

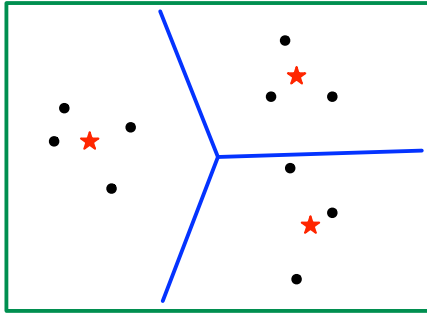
## The basics of clustering

**A:** *K*-means clustering

## The $k$ -means optimization problem

- Input: Points  $x_1, \dots, x_n \in \mathbb{R}^d$ ; integer  $k$
- Output: “Centers”, or representatives,  $\mu_1, \dots, \mu_k \in \mathbb{R}^d$
- Goal: Minimize average squared distance between points and their nearest representatives:

$$\text{cost}(\mu_1, \dots, \mu_k) = \sum_{i=1}^n \min_j \|x_i - \mu_j\|^2$$

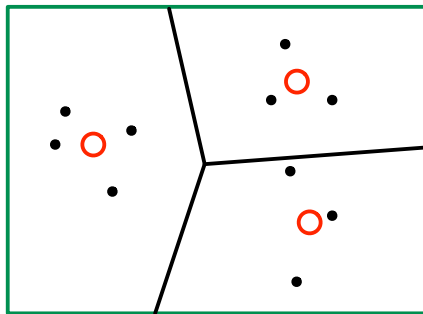


Centers carve  $\mathbb{R}^d$  into  $k$  **convex** regions:  $\mu_j$ 's region consists of points for which it is the closest center.

## Lloyd's $k$ -means algorithm

NP-hard optimization problem. Heuristic: “ $k$ -means algorithm”.

- Initialize centers  $\mu_1, \dots, \mu_k$  in some manner.
- Repeat until convergence:
  - Assign each point to its closest center.
  - Update each  $\mu_j$  to the mean of the points assigned to it.



Each iteration reduces the cost  $\Rightarrow$  convergence to a local optimum.

## Initialization matters



## Initializing the $k$ -means algorithm

Typical practice: choose  $k$  data points at random as the initial centers.

Another common trick: start with extra centers, then prune later.

A particularly good initializer:  **$k$ -means++**

- Pick a data point  $x$  at random as the first center
- Let  $C = \{x\}$  (centers chosen so far)
- Repeat until desired number of centers is attained:
  - Pick a data point  $x$  at random from the following distribution:

$$\Pr(x) \propto \text{dist}(x, C)^2,$$

where  $\text{dist}(x, C) = \min_{z \in C} \|x - z\|$

- Add  $x$  to  $C$

**B: Two uses of clustering**

## Two common uses of clustering

- **Vector quantization**  
Find a finite set of representatives that provides good coverage of a complex, possibly infinite, high-dimensional space.
- **Finding meaningful structure in data**  
Finding salient grouping in data.

## Representing images using $k$ -means codewords

How to represent a collection of images as fixed-length vectors?



- Take all  $\ell \times \ell$  patches in all images. Extract features for each.
- Run  $k$ -means on this entire collection to get  $k$  centers.
- Now associate any image patch with its nearest center.
- Represent an image by a histogram over  $\{1, 2, \dots, k\}$ .

# Looking for natural groups in data

## “Animals with attributes” data set

- 50 animals: antelope, grizzly bear, beaver, dalmatian, tiger, ...
- 85 attributes: longneck, tail, walks, swims, nocturnal, forager, desert, bush, plains, ...
- Each animal gets a score (0 – 100) along each attribute
- 50 data points in  $\mathbb{R}^{85}$

Apply  $k$ -means with  $k = 10$  and look at grouping obtained.

- 1 zebra
- 2 spider monkey, gorilla, chimpanzee
- 3 tiger, leopard, wolf, bobcat, lion
- 4 hippopotamus, elephant, rhinoceros
- 5 killer whale, blue whale, humpback whale, seal, walrus, dolphin
- 6 giant panda
- 7 skunk, mole, hamster, squirrel, rabbit, bat, rat, weasel, mouse, raccoon
- 8 antelope, horse, moose, ox, sheep, giraffe, buffalo, deer, pig, cow
- 9 beaver, otter
- 10 grizzly bear, dalmatian, persian cat, german shepherd, siamese cat, fox, chihuahua, polar bear, collie

- 1 zebra
- 2 spider monkey, gorilla, chimpanzee
- 3 tiger, leopard, fox, wolf, bobcat, lion
- 4 hippopotamus, elephant, rhinoceros, buffalo, pig
- 5 killer whale, blue whale, humpback whale, seal, otter, walrus, dolphin
- 6 dalmatian, persian cat, german shepherd, siamese cat, chihuahua, giant panda, collie
- 7 beaver, skunk, mole, squirrel, bat, rat, weasel, mouse, raccoon
- 8 antelope, horse, moose, ox, sheep, giraffe, deer, cow
- 9 hamster, rabbit
- 10 grizzly bear, polar bear

## C: Beyond $k$ -means

### $K$ -means: the good and the bad

#### The good:

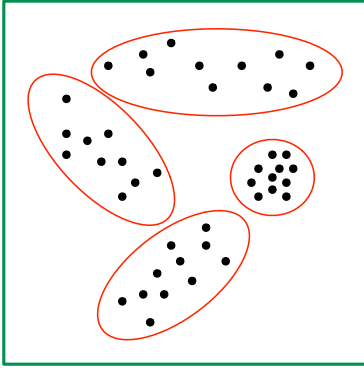
- Fast and easy.
- Effective in quantization.

#### The bad:

- Geared towards data in which the clusters are spherical, and of roughly the same radius.

Is there is a similarly-simple algorithm in which clusters of more general shape are accommodated?

## Preview: Mixtures of Gaussians



Each of the  $k$  clusters is specified by:

- a Gaussian distribution  $P_j = N(\mu_j, \Sigma_j)$
- a mixing weight  $\pi_j$

Overall distribution over  $\mathbb{R}^d$ : a **mixture of Gaussians**

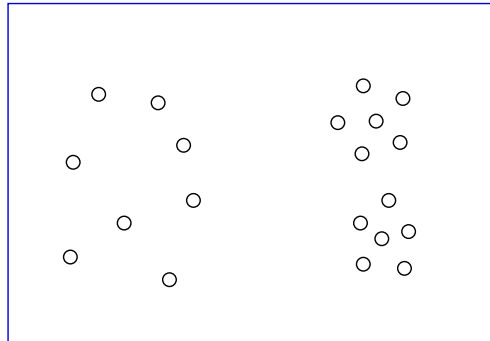
$$\Pr(x) = \pi_1 P_1(x) + \cdots + \pi_k P_k(x)$$

## D: Hierarchical clustering



# Hierarchical clustering

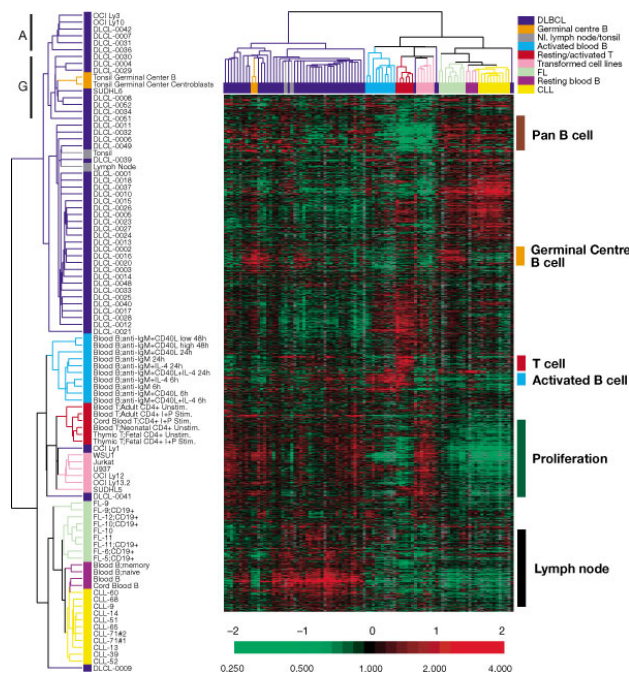
Choosing the number of clusters ( $k$ ) is difficult.



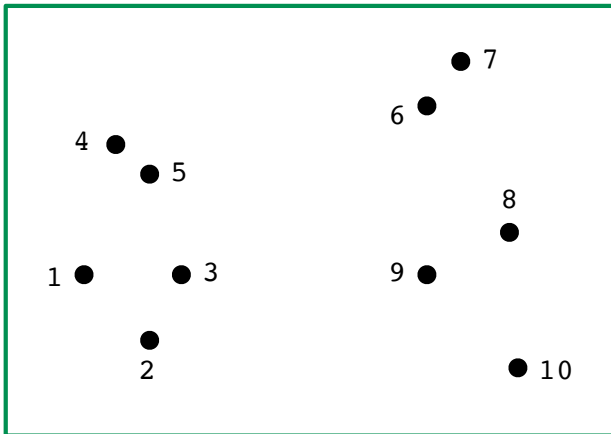
Often: no single right answer, because of multiscale structure.

Hierarchical clustering avoids these problems.

## Example: gene expression data



# The single linkage algorithm

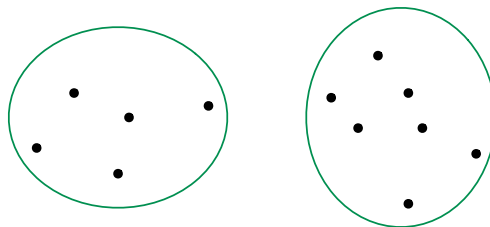


- Start with each point in its own, singleton, cluster
- Repeat until there is just one cluster:
  - Merge the two clusters with the closest pair of points
- Disregard singleton clusters

## Linkage methods

- Start with each point in its own, singleton, cluster
- Repeat until there is just one cluster:
  - Merge the two “closest” clusters

How to measure distance between two clusters  $C$  and  $C'$ ?



- Single linkage

$$\text{dist}(C, C') = \min_{x \in C, x' \in C'} \|x - x'\|$$

- Complete linkage

$$\text{dist}(C, C') = \max_{x \in C, x' \in C'} \|x - x'\|$$

## Average linkage

Three commonly-used variants:

- ① Average pairwise distance between points in the two clusters

$$\text{dist}(C, C') = \frac{1}{|C| \cdot |C'|} \sum_{x \in C} \sum_{x' \in C'} \|x - x'\|$$

- ② Distance between cluster centers

$$\text{dist}(C, C') = \|\text{mean}(C) - \text{mean}(C')\|$$

- ③ Ward's method: the increase in  $k$ -means cost occasioned by merging the two clusters

$$\text{dist}(C, C') = \frac{|C| \cdot |C'|}{|C| + |C'|} \|\text{mean}(C) - \text{mean}(C')\|^2$$