

Supplementary material for Compressive Spectral Clustering

A. Proof of Theorem 3.2

Proof. Note that $H_{\lambda_k} = U_k U_k^T$, and that $Y_k = V_k U_k$. We rewrite $\|\tilde{\mathbf{f}}_i - \tilde{\mathbf{f}}_j\|$ in a form that will let us apply the Johnson-Lindenstrauss lemma of norm conservation:

$$\begin{aligned} \|\tilde{\mathbf{f}}_i - \tilde{\mathbf{f}}_j\| &= \|\mathbf{R}^T H_{\lambda_k}^T V_k^T (\boldsymbol{\delta}_i - \boldsymbol{\delta}_j)\| \\ &= \|\mathbf{R}^T U_k U_k^T V_k^T (\boldsymbol{\delta}_i - \boldsymbol{\delta}_j)\| \\ &= \|\mathbf{R}^T U_k (\mathbf{f}_i - \mathbf{f}_j)\| \end{aligned} \quad (1)$$

where the \mathbf{f}_i are the standard SC feature vectors. Applying Theorem 1.1 of (Achlioptas, 2003) (an instance of the Johnson-Lindenstrauss lemma) to $\|\mathbf{R}^T U_k (\mathbf{f}_i - \mathbf{f}_j)\|$, the following holds. If d is larger than:

$$\frac{4 + 2\beta}{\epsilon^2/2 - \epsilon^3/3} \log N, \quad (2)$$

then with probability at least $1 - N^{-\beta}$, we have, $\forall (i, j) \in \{1, \dots, N\}^2$:

$$(1 - \epsilon) \|U_k (\mathbf{f}_i - \mathbf{f}_j)\| \leq \tilde{D}_{ij} \leq (1 + \epsilon) \|U_k (\mathbf{f}_i - \mathbf{f}_j)\|.$$

As the columns of U_k are orthonormal, we end the proof:

$$\forall (i, j) \in [1, N]^2 \quad \|U_k (\mathbf{f}_i - \mathbf{f}_j)\| = \|\mathbf{f}_i - \mathbf{f}_j\| = D_{ij}.$$

□

B. Proof of Theorem 4.1

Proof. Recall that: $\tilde{D}_{ij}^r := \|\tilde{\mathbf{f}}_{\omega_i} - \tilde{\mathbf{f}}_{\omega_j}\| = \|\mathbf{R}^T \tilde{H}_{\lambda_k}^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|$, where $\boldsymbol{\delta}_{ij}^r = \boldsymbol{\delta}_i^r - \boldsymbol{\delta}_j^r$. Given that $\tilde{H}_{\lambda_k} = H_{\lambda_k} + E$ and using the triangle inequality in the definition of \tilde{D}_{ij}^r , we obtain

$$\begin{aligned} &\|\mathbf{R}^T H_{\lambda_k}^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| - \\ &\|\mathbf{R}^T E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| \leq \tilde{D}_{ij}^r \leq \|\mathbf{R}^T E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| + \\ &\|\mathbf{R}^T H_{\lambda_k}^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|, \end{aligned} \quad (3)$$

We continue the proof by bounding $\|\mathbf{R}^T H_{\lambda_k}^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|$ and $\|\mathbf{R}^T E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|$ separately.

Let $\delta \in]0, 1]$. To bound $\|\mathbf{R}^T H_{\lambda_k}^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|$, we set $\epsilon = \delta/2$ in Theorem 3.2. This proves that if d is larger than

$$d_0 = \frac{16(2 + \beta)}{\delta^2 - \delta^3/3} \log n,$$

then with probability at least $1 - n^{-\beta}$,

$$\left(1 - \frac{\delta}{2}\right) D_{ij}^r \leq \|\mathbf{R}^T H_{\lambda_k}^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| \leq \left(1 + \frac{\delta}{2}\right) D_{ij}^r,$$

for all $(i, j) \in \{1, \dots, n\}^2$. To bound $\|\mathbf{R}^T E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|$, we use Theorem 1.1 in (Achlioptas, 2003). This theorem proves that if $d > d_0$, then with probability at least $1 - n^{-\beta}$,

$$\|\mathbf{R}^T E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| \leq \left(1 + \frac{\delta}{2}\right) \|E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|,$$

for all $(i, j) \in \{1, \dots, n\}^2$. Using the union bound and (3), we deduce that, with probability at least $1 - 2n^{-\beta}$,

$$\begin{aligned} \left(1 - \frac{\delta}{2}\right) D_{ij}^r - \left(1 + \frac{\delta}{2}\right) \|E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| \\ \leq \tilde{D}_{ij}^r \leq \left(1 + \frac{\delta}{2}\right) \|E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| + \left(1 + \frac{\delta}{2}\right) D_{ij}^r, \end{aligned} \quad (4)$$

for all $(i, j) \in \{1, \dots, n\}^2$ provided that $d > d_0$.

Then, as e is bounded by e_1 on the first k eigenvalues of the spectrum and by e_2 on the remaining ones, we have

$$\begin{aligned} \|E^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 &= \|U e(\Lambda) U^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 \\ &= \|e(\Lambda) U^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 \\ &= \sum_{l=1}^N e(\lambda_l)^2 |(M V_k \mathbf{u}_l)^T \boldsymbol{\delta}_{ij}^r|^2 \\ &\leq e_1^2 \sum_{l=1}^k |(M V_k \mathbf{u}_l)^T \boldsymbol{\delta}_{ij}^r|^2 \\ &\quad + e_2^2 \sum_{l=k+1}^N |(M V_k \mathbf{u}_l)^T \boldsymbol{\delta}_{ij}^r|^2 \\ &= e_1^2 \|U_1^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 \\ &\quad + e_2^2 \left(\|U^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 - \|U_1^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 \right) \\ &= (e_1^2 - e_2^2) \|U_1^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 \\ &\quad + e_2^2 \|U^T V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 \\ &= (e_1^2 - e_2^2) (D_{ij}^r)^2 + e_2^2 \|V_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 \\ &\leq (e_1^2 - e_2^2) (D_{ij}^r)^2 + \frac{2e_2^2}{\min_i \{v_k(i)^2\}}. \end{aligned}$$

055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083
084
085
086
087
088
089
090
091
092
093
094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109

The last step follows from the fact that

$$\begin{aligned} \|\mathbf{V}_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\|^2 &= \sum_{l=1}^N \frac{1}{v_k(l)^2} |(\mathbf{M}^T \boldsymbol{\delta}_{ij}^r)(l)|^2 \\ &= \frac{1}{v_k(\omega_i)^2} + \frac{1}{v_k(\omega_j)^2} \leq \frac{2}{\min_i \{v_k(i)\}^2} \end{aligned}$$

Define, for all $(i, j) \in \{1, \dots, n\}^2$:

$$e_{ij} := \sqrt{|e_1^2 - e_2^2|} D_{ij}^r + \frac{\sqrt{2}e_2}{\min_i \{v_k(i)\}}.$$

Thus, the above inequality may be rewritten as:

$$\|\mathbf{E}^T \mathbf{V}_k^T \mathbf{M}^T \boldsymbol{\delta}_{ij}^r\| \leq e_{ij},$$

for all $(i, j) \in \{1, \dots, n\}^2$, which combined with (4) yields

$$\begin{aligned} \left(1 - \frac{\delta}{2}\right) D_{ij}^r - \left(1 + \frac{\delta}{2}\right) e_{ij} \\ \leq \tilde{D}_{ij}^r \leq \left(1 + \frac{\delta}{2}\right) e_{ij} + \left(1 + \frac{\delta}{2}\right) D_{ij}^r, \end{aligned} \quad (5)$$

for all $(i, j) \in \{1, \dots, n\}^2$, with probability at least $1 - 2n^{-\beta}$ provided that $d > d_0$.

Let us now separate two cases. In the case where $D_{ij}^r \geq D_{min}^r > 0$, we have

$$\begin{aligned} e_{ij} &= \frac{e_{ij}}{D_{ij}^r} D_{ij}^r = \left(\sqrt{|e_1^2 - e_2^2|} + \frac{\sqrt{2}e_2}{D_{ij}^r \min_i \{v_k(i)\}} \right) D_{ij}^r \\ &\leq \left(\sqrt{|e_1^2 - e_2^2|} + \frac{\sqrt{2}e_2}{D_{min}^r \min_i \{v_k(i)\}} \right) D_{ij}^r \\ &\leq \frac{\delta}{2 + \delta} D_{ij}^r. \end{aligned}$$

provided that Eq. (7) of the main paper holds. Combining the last inequality with (5) proves the first part of the theorem.

In the case where $D_{ij}^r < D_{min}^r$, we have

$$e_{ij} < \sqrt{|e_1^2 - e_2^2|} D_{min}^r + \frac{\sqrt{2}e_2}{\min_i \{v_k(i)\}} \leq \frac{\delta}{2 + \delta} D_{min}^r.$$

provided that Eq. (7) of the main paper holds. Combining the last inequality with (5) terminates the proof. \square

C. Experiments on the SBM with heterogeneous community sizes

We perform experiments on a SBM with $N = 10^3$, $k = 20$, $s = 16$ and heterogeneous community sizes. More

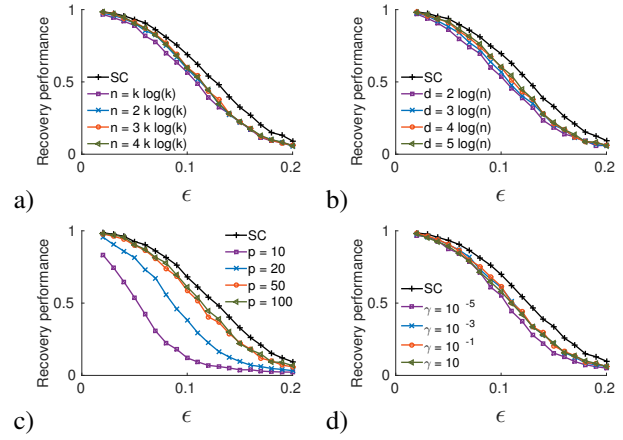


Figure 1. (a-d): recovery performance of CSC on a SBM with $N = 10^3$, $k = 20$, $s = 16$ and heterogeneous community sizes versus ϵ , for different n , d , p , γ . Default is $n = 2k \log k$, $d = 4 \log n$, $p = 50$ and $\gamma = 10^{-3}$. All results are averaged over 20 graph realisations.

specifically, the list of community sizes is chosen to be: 5, 10, 15, 20, 25, 30, 35, 40, 45, 50, 55, 60, 65, 70, 75, 80, 85, 90 and 95 nodes. In this scenario, there is no theoretical value of ϵ over which it is proven that recovery is impossible in the large N limit. Instead, we vary ϵ between 0 and 0.2 and show the recovery performance results with respect to n , d , p and γ in Fig. 1. Results are similar to the homogeneous case presented in Fig. 1(a-d) of the main paper.

References

Achlioptas, D. Database-friendly random projections: Johnson-lindenstrauss with binary coins. *Journal of Computer and System Sciences*, 66(4):671 – 687, 2003.