in **few docs**
 - Equivalently, discount word **w** based on
of docs containing w in corpus

Important words

- Do we want only rare words to dominate???
- What characterizes an **important word**?
 - Appears frequently in document (**common locally**)
 - Appears rarely in corpus (**rare globally**)
- Trade off between **local frequency** and **global rarity**

注意這是個假設
在單一 (少數) 文章裡重複
出現且整體文章不常看到
=> **重要的字**

TF-IDF document representation

- Term frequency – inverse document frequency (tf-idf)



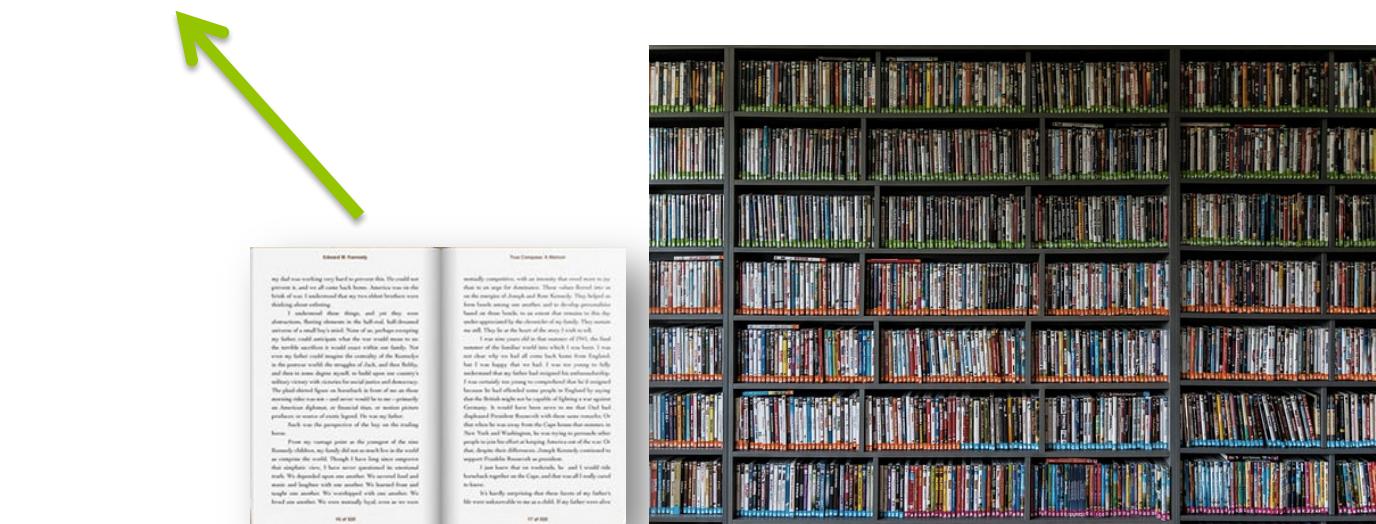
TF-IDF document representation

- Term frequency – inverse document frequency (tf-idf)
- Term frequency



- Same as word counts

其實有很多種表示法



TF-IDF document representation

- Term frequency – inverse document frequency (tf-idf)
- Term frequency



- Inverse document frequency



$$\log \frac{\# \text{ docs}}{1 + \# \text{ docs using word}}$$

避免
除以 0



TF-IDF document representation

- Term frequency – inverse document frequency (tf-idf)
- Term frequency



- Inverse document frequency



$$\log \frac{\# \text{ docs}}{1 + \# \text{ docs using word}}$$

Why log?

word in many docs
rare word

$$\log \frac{\text{large } \#}{1 + \text{small } \#} \rightarrow \text{large } \#$$

$$\frac{\text{large } \#}{1 + \text{large } \#} \approx \log 1 = 0$$



如果定義成這樣？

$$\frac{\# \text{ doc}}{1 + \# \text{ docs using word}} - 1$$

TF-IDF document representation

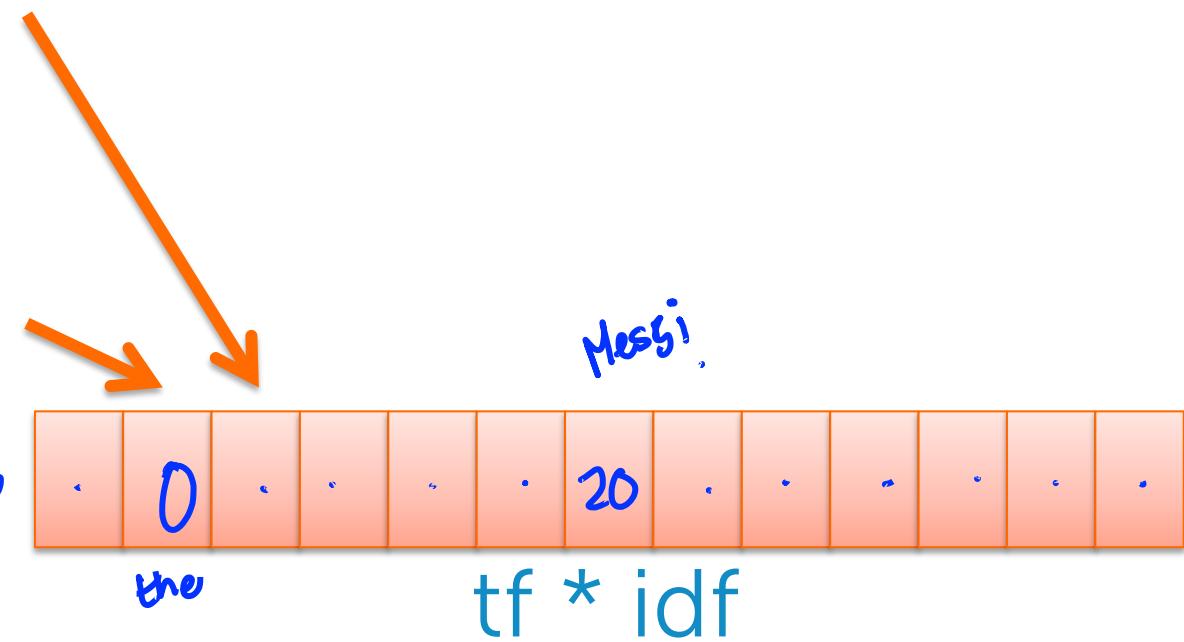
- Term frequency – inverse document frequency (tf-idf)
- Term frequency



- Inverse document frequency



$$\log \frac{64}{1+63} = 0$$
$$\log \frac{64}{1+3} = \log 16$$



Retrieving similar documents

Nearest neighbor search

- Query article:



- Corpus:



- Specify: Distance metric
- Output: Set of most similar articles

距離公式很可能影響結果



1 – Nearest neighbor

找最相近的

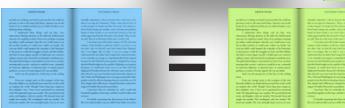
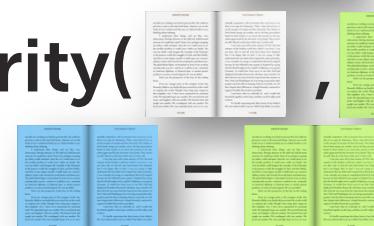
假設?



- **Input:** Query article
- **Output:** *Most* similar article



- Algorithm:
 - Search over each article in corpus
 - Compute $s = \text{similarity}(\text{query}, \text{article})$
 - If $s > \text{Best_s}$, record = and set $\text{Best_s} = s$
 - Return



k – Nearest neighbor

找 k 個最相近的

假設?



- **Input:** Query article
- **Output:** *List of k* similar articles



Clustering documents

參照 Andrew 的 machine learning , week 8 前半段



Structure documents by topic

- Discover groups (*clusters*) of related articles



SPORTS

WORLD NEWS

What if some of the labels are known?

- Training set of labeled docs



SPORTS



WORLD NEWS

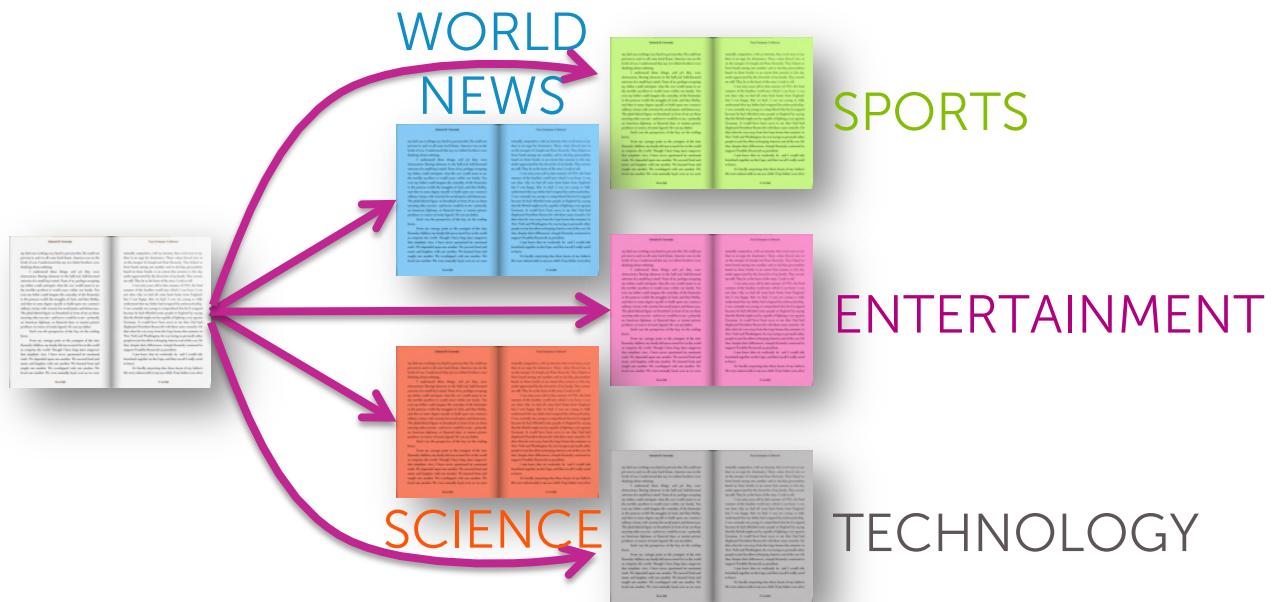


ENTERTAINMENT



SCIENCE

Multiclass classification problem



分群 (clustering) 的問題
在 training data label 已知
的情形下可視為多類別的
supervised learning

?

Example of
supervised learning

Clustering

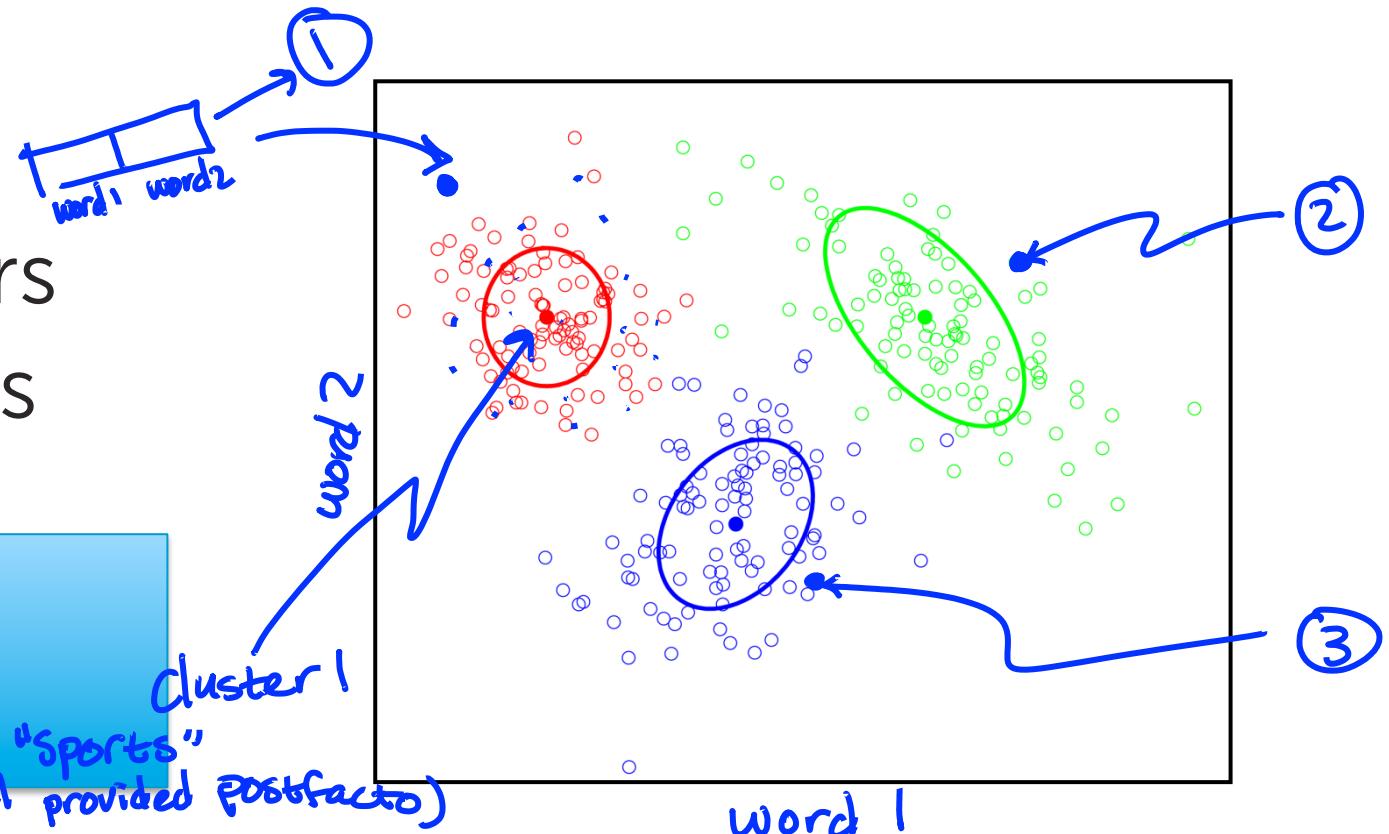
實際上沒有 label
需要找出結構

- No labels provided
- Want to uncover cluster structure
- **Input:** docs as vectors
- **Output:** cluster labels

找出結構後，人為地
在事後給定 label

An **unsupervised**
learning task

“Sports”
(label provided postfacto)

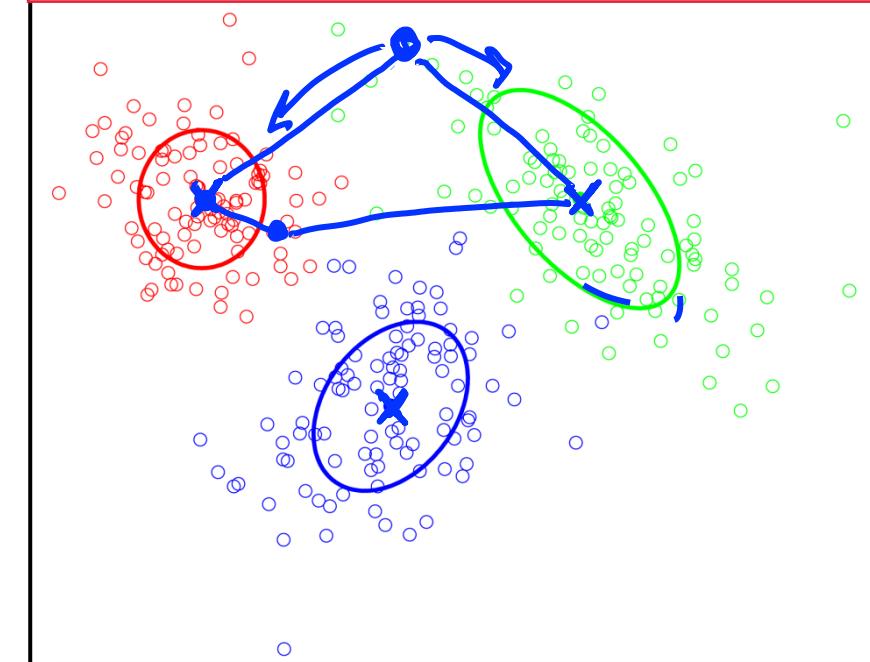


What defines a cluster?

- Cluster defined by center & shape/spread
- Assign observation (doc) to cluster (topic label)
 - Score under cluster is higher than others
 - Often, just more similar to assigned cluster center than other cluster centers

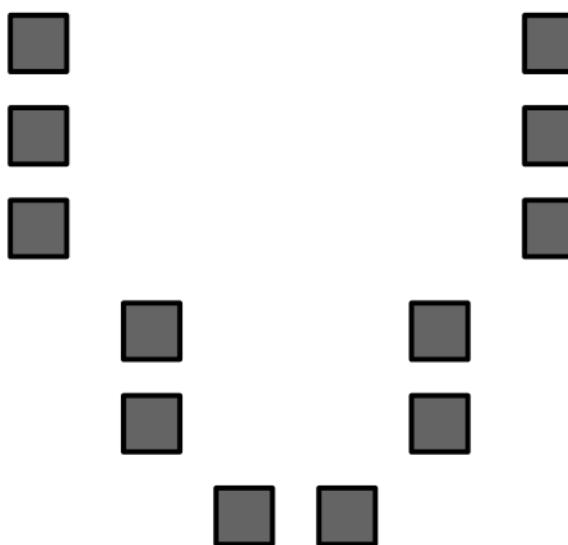
別忽視形狀、分布

只用跟中心點的距離，此點可能會被視為紅色或綠色



k-means

- Assume
 - Similarity metric =
distance to cluster
center
(smaller better)

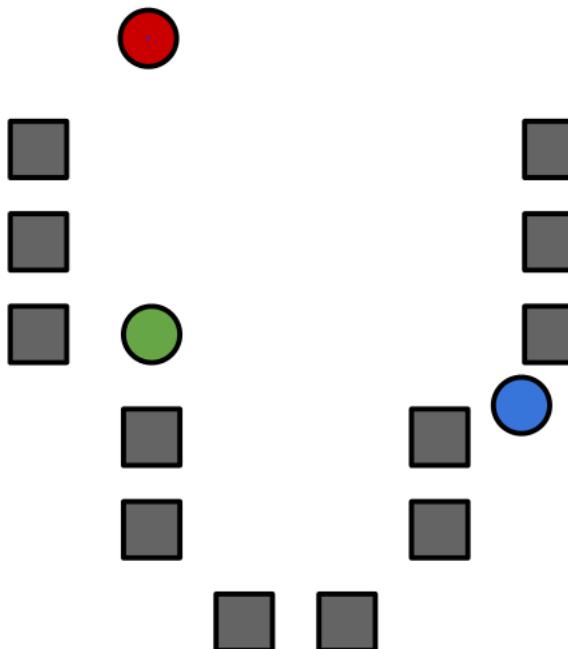


DATA
to
CLUSTER

k-means algorithm

0. Initialize cluster centers

此例子看起來是隨機
生成中心點



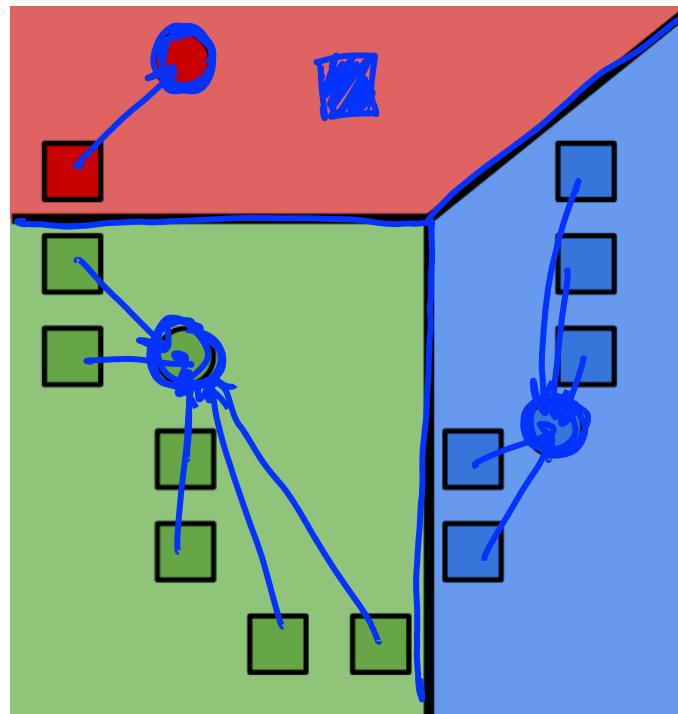
k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center

還記得假設嗎？

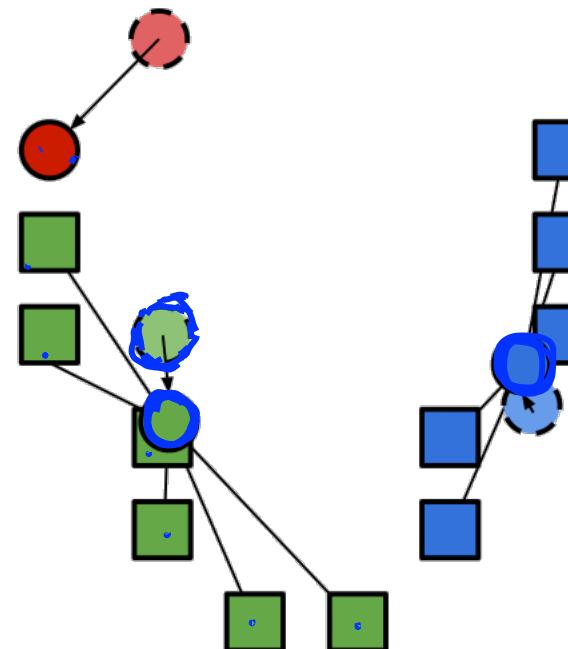


這個示意圖有什麼問題？



k-means algorithm

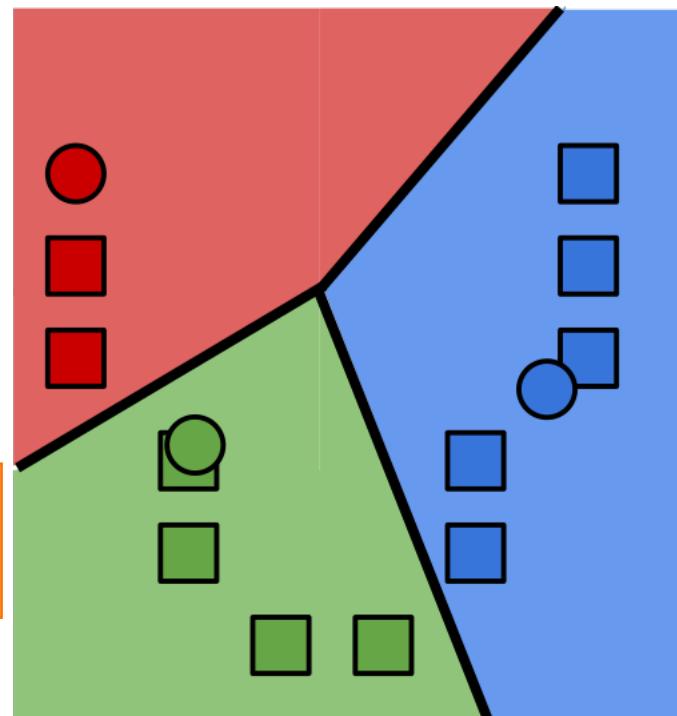
0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations



k-means algorithm

0. Initialize cluster centers
1. Assign observations to closest cluster center
2. Revise cluster centers as mean of assigned observations
3. Repeat 1.+2. until convergence

這次就忽視
小瑕疵吧

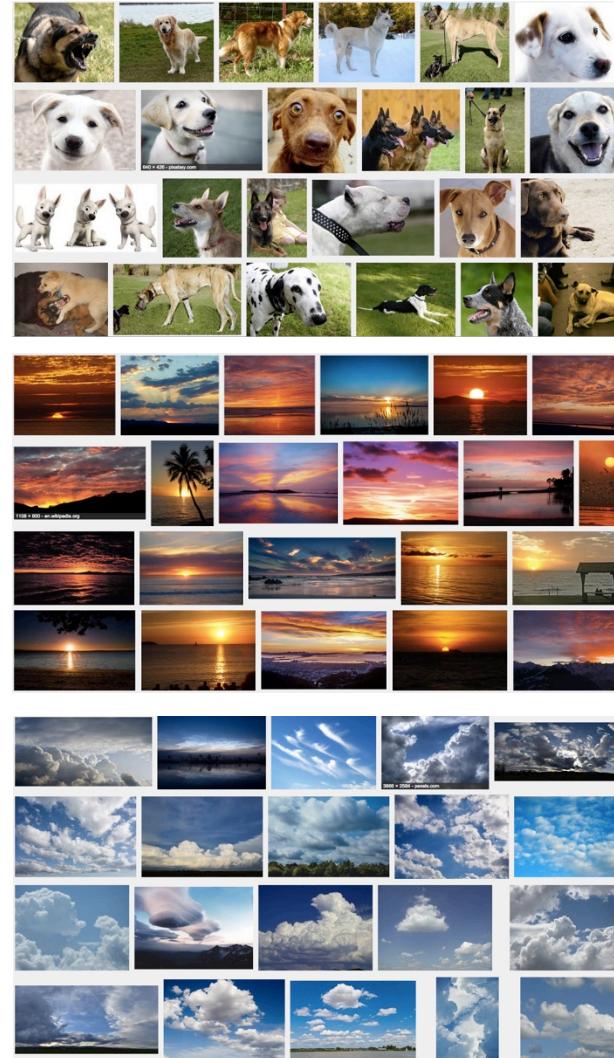
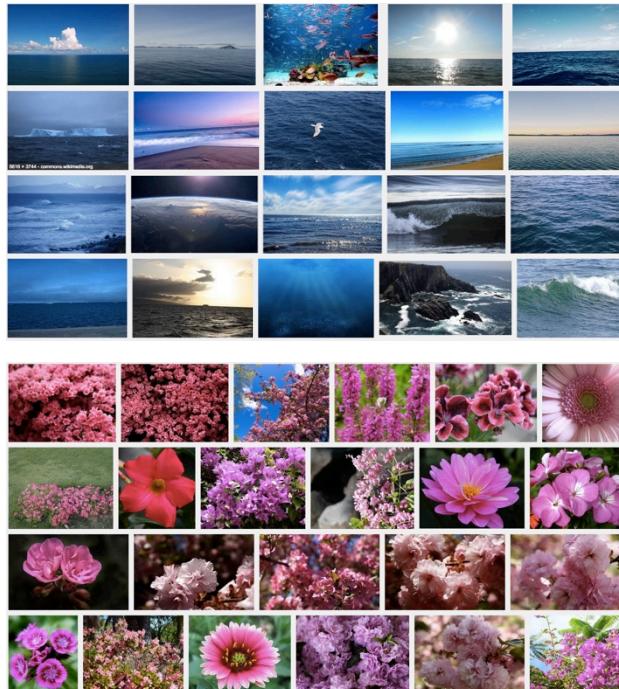


Other examples

後續無 note

Clustering images

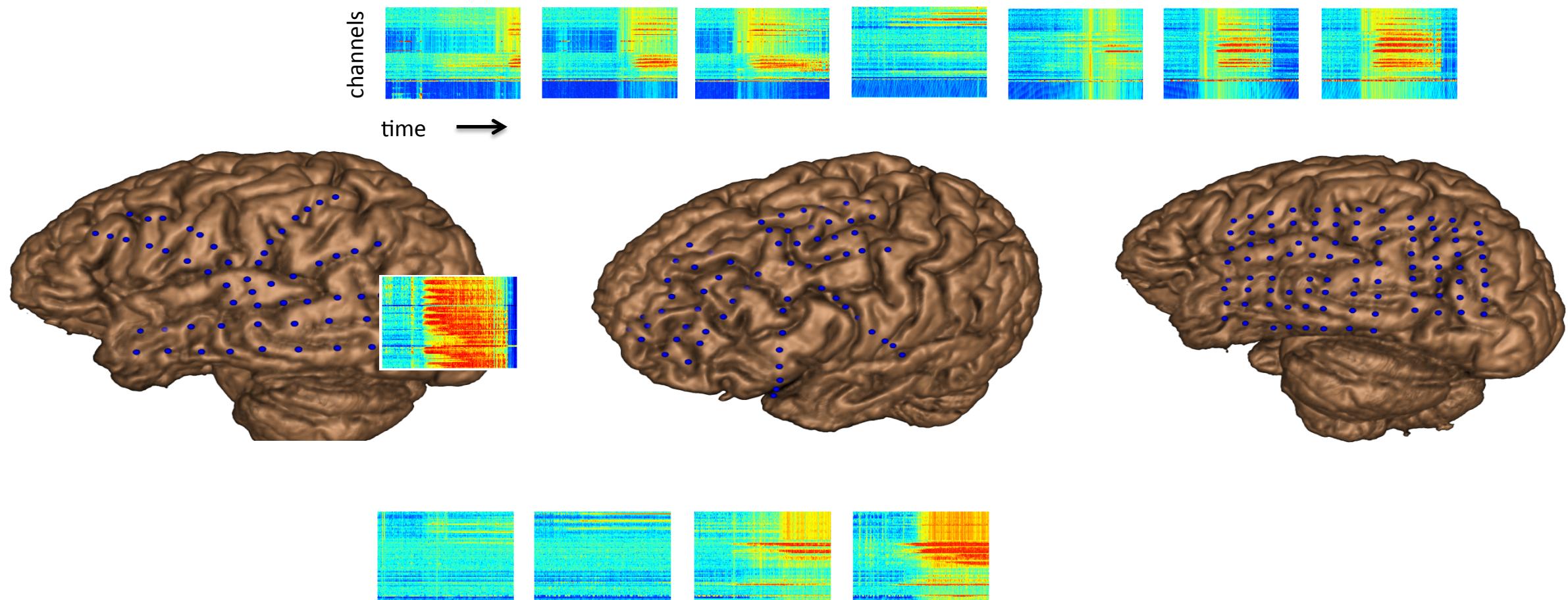
- For search, group as:
 - Ocean
 - Pink flower
 - Dog
 - Sunset
 - Clouds
 - ...



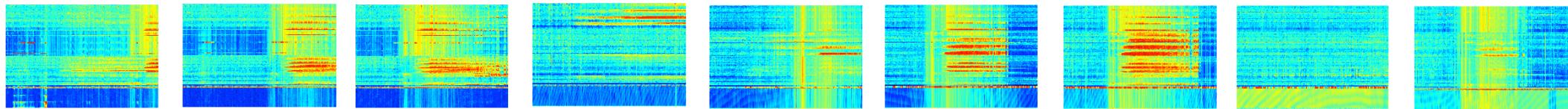
Grouping patients by medical condition

- Better characterize subpopulations and diseases

Example: Patients and seizures are diverse



Cluster seizures by observed time courses



Products on Amazon

- Discover product categories from purchase histories



~~"furniture"~~
"baby"



- Or discovering groups of **users**

Structuring web search results

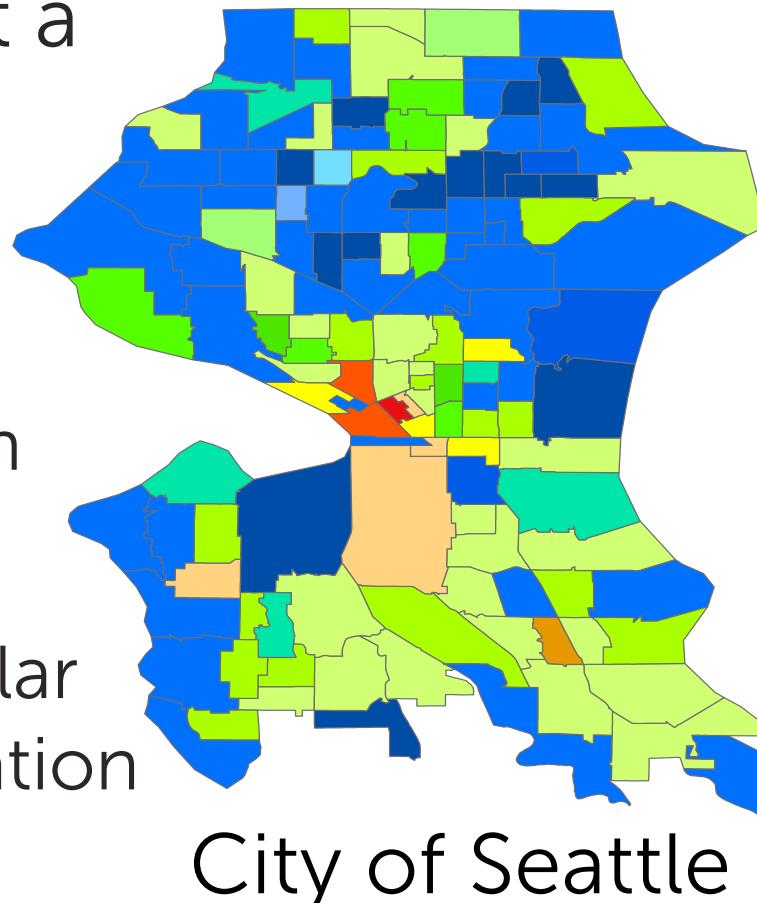
- Search terms can have multiple meanings
- Example: “**cardinal**”



- Use clustering to **structure output**

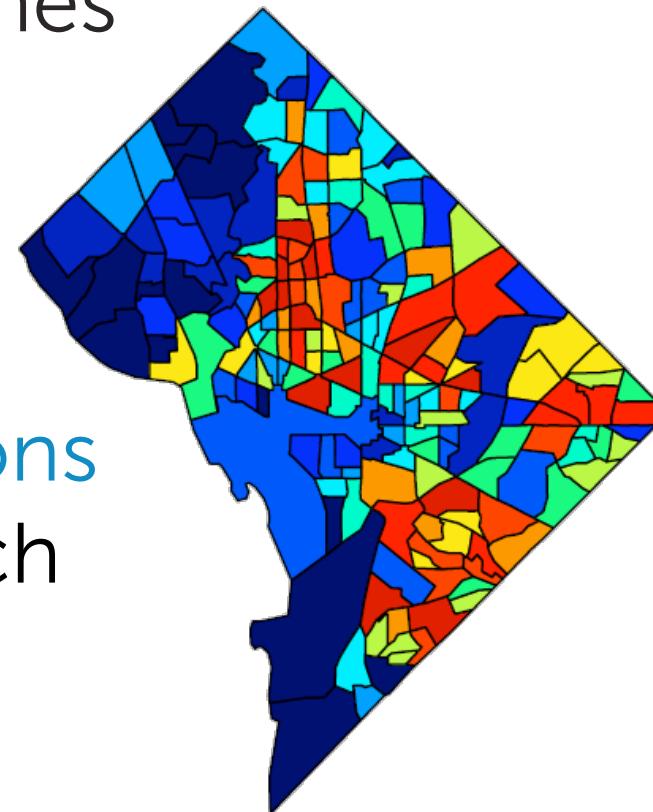
Discovering similar neighborhoods

- **Task 1:** Estimate price at a small regional level
- **Challenge:**
 - Only a few (or no!) sales in each region per month
- **Solution:**
 - Cluster regions with similar trends and share information within a cluster



Discovering similar neighborhoods

- **Task 2:** Forecast violent crimes to better task police
- Again, **cluster regions** and **share information!**
- Leads to **improved predictions** compared to examining each region independently



Washington, DC

Summary for clustering and similarity

What you can do now...

- Describe ways to represent a document (e.g., raw word counts, tf-idf,...)
- Measure the similarity between two documents
- Discuss issues related to using raw word counts
 - Normalize counts to adjust for document length
 - Emphasize important words using tf-idf
- Implement a nearest neighbor search for document retrieval
- Describe the input (unlabeled observations) and output (labels) of a clustering algorithm
- Determine whether a task is supervised or unsupervised
- Cluster documents using k-means (algorithmic details to come...)
- Describe other applications of clustering