

Lecture 11:  $k$ -means clustering*Lecturer: Edo Liberty*

**Warning:** This note may contain typos and other inaccuracies which are usually discussed during class. Please do not cite this note as a reliable source. If you find mistakes, please inform me.

## 1 Introduction

**Definition 1.1** ( $k$ -means). Given  $n$  vectors  $x_1, \dots, x_n \in \mathbb{R}^d$ , and an integer  $k$ , find  $k$  points  $\mu_1, \dots, \mu_k \in \mathbb{R}^d$  which minimize the expression:

$$f_{k\text{-means}} = \sum_{i \in [n]} \min_{j \in [k]} \|x_i - \mu_j\|^2$$

In words, we aim to find  $k$  cluster centers. The cost is the squared distance between all the points to their closest cluster center.  $k$ -means clustering and Lloyd's algorithm [6] are probably the most widely used clustering procedure. This is for three main reasons:

- The objective function is simple and natural.
- Lloyd's algorithm (which we see below) is simple, efficient and often results in the optimal solution.
- The results are easily interpretable and are often quite descriptive for real data sets.

In 1957 Stuart Lloyd suggested a simple iterative algorithm which efficiently finds a local minimum for this problem. This algorithm (a.k.a. Lloyd's algorithm) seems to work so well in practice that it is sometimes referred to as  $k$ -means or the  $k$ -means algorithm.

---

### Algorithm 1 Lloyd's Algorithm

---

```

 $\mu_1, \dots, \mu_k \leftarrow$  randomly chosen centers
while Objective function still improves do
   $S_1, \dots, S_k \leftarrow \emptyset$ 
  for  $i \in 1, \dots, n$  do
     $j \leftarrow \arg \min_{j'} \|x_i - \mu_{j'}\|^2$ 
    add  $i$  to  $S_j$ 
  end for
  for  $j \in 1, \dots, k$  do
     $\mu_j = \frac{1}{|S_j|} \sum_{i \in S_j} x_i$ 
  end for
end while

```

---

This algorithm can be thought of as a potential function reducing algorithm. The potential function is

$$f_{k\text{-means}} = \sum_{j \in [k]} \sum_{i \in S_j} \|x_i - \mu_j\|^2.$$

The sets  $S_j$  are the sets of points to which  $\mu_j$  is the closest center. In each step of the algorithm the potential function is reduced. Let's examine that. First, if the set of centers  $\mu_j$  are fixed, the best assignment is clearly

the one which assigns each data point to its closest center. Also, assume that  $\mu$  is the center of a set of points  $S$ . Then, if we move  $\mu$  to  $\frac{1}{|S|} \sum_{i \in S} x_i$  then we only reduce the potential. This is because  $\frac{1}{|S|} \sum_{i \in S} x_i$  is the best possible value for  $\mu$  (can easily be seen by derivation of the cost function).

The algorithm therefore terminates in a local minimum. The question of course is whether we can guaranty that the solution is close to optimal and under what computational cost.

## 2 k-means and PCA

This section will present a simple connection between  $k$ -means and PCA (similar ideas given here [3]).

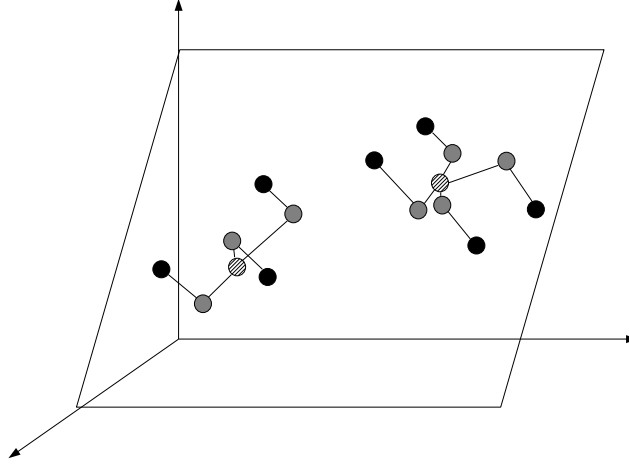


Figure 1: Example of  $k$ -means ( $k = 2$ ) cost broken into a PCA cost and a  $k$ -means cost in dimension  $k$ .

First, consider the similarity between the  $k$ -means cost function

$$f_{k\text{-means}} = \min_{\mu_1, \dots, \mu_k} \sum_{i \in [n]} \min_{j \in [k]} \|x_i - \mu_j\|^2$$

and that of PCA

$$f_{PCA} = \min_{P_k} \sum_{i \in [n]} \|x_i - P_k x_i\|^2 = \min_{P_k} \sum_{i \in [n]} \min_{y_i \in P_k} \|x_i - y_i\|^2$$

where  $P_k$  is a projection into dimension  $k$  and  $y \in P_k$  means that  $P_k y = y$ . The equality stems from the fact that for any point  $x$  and projection matrix  $P$  we have that  $\arg \min_{y \in P} \|x - y\| = Px$ .

Now, think about the subspace  $P_k^*$  which contains the  $k$  optimal centers. Since  $\mu_j^* \in P_k^*$  we have that:

$$f_{k\text{-means}} = \sum_{i \in [n]} \min_{j \in [k]} \|x_i - \mu_j^*\|^2 \tag{1}$$

$$\geq \sum_{i \in [n]} \min_{y_i \in P_k^*} \|x_i - y_i\|^2 \tag{2}$$

$$\geq \min_{P_k} \sum_{i \in [n]} \min_{y_i \in P_k} \|x_i - y_i\|^2 \tag{3}$$

$$= \min_{P_k} \sum_{i \in [n]} \|x_i - P_k x_i\|^2 = f_{PCA} \tag{4}$$

Now, consider solving  $k$ -means on the points  $y_i$  instead. This intuitively will be an easier task because they are isometrically embedded into dimension exactly  $k$  (by the projection  $P_k$ ). Before we do that though, we should argue that a good clustering for  $y_i$  results in a good clustering to  $x_i$ . Let  $P$  be any projection and let  $y_i = Px_i$  and  $\hat{\mu}_j = P\mu_j$ . We have that:

$$\sum_{j \in [k]} \sum_{i \in S_j} \|x_i - \mu_j\|^2 \geq \sum_{j \in [k]} \sum_{i \in S_j} \|Px_i - P\mu_j\|^2 \quad (5)$$

$$\geq \sum_{j \in [k]} \sum_{i \in S_j} \|y_i - \hat{\mu}_j\|^2 \quad (6)$$

$$\geq \sum_{j \in [k]} \sum_{i \in \hat{S}_j} \|y_i - \hat{\mu}_j\|^2 = \hat{f}_{k\text{-means}} \quad (7)$$

where  $\hat{S}$  and  $\hat{\mu}$  are the assignments and centers of the projected points  $y_i$ .

The following gives us a simple algorithm. Compute the *PCA* of the points  $x_i$  into dimension  $k$ . Solve  $k$ -means on the points  $y_i$  in dimension  $k$ . Output the resulting clusters and centers.

$$f_{alg} = \sum_{j \in [k]} \sum_{i \in S_j} \|x_i - \hat{\mu}_j\|^2 \quad (8)$$

$$= \sum_{j \in [k]} \sum_{i \in S_j} \|x_i - y_i\|^2 + \|y_i - \hat{\mu}_j\|^2 \quad (9)$$

$$= \sum_{i \in [n]} \|x_i - y_i\|^2 + \sum_{j \in [k]} \sum_{i \in \hat{S}_j} \|y_i - \hat{\mu}_j\|^2 \quad (10)$$

$$= f_{PCA} + \hat{f}_{k\text{-means}} \leq 2f_{k\text{-means}} \quad (11)$$

## 2.1 $\varepsilon$ -net argument for fixed dimensions

Since computing the SVD of a matrix (and hence PCA) is well known. We get that computing a 2-approximation to the  $k$ -means problem in dimension  $d$  is possible if it can be done in dimension  $k$ .

To solve this problem we adopt a brut force approach. Let  $Q_\varepsilon$  be a set of points inside the unit ball  $B_1^k$  such that:

$$\forall z \in B_1^k \exists q \in Q_\varepsilon \text{ s.t. } \|z - q\| \leq \varepsilon$$

Such sets of points exist such that  $|Q_\varepsilon| \leq c(\frac{1}{\varepsilon})^k$ . There are probabilistic constructions for such sets as well but we will not go into that. Assuming w.l.o.g. that  $\|x_i\| \leq 1$  we can constrain the centers of the clusters to one of the points in the  $\varepsilon$ -net  $Q_\varepsilon$ . Let  $q_j$  be the closes point in  $Q_\varepsilon$  to  $\mu_j$  (so  $\|\mu_j - q_j\| \leq \varepsilon$ ). From a simple calculation we have that:

$$\sum_{j \in [k]} \sum_{i \in S_j} \|x_i - q_j\|^2 \leq \sum_{j \in [k]} \sum_{i \in S_j} \|x_i - \mu_j\|^2 + 5\varepsilon.$$

To find the best clustering we can exhaustively search through every set of  $k$  points from  $Q_\varepsilon$ . For each such set, compute the cost of this assignment on the original points and return the one minimizing the cost. That will require  $\binom{c(\frac{1}{\varepsilon})^k}{k}$  iterations over candidate solutions each of which requires  $O(ndk)$  time. The final running time we achieve is  $2^{O(k^2 \log(1/\varepsilon))} nd$ .

## 3 Sampling

Another simple idea is to sample sufficiently many points from the input as candidate centers. Ideas similar to the ones described here can be found here [7].

First, assume we have only one set of points  $S$ . Also, denote by  $\mu$  the centroid of  $S$ ,  $\mu = \frac{1}{\|S\|} \sum_{i \in S} x_i$ . We will claim that picking one of the members of  $S$  as a centroid is not much worse than picking  $\mu$ . Let  $q$  be a member of  $S$  chosen uniformly at random. Let us compute the expectation of the cost function.

$$\mathbb{E}[\sum_{i \in S} \|x_i - q\|^2] = \sum_{i \in S} \sum_{j \in S} \frac{1}{n} \|x_i - x_j\|^2 \quad (12)$$

$$\leq \sum_{i \in S} \sum_{j \in S} \frac{1}{n} \cdot 2(\|x_j - \mu\|^2 + \|x_i - \mu\|^2) \quad (13)$$

$$\leq 4 \sum_{i \in S} \|x_i - \mu\|^2. \quad (14)$$

Using Markov's inequality we get that

$$\Pr[\sum_{i \in S} \|x_i - q\|^2 \leq 8 \sum_{i \in S} \|x_i - \mu\|^2] \geq 1/2$$

If this happens we say that  $q$  is a good representative for  $S$ . Now consider again the situation where we have  $k$  clusters  $S_1, \dots, S_k$ . If we are given a set  $Q$  which contains a good candidate for each of the sets. Then, restricting ourselves to pick centers from  $Q$  will result in at most a multiplicative factor of 8 to the cost.

The set  $Q$  can be quite small if the set are roughly balanced. Let the smallest set contain  $n_s$  points. We therefore succeed in finding a good representative for any set with probability at least  $\frac{1}{2} \frac{n_s}{n}$ . The probability of failure for any set is thus bounded by  $k(1 - \frac{n_s}{2n})^{|Q|}$ . Therefore  $|Q| = O(k \log(k))$  if  $n_s \in \Omega(n/k)$ .

Again, iterating over all subsets of  $Q$  of size  $k$  we can find an approximate solution in time  $O(\binom{ck \log(k)}{k} knd) = 2^{O(k \log(k))} nd$ .

## 4 Advanced reading

In the above, we gave approximation algorithms to the  $k$ -means problem. Alas, any solution can be improved by performing Lloyd's algorithm on its output. Therefore, such algorithms can be considered as 'seeding' algorithms which give initial assignments to Lloyd's algorithm. A well known seeding procedure [2] is called  $k$ -means++. In each iteration, the next center is chosen randomly from the input points. The distribution over

---

**Algorithm 2**  $k$ -means++ algorithm [2]

---

```

 $C \leftarrow \{x_i\}$  where  $x_i$  is a uniformly chosen from  $[n]$ .
for  $j \in [k]$  do
    Pick node  $x$  with probability proportional to  $\min_{\mu \in C} \|x - \mu\|^2$ 
    Add  $x$  to  $C$ 
end for
return:  $C$ 

```

---

the points is not uniform. Each point is picked with probability proportional to the minimal square distance from it to a picked center. Surprisingly, this simple and practical approach already gives an  $O(\log(k))$  approximation guarantee. More precisely, let  $f_{k\text{-means}}(C)$  denote the cost of  $k$ -means with the set of centers  $C$ . Also, denote by  $C^*$  the optimal set of centers. Then

$$\mathbb{E}[f_{k\text{-means}}(C)] \leq 8(\log(k) + 2).$$

In [1] the authors give a streaming algorithm for this problem. They manipulate ideas from [2] and combine them with a hierarchical divide and conquer methodology. See also [4] for a thorough survey and new techniques for clustering in streams.

Another problem which is very related to  $k$ -means is the  $k$ -medians problem. Given a set of points  $x_1, \dots, x_n$  the aim is to find centers  $\mu_1, \dots, \mu_k$  which minimize:

$$f_{k\text{-medians}} = \sum_{i \in [n]} \min_{j \in [k]} \|x_i - \mu_j\|$$

Both  $k$ -means and the  $k$ -median problem admit  $1 + \varepsilon$  multiplicative approximation algorithms but these are far from being simple. See [5] for more details, related work, and a new core set based solution.

## References

- [1] Nir Ailon, Ragesh Jaiswal, and Claire Monteleoni. Streaming  $k$ -means approximation. In *NIPS*, pages 10–18, 2009.
- [2] David Arthur and Sergei Vassilvitskii.  $k$ -means++: the advantages of careful seeding. In *SODA*, pages 1027–1035, 2007.
- [3] Chris H. Q. Ding and Xiaofeng He.  $K$ -means clustering via principal component analysis. In *ICML*, 2004.
- [4] Sudipto Guha, Adam Meyerson, Nina Mishra, Rajeev Motwani, and Liadan O’Callaghan. Clustering data streams: Theory and practice. *IEEE Trans. Knowl. Data Eng.*, 15(3):515–528, 2003.
- [5] S. Har-Peled and A. Kushal. Smaller coresets for  $k$ -median and  $k$ -means clustering. pages 126–134, 2005.
- [6] Stuart P. Lloyd. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 28:129–137, 1982.
- [7] Hongyuan Zha, Xiaofeng He, Chris H. Q. Ding, Ming Gu, and Horst D. Simon. Spectral relaxation for  $k$ -means clustering. In *NIPS*, pages 1057–1064, 2001.