

# Assignment-Based Clustering

Note Title

2/10/2016

(class starting at 3:05)

k-means  
 $X \subset \mathbb{R}^d$

$X$  set of data points  
 $d: X \times X \rightarrow \mathbb{R}_+$

$$d(a,b) = \|a-b\|_2$$

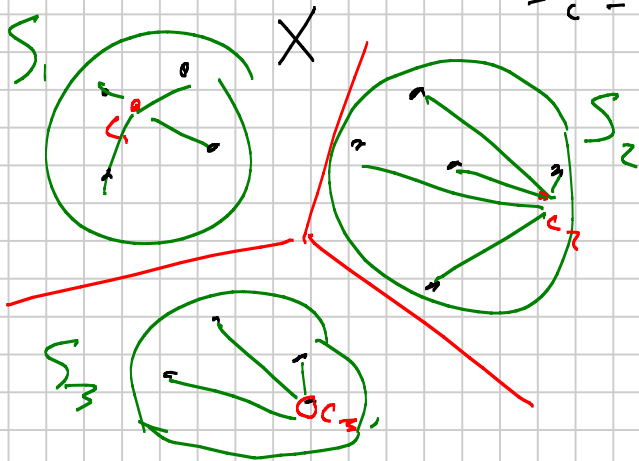
Goal: Find a set of k centers  $C$

Minimize a cost  $(C, X) = f(\text{cost}(c_i, S_i))$

$$\phi_c(x) = \arg \min_{c_i \in C} d(x, c_i)$$

k clusters  $S_1, \dots, S_k$

$$S_i = \{x \in X \mid \phi_c(x) = c_i\}$$



k-means

$$\begin{aligned} \text{cost}_2(C, X) &= \sum_{x \in X} d(x, \phi_c(x))^2 \\ &= \sum_{c_i \in C} \sum_{x \in S_i} d(x, c_i)^2 \end{aligned}$$

k-center

$$\text{cost}_\infty(C, X) = \max_{x \in X} d(x, \phi_c(x))$$

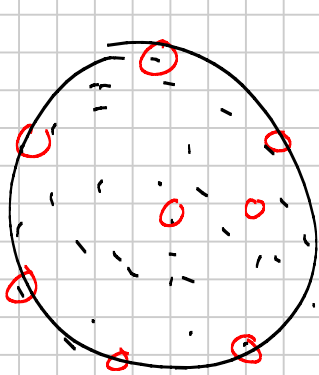
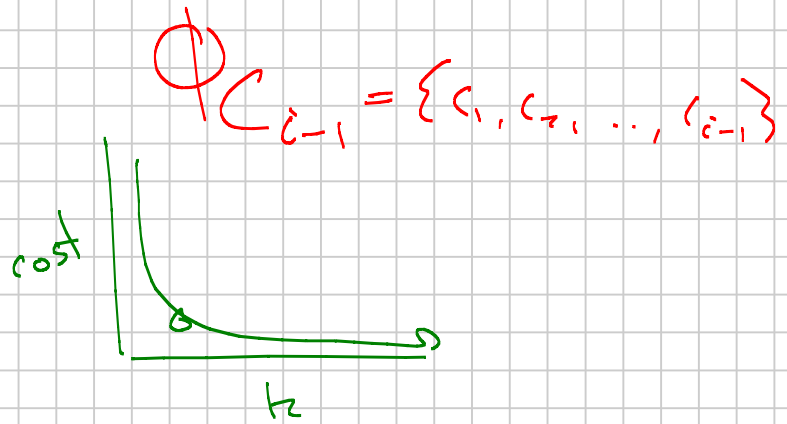
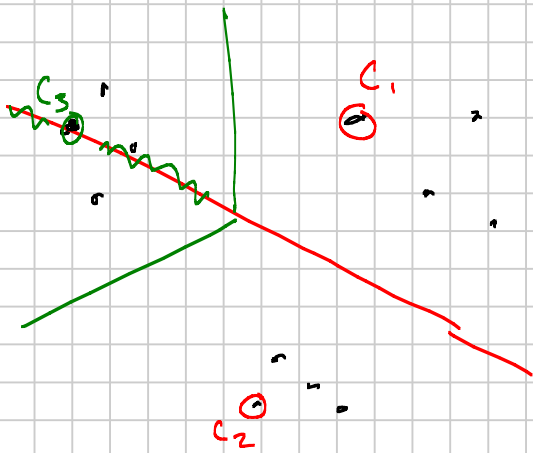
k-median

$$\text{cost}_1(C, X) = \sum_{x \in X} d(x, \phi_c(x))$$

k-medoid  $\text{cost}_1(C, X)$  s.t.  $c \in X$

# Algorithm for $k$ -center (Gonzalez)

1. Choose  $c_1 \in X$  arbitrarily  $C_1 = \{c_1\}$
2. for  $i=2$  to  $k$  do  $C_2 = \{c_1, c_2\}$
3. Set  $c_i = \arg \max_{x \in X} d(x, \phi_{C_{i-1}}(x))$   $C_i = \{c_1, c_2, \dots, c_i\}$



$k$ -approximation  
 $cost_{\infty}(X, c) \leq$

$2 cost_{\infty}(X, c^*)$   
 $\uparrow$   
 optimal center  
 NP-hard to  $k$ -approx

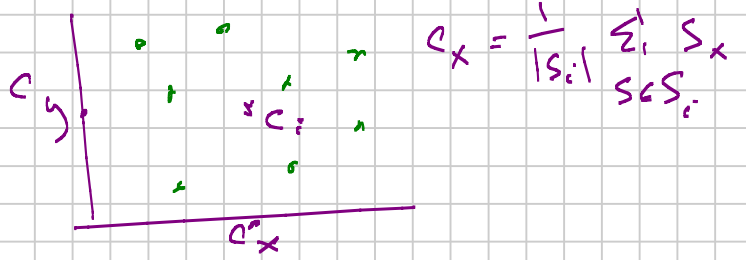
## Lloyd's Algorithm for $k$ -means

1957  $\rightarrow$  publish 1982

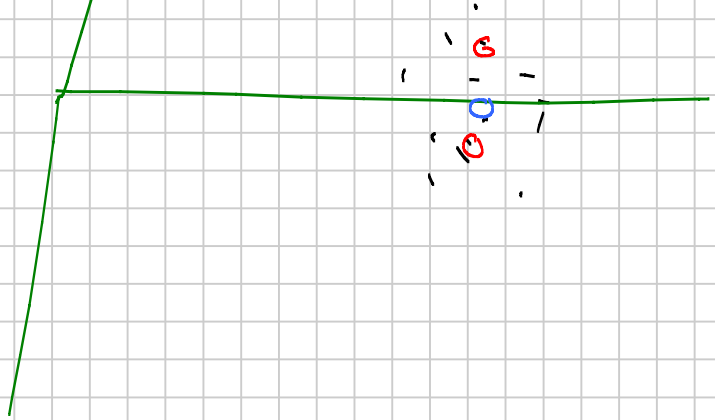
1. Choose  $k$  points  $C \subset X$  [arbitrarily?]
2. repeat
  - a) for  $x \in X$ , find  $\phi_C(x)$   $O(k \cdot n)$   $\leftarrow$  reassign to closest center
  - b) for  $i \in [k]$ ,  $c_i = \text{average}(S_i) = \{x \in X \mid \phi_C(x) = c_i\}$   $O(n)$
3. until (converged)  $\leftarrow$  if cost change  $\leq \tau$   $\leftarrow$  threshold

$$c_i = \text{average}\{S_i\}$$

$$\iff c_i = \operatorname{argmin}_{z \in \mathbb{R}^d} \sum_{S \in S_i} \|S - z\|^2$$



$$\text{cost}(X, C) = \sum_{x \in X} \|x - \phi_c(x)\|^2$$



$k^n$  possible clusters  
 $n$  data points

↑  
 $\{1, 2, 3, \dots, k\}$

How to choose initial  $C$  centers

- random (keep lowest cost solution after running Lloyd's)

• Coupon Collectors

Run Lloyd's w/  $k$  best centers

Regroup. w/ H/A

• Gonzalez  $k$ -center Algo

# k-means++

1. Choose  $c_1 \in X$  arbitrarily ( $C = \{c_1, c_2, \dots\}$ )
2. for  $i = 2$  to  $k$   
| Choose  $c_i$  from  $X$  with probability proportional to  $d(x, c_{i-1})^2$

$$M = \sum_{x \in X} (d(x, c_{i-1})^2 = m_x)$$

$$\text{Prob}[\text{choosing } x] = \frac{m_x}{M}$$

