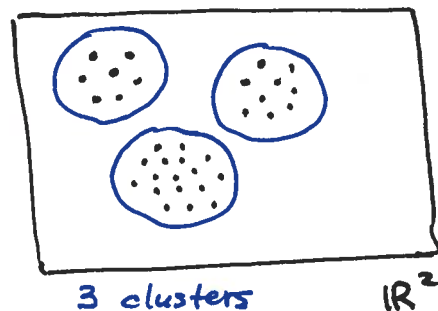


## Lecture 5 The K-means Algorithm for Clustering Vectors (26.1)

Last time we introduced the notion of an innerproduct, discussed how to measure the length of a vector, and how to compute angles and distances between vectors. In this lecture, we will discuss the clustering problem, which often appears in applications whenever we have many vectors in a normed vector space (so that we can compute distances between them).

The clustering problem is then to divide vectors into groups (clusters) of vectors that are close to each other. In particular, we will discuss a famous clustering method, called the K-means algorithm, which is widely used in applications.



Remark 1 K-means will also illustrate one old and very useful idea often used in numerical optimization

Remark 2 Clustering can be discussed in general normed vector spaces, but for the sake of simplicity, we will focus on  $(\mathbb{R}^d, \|\cdot\|_2)$ .

### Clustering Problem

Suppose we have  $n$  vectors  $v_1, \dots, v_n \in \mathbb{R}^d$ .

Goal: partition  $v_1, \dots, v_n$  into  $K$  clusters ( $K \leq n$ ), with the vectors in each cluster close to each other.

#### Applications:

① Topic discovery

We have  $n$  documents (web, text, etc) and a list of  $d$  "key-words".  $(v_i)_j = \# \text{ word } j \text{ in document } i$ .  
Then clusters are groups of <sup>similar</sup> documents.

$$n \sim 10^2 - 10^6$$

$$K \sim 2, 3 - 10^2$$

② Patient clustering

~~Coordinates~~ Coordinates of  $v_i$  are features (height, weight, age, etc) associated with patients admitted to a hospital.  
Then clusters are groups of similar patients.

③ ~~Market~~ Market segmentation

We have  $n$  customers and  $d$  products.  
 $(v_i)_j = \# \$ \text{ customer } i \text{ spent on product } j$

Then clusters are market segments.

# Why is clustering difficult?

(26.2)

- ①  $d \gg 1$  (if  $d=2 \Rightarrow$  easy, use your eyes  $\therefore$ )
- ② Not clear what the best value of  $K$  is
- ③ Real data are not "cleanly" clustered.

## Problem Formalization

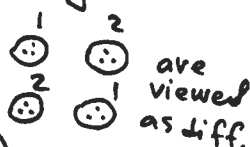
Assume  $K$  is given (how to choose  $K$  we will discuss at the end).  
Let's describe any clustering by a vector  $c \in \mathbb{R}^n$ , where

$c_i$  = cluster number that  $\vec{v}_i$  is assigned to ( $\in \{1, \dots, K\}$ )

Q: How many possible clusterings?

Anyway,  
the number  
is huge.

- $K^n$  (if clusters are "labeled")



- $\left\{ \begin{matrix} n \\ K \end{matrix} \right\} = \frac{1}{K!} \sum_{j=0}^K (-1)^{K-j} \binom{K}{j} j^n$  (if clusters are "unlabeled")

Stirling number of the second kind.

Example

$n=5$   
 $K=2$



$c = (1, 1, 2, 2, 1)$

~~For each clustering, we need to introduce a natural measure of the quality of the clustering.~~

To find a good clustering, we should be able to compare different clusterings. To compare clusterings, we need to introduce a natural measure of the clustering quality.

Idea: Let's associate with each cluster  $k$  a cluster representative  $r_k \in \mathbb{R}^d$ .  
(K-means) We want all vectors in cluster  $k$  to be close to the representative  $r_k$ :  
 $\downarrow$  can be any vector in  $\mathbb{R}^d$   
(not necessarily one of  $\vec{v}_1, \dots, \vec{v}_n$ )

we want  $\|\vec{v}_i - r_{c_i}\|$  to be small for all  $i=1, \dots, n$ .

Given clustering  $c$  and representatives  $r_1, \dots, r_K$ , we can measure the quality of  $(c; r_1, \dots, r_K)$  by the following objective function:

$$p(c; r_1, \dots, r_K) = \frac{1}{n} \sum_{i=1}^n \|\vec{v}_i - r_{c_i}\|^2$$

The smaller  $p$  is, the better the clustering.

Thus, we formalized the clustering problem to the following optimization problem.  
(reduced)

(\*) 
$$\begin{cases} p(c; r_1 \dots r_K) \rightarrow \min \\ r_1, \dots, r_K \in \mathbb{R}^d \\ c \in \mathbb{R}^n, c_i \in \{1, \dots, K\} \end{cases}$$

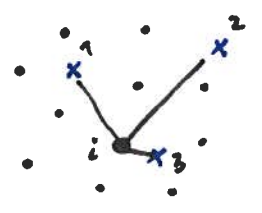
The solution of this problem  $c^*$  is called optimal clustering.

In practice, it is impossible to find the optimal clustering: the optimization problem is too difficult for large  $n$ .

The K-means algorithm allows to find a suboptimal clustering, i.e. clustering which is not optimal, but close to optimal,  $p_{K\text{-means}} \approx p_{\text{opt}}$   
K-means is based on the following idea which is very often used in numerical optimization.

① We can solve (\*) exactly if  $r_1, \dots, r_K$  are fixed.

$$p(c; r_1 \dots r_K) = \frac{1}{n} \sum_{i=1}^n \|v_i - r_{c_i}\|^2 \rightarrow \min$$



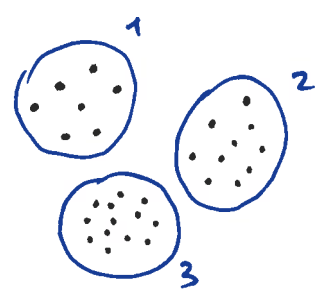
We need to find  $c = (c_1, \dots, c_n)$ , which minimizes the objective function. Key observation:  $c_i$  affects only  $i^{\text{th}}$  term.

$\Rightarrow$  to minimize  $p$ , we need to minimize each term independently:

$$f(c_i) = \|v_i - r_{c_i}\|^2 \rightarrow \min \Rightarrow \boxed{c_i = \arg \min_k \|v_i - r_k\|^2} \quad k=1 \dots K$$

② We can solve (\*) exactly if  $c$  is fixed.

$$p(c; r_1 \dots r_K) = \frac{1}{n} \sum_{i=1}^n \|v_i - r_{c_i}\|^2$$



$$= p_1 + p_2 + \dots + p_K$$

vectors associated with cluster  $k$

contribution from cluster  $k$

$$p_k = \frac{1}{n} \sum_{i: c_i=k} \|v_i - r_k\|^2$$

that is why squared distance is convenient

Representative  $r_k$  affects only  $p_k$   
 $\Rightarrow$  to minimize  $p$ , we need to minimize each  $p_k$  by choosing  $r_k$ :

$$p_k(r_k) \rightarrow \min$$

The solution is simple:

$r_k$  is the centroid (average) of vectors in cluster  $k$

$$\boxed{r_k = \frac{1}{n_k} \sum_{i: c_i=k} v_i}$$

$$n_k = \# \text{ vectors in cluster } k = |\{i: c_i=k\}|$$

# The K-means algorithm

26.4

So, we can minimize  $p(c, \{r\})$  over  $c$  for fixed  $\{r\}$  and over  $\{r\}$  for fixed  $c$ , but we can't minimize it over  $c$  and  $\{r\}$  simultaneously.

Standard solution: iterate between the two minimizations (popular strategy)



Algorithm:

- ① Pick initial representatives  $r_1, \dots, r_K$  randomly from the original vectors  $v_1 \dots v_n$

Repeat until convergence

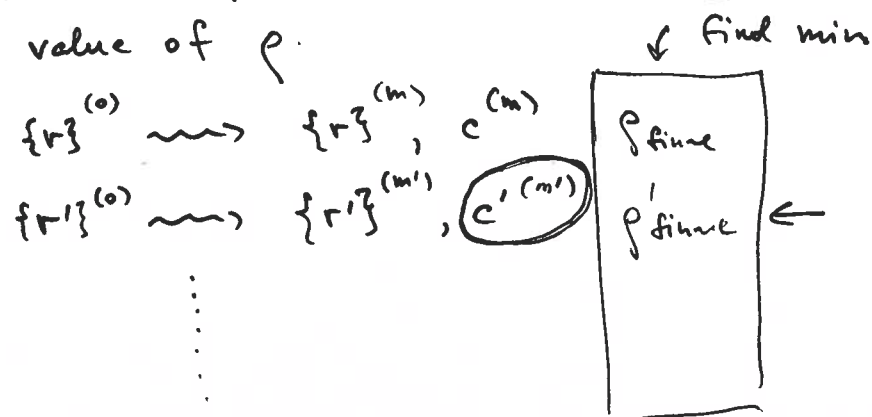
- ② Partition  $v_1 \dots v_n$  into  $K$  clusters  
Assign  $v_i$  to the cluster associated with the nearest  $r_k$

- ③ Update representatives  $r_1 \dots r_K$   
 $r_k = \text{mean of } v_i \text{ in cluster } k.$

There are more sophisticated methods, but this is beyond the scope of acm 104.

Remark 1 Let  $c^{(1)}, c^{(2)}, \dots$  be clusterings obtained in these iterations.  
If  $c^{(m)} = c^{(m-1)} \Rightarrow \{r\}^{(m)} = \{r\}^{(m-1)}$  and both  $c$  and  $\{r\}$  will not change in future ( $> m$ ) iterations.  
"Convergence" =  $\{c^{(m)} = c^{(m-1)}\}$

Remark 2 At each iteration, the value of  $p$  decreases  $\Rightarrow$  the k-means alg converges in a finite number of steps.  
Depending on the initial choice of  $r_1 \dots r_K$ , it may converge to different final clusterings with different values of  $p$ .  
 $\Rightarrow$  it is common to run the algorithm several times with different  $r_1 \dots r_K$  and choose clustering with the smallest value of  $p$ .



Remark 3 how to choose  $K$  ?

It is difficult to choose the right value of  $K$  in advance.

Common Strategy:

Run  $K$ -means with different values of  $K$ ,  
for each  $K$ , compute the corresponding minimized value  
of the objective function  $J_{\min}(K)$ . This is a  
decreasing function ( $J_{\min}(n) = 0$ )

Look for the following pattern

