

These lecture notes are based on a set of slides written by Marco Bressan in 2023.

Recall the k -means problem: given a set $\mathcal{X} \subset \mathbb{R}^d$ of size n and $1 < k < n$, find

$$\mathcal{C}^* \in \operatorname{argmin}_{\mathcal{C}_1, \dots, \mathcal{C}_k \in \mathbb{R}^d} \Phi(\mathcal{C}_1, \dots, \mathcal{C}_k)$$

where, for any $\mathcal{C} \subset \mathbb{R}^d$,

$$\Phi(\mathcal{C}) = \sum_{\mathbf{x} \in \mathcal{X}} \phi(\mathcal{C}, \mathbf{x}) = \sum_{\mathbf{x} \in \mathcal{X}} \min_{\mathbf{c}_i \in \mathcal{C}} \|\mathbf{x} - \mathbf{c}_i\|^2$$

Let $\text{OPT} = \Phi(\mathcal{C}^*)$ and, for any $\mathcal{C} \subset \mathbb{R}^d$ and $A \subseteq \mathcal{X}$, let

$$\phi(\mathcal{C}, A) = \sum_{\mathbf{x} \in A} \phi(\mathcal{C}, \mathbf{x})$$

We identify a clustering \mathcal{C} via its centers $\{\mathbf{c}_1, \dots, \mathbf{c}_k\}$ or with its clusters $\{C_1, \dots, C_k\}$. Note that, for any clustering \mathcal{C} output by Lloyd's algorithm, including the optimal clustering \mathcal{C}^* ,

$$\phi(\mathcal{C}, C) = \sum_{\mathbf{x} \in C} \|\mathbf{x} - \boldsymbol{\mu}_C\|^2 \quad \text{for all } C \in \mathcal{C}, \text{ where } \boldsymbol{\mu}_C \text{ is the centroid of } C. \quad (1)$$

We proved that Lloyd's algorithm does not have any approximation guarantee because, while outliers can contribute a lot to the overall cost, they are not favored in the initial random draw of centers.

Algoritmo 1 **k-means++**

Input: Finite set of points $\mathcal{X} \subset \mathbb{R}^d$, integer $1 < k < |\mathcal{X}|$.

1: Draw a center \mathbf{c}_1 u.a.r. from \mathcal{X} and let $\mathcal{C}_1 = \{\mathbf{c}_1\}$

2: **for** $i = 2, \dots, k$ **do**

3: draw \mathbf{c}_i from \mathcal{X} according to the distribution $\mathbb{P}(\mathbf{c}_i = \mathbf{x} \mid \mathcal{C}_{i-1}) = \frac{\phi(\mathcal{C}_{i-1}, \mathbf{x})}{\Phi(\mathcal{C}_{i-1})}$

4: $\mathcal{C}_i = \mathcal{C}_{i-1} \cup \{\mathbf{c}_i\}$

5: **end for**

Output: The output of Lloyd's algorithm run with initial centers $\mathbf{c}_1, \dots, \mathbf{c}_k$

We prove a simplified version of the following theorem.

Teorema 1 *The clustering \mathcal{C} found by **k-means++** satisfies $\mathbb{E}[\Phi(\mathcal{C})] \leq 8(\ln k + 2)\text{OPT}$.*

Note that the currently best known approximation algorithms for k -means is based on a linear programming rounding approach and produces a clustering with a cost $c \times \text{OPT}$ where $c \in [6, 7]$.

Consider any optimal clustering $\mathcal{C}^* = (A_1, \dots, A_k)$ and let \mathcal{C}_i be the clustering of **k-means++** after drawing the first i centers in Line 3.

Lemma 2 For any $A \in \mathcal{C}^*$ and for any $i \in [k]$,

$$\mathbb{E}\left[\phi(\mathcal{C}_i, A) \mid \mathbf{c}_i \in A, \mathcal{C}_{i-1}\right] \leq 8\phi(\mathcal{C}^*, A)$$

DIMOSTRAZIONE. Consider first $i = 1$. Then $\mathcal{C}_{i-1} = \mathcal{C}_0 = \emptyset$ and c_i is drawn according to the uniform distribution over \mathcal{X} , and we can write

$$\begin{aligned} \mathbb{E}[\phi(\mathcal{C}_1, A) \mid \mathbf{c}_1 \in A] &= \frac{1}{|A|} \sum_{\mathbf{a} \in A} \left(\sum_{\mathbf{x} \in A} \|\mathbf{x} - \mathbf{a}\|^2 \right) && (\mathcal{C}_1 = \{\mathbf{c}_1\}) \\ &\leq \frac{1}{|A|} \sum_{\mathbf{a} \in A} \left(\sum_{\mathbf{x} \in A} \|\mathbf{x} - \boldsymbol{\mu}\|^2 + |A| \|\mathbf{a} - \boldsymbol{\mu}\|^2 \right) && (\boldsymbol{\mu} \text{ is the centroid of } A) \\ &= \sum_{\mathbf{x} \in A} \|\mathbf{x} - \boldsymbol{\mu}\|^2 + \sum_{\mathbf{a} \in A} \|\mathbf{a} - \boldsymbol{\mu}\|^2 \\ &= 2 \sum_{\mathbf{x} \in A} \|\mathbf{x} - \boldsymbol{\mu}\|^2 = 2\phi(\mathcal{C}^*, A) && (\text{because of (1).}) \end{aligned}$$

In particular, note that

$$\frac{1}{|A|} \sum_{\mathbf{a} \in A} \sum_{\mathbf{x} \in A} \|\mathbf{x} - \mathbf{a}\|^2 \leq 2\phi(\mathcal{C}^*, A) \quad (2)$$

Now assume $i > 1$. Then

$$\mathbb{P}(\mathbf{c}_i = \mathbf{a} \mid \mathbf{a} \in A, \mathcal{C}_{i-1}) = \frac{\phi(\mathcal{C}_{i-1}, \mathbf{a})}{\sum_{\mathbf{x} \in A} \phi(\mathcal{C}_{i-1}, \mathbf{x})}$$

For any $\mathbf{x}, \mathbf{a} \in A$, let \mathbf{c} be the center of \mathcal{C}_{i-1} closest to \mathbf{x} . Then

$$\begin{aligned} \min_{j=1, \dots, i-1} \|\mathbf{a} - \mathbf{c}_j\| &\leq \|\mathbf{a} - \mathbf{c}\| \\ &\leq \|\mathbf{x} - \mathbf{c}\| + \|\mathbf{a} - \mathbf{x}\| && (\text{by the triangular inequality.}) \end{aligned}$$

Using $(a + b)^2 \leq 2(a^2 + b^2)$ for all $a, b \in \mathbb{R}$ and $\|\mathbf{x} - \mathbf{c}\|^2 = \phi(\mathcal{C}_{i-1}, \mathbf{x})$ we get

$$\phi(\mathcal{C}_{i-1}, \mathbf{a}) \leq 2\left(\phi(\mathcal{C}_{i-1}, \mathbf{x}) + \|\mathbf{a} - \mathbf{x}\|^2\right)$$

By averaging the above inequality over all $\mathbf{x} \in A$, we get

$$\phi(\mathcal{C}_{i-1}, \mathbf{a}) \leq \frac{2}{|A|} \sum_{\mathbf{x} \in A} \left(\phi(\mathcal{C}_{i-1}, \mathbf{x}) + \|\mathbf{a} - \mathbf{x}\|^2 \right)$$

Note also that, for any $\mathbf{x} \in \mathcal{X}$,

$$\phi(\mathcal{C}_i, \mathbf{x}) = \min \left\{ \phi(\mathcal{C}_{i-1}, \mathbf{x}), \|\mathbf{x} - \mathbf{c}_i\|^2 \right\}$$

Therefore, for $\mathcal{C}_i = \mathcal{C}_{i-1} \cup \{\mathbf{c}_i\}$,

$$\begin{aligned}
\mathbb{E}[\phi(\mathcal{C}_i, A) \mid \mathbf{c}_i \in A, \mathcal{C}_{i-1}] &= \sum_{\mathbf{a} \in A} \frac{\phi(\mathcal{C}_{i-1}, \mathbf{a})}{\sum_{\mathbf{x} \in A} \phi(\mathcal{C}_{i-1}, \mathbf{x})} \phi(\mathcal{C}_i, A) \\
&\leq \frac{2}{|A|} \sum_{\mathbf{a} \in A} \sum_{\mathbf{x} \in A} \frac{\phi(\mathcal{C}_{i-1}, \mathbf{x}) + \|\mathbf{a} - \mathbf{x}\|^2}{\sum_{\mathbf{x}' \in A} \phi(\mathcal{C}_{i-1}, \mathbf{x}')} \sum_{\mathbf{a}' \in A} \min \left\{ \phi(\mathcal{C}_{i-1}, \mathbf{a}'), \|\mathbf{a}' - \mathbf{a}\|^2 \right\} \\
&= \frac{2}{|A|} \sum_{\mathbf{a} \in A} \frac{\sum_{\mathbf{x} \in A} \phi(\mathcal{C}_{i-1}, \mathbf{x})}{\sum_{\mathbf{x}' \in A} \phi(\mathcal{C}_{i-1}, \mathbf{x}')} \sum_{\mathbf{a}' \in A} \min \left\{ \phi(\mathcal{C}_{i-1}, \mathbf{a}'), \|\mathbf{a}' - \mathbf{a}\|^2 \right\} \\
&\quad + \frac{2}{|A|} \sum_{\mathbf{a} \in A} \sum_{\mathbf{x} \in A} \frac{\|\mathbf{a} - \mathbf{x}\|^2}{\sum_{\mathbf{x}' \in A} \phi(\mathcal{C}_{i-1}, \mathbf{x}')} \sum_{\mathbf{a}' \in A} \min \left\{ \phi(\mathcal{C}_{i-1}, \mathbf{a}'), \|\mathbf{a}' - \mathbf{a}\|^2 \right\} \\
&\leq \frac{2}{|A|} \sum_{\mathbf{a} \in A} \sum_{\mathbf{a}' \in A} \|\mathbf{a}' - \mathbf{a}\|^2 + \frac{2}{|A|} \sum_{\mathbf{a} \in A} \sum_{\mathbf{x} \in A} \|\mathbf{a} - \mathbf{x}\|^2 \\
&= \frac{4}{|A|} \sum_{\mathbf{a} \in A} \sum_{\mathbf{x} \in A} \|\mathbf{x} - \mathbf{a}\|^2 \\
&\leq 8 \phi(\mathcal{C}^*, A) \tag{because of (2).}
\end{aligned}$$

concluding the proof. \square

A cluster $A \in \mathcal{C}^*$ is uncovered in \mathcal{C}_i if $A \cap \{\mathbf{c}_1, \dots, \mathbf{c}_i\} = \emptyset$. Lemma 2 shows that we pay $\mathcal{O}(\text{OPT})$ for every optimal cluster that we cover. This justifies the following simplifying assumptions, stating that the cost of each optimal cluster is set to 1, and we pay 1 for each optimal cluster that is covered and L for each optimal cluster that remains uncovered.

Assunzione 3 *For all $A \in \mathcal{C}^*$:*

1. $\phi(\mathcal{C}^*, A) = 1$
2. *for all $i \in [k]$, if A is covered in \mathcal{C}_i , then $\phi(\mathcal{C}_i, A) = 1$; otherwise, $\phi(\mathcal{C}_i, A) = L$.*

Lemma 4 *Under the above assumptions, $\mathbb{E}[\Phi(\mathcal{C})] \leq (2 + \ln k)\text{OPT}$.*

DIMOSTRAZIONE. Let $\mathcal{C}_i = (\mathbf{c}_1, \dots, \mathbf{c}_i)$. Conventionally, $\mathcal{C}_0 = \emptyset$ and $\Phi(\mathcal{C}_0) = kL$ (as if there were a default faraway center). Now, observing that $\mathcal{C} = \mathcal{C}_k$,

$$\Phi(\mathcal{C}_k) = \Phi(\mathcal{C}_0) + \sum_{i=0}^{k-1} (\Phi(\mathcal{C}_{i+1}) - \Phi(\mathcal{C}_i))$$

Taking expectations,

$$\begin{aligned}
\mathbb{E}[\Phi(\mathcal{C}_k)] &= \Phi(\mathcal{C}_0) + \sum_{i=0}^{k-1} (\mathbb{E}[\Phi(\mathcal{C}_{i+1})] - \mathbb{E}[\Phi(\mathcal{C}_i)]) \\
&= kL + \sum_{i=0}^{k-1} (\mathbb{E}[\Phi(\mathcal{C}_{i+1})] - \mathbb{E}[\Phi(\mathcal{C}_i)]) \\
&= k + \sum_{i=0}^{k-1} ((L-1) + \mathbb{E}[\Phi(\mathcal{C}_{i+1})] - \mathbb{E}[\Phi(\mathcal{C}_i)])
\end{aligned}$$

Let N_i the number of uncovered clusters in \mathcal{C}_i . Because of our assumptions, $\Phi(\mathcal{C}_i) = N_iL + (k - N_i)$.

For any uncovered A , the probability that at round $i+1$ we choose a center from A is

$$\mathbb{P}(\mathbf{c}_{i+1} \in A \mid \mathcal{C}_i) = \frac{\phi(\mathcal{C}_i, A)}{\Phi(\mathcal{C}_i)} = \frac{L}{N_iL + (k - N_i)}$$

So the probability p_{i+1} that we choose a center from some uncovered cluster is:

$$\mathbb{P}(\exists A \in \mathcal{C}^* : \mathbf{c}_{i+1} \in A \wedge A \cap \{\mathbf{c}_1, \dots, \mathbf{c}_i\} = \emptyset \mid \mathcal{C}_i) = \frac{N_iL}{N_iL + (k - N_i)} \geq \frac{(k-i)L}{(k-i)L + i}$$

where in the last inequality we used $N_i \geq k - i$.

Now, if \mathbf{c}_{i+1} does not cover any A that was uncovered in \mathcal{C}_i (which happens with probability $1 - p_{i+1}$), then $\Phi(\mathcal{C}_{i+1}) \leq \Phi(\mathcal{C}_i)$. On the other hand, if \mathbf{c}_{i+1} covers some A that was uncovered in \mathcal{C}_i (which happens with probability p_{i+1}), then $\Phi(\mathcal{C}_{i+1}) = \Phi(\mathcal{C}_i) - L + 1 = \Phi(\mathcal{C}_i) - (L-1)$. Therefore

$$\begin{aligned}
(L-1) + \mathbb{E}[\Phi(\mathcal{C}_{i+1}) \mid \mathcal{C}_i] - \mathbb{E}[\Phi(\mathcal{C}_i) \mid \mathcal{C}_i] &\leq (L-1) + 0 \times (1 - p_{i+1}) - (L-1)p_{i+1} \\
&\leq (L-1) - (L-1) \frac{(k-i)L}{(k-i)L + i} \\
&= (L-1) \left(\frac{i}{(k-i)L + i} \right) \\
&< L \frac{i}{(k-i)L + i} \\
&< L \frac{k}{(k-i)L} = \frac{k}{k-i}
\end{aligned}$$

Therefore,

$$\begin{aligned}
\mathbb{E}[\Phi(\mathcal{C}_k)] &= k + \sum_{i=0}^{k-1} ((L-1) + \mathbb{E}[\Phi(\mathcal{C}_{i+1})] - \mathbb{E}[\Phi(\mathcal{C}_i)]) \\
&= k + \sum_{i=0}^{k-1} \mathbb{E}[(L-1) + \mathbb{E}[\Phi(\mathcal{C}_{i+1}) \mid \mathcal{C}_i] - \mathbb{E}[\Phi(\mathcal{C}_i) \mid \mathcal{C}_i]] \\
&\leq k + \sum_{i=0}^{k-1} \frac{k}{k-i} \\
&= k + k \sum_{i=1}^k \frac{1}{i} \leq k(2 + \ln k)
\end{aligned}$$

where we used the bound on the harmonic sum $1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{k} \leq 1 + \ln k$. The proof is concluded by noticing that, under our assumptions, $\text{OPT} = \Phi(\mathcal{C}^*) = k$. \square