# Assignment-based Clustering

## Feb 8, 2018

- K-center
- K-means
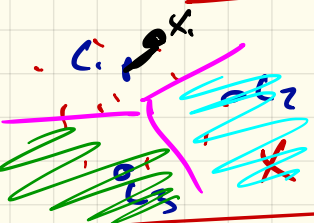- K-median/mediod
- K-means++
- Mixture of Gaussians (EM)

# Homework #1

| Largest setting | Q1D Birthday | Q2D coupon coll |
|---|---|---|
| min | 0.45 | 18 |
| mean | 300 ~5 minutes | 3000 ~1 hour |
| median | 30 | 1100 |
| max (seconds) | 15,000 ~4 hours | 80,000 ~22 hours |

$\|c_i - x_i\|$

**Input** • $X \subset \mathbb{R}^d = \{x_1, x_2 \ldots x_n\}$

• distance $d : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$

• $k = \#$ clusters

$$d(x_1, x_2) = \|x_1 - x_2\|$$

① **Outputs :** Centers $C = \{c_1, c_2 \ldots c_k\} \in \mathbb{R}^\theta$

Mapping $\phi_C : \mathbb{R}^d \to C$

• **k-means** formulation

$$\text{Cost}_2(X, C) = \sum_{x \in X} d(x, \phi_C(x))$$

② $\phi_C(x) = \arg\min_{c_i \in C} \|x - c_i\|$

• **k-median** formulation

$$\text{Cost}_1(X, C) = \sum_{x \in X}^{'} d(x, \phi_C(x))$$

k-mediod ← same but $C \subset X$

**k-center**

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

# Gonzalez Alg. for $k$-center Clustering

- NP-hard to solve w/in factor 2 of OPT

- Gonz Alg : 2-apx OPT, in metric $d$.
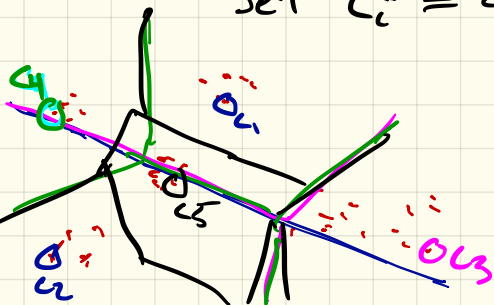
1. Choose center $c_i \in X$ arbitrarily
   Let $C_1 = \{c_1\}$      $C_i = \{c_1, c_2, \dots c_i\}$

2. for $i = 2$ to $k$ do

   Set $c_i = \arg\max_{x \in X} d\left(x, \phi_{C_{i-1}}(x)\right)$



$\boxed{C_1 \; C_2 \dots C_k}$

| data | | $x_1$ | $x_2$ | | | $x_n$ |
|---|---|---|---|---|---|---|
| | | | | | | |

| assign | 1 | 2 | | 3 | | 2 |
|---|---|---|---|---|---|---|

# Lloyd's Alg. for k-means

$D =$ Euclidean $X \subset \mathbb{R}^d$

1. **Choose** k centers $C \subset X$ (arbitrarily)

2. **repeat**

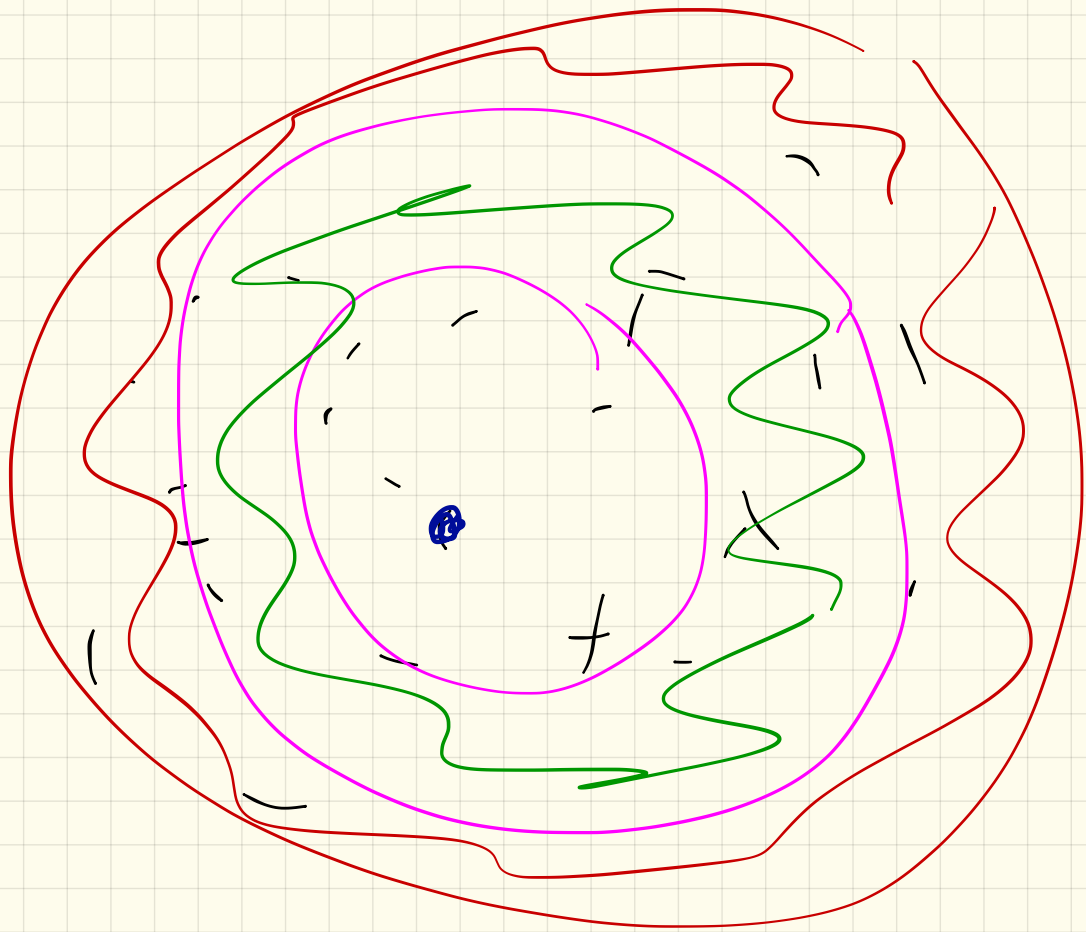both a,b    a. For all $x \in X$, set $x \to c_i = \phi_C(x) = \underset{c_i \in C}{\arg\min} \|x - c_i\|_d$

$\|x - c_i\|$

Cost decreases    b. For all $c_i \in C$, set $c_i = \text{average} \left\{ x \in X \mid \phi_C(x) = c_i \right\}$

$= S_i$

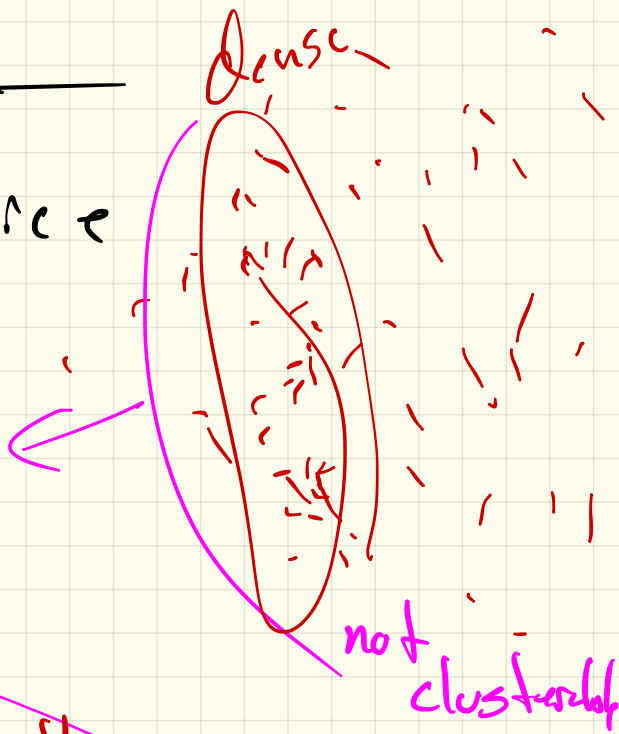$c_i = \frac{1}{|S_i|} \underset{x \in S_i}{\Sigma} x$    "such that"

**until** "converged"

Voronoi diagram

# Choosing $k$

- Modeling choice

$Var(X)$

$Cost_k(X,C)$
of
optimal $C$

elbow

well
clusterable

$k=1$   $k^*$   $k$   $k=n$

Dense

not
clusterable

$C_5$   $C_6$   $C_7$

$C_1$

local minimum
of
Cost2

if happens

↳ random
restart.

$C_2$   $C_4$
$C_3$   $C_1$
$C_3$

# How to choose initial centers?

- k-center (Gonz Alg)

- Choose $O(k \log k)$ centers.
    "coupon collectors"
    $\hookrightarrow$ then cluster
        then merge.

    ( mil

- k-means ++

$c_7$ fpt

far

$c_2$

$c_i$
$c_i$

1 mil

# k-means++ Algorithm "$D^2$-sampling"

1. Choosing $c_i \in X$ arbitrarily

$$\zeta_i = \{c_1, c_2, \dots c_i\}$$

2. for $i = 2$ to $k$

Choose $c_i \in X$ w/ probability proportional

$$w_x = D(x, \phi_{c_i}(x))^2$$

$W = \sum_{x \in X} w_x$

$P_x = \frac{w_x}{W}$

| | 0 | 0.09 | 0.18 | 0.23 | |
|---|---|---|---|---|---|
| | $P_1$ | $P_2$ | $P_3$ | $\cdots$ | |

0          ↑ 0.21          1