
L10: k -Means Clustering

Probably the most famous clustering formulation is k -means. This is the focus today. Note: k -means is not an algorithm, it is a problem formulation.

k -Means is in the family of *assignment based clustering*. Each cluster is represented by a single point, to which all other points in the cluster are “assigned.” Consider a set X , and distance $\mathbf{d} : X \times X \rightarrow \mathbb{R}_+$, and the output is a set $C = \{c_1, c_2, \dots, c_k\}$. This implicitly defines a set of clusters where $\phi_C(x) = \arg \min_{c \in C} \mathbf{d}(x, c)$. Then the *k -means clustering problem* is to find the set C of k clusters (often, but always as a subset of X) to

$$\text{minimize } \sum_{x \in X} \mathbf{d}(\phi_C(x), x)^2.$$

So we want every point assigned to the closest center, and want to minimize the sum of the squared distance of all such assignments.

Recall, there are other variants:

- the *k -center clustering problem*: minimize $\max_{x \in X} \mathbf{d}(\phi_C(x), x)$
This will be the topic of next lecture
- the *k -median clustering problem*: minimize $\sum_{x \in X} \mathbf{d}(\phi_C(x), x)$

10.1 Lloyd’s Algorithm

When people think of k -means, they usually think of the following algorithm. It is usually attributed to Lloyd from a document in 1957, although it was not published until 1982.

Algorithm 10.1.1 Lloyd’s Algorithm for k -Means Clustering

Choose k points $C \subset X$ (arbitrarily?)

repeat

For all $x \in X$, find $\phi_C(x)$ (closest center $c \in C$ to x)

For all $i \in [k]$ let $c_i = \text{average}\{x \in X \mid \phi_C(x) = c_i\}$

until The set C is unchanged

If the main loop has R rounds, then this take roughly Rnk steps (and can be made closer to $Rn \log k$ with faster nearest neighbor search in some cases).

But what is R ?

- It is finite. The cost ($\sum_{x \in X} (\mathbf{d}(x, \phi_C(x))^2)$) is always decreasing, and there are a finite (precisely, $\binom{n}{k} = O(n^k)$) number of possible distinct cluster centers. But it could be exponential in k and d (the dimension when Euclidean distance used).
- However, usually $R = 10$ is fine.
- Smoothed analysis: if data perturbed randomly slightly, then $R = O(n^{35} k^{34} d^8)$. This is “polynomial,” but still ridiculous.
- If all points are on a grid of length M , then $R = O(dn^4 M^2)$. But that’s still way too big.

Lesson: there are crazy special cases that can take a long time, but usually it works. Recall:

When data is easily cluster-able, most clustering algorithms work quickly and well.

When is not easily cluster-able, then no algorithm will find good clusters.

Sometimes there is a good k -means clustering, but it is not found by Lloyd's algorithm. Then we can choose new centers again (with randomness), and try again.

How do we initialize C ? The goal is to get one point from each final cluster. Then it will converge quickly.

- Random set of k points. By coupon collectors, we know that we need about $k \log k$ to get one in each cluster.
- Randomly partition $X = \{X_1, X_2, \dots, X_k\}$ and take $c_i = \text{average}(X_i)$. This biases towards "center" of X (by Chernoff-Hoeffding).
- Gonzalez algorithm (for k -center). This biases too much to outlier points.

Recent algorithm (Arthur + Vassilvitskii) called k -means++.

Algorithm 10.1.2 k -Means++ Algorithm

Choose $c_1 \in X$ arbitrarily. Let $C_1 = \{c_1\}$.

(In general let $C_i = \{c_1, \dots, c_i\}$.)

for $i = 2$ to k **do**

 Choose c_i from X with probability proportional to $\mathbf{d}(x, \phi_{C_{i-1}}(x))^2$.

As Algorithm 10.1.2 describes, the algorithm is like Gonzalez algorithm, but is not completely greedy.

How accurate is Lloyd's algorithm for k -means? It can be arbitrarily bad.

Theory algorithm: Gets $(1 + \varepsilon)$ -approximation for k -means in $2^{(k/\varepsilon)^{O(1)}} nd$ time. (Kumar, Sabharwal, Sen 2004).

But k -means++ is $O(\log n)$ -approximate (or 8-approximate if data is well-spaced). Can then be refined with k -means, if desired.

10.2 Problems with k -Means

- The key step that makes Lloyd's algorithm so cool is $\text{average}\{x \in X\} = \arg \min_{c \in \mathbb{R}^d} \sum_{x \in X} \|c - x\|^2$. But this only works with $\mathbf{d}(x, c) = \|x - c\|_2$.

As an alternative, can enforce that $C \subset X$. Then choose each c_i from $\{x \in X \mid \phi_C(x) = c_i\}$ that minimizes distance. But slower.

- Is effected by outliers more than k -median clustering. Can adapt Lloyd's algorithm, but then step two (recentering) is harder: Called "Fermat-Weber problem," and can be approximated with gradient descent.
- Enforces equal-sized clusters. Based on distance to cluster centers, not density.

One adaptation that perhaps has better modeling is the EM formulation: Expectation-Maximization. It models each cluster as a Gaussian distribution G_i centered at c_i .

- For each point $x \in X$, find cluster c_i with largest probability of containing that point.
- For each cluster, find best fit Gaussian G_i with $c_i = \text{average}\{x \in X \mid \phi_C(x) = c_i\}$, but estimated variance from data.

This can also allow for non-uniform Gaussians, but first taking PCA of data in cluster, and then estimating variance along each PCA axis. Can be made more robust with regularization.

10.3 Speeding-Up k -Means

- First run Lloyds (or k -means++) on random sample of points (of size $n' \ll n$). Then given good estimate of centers, run on full set (will hopefully be close to converged).
- Run a one-pass algorithm (streaming, covered later) getting $O(k \log k)$ clusters. Reduce to k clusters at end, but merging extra clusters.

Can use another streaming trick where there are a hierarchy of clusters of recent subsets representing geometrically increasing size.

- BFR algorithm: Process data in batches. For each batch find some good clusters (represented as Gaussians), plus some leftover points that don't fit in good clusters. Then merge representatives of all batches.

This increase error on the second merge. An open question is to do this so it is *mergeable* and does not increase error on merge. May need a little extra information.