

Algorithm

1 Sketch Version

Input: $x_1 \dots x_n \in \mathbb{R}^d$ ordered in decreasing frequency, number of clusters m

Output: hierarchical clustering of $x_1 \dots x_n$

1. $A \leftarrow \{1, \dots, m\}$
2. For $i = 1 \dots m$
 - (a) $\text{twin}(i) \leftarrow \arg \min_{j \in A: j \neq i} \text{cost}(x_i, x_j)$
 - (b) $\text{lb}(i) \leftarrow \min_{j \in A: j \neq i} \text{cost}(x_i, x_j)$
 - (c) $\text{tight}(i) \leftarrow 1$
3. For $i = m + 1 \dots 2n - 1$
 - (a) If $i \leq n$
 - i. $A \leftarrow A \cup \{i\}$
 - ii. $\text{twin}(i) \leftarrow \arg \min_{j \in A: j \neq i} \text{cost}(x_i, x_j)$
 - iii. $\text{lb}(i) \leftarrow \min_{j \in A: j \neq i} \text{cost}(x_i, x_j)$
 - iv. $\text{tight}(i) \leftarrow 1$
 - v. For $j \in A$, if $\text{cost}(x_i, x_j) < \text{lb}(j)$, set $\text{twin}(j) \leftarrow i$, $\text{lb}(i) \leftarrow \text{cost}(x_i, x_j)$, $\text{tight}(j) \leftarrow 1$
 - (b) $a \leftarrow \arg \min_{j \in A} \text{lb}(j)$
 - (c) While $\text{tight}(a) \neq 1$
 - i. $\text{twin}(a) \leftarrow \arg \min_{j \in A: j \neq a} \text{cost}(x_a, x_j)$
 - ii. $\text{lb}(a) \leftarrow \min_{j \in A: j \neq a} \text{cost}(x_a, x_j)$
 - iii. $\text{tight}(a) \leftarrow 1$
 - iv. $a \leftarrow \arg \min_{j \in A} \text{lb}(j)$
 - (d) $b \leftarrow \text{twin}(a)$, $c \leftarrow i - m + n$
 - (e) $x_c \leftarrow \text{merge}(x_a, x_b)$
 - (f) $A \leftarrow A \cup \{c\} \setminus \{a, b\}$
 - (g) $\text{twin}(c) \leftarrow \arg \min_{j \in A: j \neq c} \text{cost}(x_c, x_j)$
 - (h) $\text{lb}(c) \leftarrow \min_{j \in A: j \neq c} \text{cost}(x_c, x_j)$
 - (i) $\text{tight}(c) \leftarrow 1$
 - (j) For $j \in A$, if $j \neq c$ and $\text{twin}(j) \in \{a, b\}$, set $\text{tight}(j) \leftarrow 0$

2 Implementation Version

We will have total $2n - 1$ clusters, referred to by their “**ID**”s $1 \dots 2n - 1$:

- **ID** $\leq n$: clusters corresponding to words $1 \dots n$
- **ID** $> n$: clusters corresponding to the **ID**th merge

But we only need to keep around at most $m + 1$ clusters at any point. To achieve this, we will use the following vectors of length $m + 1$ which we call *holders*. The elements holders contain will change dynamically.

- **I**: $\mathbf{I}[i] \in \{1 \dots 2n - 1\}$ is the cluster **ID** at position $i \in [m + 1]$
- **C**: $\mathbf{C}[i] \in \mathbb{R}^d$ is the center of cluster $\mathbf{I}[i]$
- **lb**: $\mathbf{lb}[i] \in \mathbb{R}$ is the lower bound of cluster $\mathbf{I}[i]$
- **twin**: $\mathbf{twin}[i] \in \{1 \dots m + 1\}$ and $\mathbf{I}[\mathbf{twin}[i]]$ is the twin cluster **ID** of cluster $\mathbf{I}[i]$
- **tight**: $\mathbf{tight}[i] \in \{0, 1\}$ is the tightness of cluster $\mathbf{I}[i]$

We also keep a vector **S** of length $2n - 1$ for recording the sizes of the clusters. In contrast to holders, it has static elements—once stored, a value in this vector does not change.

- **S**: $\mathbf{S}[\mathbf{ID}] \in \{1 \dots n\}$ is the size of cluster **ID** = $1 \dots 2n - 1$

Finally, the output of the algorithm will be recorded in matlab style as $\mathbf{Z} \in \mathbb{R}^{(n-1) \times 3}$ where $\mathbf{Z}(\mathbf{ID} - n) \in \mathbb{R}^3$ for **ID** = $n + 1 \dots 2n - 1$ corresponds to cluster **ID** and

- $\mathbf{Z}(\mathbf{ID} - n)_1$: left child of cluster **ID**
- $\mathbf{Z}(\mathbf{ID} - n)_2$: right child of cluster **ID**
- $\mathbf{Z}(\mathbf{ID} - n)_3$: merge cost of cluster **ID**

For $i, j \leq m + 1$, we use function $\text{cost}(i, j)$ to compute the distance between clusters $\mathbf{I}[i]$ and $\mathbf{I}[j]$. One call to cost has runtime $O(d)$.

$$\text{cost}(i, j) = \frac{\mathbf{S}[\mathbf{I}[i]] \times \mathbf{S}[\mathbf{I}[j]]}{\mathbf{S}[\mathbf{I}[i]] + \mathbf{S}[\mathbf{I}[j]]} \|\mathbf{C}[i] - \mathbf{C}[j]\|_2^2$$

Also, use function $\text{merge}(i, j)$ to compute the new center of the result of merging clusters $\mathbf{I}[i]$ and $\mathbf{I}[j]$. One call to merge has runtime $O(d)$.

$$\text{merge}(i, j) = \frac{\mathbf{S}[\mathbf{I}[i]]}{\mathbf{S}[\mathbf{I}[i]] + \mathbf{S}[\mathbf{I}[j]]} \mathbf{C}[i] + \frac{\mathbf{S}[\mathbf{I}[j]]}{\mathbf{S}[\mathbf{I}[i]] + \mathbf{S}[\mathbf{I}[j]]} \mathbf{C}[j]$$

2.1 Initialize the first m clusters: for $i = 1 \dots m$

1. $\mathbf{S}[i] = 1$
2. $\mathbf{I}[i] = i$
3. $\mathbf{C}[i] = x_i$
4. $\mathbf{lb}[i] = \min_{j \in [m]: j \neq i} \text{cost}(i, j)$
5. $\mathbf{twin}[i] = \arg \min_{j \in [m]: j \neq i} \text{cost}(i, j)$
6. $\mathbf{tight}[i] = 1$

Note that this takes $O(dm^2)$ operations. Each holder has an empty slot at position $m + 1$.

2.2 Perform $n - 1$ merges: for $i = n + 1 \dots 2n - 1$

Let $t = i - n + m$. Thus

- $t \in [m + 1, n]$ is the cluster **IDs** corresponding to the last $n - m$ words
- $i \in [n + 1, 2n - 1]$ is the cluster **IDs** corresponding to the $n - 1$ merges

2.2.1 If $t \leq n$, do this extra step

1. $S[t] = 1$
2. $I[m + 1] = t$
3. $C[m + 1] = x_t$
4. $lb[m + 1] = \min_{j \in [m]} \text{cost}(m + 1, j)$
5. $twin[m + 1] = \arg \min_{j \in [m]} \text{cost}(m + 1, j)$
6. $tight[m + 1] = 1$
7. For $j \in [m]$, if $\text{cost}(m + 1, j) < lb[j]$, set $lb[j] \leftarrow \text{cost}(m + 1, j)$, $twin[j] \leftarrow m + 1$, $tight[j] \leftarrow 1$

This extra step takes $O(dm)$ operations. Now each holder has full $m + 1$ elements.

2.2.2 Find the positions a and b of the closest pair

The clusters in consideration are the first $r := \min\{m + 1, 2n - i + 1\}$ elements in holders.¹

- $a \leftarrow \arg \min_{j \in [r]} lb[j]$
- While $tight[a] \neq 1$
 1. $lb[a] = \min_{j \in [r]: j \neq a} \text{cost}(a, j)$
 2. $twin[a] = \arg \min_{j \in [r]: j \neq a} \text{cost}(a, j)$
 3. $tight[a] = 1$
 4. For $j \in [r]$, if $\text{cost}(a, j) < lb[j]$, set $lb[j] \leftarrow \text{cost}(a, j)$, $twin[j] \leftarrow a$, $tight[j] \leftarrow 1$
 5. $a \leftarrow \arg \min_{j \in [r]} lb[j]$
- $b \leftarrow twin[a]$
- Make sure $a < b$

This step takes $O(m)$ in the best case and $O(dm^2)$ in the worst case.

2.2.3 Record the merge between $I[a]$ and $I[b]$ as cluster i

1. $Z[i - n]_1 = I[a]$
2. $Z[i - n]_2 = I[b]$
3. $Z[i - n]_3 = \text{cost}(a, b)$

¹At the $(n - m + 1)^{th}$ merge (i.e., the first merge without adding a cluster), we will have $i = n + (n - m + 1) = 2n - m + 1$ so that $r = \min\{m + 1, m\} = m$.

2.2.4 Put cluster i at position a in holders

1. $S[i] = S[I[a]] + S[I[b]]$
2. $C[a] = \text{merge}(a, b)$ // overwriting: now $C[a]$ is the center of cluster i
3. $I[a] = i$
4. $lb[a] = \min_{j \in [r]: j \neq a, b} \text{cost}(a, j)$
5. $\text{twin}[a] = \arg \min_{j \in [r]: j \neq a, b} \text{cost}(a, j)$
6. $\text{tight}[a] = 1$

2.2.5 Push the cluster at position b out of consideration: for $j = 1 \dots r - 1$

- If $j < b$
 - If $\text{twin}[j] > b$, set $\text{twin}[j] = \text{twin}[j] - 1$
- If $j \geq b$
 1. $I[j] = I[j + 1]$
 2. $C[j] = C[j + 1]$
 3. $lb[j] = lb[j + 1]$
 4. $\text{twin}[j] = \text{twin}[j + 1]$ if $\text{twin}[j + 1] < b$, $\text{twin}[j] = \text{twin}[j + 1] - 1$ if $\text{twin}[j + 1] \geq b$!
 5. $\text{tight}[j] = \text{tight}[j + 1]$

Now the values at position r in our vectors are meaningless.

2.2.6 Loosen the bounds of twins of a and b : for $j = 1 \dots r - 1$

- If $\text{twin}[j] \in \{a, b\}$, then set $\text{tight}[j] \leftarrow 0$