

# Hierarchical clustering

## Hierarchical clustering

Use for retrieval

## Our investigation

### Algorithm complexity analysis

Let  $n$  be the number of data points,  $k$  be the number of clusters (number of means),  $d$  be the dimensionality of the data, and  $i$  be the number of iterations until convergence. The runtime of the  $k$ -means algorithm (including spherical  $k$ -means) is  $O(nkdi)$  (TODO: CITE).

Let us assume that the cost of initializing the means and splitting empty clusters is negligible, i.e.  $O(1)$ . In practice, the two are related: we find that explicitly picking good starting means for the agglomerate algorithm (e.g. randomly, or by maximizing distance between the means) makes the probability of empty clusters small, whereas it becomes excessive and unnecessary for the divisive algorithm. TODO: We explain this in a later section.

We will also assume that  $d$  and  $i$  are constants. Though in degenerate data sets  $i$  can be an exponential function of  $n$ , in practice  $i$  is a small constant on datasets with clustering structuring (CITE: WIKI?).

Let us also assume that we choose a branching factor of 2, i.e. a binary hierarchical tree structure. The following analysis does not depend on the branching factor, assuming it is constant.

### Theorem. *Agglomerative complexity*

$$T_{aggl}(n) = \Theta(n^2 di)$$

**Proof.** Our agglomerate algorithm has a recurrence relation of the following form:

$$T_{aggl}(n) = T\left(\frac{n}{2}\right) + nkdi$$

for some constant  $c$ , since we divide the problem size into two for every recursive call. But we choose  $k = \frac{n}{2}$  each time, since the number of means we choose to cluster with is also half of  $n$ .

We use the master theorem (TODO: CITE) to solve this recurrence relation: This satisfies case 3 of the master theorem:

$$f(n) \in \Omega(n^c) \text{ s.t. } c > \log_b(a)$$

since  $f(n) \in \Omega(n^2)$  s.t.  $c = 2 > \log_b(a) = 0$ , and

$$af\left(\frac{n}{b}\right) \leq k_0 f(n) \text{ for some } k_0 < 1 \text{ and sufficiently large } n$$

Master theorem variable	Value
$a$	1
$b$	2
$f(n)$	$\frac{n^2 di}{2}$
$c$	2
$\log_b(a)$	0

since when  $k_0 = \frac{1}{4}$ ,  $\left\{f\left(\frac{n}{2}\right) = \left(\frac{1}{4}\right)\frac{n^2 di}{2}\right\} \leq \left\{k_0 f(n) = \left(\frac{1}{4}\right)\frac{n^2 di}{2}\right\}$   
and so, by the master theorem:  $T_{aggl}(n) = \Theta(n^2 di)$ .  $\square$

**Theorem. *Divisive complexity***

$$T_{div}(n) = \Theta(n \log(n) di)$$

**Proof.** Our divisive algorithm has a recurrence relation of the following form:

$$T_{div}(n) = 2T\left(\frac{n}{2}\right) + nk di$$

for some constant  $c$ , since we divide the problem size into two for every recursive call and we have two subproblems each time. But we choose  $k = 2$  each time, since the number of means we choose to cluster with is always the branching factor.

We use the master theorem (TODO: CITE) to solve this recurrence relation:

Master theorem variable	Value
$a$	2
$b$	2
$f(n)$	$n(2) di$
$c$	1
$\log_b(a)$	1

This satisfies case 2 of the master theorem:

$$f(n) \in \Theta(n^c \log_0^k n) \text{ s.t. } c = \log_b(a)$$

by letting  $k_0 = 0$ , since  $f(n \log^0(n)) = f(n) \in \Theta(n)$  s.t.  $c = 1 = \log_b(a) = 1$ .  
and so, by the master theorem:  $T_{div}(n) = \Theta(n \log(n) di)$ .  $\square$