

# Graph sketching-based Space-efficient Data Clustering (Supplementary material)

Anne Morvan<sup>\*†‡</sup>  
anne.morvan@cea.fr

Krzysztof Choromanski<sup>§</sup>  
kchoro@google.com

Cédric Gouy-Pailler<sup>\*</sup>  
cedric.gouy-pailler@cea.fr

Jamal Atif<sup>†</sup>  
jamal.atif@dauphine.fr

## 1 Proofs.

### 1.1 Proof of Prop. 4.1.

PROPOSITION 4.1. WHEN THE FIRST CUT IS NOT THE HEAVIEST *Let  $\mathcal{T}$  be an MST of the dissimilarity data graph with  $N$  nodes. Let us consider this specific case: all edges have a weight equal to  $w$  except two edges  $e_1$  and  $e_2$  resp. with weight  $w_1$  and  $w_2$  s.t.  $w_1 > w_2 > w > 0$ . DBMSTClu does not cut any edge with weight  $w$  and cuts  $e_2$  instead of  $e_1$  as a first cut iff:*

$$w_2 > \frac{2n_2w_1 - n_1 + \sqrt{n_1^2 + 4w_1(n_2^2w_1 + N^2 - Nn_1 - n_2^2)}}{2(N - n_1 + n_2)}$$

where  $n_1$  (resp.  $n_2$ ) is the number of nodes in the first cluster resulting from the cut of  $e_1$  (resp.  $e_2$ ). Otherwise,  $e_1$  gets cut.

*Proof.* Let  $DBCVI_1$  (resp.  $DBCVI_2$ ) be the DBCVI after cut of  $e_1$  (resp.  $e_2$ ). As  $w$  (resp.  $w_1$ ) is the minimum (resp. maximal) weight, the algorithm does not cut  $e$  since the resulting DBCVI would be negative (cf. Lemma 4.2) while  $DBCVI_1$  is guaranteed to be positive (cf. Lemma 4.1). So, the choice will be between  $e_1$  and  $e_2$  but  $e_2$  gets cut iff  $DBCVI_2 > DBCVI_1$ .  $DBCVI_1$  and  $DBCVI_2$  expressions are simplified w.l.o.g. by scaling the weights by  $w$  s.t.  $w \leftarrow 1$ ,  $w_1 \leftarrow w_1/w$ ,  $w_2 \leftarrow w_2/w$ , hence  $w_1 > w_2 > 1$ .

Then,

$$\begin{aligned} DBCVI_2 > DBCVI_1 > 0 \\ \iff \frac{n_2}{N}(\frac{w_2}{w_1} - 1) + (1 - \frac{n_2}{N})(1 - \frac{1}{w_2}) \\ \quad - \frac{n_1}{N}(1 - \frac{1}{w_1}) + (1 - \frac{n_1}{N})(1 - \frac{w_2}{w_1}) > 0 \\ \iff w_2^2 \underbrace{(N + n_2 - n_1)}_a + w_2 \underbrace{(n_1 - 2n_2w_1)}_b \\ \quad + \underbrace{(n_2 - N)w_1}_{c < 0} > 0. \end{aligned}$$

Clearly,  $\Delta = b^2 - 4ac$  is positive and  $c/a$  is negative. But  $w_2 > 0$ , then  $w_2 > \frac{-b + \sqrt{b^2 - 4ac}}{2a}$  which gives the final result after some simplifications.

### 1.2 Proof of Prop. 4.2.

PROPOSITION 4.2. FIRST CUT ON THE HEAVIEST EDGE IN THE MIDDLE *Let  $\mathcal{T}$  be an MST of the dissimilarity data graph with  $N$  nodes. Let us consider this specific case: all edges have a weight equal to  $w$  except two edges  $e_1$  and  $e_2$  resp. with weight  $w_1$  and  $w_2$  s.t.  $w_1 > w_2 > w > 0$ . Denote  $n_1$  (resp.  $n_2$ ) the number of nodes in the first cluster resulting from the cut of  $e_1$  (resp.  $e_2$ ). In the particular case where edge  $e_1$  with maximal weight  $w_1$  stands between two subtrees with the same number of points, i.e.  $n_1 = N/2$ ,  $e_1$  is always preferred over  $e_2$  as the first optimal cut.*

*Proof.* A reductio ad absurdum is made by showing that cutting edge  $e_2$  i.e.  $DBCVI_2 > DBCVI_1$  leads to the contradiction  $w_1/w < 1$ . With the scaling process from

<sup>\*</sup>CEA, LIST, 91191 Gif-sur-Yvette, France.

<sup>†</sup>Université Paris-Dauphine, PSL Research University, CNRS, LAMSADE, 75016 Paris, France.

<sup>‡</sup>Partly supported by the DGA (French Ministry of Defense).

<sup>§</sup>Google Brain Robotics, New York, USA.

Prop. 4.1'proof:

$$\begin{aligned}
DBCVI_1 &= \frac{1}{2}\left(1 - \frac{1}{w_1}\right) + \frac{1}{2}\left(1 - \frac{w_2}{w_1}\right) = 1 - \frac{1}{2w_1} - \frac{w_2}{2w_1} \\
DBCVI_2 &= \frac{n_2}{N}\left(\frac{w_2}{w_1} - 1\right) + \left(1 - \frac{n_2}{N}\right)\left(1 - \frac{1}{w_2}\right) \\
&= 1 - \frac{1}{w_2} + \frac{n_2}{N}\underbrace{\left(\frac{w_2}{w_1} + \frac{1}{w_2} - 2\right)}_{=A}
\end{aligned}$$

There is  $w_2 > w = 1$ , so  $\frac{1}{w_2} < 1$ . Besides  $w_2 < w_1$  so  $\frac{w_2}{w_1} < 1$  thus,  $A < 0$ . Let now consider w.l.o.g. that edge  $e_2$  is on the "right side" (right cluster/subtree) of  $e_1$  (similar proof if  $e_2$  is on the left side of  $e_1$ ). Hence, it is clear that for maximizing  $DBCVI_2$  as a function of  $n_2$ , we need  $n_2 = n_1 + 1$ . Then,

$$\begin{aligned}
DBCVI_2 &> DBCVI_1 \\
\iff -\frac{1}{w_2} + \left(\frac{1}{2} + \frac{1}{N}\right)\left(\frac{w_2}{w_1} - 2 + \frac{1}{w_2}\right) &> -\frac{1}{w_1} - \frac{w_2}{w_1} \\
\iff \left(\frac{1}{2w_1} + \frac{1}{Nw_1} + \frac{1}{2w_1}\right)w_2 - 1 - \frac{2}{N} + \frac{1}{2w_1} &+ \left(-1 + \frac{1}{2} + \frac{1}{N}\right)\frac{1}{w_2} > 0 \\
\iff \underbrace{\left(1 + \frac{1}{N}\right)w_2^2 + w_2}_{a>0} \underbrace{\left(\frac{1}{2} - w_1\left(1 + \frac{2}{N}\right)\right)}_{b<0} &+ w_1 \underbrace{\left(\frac{1}{N} - \frac{1}{2}\right)}_{c<0} > 0
\end{aligned}$$

As  $c/a < 0$  and  $w_2 > 0$ ,  $w_2 > \frac{N}{2(N+1)} \left[ w_1\left(1 + \frac{2}{N}\right) - \frac{1}{2} + \sqrt{\Delta} \right]$  with  $\Delta = \left(w_1\left(1 + \frac{2}{N}\right) - \frac{1}{2}\right)^2 + 4\left(1 + \frac{1}{N}\right)\left(\frac{1}{2} - \frac{1}{N}\right)w_1$ . This inequality is incompatible with  $w_1 > w_2$  since:

$$\begin{aligned}
w_1 > w_2 &\iff w_1 > \frac{N}{2(N+1)} \left[ w_1\left(1 + \frac{2}{N}\right) - \frac{1}{2} + \sqrt{\Delta} \right] \\
&\iff w_1 + \frac{1}{2} > \sqrt{\Delta} \\
&\iff \frac{4}{N} w_1^2 \left(1 + \frac{1}{N}\right) + \frac{4}{N} w_1 \left(-1 - \frac{1}{N}\right) < 0 \\
&\iff w_1 < 1 : \text{ILLICIT}
\end{aligned}$$

Indeed, after the scaling process,  $w_1 < 1 = w$  is not possible since by hypothesis,  $w_1 > w$ . Finally, it is not allowed to cut  $e_2$ , the only remaining possible edge to cut is  $e_1$ .

### 1.3 Proof of Prop. 4.3.

**PROPOSITION 4.3.** FATE OF NEGATIVE  $V_C$  CLUSTER  
Let  $K = t + 1$  be the number of clusters in the clustering partition at iteration  $t$ . If for some  $i \in [K]$ ,  $V_C(C_i) < 0$ , then  $DBMSTClu$  will cut an edge at this stage.

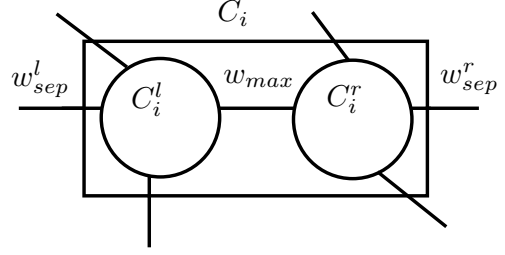


Figure 1: Generic example for proof of Prop. 4.3 and 4.5.

*Proof.* Let  $i \in [K]$  s.t.  $V_C(C_i) < 0$  i.e.  $SEP(C_i) < DISP(C_i)$ . We denote  $w_{sep}^l$  the minimal weight outting cluster  $C_i$  and  $w_{max}$  the maximal weight in subtree  $S_i$  of  $C_i$  i.e.  $SEP(C_i) \stackrel{def}{=} w_{sep}^l$  and  $DISP(C_i) \stackrel{def}{=} w_{max}$ . Hence,  $w_{sep}^l < w_{max}$ . By cutting the cluster  $C_i$  on the edge with weight  $w_{max}$ , we define  $C_i^l$  and  $C_i^r$  resp. the left and right resulting clusters.

Let us look at  $V_C(C_i^l)$ . If  $SEP(C_i^l) \geq DISP(C_i^l)$  then  $V_C(C_i^l) \geq 0 \geq V_C(C_i)$  else  $V_C(C_i^l) = \frac{SEP(C_i^l)}{DISP(C_i^l)} - 1$ . The definition of the Separation as a minimum and our cut imply that

$$SEP(C_i^l) \geq \min(SEP(C_i), w_{max}) \geq SEP(C_i).$$

Also the definition of the Dispersion as a maximum implies that  $DISP(C_i^l) \leq DISP(C_i)$ . Hence we get that  $\frac{SEP(C_i^l)}{DISP(C_i^l)} - 1 \geq \frac{SEP(C_i)}{DISP(C_i)} - 1$  i.e.  $V_C(C_i^l) \geq V_C(C_i)$  in this case too. The same reasoning holds for  $C_i^r$  showing that  $V_C(C_i^r) \geq V_C(C_i)$ . Finally,

$$\begin{aligned}
DBCVI_{aftercut} &= \sum_{j \neq i} \frac{n_j}{N} V_C(C_j) + \frac{n_i^l}{N} V_C(C_i^l) + \frac{n_i^r}{N} V_C(C_i^r) \\
&\geq \sum_{j \neq i} \frac{n_j}{N} V_C(C_j) + \frac{n_i^l}{N} V_C(C_i) + \frac{n_i^r}{N} V_C(C_i) \\
&= DBCVI_{beforecut}.
\end{aligned}$$

Hence cutting the edge with maximal weight in  $C_i$  improves the resulting DBCVI.

### 1.4 Proof of Prop. 4.4.

**PROPOSITION 4.4.** FATE OF POSITIVE  $V_C$  CLUSTER I  
Let  $\mathcal{T}$  be an MST of the dissimilarity data graph and  $C$  a cluster s.t.  $V_C(C) > 0$  and  $SEP(C) = s$ .  $DBMSTClu$  does not cut an edge  $e$  of  $C$  with weight  $w < s$  if both resulting clusters have at least one edge with weight greater than  $w$ .

*Proof.* Let us consider clusters  $C_1$  and  $C_2$  resulting from the cut of edge  $e$ . Assume that in the associated subtree of  $C_1$  (resp.  $C_2$ ), there is an edge  $e_1$  (resp.  $e_2$ ) with a weight  $w_1$  (resp.  $w_2$ ) higher than  $w$  s.t. without loss of generality,  $w_1 > w_2$ . Since  $V_C(C) > 0$ ,  $s > w_1 > w_2 > w$ . But cutting edge  $e$  implies that for  $i \in \{1, 2\}$ ,  $\text{DISP}(C_i) > \text{SEP}(C_i) = w$ , and thus  $V_C(C_i) < 0$ . Cutting edge  $e$  would therefore mean to replace a cluster  $C$  s.t.  $V_C(C) > 0$  by two clusters s.t. for  $i \in \{1, 2\}$ ,  $V_C(C_i) < 0$  which obviously decreases the current DBCVI. Thus,  $e$  does not get cut at this step of the algorithm.

### 1.5 Proof of Prop. 4.5.

**PROPOSITION 4.5. FATE OF POSITIVE  $V_C$  CLUSTER II**  
*Consider a partition with  $K$  clusters s.t. some cluster  $C_i$ ,  $i \in [K]$  with  $V_C(C_i) > 0$  is in the setting of Fig. 1 i.e. cutting the heaviest edge  $e$  with weight  $w_{\max}$  results in two clusters: the left (resp. right) cluster  $C_i^l$  (resp.  $C_i^r$ ) with  $n_1$  points (resp.  $n_2$ ) s.t.  $\text{DISP}(C_i^l) = d_1$ ,  $\text{SEP}(C_i^l) = w_{\text{sep}}^l$ ,  $\text{DISP}(C_i^r) = d_2$  and  $\text{SEP}(C_i^r) = w_{\text{sep}}^r$ . Assuming w.l.o.g.  $w_{\text{sep}}^l > w_{\text{sep}}^r$ ,  $e$  gets cut iff:*

$$\frac{\left( \frac{n_1 d_1 + n_2 d_2}{n_1 + n_2} \right)}{w_{\max}} \leq \frac{w_{\max}}{w_{\text{sep}}^r}.$$

*Proof.* As  $V_C(C_i) > 0$ , there is  $\text{SEP}(C_i) = w_{\text{sep}}^r > w_{\max}$ . Then, the DBCVI before ( $K$  clusters) and after cut of  $w_{\max}$  ( $K+1$  clusters) are:

$$\begin{aligned} \text{DBCVI}_K &= \sum_{j \neq i}^K V_C(C_j) + \frac{n_1 + n_2}{N} \left( 1 - \frac{w_{\max}}{w_{\text{sep}}^r} \right) \\ \text{DBCVI}_{K+1} &= \sum_{j \neq i}^K V_C(C_j) + \frac{n_1}{N} \left( 1 - \frac{d_1}{w_{\max}} \right) \\ &\quad + \frac{n_2}{N} \left( 1 - \frac{d_2}{w_{\max}} \right) \end{aligned}$$

DBMSTClu cuts  $w_{\max}$  iff  $\text{DBCVI}_{K+1} \geq \text{DBCVI}_K$ . So the result after simplification.

## 2 Complements on experiments.

Experiments were conducted using Python and scikit-learn library [1] on a single-thread process on an intel processor based node.

**2.1 Safety of the sketching.** Fig. 2 shows another result on a synthetic dataset: three blobs generated from three Gaussian distributions. With the three blobs, each method SEMST, DBSCAN and DBMSTClu performs well: they all manage to retrieve three clusters.

Quantitative results for the three synthetic datasets are shown in Table 1: the achieved silhouette coefficient, Adjusted Rand Index (ARI) and DBCVI. For all the indices, the higher, the better. Silhouette coefficient (between  $-1$  and  $1$ ) is used to measure a clustering partition without any external information. For DBSCAN it is computed by considering noise points as singletons. We see that this measure is not very suitable for nonconvex clusters like noisy circles or moons. The ARI (between  $0$  and  $1$ ) measures the similarity between the experimental clustering partition and the known groundtruth. DBSCAN and DBMSTClu give similar almost optimal results. Finally, the obtained DBCVIs are consistent, since the best ones are reached for DBMSTClu.

## References

- [1] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot and E. Duchesnay, *Scikit-learn: Machine Learning in Python*, Journal of Machine Learning Research, 12 (2011), pp. 2825–2830.

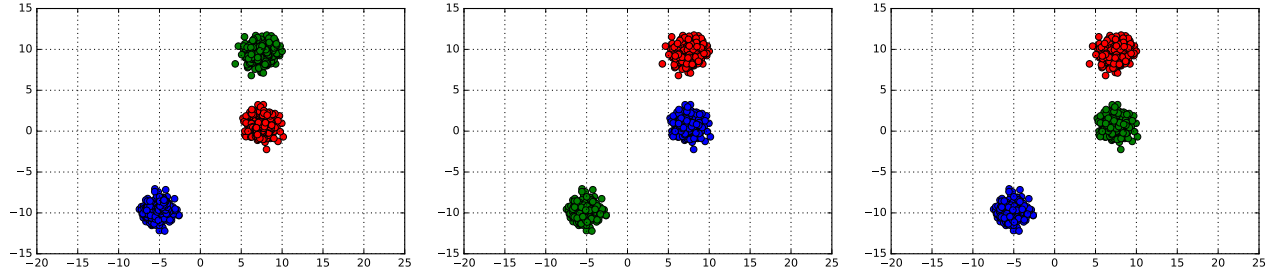


Figure 2: Three blobs: SEMST, DBSCAN ( $\epsilon = 1.4$ ,  $minPts = 5$ ), DBMSTClu with an approximate MST.

	Silhouette coeff.			Adjusted Rand Index			DBCVI		
SEMST	<b>0.84</b>	<b>0.16</b>	-0.12	<b>1</b>	0	0	<b>0.84</b>	0.001	0.06
DBSCAN	<b>0.84</b>	0.02	<b>0.26</b>	<b>1</b>	<b>0.99</b>	<b>0.99</b>	<b>0.84</b>	-0.26	<b>0.15</b>
DBMSTClu	<b>0.84</b>	-0.26	<b>0.26</b>	<b>1</b>	<b>0.99</b>	<b>0.99</b>	<b>0.84</b>	<b>0.18</b>	<b>0.15</b>

Table 1: Silhouette coefficients, Adjusted Rand Index and DBCVI for the blobs, noisy circles and noisy moons datasets with SEMST, DBSCAN and DBMSTClu.