Complementi di Algoritmi e Strutture Dati

k-Means++

Docente: Nicolò Cesa-Bianchi

versione 25 maggio 2025

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

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

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

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

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

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

We identify a clustering C via its centers $\{c_1, \ldots, c_k\}$ or with its clusters $\{C_1, \ldots, C_k\}$. Note that, for any clustering C output by Lloyd's algorithm, including the optimal clustering C^* ,

$$\phi(\mathcal{C}, C) = \sum_{\boldsymbol{x} \in C} \|\boldsymbol{x} - \boldsymbol{\mu}_C\|^2 \quad \text{for all } C \in \mathcal{C}, \text{ where } \boldsymbol{\mu}_C \text{ is the centroid of } C.$$
(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 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 c_1 u.a.r. from \mathcal{X} and let $\mathcal{C}_1 = \{c_1\}$ 2: for i = 2, ..., k do 3: draw c_i from \mathcal{X} according to the distribution $\mathbb{P}(c_i = \boldsymbol{x} \mid \mathcal{C}_{i-1}) = \frac{\phi(\mathcal{C}_{i-1}, \boldsymbol{x})}{\Phi(\mathcal{C}_{i-1})}$ 4: $\mathcal{C}_i = \mathcal{C}_{i-1} \cup \{c_i\}$ 5: end for Output: The output of Lloyd's algorithm run with initial centers $c_1, ..., c_k$

We prove a simplified version of the following theorem.

Teorema 1 The clustering C found by *k*-means++ satisfies $\mathbb{E}[\Phi(C)] \leq 8(\ln k + 2)$ 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 $C^* = (A_1, \ldots, A_k)$ and for let C_i be the clustering of k-means++ after drawing the first *i* centers in Line 3.

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

$$\mathbb{E}\Big[\phi(\mathcal{C}_i, A) \,\Big|\, \boldsymbol{c}_i \in A, \mathcal{C}_{i-1}\Big] \leq 8\,\phi(\mathcal{C}^*, A)$$

DIMOSTRAZIONE. Consider first i = 1. Then $C_{i-1} = C_0 = \emptyset$ and c_i is drawn according to the uniform distribution over \mathcal{X} , and we can write

$$\mathbb{E}[\phi(\mathcal{C}_{1}, A) \mid c_{1} \in A] = \frac{1}{|A|} \sum_{a \in A} \left(\sum_{x \in A} \|x - a\|^{2} \right) \qquad (\mathcal{C}_{1} = \{c_{1}\})$$

$$\leq \frac{1}{|A|} \sum_{a \in A} \left(\sum_{x \in A} \|x - \mu\|^{2} + |A| \|a - \mu\|^{2} \right) \qquad (\mu \text{ is the centroid of } A)$$

$$= \sum_{x \in A} \|x - \mu\|^{2} + \sum_{a \in A} \|a - \mu\|^{2}$$

$$= 2 \sum_{x \in A} \|x - \mu\|^{2} = 2 \phi(\mathcal{C}^{*}, A) \qquad (\text{because of } (1).)$$

In particular, note that

$$\frac{1}{|A|} \sum_{\boldsymbol{a} \in A} \sum_{\boldsymbol{x} \in A} \|\boldsymbol{x} - \boldsymbol{a}\|^2 \le 2 \,\phi(\mathcal{C}^*, A) \tag{2}$$

Now assume i > 1. Then

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

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

$$\min_{j=1,...,i-1} \| oldsymbol{a} - oldsymbol{c}_j \| \le \| oldsymbol{a} - oldsymbol{c} \|$$

 $\le \| oldsymbol{x} - oldsymbol{c} \| + \| oldsymbol{a} - oldsymbol{x} \|$ (by the triangular inequality.)

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

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

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

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

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

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

Therefore, for $C_i = C_{i-1} \cup \{c_i\},\$

$$\mathbb{E}\left[\phi(\mathcal{C}_{i},A) \mid \mathbf{c}_{i} \in A, \mathcal{C}_{i-1}\right] = \sum_{a \in A} \frac{\phi(\mathcal{C}_{i-1},a)}{\sum_{x \in A} \phi(\mathcal{C}_{i-1},x)} \phi(\mathcal{C}_{i},A)$$

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

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

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

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

$$\leq \frac{4}{|A|} \sum_{a \in A} \sum_{x \in A} \|a' - a\|^{2} + \frac{2}{|A|} \sum_{a \in A} \sum_{x \in A} \|a - x\|^{2}$$

$$\leq 8 \phi(\mathcal{C}^{*}, A) \qquad (\text{because of } (2).)$$

concluding the proof.

A cluster $A \in \mathcal{C}^*$ is uncovered in \mathcal{C}_i if $A \cap \{c_1, \ldots, c_i\} = \emptyset$. 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.

Assumption 3 For all $A \in C^*$:

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

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

DIMOSTRAZIONE. Let $C_i = (c_1, \ldots, c_i)$. Conventionally, $C_0 = \emptyset$ and $\Phi(C_0) = kL$ (as if there were a default faraway center). Now, observing that $C = C_k$,

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

Taking expectations,

$$\mathbb{E}\left[\Phi(\mathcal{C}_{k})\right] = \Phi(\mathcal{C}_{0}) + \sum_{i=0}^{k-1} \left(\mathbb{E}\left[\Phi(\mathcal{C}_{i+1})\right] - \mathbb{E}\left[\Phi(\mathcal{C}_{i})\right]\right)$$
$$= kL + \sum_{i=0}^{k-1} \left(\mathbb{E}\left[\Phi(\mathcal{C}_{i+1})\right] - \mathbb{E}\left[\Phi(\mathcal{C}_{i})\right]\right)$$
$$= k + \sum_{i=0}^{k-1} \left((L-1) + \mathbb{E}\left[\Phi(\mathcal{C}_{i+1})\right] - \mathbb{E}\left[\Phi(\mathcal{C}_{i})\right]\right)$$

Let N_i the number of uncovered clusters in C_i . Because of our assumptions, $\Phi(C_i) = N_i L + (k - N_i)$. For any uncovered A, the probability that at round i + 1 we choose a center from A is

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

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

$$\mathbb{P}\Big(\exists A \in \mathcal{C}^* : c_{i+1} \in A \land A \cap \{c_1, \dots, c_i\} = \emptyset \mid \mathcal{C}_i\Big) = \frac{N_i L}{N_i L + (k - N_i)} \ge \frac{(k - i)L}{(k - i)L + i}$$

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

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

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

Therefore,

$$\mathbb{E}\left[\Phi(\mathcal{C}_{k})\right] = k + \sum_{i=0}^{k-1} \left((L-1) + \mathbb{E}\left[\Phi(\mathcal{C}_{i+1})\right] - \mathbb{E}\left[\Phi(\mathcal{C}_{i})\right] \right)$$
$$= k + \sum_{i=0}^{k-1} \mathbb{E}\left[(L-1) + \mathbb{E}\left[\Phi(\mathcal{C}_{i+1}) \mid \mathcal{C}_{i}\right] - \mathbb{E}\left[\Phi(\mathcal{C}_{i}) \mid \mathcal{C}_{i}\right] \right]$$
$$\leq k + \sum_{i=0}^{k-1} \frac{k}{k-i}$$
$$= k + k \sum_{i=1}^{k} \frac{1}{i} \leq k(2 + \ln k)$$

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