

Clustering is a central problem in unsupervised learning. A clustering problem is typically represented by a set of elements together with a notion of similarity (or dissimilarity) between them. When the elements are points in a metric space, dissimilarity can be measured via a distance function. In more general settings, when the elements to be clustered are members of an abstract set V , similarity is defined by an arbitrary symmetric function σ defined on pairs of distinct elements in V .

Correlation Clustering (CC) is a well-known special case where σ is a $\{-1, +1\}$ -valued function establishing whether any two distinct elements of V are similar or not. The objective of CC is to cluster the points in V so to minimize the number of errors, where an error is given by any pair of elements having similarity -1 and belonging to the same cluster, or having similarity $+1$ and belonging to different clusters. Importantly, there are no a priori limitations on the number of clusters or their sizes: all partitions of V , including the trivial ones, are valid. Given V and σ , the error achieved by an optimal clustering is known as the *Correlation Clustering index*, denoted by OPT.

Note that $\text{OPT} = 0$ implies that V can be perfectly clustered: any two elements in the same cluster have similarity $+1$ and any two elements in different clusters have similarity -1 . Since its introduction, CC has attracted a lot of interest and has found numerous applications in entity resolution, image analysis, and social media analysis.

Minimizing the correlation clustering error is hard, and the best efficient algorithm found so far achieves 2OPT (the actual coefficient multiplying OPT is slightly smaller than 2). A very simple and elegant algorithm for approximating CC is KwikCluster. At each round, KwikCluster draws a random pivot π_r from V , queries the similarities between π_r and every other node in V , and creates a cluster C containing π_r and all points u such that $\sigma(\pi_r, u) = +1$. The algorithm then recursively invokes itself on $V \setminus C$. On any instance of CC, KwikCluster achieves an expected error bounded by 3OPT .

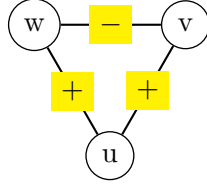
Algoritmo 1 KwikCluster

Parameters: residual node set V_r , round index r

- 1: **if** $|V_r| = 0$ **then** RETURN
 - 2: **end if**
 - 3: **if** $|V_r| = 1$ **then** output singleton cluster V_r and RETURN
 - 4: **end if**
 - 5: Draw pivot π_r u.a.r. from V_r
 - 6: $C_r \leftarrow \{\pi_r\}$ ▷ Create new cluster and add the pivot to it
 - 7: $C_r \leftarrow C_r \cup \{u \in V_r : \sigma(\pi_r, u) = +1\}$ ▷ Populate cluster
 - 8: Output cluster C_r
 - 9: KwikCluster($V_r \setminus C_r, r + 1$) ▷ Recursive call on the remaining nodes
-

We denote by $V \equiv \{1, \dots, n\}$ the set of input nodes, by $\mathcal{E} \equiv \binom{V}{2}$ the set of all pairs $\{u, v\}$ of distinct nodes in V , and by $\sigma : \mathcal{E} \rightarrow \{-1, +1\}$ the binary similarity function. A clustering \mathcal{C} is a partition of V in disjoint clusters $C_i : i = 1, \dots, k$. Given \mathcal{C} and σ , the set $\Gamma_{\mathcal{C}}$ of mistaken edges contains all pairs $\{u, v\}$ such that $\sigma(u, v) = -1$ and u, v belong to same cluster of \mathcal{C} and all pairs $\{u, v\}$ such that $\sigma(u, v) = +1$ and u, v belong to different clusters of \mathcal{C} . The cost $\Delta_{\mathcal{C}}$ of \mathcal{C} is $|\Gamma_{\mathcal{C}}|$. The correlation clustering index is $\text{OPT} = \min_{\mathcal{C}} \Delta_{\mathcal{C}}$, where the minimum is over all clusterings \mathcal{C} .

A triangle is any unordered triple $T = \{u, v, w\} \subseteq V$. We denote by $e = \{u, w\}$ a generic triangle edge; we write $e \subset T$ and $v = T \setminus e$. We say T is a *bad triangle* if the labels $\sigma(u, v), \sigma(u, w), \sigma(v, w)$ are $\{+, +, -\}$ (the order is irrelevant), see the figure below here.



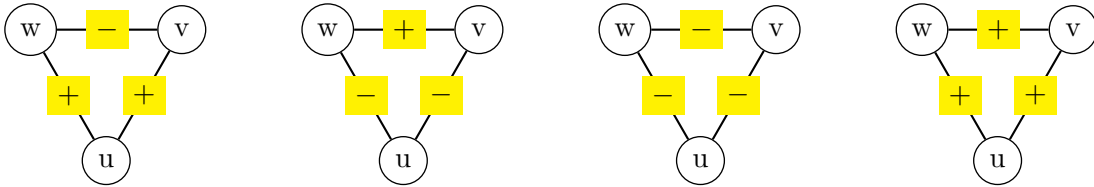
We denote by \mathcal{T} the set of all bad triangles in V and also define $\mathcal{T}(e) \equiv \{T \in \mathcal{T} : e \subset T\}$. It is easy to see that the number of edge-disjoint bad triangles is a lower bound on OPT : no matter how we cluster its nodes, a bad triangle contributes by at least 1 to the total cost of the partition. The following lemma (which we state without proof) shows that the weighted sum of all bad triangles is also a lower bound on OPT provided the sum of the weights of the bad triangles insisting on any single edge e is at most 1.

Lemma 1 *If $\{\beta_T \geq 0 : T \in \mathcal{T}\}$ is a set of weights on the bad triangles such that $\sum_{T \in \mathcal{T}(e)} \beta_T \leq 1$ for all $e \in \mathcal{E}$, then $\sum_{T \in \mathcal{T}} \beta_T \leq \text{OPT}$.*

We now bound the expected error of KwikCluster. We use V_r to denote the set of remaining nodes at the beginning of the r -th recursive call.

Let Γ_A be the set of mistaken edges for the clustering output by KwikCluster and let $\Delta_A = |\Gamma_A|$ be the cost of this clustering. The expected cost of the clustering is therefore:

$$\mathbb{E}[\Delta_A] = \sum_{e \in \mathcal{E}} \mathbb{P}(e \in \Gamma_A)$$



Lemma 2 *An edge e is mistaken by Kwik in a recursive call r if and only if there exists a bad triangle T such that $T \subseteq V_r$, $T \in \mathcal{T}(e)$, and $\pi_r = T \setminus e$.*

DIMOSTRAZIONE. Pick a round r with pivot π_r and pick an arbitrary edge $e \in V_r$. Consider the triangle $T = e \cup \{\pi_r\}$. If e is mistaken in round r , then: $\pi_r = T \setminus e$ and $T \in \mathcal{T}(e)$ (see the pictures

above here). Moreover, $T \subseteq V_r$ by construction. We prove the rest of the lemma via a case analysis with $e = \{u, w\}$. Assume $T = \{u, \pi_r, w\} \subseteq V_r$, $T \in \mathcal{T}(e)$, and $\pi_r = T \setminus e$.

Case 1: $\sigma(u, w) = +1$. If $\sigma(u, w) = +1$, $\pi_r = T \setminus e$, and T is a bad triangle, then $\sigma(\pi_r, w) \neq \sigma(\pi_r, u)$. But then u and w must end up in different clusters, which implies that e is mistaken.

Case 2: $\sigma(u, w) = -1$. If $\sigma(u, w) = -1$, $\pi_r = T \setminus e$, and T is a bad triangle, then $\sigma(\pi_r, u) = \sigma(\pi_r, w) = +1$. But then u and w end up in the same cluster, which implies that e is mistaken. \square

Since e can be mistaken only once (because when $e = \{u, w\}$ is mistaken at least one between u and w is removed from V), for each $e \in \Gamma_A$ there is a unique pair (r, T) satisfying the conditions of Lemma 2. We may thus write

$$|\Gamma_A| = \sum_{e \in \mathcal{E}} \mathbb{I}\{e \in \Gamma\} = \sum_{T \in \mathcal{T}} \sum_r \mathbb{I}\{T \subseteq V_r \wedge \pi_r \in T\} = \sum_{T \in \mathcal{T}} \mathbb{I}\{A_T\}$$

where A_T is the event $\{(\exists r) : T \subseteq V_r \wedge \pi_r \in T\}$.

Note that for any $e \in \Gamma_A$ and for any two distinct $T, T' \in \mathcal{T}(e)$, A_T and $A_{T'}$ can not both occur because, due to Lemma 2, for each $e \in \Gamma_A$ there is a unique pair (r, T) such that $T \subseteq V_r$ and $\pi_r = T \setminus e$. Thus we can write

$$1 = \sum_{T \in \mathcal{T}(e)} \mathbb{P}(A_T \wedge e \in \Gamma_A) = \sum_{T \in \mathcal{T}(e)} \mathbb{P}(e \in \Gamma_A \mid A_T) \mathbb{P}(A_T) = \sum_{T \in \mathcal{T}(e)} \frac{1}{3} \mathbb{P}(A_T). \quad (1)$$

Applying Lemma 1 with $\beta_T = \frac{1}{3} \mathbb{P}(A_T)$, we get $\sum_{T \in \mathcal{T}} \mathbb{P}(A_T) \leq 3\text{OPT}$.