Scalable K-Means++

Steven Wu

University of Pennsylvania

October 3, 2014

Paper by Bahmani, Moseley, Vattani, Kumar and Vassilvitskii

Clustering is very important!



Clustering is very important!



- Fundamental problem in data analysis and machine learning
- "Frequently asked interview question by Big Data tech firms"
 - Yaroslavtsev

Steven Wu (Penn)

Clustering is very important!



- Fundamental problem in data analysis and machine learning
- "Frequently asked interview question by Big Data tech firms"
 - Yaroslavtsev

Steven Wu (Penn)

k-means Objective

Let $X = \{x_1, \ldots, x_n\}$ be a set of points in *d*-dimensional space, find a set of centers $C = \{c_1, \ldots, c_k\}$ to minimize

$$\sum_{k \in X} \min_{i \in [k]} \|x - c_i\|^2$$

3

k-means Objective

Let $X = \{x_1, \ldots, x_n\}$ be a set of points in *d*-dimensional space, find a set of centers $C = \{c_1, \ldots, c_k\}$ to minimize

$$\sum_{i \in X} \min_{i \in [k]} \|x - c_i\|^2$$

r

NP-Hard Problem

k-means (sometimes called Lloyd's algorithm)

- Initialization: Start with a set of randomly chosen initial centers
- Repeat:
 - Assign each point to its nearest center;
 - Recompute the center given the point assignment
- Until convergence

 $k\mbox{-means}$ algorithm not very appealing

- Efficiency: run time can be exponential in the worst case
- Quality: final solution is locally optimal, but far away from global optimum

In practice

wetit

A scalable k-means algorithm with

- theoretical guarantee
- good practical performance

A scalable k-means algorithm with

- theoretical guarantee
- good practical performance



A better way to initialize the clustering dramatically changes the performance of the algorithm



Potential Problem: Sensitive to Initialization



Potential Problem: Sensitive to Initialization



Potential Problem: Sensitive to Initialization



Stuck in local optimum Photo credited to David Arthur

Intuition: Spread Out





- First center is selected uniformly at random from the data
- Subsequent centers: each point is selected with probability

$$\frac{d^2(x,\mathcal{C})}{\sum_x d^2(x,\mathcal{C})}$$

- First center is selected uniformly at random from the data
- Subsequent centers: each point is selected with probability

$$\frac{d^2(x,\mathcal{C})}{\sum_x d^2(x,\mathcal{C})}$$

proportional to its contribution to the overall error given the previous selections:

The initialization step itself already obtains an $8\log k$ approximation to OPT in expectation.

The initialization step itself already obtains an $8\log k$ approximation to OPT in expectation.

(The Lloyd iterations would only make it better)

The initialization step itself already obtains an $8\log k$ approximation to OPT in expectation.

(The Lloyd iterations would only make it better)

Disadvantage

Not scalable!

The initialization step itself already obtains an $8 \log k$ approximation to OPT in expectation.

(The Lloyd iterations would only make it better)

Disadvantage

Not scalable!

Sequential nature: the choice of the next center depends on the current set of centers.

The initialization step itself already obtains an $8 \log k$ approximation to OPT in expectation. (The Lloyd iterations would only make it better)

Disadvantage

Not scalable!

Sequential nature: the choice of the next center depends on the current set of centers. k passes over the data (think of $k=1\,000$)

- Fewer number of iterations (sample more than 1 points each round)
- Provable approximation guarantee

- 2 Initial cost $\psi = \sum_x d^2(x, \mathcal{C})$

- **①** First center C: sample a point uniformly at random
- 2 Initial cost $\psi = \sum_x d^2(x, \mathcal{C})$
- 3 for $O(\log \psi)$ times do
 - $\mathcal{C}' \leftarrow$ sample each point $x \in X$ independently with probability

$$p_x = \frac{\ell \cdot d^2(x, \mathcal{C})}{\sum_x d^2(x, \mathcal{C})}$$

 $\blacktriangleright \ \mathcal{C} \leftarrow \mathcal{C} \cup \mathcal{C}'$

- 2 Initial cost $\psi = \sum_x d^2(x, \mathcal{C})$
- 3 for $O(\log \psi)$ times do
 - $\mathcal{C}' \leftarrow$ sample each point $x \in X$ independently with probability

$$p_x = \frac{\ell \cdot d^2(x, \mathcal{C})}{\sum_x d^2(x, \mathcal{C})}$$

- $\blacktriangleright \ \mathcal{C} \leftarrow \mathcal{C} \cup \mathcal{C}'$
- **④** For $x \in C$, let w_x be the number of points belonging to this center
- **(3)** Recluster the *weighted* points in C into k clusters

Number of intermediate centers?

- Oversampling factor $\ell = \Theta(k)$.
- Expected number of points in $\mathcal{C} \colon \ell \log \psi$



Photo credited to Bahmani

Steven Wu (Penn)

Theorem

If an α -approximation is used in the last step, then k-means|| obtains a solution that is an $O(\alpha)$ -approximation to k-means.

Theorem

If an α -approximation is used in the last step, then k-means \parallel obtains a solution that is an $O(\alpha)$ -approximation to k-means.

For example, we could use k-means++ and get $O(\log k)$ -approximation.

Theorem

If Ψ and Ψ' are the costs of the clustering at the beginning and end of an iteration, and OPT is the cost of the optimum clustering

$$\mathbb{E}[\Psi'] \le O(OPT) + \frac{k}{e\ell}\Psi.$$

Consider a cluster $A \mbox{ in } \mbox{OPT}$

$$A = \{a_1, \dots, a_T\}$$

Centroid $c_A = \frac{1}{|T|} \sum a_t$
Increasing order of their distance to c_A



$$A = \{a_1, \dots, a_T\}$$

Centroid $c_A = \frac{1}{|T|} \sum a_t$
Increasing order of their distance to c_A

$$\phi(\mathcal{C}) = \sum_{x} d^{2}(x, \mathcal{C}); \phi_{A}(\mathcal{C}) = \sum_{a} d^{2}(a, \mathcal{C})$$

$$A = \{a_1, \dots, a_T\}$$

Centroid $c_A = \frac{1}{|T|} \sum a_t$
Increasing order of their distance to c_A

$$\phi(\mathcal{C}) = \sum_{x} d^{2}(x, \mathcal{C}); \phi_{A}(\mathcal{C}) = \sum_{a} d^{2}(a, \mathcal{C})$$

Let $p_t = \ell d^2(a_t, \mathcal{C}) / \phi(\mathcal{C})$ be the probability of selecting a_t

$$A = \{a_1, \dots, a_T\}$$

Centroid $c_A = \frac{1}{|T|} \sum a_t$
Increasing order of their distance to c_A

$$\phi(\mathcal{C}) = \sum_{x} d^{2}(x, \mathcal{C}); \phi_{A}(\mathcal{C}) = \sum_{a} d^{2}(a, \mathcal{C})$$

Let $p_t = \ell d^2(a_t, C) / \phi(C)$ be the probability of selecting a_t . For any $1 \le t \le T$,

$$q_t = p_t \prod_{j=1}^{t-1} (1 - p_j)$$

$$A = \{a_1, \dots, a_T\}$$

Centroid $c_A = \frac{1}{|T|} \sum a_t$
Increasing order of their distance to c_A

$$\phi(\mathcal{C}) = \sum_{x} d^{2}(x, \mathcal{C}); \phi_{A}(\mathcal{C}) = \sum_{a} d^{2}(a, \mathcal{C})$$

Let $p_t = \ell d^2(a_t, C) / \phi(C)$ be the probability of selecting a_t . For any $1 \le t \le T$,

$$q_t = p_t \prod_{j=1}^{t-1} (1 - p_j)$$

$$s_t = \min\left\{\phi_A, \sum_{a \in A} \|a - a_t\|^2\right\}$$

$$s_t = \min\left\{\phi_A, \sum_{a \in A} \|a - a_t\|^2\right\}$$
$$\mathbb{E}\left[\phi_A(\mathcal{C} \cup \mathcal{C}')\right] \le \sum q_t s_t + q_{T+1}\phi_A(\mathcal{C})$$

t

where q_{T+1} is the probability no point in A is selected.

$$s_t = \min\left\{\phi_A, \sum_{a \in A} \|a - a_t\|^2\right\}$$

$$\mathbb{E}\left[\phi_A(\mathcal{C}\cup\mathcal{C}')\right] \le \sum_t q_t s_t + q_{T+1}\phi_A(\mathcal{C})$$

where q_{T+1} is the probability no point in A is selected. Plug in $p_t = p$ (the case in which all points are far from C and they are tightly clustered)

$$s_t = \min\left\{\phi_A, \sum_{a \in A} \|a - a_t\|^2\right\}$$

$$\mathbb{E}\left[\phi_A(\mathcal{C}\cup\mathcal{C}')\right] \le \sum_t q_t s_t + q_{T+1}\phi_A(\mathcal{C})$$

where q_{T+1} is the probability no point in A is selected. Plug in $p_t = p$ (the case in which all points are far from C and they are tightly clustered) $q_t = p(1-p)^t$

$$s_t' = \sum_{a \in A} \|a - a_t\|^2$$

 $\{s_t'\}$ is an increasing sequence.

$$s'_t = \sum_{a \in A} \|a - a_t\|^2$$

 $\{s_t'\}$ is an increasing sequence.

$$\sum_{t} q_t s_t \leq \sum_{t} q_t s'_t$$
$$\leq 1/T \left(\sum_{t} q_t \cdot \sum_{t} s'_t \right)$$
$$= \left(\sum_{t} q_t \cdot 1/T \sum_{t} s'_t \right)$$
$$= \left(\sum_{t} q_t \right) 2\phi_A^*$$

$E[\phi_A(\mathcal{C}\cup\mathcal{C}')] \le (1-q_{T+1})2\phi_A^* + q_{T+1}\phi_A(\mathcal{C})$

Lloyd's iteration: easy to implement as long as we can store the set $\ensuremath{\mathcal{C}}$ among all mappers

 $\textbf{ I First center } \mathcal{C}: \text{ sample a point uniformly at random }$

- **①** First center C: sample a point uniformly at random
- lnitial cost $\psi = \sum_{x} d^{2}(x, C)$ (reducer simply adds)

- First center C: sample a point uniformly at random
- 2 Initial cost $\psi = \sum_{x} d^{2}(x, C)$ (reducer simply adds)
- **3** for $O(\log \psi)$ times do
 - $C' \leftarrow$ sample each point $x \in X$ independently with probability (mapper independently sample)

$$p_x = \frac{\ell \cdot d^2(x, \mathcal{C})}{\sum_x d^2(x, \mathcal{C})}$$

$$\blacktriangleright \ \mathcal{C} \leftarrow \mathcal{C} \cup \mathcal{C}'$$

	Clustering Cost Right After Initialization	Clustering Cost After Lloyd Convergence
Random	NA	22,000
K-means++	430	65
K-means	16	14

GAUSSMIXTURE: 10,000 points in 15 dimensions

K=50

Costs scaled down by 10⁴

Photo credited to Bahmani

Scalable K-Means++

Steven Wu

University of Pennsylvania

October 3, 2014

Paper by Bahmani, Moseley, Vattani, Kumar and Vassilvitskii