

Compressibility Barriers to Neighborhood-Preserving Data Visualizations

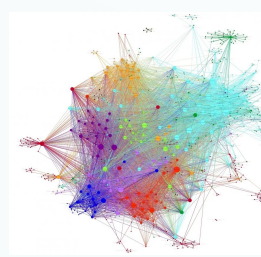
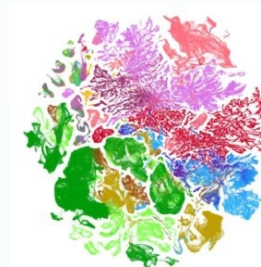


Noah Bergam (Columbia University, CS Dept.)
Joint work w. Szymon Snoeck and Nakul Verma

PROBLEM

Many scientists want to visualize big (say, 100k points, 10k dimensions) datasets in 2D or 3D plots, and often resort to (unreliable) tools like t-SNE and UMAP.

Is it even possible to accurately visualize neighborhoods in large, structured datasets?



FINDINGS

Most graphs are maximally difficult.

For most* n-vertex G,

$$\dim_{\alpha}(G) \gtrsim \frac{\log n}{\log\left(\frac{8}{\alpha}\right)}$$

(Note: $\dim_{\alpha}(G) \leq \log(n)$ for all G)

Sparsity hardly helps in general metrics.

For most n-vertex, k-regular G, $k = O(1)$,

$$\dim_{\alpha}(G) \gtrsim \frac{\log n}{\log\left(\frac{\log n}{\alpha}\right)}$$

*most = for all but $2^{-\Omega(n)}$ fraction

Sparsity helps in Euclidean space.

For most n-vertex graphs:

$$\dim_{\alpha=1}(G, \ell_2) = \Theta(n)$$

Meanwhile, for *any* degree-k graph:

$$\dim_{\alpha \leq 1+1/\sqrt{k}}(G, \ell_2) \lesssim k^2 \log n$$

(If α exceeds this threshold, ℓ_2 preservation can become impossible!)

Cluster structure rarely helps.

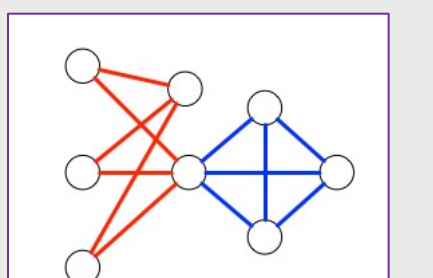
For G sampled from $(p > q)$ planted partition model on n nodes and k components, w. h. p. ,

If $p = 1$ and $\alpha \leq 1$, $\dim_{\alpha}(G) = \Theta(\log k)$,
otherwise... $\dim_{\alpha}(G) = \Theta(\log n)$.

$(\alpha \leq 1)$ - versus $(\alpha > 1)$ -preservation .

An interesting phase-change occurs.

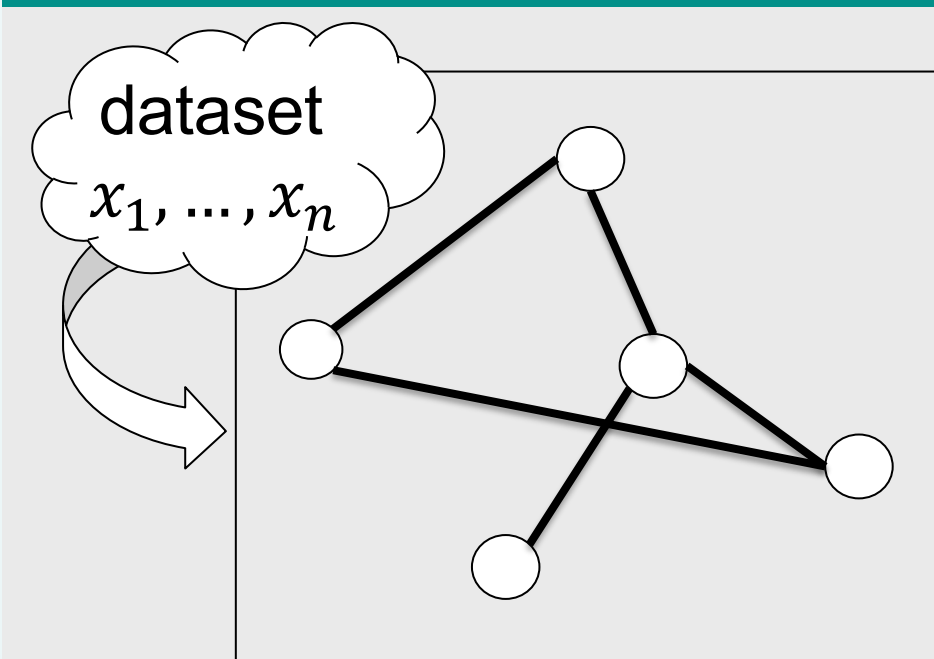
$\dim_{\alpha \leq 1}(G)$ depends on the **minimal clique partition of G** (NP-hard to compute).



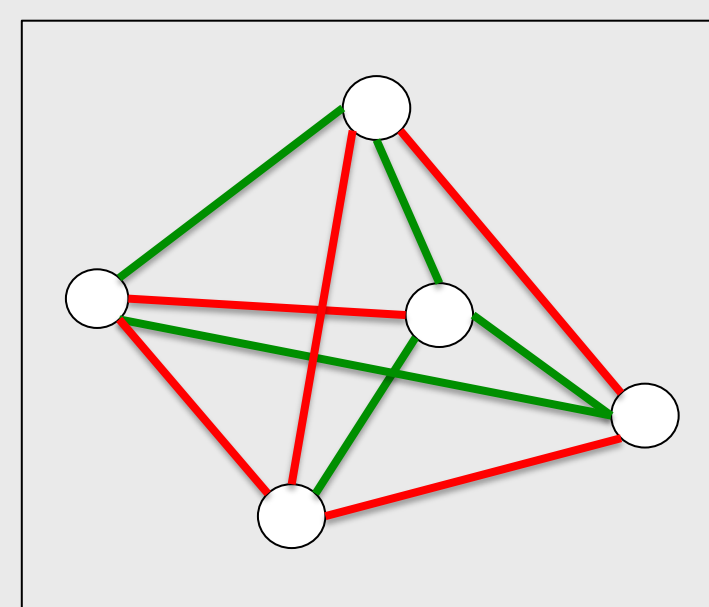
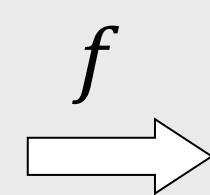
a size-2 clique partition

$\dim_{\alpha > 1}(G)$ depends on the **# of distinct neighborhoods in G** (often much larger, but easy to compute).

A simple model: α -preservation



ground-truth graph $G = (V, E)$
(representing neighborhoods)



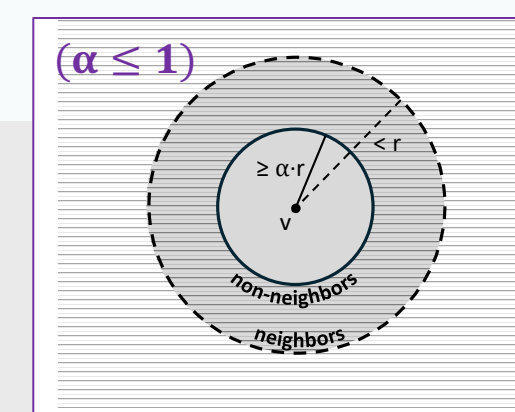
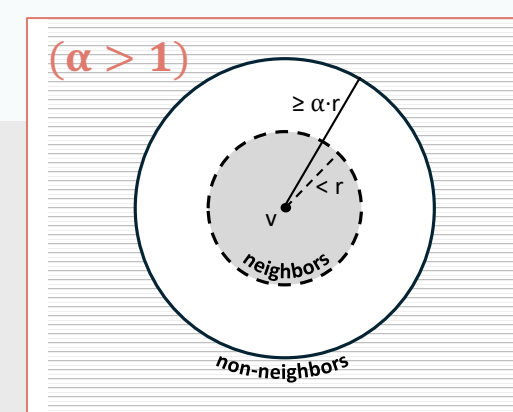
embedding in
metric space (X, ρ)

Definition: A map $f: V \rightarrow X$ is an **α -preservation** of $G = (V, E)$ if

$$(u, v) \in E \Rightarrow d(f(u), f(v)) < 1$$

$$(u, v) \notin E \Rightarrow d(f(u), f(v)) \geq \alpha$$

As α increases, the problem gets harder.
 $(\alpha > 1)$ is significantly more desirable than $(\alpha \leq 1)$.



DISCUSSION

- Special cases of α -preservation dimension have been studied before, e.g. $\dim_{\alpha=1}(G, \ell_2) = \text{"sphericity"}$.
- $(\alpha \leq 1)$ -preservation is a strict generalization of $(1/\alpha)$ -distortion embedding. Low-distortion embeddings are well-studied in TCS, but too strict a notion for visualization.
- Future work: **algorithms and complexity** of α -preservation; and **approximate** α -preservation: how many edges do need to "sacrifice" to preserve in 2D or 3D?

GOAL: characterize preservation dimension for "realistic" graphs (sparse, clustered, etc.).

IDEA: constant preservation dimension suggests that 2D/3D visualization is possible!