

Seeded Center Initialization In k-means Clustering

Masum Billal

Farhad Naeem

Nobel Khandaker

Walid Khan

Data Science Department

Shohoz

Bangladesh

BILLALMASUM93@GMAIL.COM

EMAIL

EMAIL

EMAIL

Editor: Not assigned

Abstract

k-means is a popular clustering algorithm that aims to reduce the sum of minimum squared distances from data points to the centers. Now a days most of the times seeded centers are used instead of choosing all uniformly at random. We will talk about the most popular seeding technique of **k-means** known as **k-means++** in details. We take a second look at theoretical limits and experimental results presented in the **k-means++** paper. Experimental results will be shown for five different type of seeding. Therefore, this paper also serves as an updated survey of **k-means** initial center seeding methods. Our conclusion is that contrary to popular belief, usually other initialization methods than **k-means++** perform better. We also prove a theorem on upper bound of expected value of inertia when **k-means** is seeded with a probability.

1. Introduction

Widely regarded as the most popular clustering techniques, **k-means** remains a humble interesting topic in machine learning as well as computational geometry. Roughly the problem is: given a set of points \mathbf{X} in \mathbb{R}^d . Find a set of centers \mathcal{C} such that the function *inertia*

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

is minimum where $\|\cdot\|$ is the L_2 norm¹.

Default **k-means** algorithm starts with random centers and then converge based on minimum distances of the centers from the data points. New centers are calculated based on the centroid. This is known as Lloyd's algorithm (Lloyd, 1982). We repeat this process until no more change is possible. Ostrovsky et al. (2006) and Arthur et al. (2007) take it one step further by choosing the initial centers with a probability. We intend to introduce other ways of initialization and compare their performances in terms of inertia, convergence speed and CPU time taken.

1. $\|c - x\|$ or L_2 norm of $\mathbf{c} - \mathbf{x}$ is the distance between the center \mathbf{c} and point \mathbf{x} or the magnitude of the vector $\mathbf{c} - \mathbf{x}$.

Before proceeding on to the main discussion, we would like to discuss a sensitive issue. The primary motivation behind our work is not simply improving **k-means**. **k-means++** is the most popular center initialization seeding method. It seems very justified given how simple it is and the theoretical as well as experimental results given in the original **k-means++** paper (Arthur et al., 2007) are very impressive. Even in the very popular scientific python package **scikit-learn**, **k-means++** initialization is supported by default. Unfortunately, we believe neither theoretical nor experimental results shown in that paper are accurate. We will discuss our reasoning in section (3).

2. Related Work

There has been multiple surveys on **k-means** in the literature. Probably the most relevant work in this regard is done by Celebi et al. (2013). However, most of the algorithms used in that paper are not used practically very much. It was also noted by the authors themselves that **k-means++** and its greedy version work better than most. They also mention that probabilistic algorithms perform better than deterministic ones. Moreover, the not so popular algorithm (Ostrovsky et al., 2006) (ORSS) was not considered in their experiments.

We would like to point out that to our knowledge no surveys were done after removing linear dependency prior to running the experiments. This is a very important step if we are to get a meaningful clustering out of **k-means** algorithm. We can achieve this using some known algorithm such as **principle component analysis** or **singular value decomposition**. We opted for PCA in this paper. PCA decomposes the existing data points into orthogonal² ones. If we use PCA (principal component analysis) before running a clustering algorithm, we can redefine the variables into linearly independent ones. For our experiment, we have used PCA on every data set before running cluster algorithm. However, we did not reduce dimensions in order to preserve the originality of the data set. We only used PCA for removing linear dependency among variables. We strongly believe this strengthens our result over other available results.

3. How Good Is k-means++?

There is no mention of Ostrovsky’s algorithm in their paper whatsoever. Surprisingly, there is no mention of Ostrovsky’s algorithm in **k-means++** paper either even though it was published in 2007 whereas Ostrovsky’s algorithm was published in 2006. In fact, one can argue that **k-means++** algorithm is actually a special case of ORSS algorithm. Including ORSS makes the results very interesting. ORSS is like a dark horse. There are cases where ORSS is better than both **k-means++** and our proposed algorithms.

Our results conflict with some results shown in **k-means++** paper Arthur et al. (2007). For this reason, we used $k \in \{10, 25, 50\}$ to check our results against **k-means++**.

4. Proposed Initialization Methods

First of all, we would like to say that our proposed method of initializing first center in **k-means++** is not novel. In fact, this was already discussed in (Ostrovsky et al., 2006)

2. It is well known that orthogonal vectors are linearly independent.

already. However, we did not know that at the time because we did not read the full paper back then. It was only after we discovered that **k-means++** results could be wrong, we went through ORSS paper. We would really like to commend Ostrovsky et al. (2006) for such an impressive work and it is indeed very unlucky of them not getting the recognition they deserve. When we read their paper for the second time, we found out that they had already mentioned a similar result (Ostrovsky et al., 2006, Page 4, section 3, Paragraph: Running Time) in their paper. We have decided to discuss our reasoning in this paper because we believe our motivation was different than theirs. We talk about this in section (5).

To summarize, we first tried to improve **k-means++** by choosing even the first center with a probability rather than uniformly at random. Then we found out ORSS algorithm does exactly that, in a better fashion. Therefore, we consider ORSS method to be an improvement over **k-means++**. Although, to be fair, we should mention that ORSS was published before **k-means++**. So, it was a better algorithm than **k-means++** all along. It's just that we do know about it as much.

For a set of points S and a point x , we use $\min(\|x - S\|)$ to denote the minimum of distances from x to the points of S that is $\min(\|x - S\|) = \min_{a \in S}(\|x - a\|)$. Set $D(\mathbf{x}) = \min(\|\mathbf{x} - \mathcal{C}\|)$ for a point x and a set of centers \mathcal{C} . Let us denote the centroid of S by μ_S that is $\mu_S = \frac{1}{|S|} \sum_{x \in S} x$.

4.1 First center for **k-means++**

In **k-means++** algorithm, the first center is chosen uniformly at random. However, not all points have the same contribution to inertia. We choose x as a first center in a way that is equivalent to the variance explained by x .

- i Choose \mathbf{x} with probability $\frac{\|\mathbf{x} - \mu_{\mathbf{X}}\|^2}{\sum_{\mathbf{x} \in \mathbf{X}} \|\mathbf{x} - \mu_{\mathbf{X}}\|^2}$. Set $\mathcal{C}_1 = \{\mathbf{x}\}$.
- ii Repeat the remaining steps in **k-means++** (Arthur et al., 2007, Section 2.2, Page 3).

4.2 Centroid of Centers Based Seeding

We want to choose x with probability proportional to squared distance from centroid of the cluster centers. Our motivation for doing so is the following. In **k-means++**, probability is considered proportional to squared minimum distance from the centers. Therefore, the larger this minimum squared distance is, the higher the probability is for x to be chosen as a center. So, in a sense this can be thought of maximizing the minimum squared distance from the centers to the point in consideration. We intend to check the case where we choose the probability proportional to the total sum of squared distances rather than just the minimum one. As we will show later, sum of all squared distances from centers to the point in discussion is dependent on the distance from centroid of those cluster centers to that point. Note that this opens up the discussion where we need to check other ways of taking this probability. For example, minimizing variance of distances from centers, etc. However, those ideas are out of our scope for this paper.

- i Choose a point \mathbf{x} as stated in step (i) of section (4.1). Set $\mathcal{C}_1 = \{\mathbf{x}\}$.

- ii For an already existing set of i centers $\mathcal{C}_i = \{c_1, \dots, c_i\}$, choose a new center $\mathbf{x} \in \mathbf{X}$ with probability proportional to $\|\mathbf{x} - \mu_{\mathcal{C}_i}\|^2$.
- iii Repeat step (ii) until $i = k$.
- iv For each $1 \leq i \leq k$, set $\mathcal{C}_i = \{\mathbf{x} \in \mathbf{X} : \|\mathbf{x} - c_i\| = \min(x - \mathcal{C})\}$.
- v Set $c_i = \mu_{\mathcal{C}_i}$.
- vi Repeat (iv) and (v) until convergence or number of iteration is reached.

5. Motivation

In this section, we discuss the motivation behind our algorithm and why they make sense mathematically.

5.1 k-means++ Improved

Consider a set of n points \mathbf{X} and that the probability of $x \in \mathbf{X}$ being chosen as a center as $p(x)$. Following the definition of variance, for a set of points S , we define

$$\begin{aligned}\sigma^2(S) &= \frac{\sum_{x \in S} \|x - \mu_S\|^2}{|S|} \\ \sum_{x \in S} \|x - \mu_S\|^2 &= |S| \sigma^2\end{aligned}\tag{1}$$

For any arbitrary point a and μ as the centroid of \mathbf{X} ,

$$\begin{aligned}\sum_{\mathbf{x} \in \mathbf{X}} \|\mathbf{x} - a\|^2 &= \sum_{\mathbf{x} \in \mathbf{X}} \|\mathbf{x} - \mu + \mu - a\|^2 \\ &= \sum_{\mathbf{x} \in \mathbf{X}} (\|\mathbf{x} - \mu\|^2 + 2\langle \mathbf{x} - \mu, \mu - a \rangle + \|\mu - a\|^2) \\ &= \sum_{\mathbf{x} \in \mathbf{X}} \|\mathbf{x} - \mu\|^2 + 2 \left\langle \sum_{\mathbf{x} \in \mathbf{X}} \mathbf{x} - n\mu, \mu - a \right\rangle + n\|\mu - a\|^2 \\ &= n\sigma^2 + 2\langle n\mu - n\mu, \mu - a \rangle + n\|\mu - a\|^2 \\ &= n(\sigma^2 + \|\mu - a\|^2)\end{aligned}\tag{2}$$

Here $\langle a, b \rangle$ is the dot product of vectors a and b . Using equation (2), we have the following.

$$\begin{aligned}\sum_{\mathbf{x} \in \mathbf{X}} \sum_{\mathbf{y} \in \mathbf{X}} \|\mathbf{x} - \mathbf{y}\|^2 &= \sum_{\mathbf{x} \in \mathbf{X}} n(\sigma^2 + \|\mu - \mathbf{x}\|^2) \\ &= n(n\sigma^2 + \sum_{\mathbf{x} \in \mathbf{X}} \|\mu - \mathbf{x}\|^2) \\ &= n(n\sigma^2 + n\sigma^2) \\ &= 2n^2\sigma^2\end{aligned}\tag{3}$$

Using equation (3), the probability becomes

$$\begin{aligned}
p(x) &= \frac{\sum_{\mathbf{y} \in \mathbf{X}} \|\mathbf{x} - \mathbf{y}\|^2}{\sum_{\mathbf{y} \in \mathbf{X}} \sum_{\mathbf{x}' \in \mathbf{X}} \|\mathbf{x}' - \mathbf{y}\|^2} \\
&= \frac{n(\sigma^2 + \|\mu - \mathbf{x}\|^2)}{2n^2\sigma^2} \\
&= \frac{\sigma^2 + \|\mu - \mathbf{x}\|^2}{2n\sigma^2} \\
&= \frac{1}{2} + \frac{\|\mu - \mathbf{x}\|^2}{2n\sigma^2}
\end{aligned}$$

However, to make things smoother, one can also choose to use the following as the probability of x being chosen as the first center.

$$p(x) = \frac{\|\mu - \mathbf{x}\|^2}{2n\sigma^2}$$

Notice that the variance of S in the denominator. We can consider this as the amount of variance explained by x . Also, notice that computationally this version of $p(x)$ is much cheaper than using $\sum_{\mathbf{y} \in \mathbf{X}} \|\mathbf{x} - \mathbf{y}\|^2$. Therefore, we considered $p(x)$ proportional to $\|x - \mu\|^2$ for choosing x as the first center in our experiments.

5.2 Seeded Inertia

Since the algorithms in discussion are probabilistic, we should look at the expected value of inertia. For a set of points \mathbf{X} and a set of centers \mathcal{C} , we denote the inertia by $\mathcal{I}_{\mathcal{C}}(\mathbf{X})$.

$$\mathcal{I}_{\mathcal{C}}(\mathbf{X}) = \sum_{\mathbf{x} \in \mathbf{X}} \min_{\mathbf{c} \in \mathcal{C}} (\|\mathbf{c} - \mathbf{x}\|^2)$$

If the context is clear, we may omit \mathcal{C} and \mathbf{X} . For a fixed set of points \mathbf{X} , if the probability of $\mathbf{x} \in \mathbf{X}$ being chosen to be a center is $p(x)$ with respect to a set of centers \mathcal{C} , then the expected value of inertia $E[\mathcal{I}(\mathbf{X})]$ is

$$E[\mathcal{I}(\mathbf{X})] = \sum_{\mathbf{x} \in \mathbf{X}} p(x) \sum_{\mathbf{y} \in \mathbf{X}} \min(D(x), \|\mathbf{y} - \mathbf{x}\|^2)$$

If we take $p(x)$ proportional to $f(x) = \|x - \mu_{\mathcal{C}}\|^2$, then using equation (2),

$$\begin{aligned}
f(x) &= \|x - \mu_{\mathcal{C}}\|^2 \\
\sum_{x \in S} f(x) &= \sum_{x \in S} \|x - \mu_{\mathcal{C}}\|^2 \\
&= |S|(\sigma^2 + \|\mu_S - \mu_{\mathcal{C}}\|^2)
\end{aligned}$$

Then the probability of x being chosen as a center is

$$\begin{aligned}
p(x) &= \frac{f(x)}{\sum_{y \in S} f(y)} \\
&= \frac{\|x - \mu_{\mathcal{C}}\|^2}{|S|(\sigma^2 + \|\mu_S - \mu_{\mathcal{C}}\|^2)}
\end{aligned}$$

Note the following about minimum distance from x to the set of points \mathcal{C} .

$$\begin{aligned}
D(x)^2 &\leq \frac{1}{k} \sum_{c \in \mathcal{C}} \|x - c\|^2 \\
&= \frac{1}{k} \sum_{c \in \mathcal{C}} \|x - \mu_{\mathcal{C}} + \mu_{\mathcal{C}} - c\|^2 \\
&= \frac{1}{k} \sum_{c \in \mathcal{C}} (\|x - \mu_{\mathcal{C}}\|^2 + 2\langle x - \mu_{\mathcal{C}}, \mu_{\mathcal{C}} - c \rangle + \|\mu_{\mathcal{C}} - c\|^2) \\
&= \frac{1}{k} \left(k\|x - \mu_{\mathcal{C}}\|^2 + \sum_{c \in \mathcal{C}} \langle x - \mu_{\mathcal{C}}, \mu_{\mathcal{C}} - c \rangle + \|\mu_{\mathcal{C}} - c\|^2 \right) \\
&= \|x - \mu_{\mathcal{C}}\|^2 + \frac{1}{k} \sum_{c \in \mathcal{C}} \|\mu_{\mathcal{C}} - c\|^2
\end{aligned}$$

Thus, we can write inertia as

$$\begin{aligned}
E[\mathcal{I}(S)] &= \sum_{x \in S} p(x) \sum_{y \in S} \min(D(y), \|x - y\|)^2 \\
&\leq \sum_{x \in S} p(x) \sum_{y \in S} \left(\|y - \mu_{\mathcal{C}}\|^2 + \frac{1}{k} \sum_{c \in \mathcal{C}} \|\mu_{\mathcal{C}} - c\|^2 \right) \\
&= \sum_{x \in S} p(x) \sum_{y \in S} \|y - \mu_{\mathcal{C}}\|^2 + \frac{n}{k} \sum_{c \in \mathcal{C}} \|\mu_{\mathcal{C}} - c\|^2 \\
&= 1 \cdot \left(\sum_{y \in S} \|y - \mu\|^2 + n\|\mu - \mu_{\mathcal{C}}\|^2 \right) + \frac{n}{k} \mathcal{I}_{opt}(\mathcal{C}) \\
&= n(\sigma^2 + \|\mu - \mu_{\mathcal{C}}\|^2) + \frac{n}{k} \mathcal{I}_{opt}(\mathcal{C})
\end{aligned}$$

This holds for any seeding technique.

6. Experiment Setup

We ensure that all algorithms are run under the same conditions. All of them share the same environment and no special optimizations were made for any particular algorithm. Only CPU was used to determine the values we are interested in and no parallelism mechanism was in place for speeding up the process. This way, we can get an idea about the raw performances of the algorithms involved.

Python is used as the programming language to write necessary codes. Some common auxiliary packages such as *scipy*, *scikit-learn*, *numpy* etc are used to help with the code. The algorithms are simply different methods of the same class, so they share the same fitting and prediction function. Only the initialization differs for different algorithm. It should be mentioned that even though some packages have native support for **k-means** implementation, we did not use them to run the experiments. Not all algorithms we want to test are available in those packages. Therefore, in order to ensure same environment and

optimizations for every algorithm, we wrote them all from scratch so that we could be sure they are tested under the same settings.

The data sets used for the experiment are some of the popular ones.

1. Boston housing data set
2. Wine quality testing data set
3. Mall customers data set
4. Airlines cluster data set
5. Cloud data set

Miligan et al. (1988) shows that using z -score to standardize the data is not favorable for clustering because it loses between-cluster variation. Therefore, we did not use any sort of standardization or normalization lest it should lose variance or become prone to bias. Instead, we have used PCA to remove linear dependency among variables.

While experimenting on such algorithms, it is of utmost importance to run the same experiment more than once under the same parameters and conditions. For example, assume that we want to compare **k-means++** and Forgy’s algorithm (Forgy, 1965) for $k = 5$ clusters. We should run this experiment at least m times where $m > 1$ in order to eliminate bias and account for randomness. We run each experiment a total of 20 times and take the average and minimum values of inertia.

As we will see in section (7), usually every method can reach the minimum inertia at least once. Therefore, the method with lower average performs more consistently without any doubt. However, there are also cases where the difference between default **k-means** and seeded **k-means** becomes obvious.

7. Results

We present the results on the initialization procedures mentioned before. We have shown the results for $k = 10$ clusters. The ceiling of average number of iterations is taken since it may not be an integer. The time registered here is the average CPU time taken by each algorithm for those 20 iterations. As expected, default **k-means** is usually the fastest algorithm but also the worst in terms of optimizing inertia.

7.1 Comparison of Different Initialization

The comparison of inertia and number of iterations required for convergence are shown in tables 1, 3, 4.

7.2 k-means++ vs k-means++ Improved

The comparison of **k-means++** against **k-means++ improved** is shown in tables 5, 6, 7, 8, 9.

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means	2576156.1	1467415.54	9	0.09s
k-means++	1685296.83	1442170.41	7	0.1s
ORSS	1612137.72	1442170.41	7	0.1s
CoC	1604805.21	1443125.04	9	0.09s

Table 1: Results on Boston housing data set, 5 clusters

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means	2586435.45	1442170.41		2.68s
k-means++	1638992.86	1442170.41		2.66s
ORSS	1620320.77	1442170.41		2.69s
CoC	1562697.24	1442170.41		2.66s

Table 2: Results on Cloud data set, 10 clusters

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means	1011171.27	916424.19		0.93s
ORSS	1020273.2	916424.19		0.94s
k-means++	1011171.27	916424.19		0.94s
CoC	994075.55	916424.19		0.94s

Table 3: Results on Wine data set, 5 clusters

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means	2672710216978.34	2633315933118.1	75	9.33s
ORSS	2684480210031.3	2622281801735.52	40	5.52s
k-means++	2681719187656.43	2626435706106.56	40	5.41s
CoC	2667560530212.89	2627599615029.58	73	9.13s

Table 4: Results on Airlines data set, 10 clusters

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means++	6016689.97	5476422.68	24	1.01s
k-means++ Improved	5797382.38	5476422.68	19	0.84s

Table 5: **k-means++** improvement on Cloud data set, 10 clusters

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means++	41038.62	37128.01	8	0.08s
k-means++ Improved	40649.26	37663.92	6	0.07s

Table 6: **k-means++** improvement on Mall customer data set, 10 clusters

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means++	265161.75	219841.62	7	0.07s
k-means++ Improved	248798.28	217931.82	7	0.07s

Table 7: **k-means++** improvement on Wine data set, 10 clusters

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means++	2690750389197.13	2625695278621.8	41	5.57s
k-means++ Improved	2674658079797.93	2621224842855.94	35	4.8s

Table 8: **k-means++** improvement on Airlines data set, 10 clusters

8. Conclusion

We proposed methods to improve current state of **k-means** algorithms and showed their comparison in terms of convergence and inertia. **k-means++** is usually worse than both our proposed algorithm and ORSS. In most cases, our methods work best in every aspect.

References

- D. Arthur and S. Vassilvitskii. **k-means++**: The Advantages of Careful Seeding. Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms. Society for Industrial and Applied Mathematics Philadelphia, PA, USA. pp. 1027–1035.
- E. Forgy. Cluster analysis of multivariate data: Efficiency versus interpretability of classification. Biometrics, 21, 768 – 780, 1965.
- Stuart P. Lloyd. Least squares quantization in pcm. IEEE Transactions on Information Theory, 28(2) : 129–136, 1982.
- M. E. Dyer. A simple heuristic for the p-center problem. Operations Research Letters, Volume 3, February 1985, pp. 285 – 288.
- G. Milligan, M. C. Coope. A Study of Standardization of Variables in Cluster Analysis, Journal of Classification 5(2) (1988) 181–204.
- R. Ostrovsky, Y. Rabani, Leonard J. Schulman, C. Swamy. The Effectiveness of Lloyd-Type Methods for the k-Means Problem. Proceedings of the 47th Annual Symposium on Foundations of Computer Science. 2006.

Algorithm	Average \mathcal{I}	Minimum \mathcal{I}	Iterations	Time
k-means++	773980.67	677763.46	12	0.35s
k-means++ Improved	772313.3	682390.41	10	0.31s

Table 9: **k-means++** improvement on Boston housing data set, 10 clusters

- M. Emre Celebi, Hassan A. Kingravi, Patricio A. Vela. A Comparative Study of Efficient Initialization Methods for the K-Means Clustering Algorithm. *Expert Systems with Applications*, 40(1) : 200–210, 2013.
- Tom M. Apostol. *Introduction to Analytic Number Theory*. Springer, 10.1007/978 – 3 – 662 – 28579 – 4, 1976.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E. Scikit-learn: Machine Learning in Python , *JMLR* 12, pp. 2825-2830, 2011.