



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

GIRG model for realistically generating social network graphs

Master Thesis

Benjamin Dayan

August 30, 2023

Advisors: Prof. Dr. Angelika Steger, Prof. Dr. Johannes Lengler, Marc Kauffman, Ulysse Schaller

Department of Computer Science, ETH Zürich

Abstract

Many real-world graphs, in particular social networks, exhibit clustering and are hypothesized to possess an underlying geometrical generative structure. The Geometric Inhomogeneous Random Graphs (GIRG) model is one proposed generative geometric model. This thesis demonstrates the proficiency of assorted GIRG variants in generating graphs that resemble a dataset of real social network graphs across a spectrum of high-level graph features. We then investigate practical methods of fitting GIRG node geometric location parameters to a given input graph. Using these, we show that low-dimensional GIRGs can successfully replicate our social network graphs on a node and edge level. Throughout this thesis, we aim to offer guidance on generating and fitting different types of GIRGs, as well as an intuitive understanding of the GIRG model.

Contents

Contents	iii
1 Introduction	1
1.1 Thesis Introduction	1
2 GIRGs - Variants and Generation	3
2.1 GIRGs definition	3
2.1.1 Geometry	4
2.1.2 Similarity to Chung-Lu	4
2.1.3 Power Law Degree Sequence	5
2.1.4 Volume generalisation	5
2.2 Blasius C++ GIRG generation	6
2.2.1 GIRG Generation with Python	6
2.3 Cube GIRGs	6
2.3.1 Cube GIRG Formulation	7
2.3.2 Cube GIRG generation - coupling algorithm	8
2.4 Minimum Component Distance (MCD) GIRGs	8
2.4.1 Volume Formula for MCD, and Mixed Norms	9
2.4.2 MCD GIRG speedy combination of 1d GIRGs generating algorithm	9
2.5 Distorted GIRGs	11
3 Generative Graph Models (GGM) and Blasius Realism Framework	13
3.1 GGM and Framework Mathematical Outline	14
3.1.1 Fitting a GGM to a particular real graph instance	14
3.1.2 Plausibility of Fitted GGMs	14
3.1.3 Fitted GGMs as candidates for real graphs - Blasius	15
3.2 Fitting GIRG GGM to a real graph	16
3.2.1 Fitting number of nodes n	16
3.2.2 Fitting power law distribution exponent τ for node weight sampling	16
3.2.3 Power Law weight alternative: weight copying	18
3.2.4 Fitting c and α	19
3.2.5 Volume formulation torus GIRG expected average degree formula	20
3.2.6 Estimating const c in Cube GIRGs	20

3.3	Realism Framework Results	20
3.3.1	Degree distribution features	21
3.3.2	Geometric Features	22
3.3.3	Case study of distribution vs mean accuracy: closeness centrality	24
3.3.4	SVM Classifier Misclassification Rates Tables	24
4	Diffusion Maps	29
4.1	Introduction	29
4.2	Diffusion Maps Theory	30
4.2.1	Random Walk Formulation	30
4.2.2	Improving Diffusion Map Geometry	32
4.2.2.1	Expected distance approach	32
4.2.2.2	Homogeneous random walk approach	33
4.2.2.3	Bounded diffusion hop distances	34
4.3	Empirical results and post processing	35
4.3.1	Inferring dimension d	35
4.3.2	Rescaling/Shifting Diffusion Maps	37
5	Likelihood Point Estimation	41
5.1	Introduction	41
5.2	MCMC Formulation	41
5.2.1	MCMC Likelihood comparison of GIRG vs CL fit to real graph	43
5.2.2	Percent Edges Captured Metric Comparison	45
5.3	Direct Ordered Likelihood Maximisation	45
5.3.1	Analysis on GIRG to real graphs quality of fit	46
6	Graph Kernels	53
6.1	Bayes Factor Theory	53
6.1.1	Graph isomorphism problem in computing model evidence	54
6.2	Graph Kernel Introduction	55
6.2.1	Random walk kernel	55
6.2.2	Weisfeiler-Lehman kernel	55
6.2.3	Experiments	56
7	Conclusion	57
7.1	Recap	57
	Bibliography	59

Introduction

1.1 Thesis Introduction

Graphs are a useful model for describing networks of interconnected entities as an abstraction of real-life systems. Real-world graph data are abundant, and many graphs exhibit similar structural characteristics. Several graph-generating models have been proposed, each justified by its reasonable formative mechanisms which create sets of graphs similar to those observed in the wild. Studying these generative models theoretically provides insights into real networks and is useful for creating synthetic graphs for testing purposes.

Geometric Inhomogeneous Random Graphs (GIRGs) are one such graph-generating model, first proposed in [Bringmann et al., 2016] as a generalisation of *Hyperbolic Random Graphs*. GIRGs allow significant flexibility in the geometric component of their generative structure, thereby enabling better modeling of a broader range of real graphs. Their definition encapsulates the intuitive notion present in many graphs that nodes are more likely to be connected if they are “close together”. Furthermore, GIRGs exhibit inhomogeneity between nodes, resulting in a power law degree distribution, a characteristic observed in many real graphs.

This thesis evaluates the ability of GIRGs to simulate a dataset of 104 social network graphs obtained from the *Network Repository* ([Rossi and Ahmed, 2015], <https://networkrepository.com/>). These graphs represent Facebook social networks from various towns in the US, thus exhibiting structural similarities, and range from $n = 1000$ to $n = 100,000$ nodes. In chapter 3, we assess the mostly unmodified GIRG prior distribution of graphs based on their high-level similarities to the real graphs. In chapter 5, we perform a more intensive posterior style fit of GIRG to replicate real graphs at the node/edge level.

In chapter 3, the first task of assessing GIRG prior realism is conducted using a high-level feature dataset classification comparison framework proposed in [Bläsius et al., 2018]. Our results provide evidence that GIRGs are more realistic candidates than simpler generative models with respect to geometric features such as node closeness centrality, betweenness centrality, local clustering coefficient, and graph effective diameter. We compare the realism of different GIRG variants on these features and find that while max norm torus GIRGs provide a good baseline, cube GIRGs perform better on graph effective diameter, and mixed/minimum component distance GIRGs excel in closeness centrality. In most (but not all) cases,

we find that a small geometric dimensionality of $d = 1$ already provides near-peak GIRG realism, with little benefit from higher-dimensional GIRGs.

Additionally, in chapter 6, we attempt to apply the method of *graph kernels* as an alternative framework for model selection between the GIRG prior model and other generative graph models with respect to our real graph dataset. Unfortunately, we were unable to produce reasonable results.

In the second task of GIRG posterior fitting, we aim to infer the node geometric locations of a given posited GIRG generative model for an input real graph. In chapter 4, we explore the method of diffusion maps as a computationally fast initial estimate of locations with good internode distances, and as a means of proposing the underlying geometric dimensionality. In chapter 5, we describe methods to further refine node locations using either a markov chain Monte Carlo or maximum likelihood-based fitting approach. The full fitting procedure exhibits decent performance in replicating real graphs and supports the hypothesis that one dimension is mostly sufficient for this dataset of social networks.

To support all objectives, we also cover the different types of GIRGs and their unique , and demonstrate practical methods for sampling and working with them.

Chapter 2

GIRGs - Variants and Generation

Contents

2.1	GIRGs definition	3
2.1.1	Geometry	4
2.1.2	Similarity to Chung-Lu	4
2.1.3	Power Law Degree Sequence	5
2.1.4	Volume generalisation	5
2.2	Blasius C++ GIRG generation	6
2.2.1	GIRG Generation with Python	6
2.3	Cube GIRGs	6
2.3.1	Cube GIRG Formulation	7
2.3.2	Cube GIRG generation - coupling algorithm	8
2.4	Minimum Component Distance (MCD) GIRGs	8
2.4.1	Volume Formula for MCD, and Mixed Norms	9
2.4.2	MCD GIRG speedy combination of 1d GIRGs generating algorithm	9
2.5	Distorted GIRGs	11

2.1 GIRGs definition

In this section, we define GIRGs and identify their key features.

As per [Bringmann et al., 2019], the GIRG model is a random graph model characterized by the edge connection probabilities given by the formula:

$$p_{uv} = \Theta \left(\min \left\{ 1, \left(\frac{w_u w_v / W}{\|x_u - x_v\|^d} \right)^\alpha \right\} \right) \quad (2.1)$$

Here, $(w_u)_{u \in V}$ are node-specific weights in \mathbb{R}^+ , and $(x_u \in \chi)_{u \in V}$ denote positions of nodes in a particular geometric space χ , often taken to be the d-dimensional torus ¹ \mathbb{T}^n , or the d-dimensional unit cube $[0, 1]^d$. The normalising factor W^{-1} ,

¹The d-dimensional torus also resembles $[0, 1]^d$, but opposite faces of the cube are identified. Picture Pac-Man, or donuts. Just not Pac-Man eating a donut...

where $W = \sum_{u \in V} w_u$, ensures that the expected degree of a node u remains constant even as $n \rightarrow \infty$; although u will have more possible other nodes to connect with, this increase is counterbalanced by the declining probability of connecting to any one of them.

2.1.1 Geometry

The geometric space χ , along with the distance function written $r(x_u, x_v) = \|x_u - x_v\| = r_{uv}$ (which is not necessarily a norm, or even a metric), can vary a great deal without affecting key properties of GIRGs.

The edge probability p_{uv} incorporates a geometric factor from the distance r_{uv} , which inversely scales the probabilities so that nearby points (small r_{uv}) are more likely to have an edge than those further apart. r_{uv} is raised to the d th power to ensure consistency in the number of edges across different values of d . This influence in GIRGs of node proximity on p_{uv} facilitates the emergence of *clustering*, a common phenomenon in many real-world graphs, where subgroups of nodes might have more internal edges than would be expected by chance.

Using the *max norm* as the distance function, i.e., $\|x\| = \|x\|_\infty = \max_i |x_i|$, is convenient for proofs, but can be replaced by euclidean or other norms - the equivalence of norms there is not much difference. Note that on the torus, the max norm is not strictly a norm, rather we abuse notation to express:

$$\|x - y\|_\infty = \max_i |x_i - y_i|_C \quad (2.2)$$

$$|x_i - y_i|_C := \min(|x_i - y_i|, 1 - |x_i - y_i|) \quad (2.3)$$

The *minimum component distance* (MCD), defined as $\|x\|_{\text{mcd}} = \min_i |x_i|$, does not even satisfy the triangle inequality, yet it preserves most properties of GIRGs. This distance function is reasonable if we desire two nodes to have a higher edge probability if they are close in at least one dimension, rather than in every dimension.

$\alpha \in (1, \infty]$ is a parameter that impacts the affect of geometry on edge probabilities. We differentiate between *short edges* where $r_{uv} < (w_u w_v / W)^{1/d}$ and $p_{uv} = \Theta(1)$ (for these, the $\min(\cdot, 1)$ is applied to ensure $p_{uv} \leq 1$), and *long edges*, where $r_{uv} > (w_u w_v / W)^{1/d}$, and p_{uv} diminishes to zero as the distance increases. Larger values of α accelerate this decay in long edge probability, thus more severely limiting the number of long-distance edges. $\alpha = \infty$ imposes a strict cutoff, setting $p_{uv} = 0$ for long edges.

2.1.2 Similarity to Chung-Lu

A key property of the GIRG model is its nature as a geometric variant of the Chung-Lu model. In a Chung-Lu generated graph, edge probabilities between nodes for a given node weight sequence $(w_u)_{u \in V}$ are decreed as $p_{uv} = \Theta(\min\{1, \frac{w_u w_v}{W}\})$.

In a GIRG, geometry affects edge probabilities, yet when considering any two nodes u, v , marginalising over all possible locations (whether near or far) yields $E_{x_v}[p_{uv}|x_u, w_u, w_v] = \Theta(\min\{1, \frac{w_u w_v}{W}\})$, which mirrors the Chung-Lu model and disregards node locations $(xu)u \in V$.

Therefore, properties and results pertaining to graphs generated by the Chung-Lu model are directly applicable to GIRGs, such as the fact that $E[d_u] = \Theta(w_u)$ - the expected degree of a node u with weight w_u is proportional to w_u .

This might appear counterintuitive, given the presence of the $(w_u w_v / W)^\alpha$ term in the edge probabilities. However, as previously discussed, α only affects the decay rate of edge probabilities for long edges, which constitute only a constant fraction of the total edge probability, provided $\alpha > 1$. If α were smaller, it would blow up the degree d_u far beyond w_u , as the abundance of potential long edges would result in too many being realised as actual edges, so would dominate the number of short edges.

2.1.3 Power Law Degree Sequence

The weight sequence of a GIRG is usually assumed to be generated from a power law distribution with an exponent $\tau \in (2, 3)$. This can also be loosened to just include the tail (larger) weights, for instance just the top 10% largest weight nodes. This allows the GIRG to match the common real network phenomenon of a power law tail degree sequence, since node expected degrees in the GIRG are proportional to their weight.

$$\begin{aligned} x &\sim \text{powerlaw}(\tau) && \text{exact power law distribution} \\ p(w) &\propto w^{-\tau} \text{ for } w \in [x_{\min}, \infty] && \text{pdf (default } x_{\min} = 1) \\ p(x \geq w) &\propto w^{1-\tau} \end{aligned}$$

The power law distribution is heavy tailed (heavier than exponentially decaying tails). $\tau > 2$ is important to ensure that $E[w] = \Theta(1)$, which means that although we may have some very large w_i in a sequence w_1, \dots, w_n , still there is an equal balance of total weight between smaller valued w_i 's and larger weights. This produces a phenomenon known as the ‘‘Pareto principle / 80-20 rule’’, a common example being the top 20% richest individuals collectively owning roughly 80% of total private wealth. That’s certainly a heavier tail than an ideally fair society, but it also doesn’t completely skew the mean wealth of an individual more than some multiple of the median wealth.

2.1.4 Volume generalisation

The volume formulation of a GIRG presented in [Bringmann et al., 2019] provides a consistent way to handle different geometries:

$$p_{uv}(r) = \min \left\{ 1, c \left(\frac{w_u w_v / W}{\text{Vol}(r_{uv})} \right)^\alpha \right\}; \quad \rho_{uv} = \frac{w_u w_v / W}{\text{Vol}(r_{uv})} \quad (2.4)$$

Here $\text{Vol}(r_{uv}) = \text{Vol}(B_{r_{uv}})$ is the volume of the ball of radius r using the distance function $r(x_u, x_v) = r_{uv}$, which should be symmetric. For example with the max norm, $\text{Vol}(r_{uv}) = (2r_{uv})^d$ as a cube with side-length $2r_{uv}$.

Having volume in the edge probability formula is actually a generalisation of taking r_{uv} to the d th power - the point being to make $p(u \sim v | x_u, w_u, w_v) = E_r[p_{uv}(r)] = \Theta(w_u w_v / n)$, regardless of dimension. We’ll actually derive in section 3.2.5 a precise formula for expected edge probabilities / node degrees that holds across all different types of volume formulation torus GIRGs. This will however need an exact GIRG edge probability formula, without the outer $\Theta(\cdot)$, which we present in section 2.2.

2.2 Blasius C++ GIRG generation

[Bläsius et al., 2022] provides a C++ implementation of an efficient algorithm for sampling a *max torus GIRG* - with the max norm as distance function and torus geometric space. Their GIRG formulation uses a toroidal geometry $\chi = \mathbb{T}^d$, with edge probabilities

$$p_{uv} = \min \left\{ 1, c \left(\frac{w_u w_v / W}{\|x_u - x_v\|_\infty^d} \right)^\alpha \right\} \quad (2.5)$$

i.e. no outer Θ which gave us more flexibility in the general GIRG definition - just one fixed inner constant c in order to be concrete. Note that this still falls under the wider Θ GIRG definition as p_{uv} always lies in the interval

$$c \min \left\{ 1, \left(\frac{w_u w_v / W}{\|x_u - x_v\|_\infty^d} \right)^\alpha \right\} \leq p_{uv} \leq \min \left\{ 1, \left(\frac{w_u w_v / W}{\|x_u - x_v\|_\infty^d} \right)^\alpha \right\} \quad (2.6)$$

if $c \leq 1$, and with the upper and lower bounds swapped for $c > 1$.

They implement the algorithm proposed by [Bringmann et al., 2019], which has expected runtime $O(n)$ i.e. linear in the number of nodes, for a GIRG with a power law generated weight sequence. This is optimal because such a GIRG has expected number of edges $\Theta(n)$, due to the fact that nodes have finite expected weight $E[w_u] = \Theta(1)$, and expected degree proportional to their weight. I.e. just to write down / store the set of edges $E = \{(u, v) | u \sim v\}$ we need $O(n)$ time/space.

2.2.1 GIRG Generation with Python

Unfortunately this GIRG sampling algorithm scales exponentially in dimension d , as it involves breaking the Torus space down hierarchically into smaller cubes and considering edges between neighbouring cubes. We struggled to use it in experiments for $d \geq 4$. Hence we also implemented our own GIRG generation code in python for ease of modification to deal with different types of GIRG and to cope with higher dimensions d . Our algorithm is the most basic $O(n^2)$ checking of every single node pair (u, v) , and also has $O(n^2)$ memory requirement, and doesn't scale in d beyond a linear factor inherent in operating with d -dimensional points. The basic steps are:

1. sample $O(n)$ node weights w_u from a power law distribution
2. sample $O(n)$ node locations x_u from a uniform distribution on the torus
3. compute the $O(n^2)$ pairwise node distances and hence edge probabilities
4. for each of the $O(n^2)$ potential edges, sample a Bernoulli random variable with the edge probability as its parameter, to determine whether the edge exists

2.3 Cube GIRGs

Cube GIRGs are an alternative formulation to Toroidal GIRGs. Instead of the d -dimensional torus $\chi = \mathbb{T}^d$, we have the d -dimensional cube $\chi = [0, 1]^d$. Thus,

the tradeoff is to lose the symmetry and simplicity of the torus for the ability to hopefully generate more realistic graphs, as few situations in the real world are toroidal i.e. with features that wrap around in the way that a torus does.

For example in our socfb social network graphs, likely geometric features for each individual could be home address, age, athletic proclivity, political leaning, etc. In the environ of a town/city, home address is a 2d vector². Age and athleticism are clearly both non-toroidal (and hence cube-like) quantities. Political leaning could be, e.g., a 2d axis of economic left/right and authoritarian/libertarian. This might bear a toroidal-esque characterisation whereby two individuals on far opposite extremes of an axis share some common ground, e.g., a hardcore communist and a fascist might both share values such as railing against the status quo or enjoying political rallies. No model is perfect; this could potentially be captured as an extra non-toroidal feature dimension placing individuals on a scale from mainstream to unconventional.

2.3.1 Cube GIRG Formulation

The neat correspondence between volume and distance unfortunately breaks down when translating from torus GIRGs to cube GIRGs, due to the loss of spatial symmetry. With the max norm,

$$\text{in the torus } Vol(r_{uv}) = (2 \|x_u - x_v\|_\infty)^d \quad (2.7)$$

$$\text{now in the cube } Vol(B_{r_{uv}}(u)) = Vol(B_{r_{uv}}(u) \cap [0, 1]^d) \quad (2.8)$$

i.e. volume only counts within the cube itself, and hence oftentimes $Vol(B_{r_{uv}}(u)) \neq Vol(B_{r_{uv}}(v))$.

We will use a simpler cube GIRG formulation based solely on distance, i.e. the factor of r_{uv}^d in p_{uv} . This is conceptually easy to understand, and allows for a neat coupling with the original Torus geometry whereby $r_{uv}^{\text{cube}} \geq r_{uv}^{\text{T}}$, with most pairs having $r_{uv}^{\text{cube}} = r_{uv}^{\text{T}}$. We can say that a GIRG in cube geometry is stochastically dominated by its torus geometry counterpart as $p_{uv}^{\text{cube}} \leq p_{uv}^{\text{T}}$, and is hence strictly sparser.

A volume based formulation of edge probabilities is more complicated - it could look something like replacing $Vol(r_{uv}) \mapsto \sqrt{Vol(B_{r_{uv}}(u))Vol(B_{r_{uv}}(v))}$. Intuitively in the social network analogy, this is like saying that all people, no matter how extreme their geometric location (near the edge of the cube), have the same desire to make friends (modulo their inhomogeneous weights as that is the literal introversion / extroversion factor in the GIRG model). If instead of having undirected edges $u \sim v$ we had directed edges $u \rightarrow v$, we could even have a model where $p(u \rightarrow v) \propto \sqrt{Vol(B_{r_{uv}}(u))}$, such that u would in expectation have the same number of outgoing edges whatever its location x_u . Unfortunately desire to make friends doesn't equate to actual friends, as they have to want you back. So if a node u is near the edge of the cube, while it would send out the normal volume based amount of friendship invitations ($u \rightarrow v$), as most of these go to nodes v nearer the centre of the cube, fewer will reciprocate ($v \rightarrow u$). In the directed torus geometry analogy, due to symmetry, every $p(u \rightarrow v) = p(v \rightarrow u)$. This means more pairs of $u \leftrightarrow v$, and fewer pairs of "unrequited friendships" i.e. $u \rightarrow v; v \not\rightarrow u$. Hence a

²Only on the whole planet scale would home address look more torus like - strictly speaking a 3d sphere, or perhaps a cylinder if you assume that nobody lives near to the North/South Pole

volume based cube GIRG in the undirected setting would still have higher average degree than the distance based version, but lower than the torus.

2.3.2 Cube GIRG generation - coupling algorithm

We use the distance based cube GIRG formulation, which permits a simple coupling to a torus GIRG. This motivates algorithm 1 which generates a cube GIRG by appropriately randomly deleting edges from a pre-generated torus GIRG. All we need is any black box torus GIRG generation algorithm that in addition to the graph edge list, also yields the sequence of node weights and locations. algorithm 1 runs in $O(|E|)$ time, additional on top of the torus GIRG generation runtime. As GIRGs have $|E| = \Theta(n)$ whp, this means whp $O(n)$ additional time.

The coupling works for any pair of torus/cube distance functions $r_{uv}^{\text{cube}} \geq r_{uv}^T$, which results in stochastically dominated edge probabilities $p_{uv}^{\text{cube}} \leq p_{uv}^T$. Thus we can couple the events $u \sim v$ in the cube GIRG as a subset of the events $u \sim v$ in the torus GIRG. We just have to inspect each of the $\Theta(n)$ edges in the torus GIRG, and flip a **Bernoulli**($p_{uv}^{\text{cube}} / p_{uv}^T$) coin for each to decide whether to keep / delete the edge for the cube GIRG.

Algorithm 1 Generate cube GIRG from torus GIRG via coupling

Require: n, d, c, τ, α
 $(G, \{x_u\}_{u \in V}, \{w_u\}_{u \in V}) \leftarrow \text{torus-GIRG}(n, d, c, \tau, \alpha)$
for $(u, v) \in E(G)$ **do**
 $p_{uv}^T = \min\{1, c \left(\frac{w_u w_v / W}{(r_{uv}^T)^d} \right)^\alpha\}$
 $p_{uv}^{\text{cube}} = \min\{1, c \left(\frac{w_u w_v / W}{(r_{uv}^{\text{cube}})^d} \right)^\alpha\}$
 $p \leftarrow U[0, 1]$
 if $p > \frac{p_{uv}^{\text{cube}}}{p_{uv}^T}$ **then**
 delete edge (u, v) from G
 end if
end for
return G

2.4 Minimum Component Distance (MCD) GIRGs

In section 2.1 we introduced the max norm $\|\cdot\|_\infty$, euclidean norm, and the minimum component distance (MCD) $\|\cdot\|_{\text{mcd}}$ as alternative distance functions. While the max norm is commonly used in e.g. [Bringmann et al., 2019] as it is handy for proofs, the MCD can make more sense in some settings. For instance in social networks, the max norm is like stipulating that people only make friends with others that “tick all the boxes”, i.e. are similar in every dimension (this sounds more similar to the criteria in finding a long term romantic partner). The MCD is saying that you make friends with people who are similar in at least one dimension. E.g. an individual might play with a neighbourhood football group, or sing in a local choir - each a social circle formed from one shared hobby.

The MCD can be mixed with the max norm in an “and/or” fashion, e.g. on points

in 5d one grouping is

$$\|x - y\|_{\text{mix}} = \min \left(\|x_1 - y_1\|_{\infty}, \|x_{[2,3]} - y_{[2,3]}\|_{\infty}, \|x_{[4,5]} - y_{[4,5]}\|_{\infty} \right) \quad (2.9)$$

This parses as points x, y being “close” if they are close in the 1st dim, or the 2nd and 3rd dim, or the 4th and 5th dim. We call this a min/max mixed GIRG.

2.4.1 Volume Formula for MCD, and Mixed Norms

The formula for $\text{Vol}(r)$ differs according to the distance function $r(x, y)$ and geometric space χ that is in use. For the max norm in d dimensions we saw that $\text{Vol}(r) = (2r)^d$, as the ball B_r in this norm is a cube with side-length $2r$. For the euclidean norm this becomes $\text{Vol}(r) = \frac{\pi^{d/2}}{\Gamma(d/2+1)} r^d$ as the ball B_r is in this case really “ball shaped”. E.g. in $d = 1, 2, 3$ dimensions the ball volumes are respectively $2r$; πr^2 ; $\frac{4}{3}\pi r^3$.

For the MCD, we can derive a volume formula by viewing ball volumes through the lens of probabilistic events. Consider the point set $B_r(x)$, and a location sampled uniformly randomly from the torus $y \sim U(\mathbb{T}^d)$. Let event A_i be that $|y_i - x_i| < r$. Then $\text{Vol}(r) = P(\bigcup_{i=1}^d A_i)$, as the whole torus has unit volume. Since events A_i are mutually independent by orthogonality of the individual dimensions, we have that

$$\begin{aligned} \text{Vol}(r) &= P\left(\bigcup_{i=1}^d A_i\right) = 1 - P\left(\bigcap_{i=1}^d \overline{A_i}\right) \\ &= 1 - \prod_{i=1}^d P(\overline{A_i}) = 1 - (1 - 2r)^d \end{aligned}$$

This can be extended to min/max mixed GIRGs with $(I_i)_{i=1}^k$ a partition of $\{1, \dots, d\}$

$$\begin{aligned} r(x, y) &= \min \left(\|x_{I_1} - y_{I_1}\|_{\infty}, \dots, \|x_{I_k} - y_{I_k}\|_{\infty} \right) \\ \text{then } \text{Vol}(r) &= 1 - \left[\prod_{i=1}^k 1 - (2r)^{|I_i|} \right] \end{aligned}$$

2.4.2 MCD GIRG speedy combination of 1d GIRGs generating algorithm

We introduced in section 2.2 the max torus GIRG generation algorithm implemented in [Bläsius et al., 2022], which has a runtime of $O(n)$ for fixed dimension d . Unfortunately this scales exponentially in d , which hinders our experiments when generating large GIRGs (e.g. $n \geq 10,000$) with $d \geq 4$. We propose a faster algorithm that generates an MCD GIRG by first generating d 1d GIRGs each in $O(n)$ time, and combining them together into one graph in total $O(dn)$ time. This significant speedup could bias a practitioner to using MCD GIRGs when high dimensions are required.

As an example, if we want to generate an MCD GIRG of dim $d = 3$, it actually looks quite like a union of 3 one dimensional GIRGs G_1, G_2, G_3 , each sharing the same node weight sequence, but with independently 1d distributed node locations. Critically, the majority of edges in the resultant 3d MCD GIRG G occur between nodes which happen to be quite close in just one dimension, and not in the other dimensions. This means that the collective number of edges in the three 1d GIRGs is not significantly more than in the end graph, so the complexity of generating these three and handling all their edges won’t blow up our runtime too much.

Algorithm example more rigorously detailed If $\alpha = \infty$ the “union of 1d GIRGs” statement is almost true. For a node pair (u, v) , let r_1, r_2, r_3 be the distances in the 3 1d GIRGs and p_1, p_2, p_3 the corresponding edge probabilities. WLOG letting $r_1 \leq r_2 \leq r_3$, then $\alpha = \infty \implies p_i \in \{0, 1\}$, so in particular $u \overset{3}{\sim} v \implies u \overset{2}{\sim} v \implies u \overset{1}{\sim} v$. This means we could directly take the union over all edges in G_1, G_2, G_3 to produce our 3d MCD GIRG G . In the $\alpha < \infty$ case, to decide if the edge (u, v) is present in the final graph G , we rather must additionally check whether the 1d GIRG with the smallest r_i has an edge. E.g. it could be that $p_1 = 0.5, p_2 = 0.3, p_3 = 0.01$, and that $u \sim v$ only in G_2 - in this case we would take $u \not\sim v$ in G as the minimum component is $i = 1$, thus we copy edge/non-edge for (u, v) from G_1 .

There’s a subtle difference however in that the 3 1d GIRGs actually need edge probabilities with a similar formula to those in the resultant MCD GIRG. This means taking

$$p_i = \min \left\{ 1, c \left(\frac{w_u w_v / W}{1 - (1 - 2r_i)^3} \right) \right\}$$

rather than the standard 1d GIRG edge probability which just has $2r_i$ in the denominator³. This pleasingly gives us that $p_i \leq p_{uv}$, since $r_{uv} := \min_j r_j \leq r_i$ where p_{uv} and r_{uv} are the edge probability and node pair distance respectively in the final 3d MCD GIRG. This means that the expected number of edges in each 1d GIRG is fewer than in the resultant MCD GIRG, which is itself whp $O(n)$ - this matches our “union of 1d GIRGs” intuition.

The number of overcounted (and hence needed to be discarded) edges from the 1d GIRGs won’t exceed the final edge count by more than a factor of d . This depends on the point locations - e.g. in the simple $d = 3$ MCD GIRG case we expect little overcounting, however if the points all lie close to a 1d subspace like $x_1 = x_2 = x_3$, then we would get almost exactly triple counting.

Extension to max GIRGs We would like to use this higher dimension GIRG generating algorithm for max GIRGs as well. Unfortunately, although the algorithm is still correct, its runtime blows up due to the fact that each 1d GIRG will now have many more than $\Theta(n)$ edges. In short, now $r_{uv} := \max_j r_j \geq r_i$, such that $p_i \geq p_{uv} \forall i = 1, \dots, d$ and each 1d GIRG in expectation has more edges than the resultant max GIRG. This would be ok if they were only a constant factor more, but there are actually $\Theta(n^{1-1/d})$ times more edges, which means the runtime of the algorithm is now $O(d \cdot n^{2-1/d})$.

We will just give an informal intuition as to why the 1d GIRGs would now have so many more edges:

For a node u of fixed weight, consider the neighbours of u of small weight w_v , connected to u by a short edge, in both the resultant d-dim max GIRG, and one of the individual 1d GIRGs. In both GIRGs, these edges constitute on average a constant fraction of all of u ’s edges. In the resultant d-dim max GIRG, these come

³This is treacherous as with such an edge probability these “1d GIRGs” might technically no longer be actual GIRGs - however in this section we will still refer to them as such for ease of discussion. Note that having $\Theta(n)$ edges is a key GIRG property, which the “1d GIRGs” in the MCD algorithm still possess; this doesn’t hold in the case of the max norm algorithm.

from a box of side length roughly $2r = \left(\frac{\overline{\deg}}{n}\right)^{1/d}$ about u , resulting in u having about $\overline{\deg}$ neighbours. In contrast, in the 1d GIRG of dimension i , the whole strip of nodes whose projection thereon is in the range $[(x_u)_i - r, (x_u)_i + r]$ will end up as neighbours of u . This gives u a much higher degree in the 1d GIRG, of at least $\widetilde{\deg} = 2nr = \Theta(n^{1-1/d})$ in expectation.

This means that generating a d -dimensional max GIRG with our combination of 1d GIRGs algorithm has runtime at least $O(d \cdot n^{2-1/d})$. This would only beat an exponential in d , linear in n , $O(e^{kd}n)$ runtime for $d > \frac{\log n}{k}$, but at this point the naive $O(n^2)$ algorithm is already preferable.

Extension to mixed GIRGs Although the algorithm doesn't have as good a runtime on max GIRGs, it's easy to see that it does extend well to a mixed GIRG whose distance function has an outer minimum - of the form $r(\mathbf{x}, \mathbf{y}) = \min(\|\mathbf{x}_{I_1} - \mathbf{y}_{I_1}\|_1, \dots, \|\mathbf{x}_{I_k} - \mathbf{y}_{I_k}\|_k)$. We just need a routine to sample the GIRGs in each dimension subset (using the modified edge probabilities that conform to the resultant mixed GIRG), which also returns the distance in that subset for each edge. Then we can collectively consider each of these GIRGs' edges together to construct the final output GIRG.

Therefore we can efficiently sample a high dimensional mixed min/max GIRGs as long as each inner max dimension subset is small, such as $|I_j| \leq 3$.

2.5 Distorted GIRGs

A final GIRG variant that we only play with in chapter 4 is the *distorted GIRG*. This is a GIRG where the node locations' dimensions aren't all treated equally - rather they are ordered in importance $1 > 2 > \dots > d$ via some weighting scheme $1 = \lambda_1 \geq \lambda_2 \geq \dots \lambda_d > 0$. For instance a max norm distorted torus GIRG would have $r(\mathbf{x}, \mathbf{y}) = \max_i \lambda_i |x_i - y_i|_C$. Oddly, to distort a MCD GIRG you instead have to use the inverse weights λ_i^{-1} , as now the smallest dimension distance has the greatest influence on the overall node distance.

The idea of distorted GIRGs is natural - that some dimension are more important than others in influencing edges. For social networks it could be e.g. that socioeconomic status of individuals has much higher impact on their friendships than personality type. The weighted representation is easy to work with but an alternative is to literally distort the space $\chi = \mathbb{T}^d$ or $\chi = [0, 1]^d$ to be cuboidal, i.e. no longer cube shaped.

Chapter 3

Generative Graph Models (GGM) and Blasius Realism Framework

Contents

3.1	GGM and Framework Mathematical Outline	14
3.1.1	Fitting a GGM to a particular real graph instance	14
3.1.2	Plausibility of Fitted GGMs	14
3.1.3	Fitted GGMs as candidates for real graphs - Blasius	15
3.2	Fitting GIRG GGM to a real graph	16
3.2.1	Fitting number of nodes n	16
3.2.2	Fitting power law distribution exponent τ for node weight sampling	16
3.2.3	Power Law weight alternative: weight copying	18
3.2.4	Fitting c and α	19
3.2.5	Volume formulation torus GIRG expected average degree formula	20
3.2.6	Estimating const c in Cube GIRGs	20
3.3	Realism Framework Results	20
3.3.1	Degree distribution features	21
3.3.2	Geometric Features	22
3.3.3	Case study of distribution vs mean accuracy: closeness centrality	24
3.3.4	SVM Classifier Misclassification Rates Tables	24

[Bläslius et al., 2018] presents a general framework that can be used to assess the ability of a *generative graph model (GGM)* to simulate a set of real graphs of interest. They demonstrate their method on a few common GGMs: *Erdos-Renyi (ER)*, *Barabasi-Albert (BA)* preferential attachment graphs, *Chung-Lu* and *Hyperbolic Random graphs*. The real graph dataset which they attempt to simulate is a set of 219 real-world networks taken from the *Network Repository* (<https://networkrepository.com/>, [Rossi and Ahmed, 2015]).

Their framework assesses realism of a GGM by using it to produce a dataset of fake graphs. A metric of distinguishability of fake and real datasets based on a select few graph statistics is produced: training a binary classifier of graph \rightarrow fake/real

using the combined fake and real graph datasets as training data, the accuracy of this classifier signifies real/fake distinguishability.

[Bläsius et al., 2018] found poor realism of the tested models on the ~ 100 strong subset of the real-world graphs that were Facebook social networks. We however will use the framework to specifically test different types of GIRGs’ realism on this tough subset, with the hope that they will do better than the best models of [Bläsius et al., 2018] - Chung-Lu and Hyperbolic. We also think that restricting to this one subset of what should be similar types of real-world graphs will more clearly distinguish our different GIRG types.

3.1 GGM and Framework Mathematical Outline

A graph $G = (V, E)$ is said to be randomly generated from a GGM \mathcal{G} , written $G \sim \mathcal{G}$. Most of our GGMs are simple and can be described by a small number of parameters. For example for the Erdos-Renyi GGM, $G \sim \mathcal{G}_{ER}(n, p)$ means that there are n nodes: $V = [n]$, and each potential edge (u, v) exists iid with probability p . We write this as $P(u \sim v) = p$.

3.1.1 Fitting a GGM to a particular real graph instance

“Fitting” a GGM to a real graph G means finding the most plausible $\hat{\theta}$ in the hypothetical world where G was produced via $G \sim \mathcal{G}(\hat{\theta})$. Hence for our ER GGM example, choosing $\hat{n} = |V|$ is a no-brainer, and $\hat{p} = |E|/\binom{n}{2}$ follows by fitting the expected number of edges.

Fitting $\hat{\theta}$ can also be done by likelihood estimation. This is tractable for the Erdos-Renyi GGM, by solving $\arg \min_p \sum_{u \neq v} \log(p^{e_{uv}} + \log((1-p)^{1-e_{uv}}))$, which gives the same \hat{p} as above. For GIRGs however, this would be intractable: for example we cannot even easily integrate over all possible point locations in the geometric space:

$$p(G|\alpha, w, d) = \int_{x_1 \in [0,1]^d} \cdots \int_{x_n \in [0,1]^d} p(G|\alpha, (x_i)_{i=1}^n) \prod_{i=1}^n p(x_i)$$

Instead we can use a heuristic method based on Approximate Bayesian Computation (ABC). We seek the parameter $\hat{\theta}$ that minimises the expected distance $\mathbb{E}[\delta(G, G' \sim \mathcal{G}(\theta))]$ for some distance metric δ . For δ to be effective, ideally it would be something like $(f(G) - f(G'))^2$, for f a sufficient statistic of θ . In practice this looks like an algorithm where we repeatedly sample θ from a prior; use it to generate G' ; see if $f(G') \cong f(G)$, and if so keep θ as a candidate for $\hat{\theta}$.

3.1.2 Plausibility of Fitted GGMs

Similarly to the use of a statistic f of real and generated graphs as a proxy for fitting $\hat{\theta}$, we can also use other high level statistics f^{class} to evaluate the plausibility of a fitted GGM as to whether it could have produced the real graph G .

Having fit $\hat{\theta}$ of the GGM, we can then sample graphs $G' \sim \mathcal{G}(\hat{\theta})$ which are hypothetically similar to G . We can only hope for similarity on the global level,

as even in the best case scenario where both $G, G' \sim \mathcal{G}(\hat{\theta})$, due to the inherent randomness of the GGM, we cannot expect individual edges to matchup e.g. $e \in E \iff e \in E'$, let alone to be able to identify nodes in G with those in G' - see chapter 6 for more on graph similarity.

For example in the Erdos-Renyi exact fit case, in the best case scenario we have $G \sim \mathcal{G}_{ER}(n, p^*)$ and $G' \sim \mathcal{G}_{ER}(n, \hat{p} = p^* + \varepsilon)$. Although there will be some small estimation error ε in the fit parameter (which may be biased to be non-zero if e.g. G randomly has $|E|$ differing from $p^* \binom{n}{2}$), we've explicitly fit for the high level statistic of the number of edges, s.t. $|E| = \mathbb{E}[|E'|] = \hat{p} \binom{n}{2}$, and $|E'| \approx |E|$ with some small variance. For other (not fit) high level statistics f^{class} like the effective diameter of the graph, their expected value will close to that of the real graph, $E_{G' \sim \mathcal{G}(\hat{\theta})}[f^{\text{class}}(G')] \approx f^{\text{class}}(G)$, and likewise concentrated with small variance s.t. $f^{\text{class}}(G') \approx f^{\text{class}}(G)$.

Furthermore if we were to generate a whole list of graphs G'_1, \dots, G'_m from $\mathcal{G}(\hat{\theta})$, assuming the matching GGM hypothesis, we would not expect G to stand out from the crowd. Essentially $f^{\text{class}}(G'(G)) \sim f^{\text{class}}$ should follow a distribution with $f^{\text{class}}(G)$ close to its median/mean.

3.1.3 Fitted GGMs as candidates for real graphs - Blasius

[Bläslius et al., 2018] takes this method one step further. They are essentially evaluating the hypothesis that the set of real graphs G_1, \dots, G_k are generated from a particular GGM \mathcal{G} , but with differing parameters: $G_1 \sim \mathcal{G}(\theta_1^*), \dots, G_k \sim \mathcal{G}(\theta_k^*)$. Perhaps θ_i^* come from some prior distribution $p(\theta)$.

They again fit individual $\hat{\theta}_i$ with which to generate one $G'_i \sim \mathcal{G}(\hat{\theta}_i)$ per real graph G_i . This way instead of, for each real graph, comparing multiple G'_{ij} to G_i and averaging, they can compare the set of graphs $\{G_i\}_{i=1}^k$ to $\{G'_i\}_{i=1}^k$. Graph comparison is done by comparing a wide number of statistics/metrics computed for each graph. These are input as a whole feature vector $f^{\text{class}}(G)$ to an SVM classifier, which is trained to classify membership of the real or generated dataset. These features are statistics such as the number of nodes, average node degree, centrality, closeness, diameter, clustering coefficient and so on.

Aggregating the classification over the whole dataset of real graphs allows a

We suggested in the previous section that multiple G'_{ij} could be generated per real graph G_i - this would reduce the variance of the whole process. This is also necessary to be able to train a classifier on an individual real graph basis, as you cannot do binary classification on just a pair of feature vectors $f^{\text{class}}(G_i)$ and $f^{\text{class}}(G'_i)$. Blasius instead aggregates the classification over the whole dataset of real graphs, allowing sufficient datapoints to train a classifier while reducing the overall necessary computation (k generated graphs produced instead of km).

This also has the added benefit of helping to cover for GGM parameter fitting inaccuracy. If $\hat{\theta}_i = \theta_i^* + \varepsilon_i$, it may still be possible to distinguish G_i from $\{G'_{ij}\}_{i=1}^m$ due to the ε_i fitting error. If instead we compare the set $\{\theta_i^*\}_{i=1}^k$ to $\{\hat{\theta}_i\}_{i=1}^k$, Now each $\hat{\theta}_i \neq \theta_i^*$, but still well fits into the distribution of $\{\theta_i^*\}_{i=1}^k \sim p(\theta)$. Finally having one G'_i per real graph G_i also simplifies the binary classification task by having a balanced dataset.

For the SVM classifier trained on the mirrored real/fake graph dataset, it will have close to 50% accuracy if really $G_i \sim \mathcal{G}(\theta_i^*)$, and higher accuracy if instead $G_i \sim \mathcal{G}(\phi_i^*)$ some alternative GGM \mathcal{G} .

3.2 Fitting GIRG GGM to a real graph

3.2.1 Fitting number of nodes n

We follow [Bläsius et al., 2018] which first preprocesses the graph G by shrinking it to its largest connected component, i.e. $G \leftarrow \text{shrinkToGCC}(G)$ before fitting \mathcal{G} to it. Then n is just the number of remaining nodes.

Bläsius’ rational is that where some real graphs in our dataset have disconnected subgraphs, this may be due to them being a concatenation of a few distinctly generated subgraphs, which may not collectively fall under the GIRG model. All of our GGMs are capable of producing a bunch of disconnected subgraphs, however this is most likely to occur in our geometric GGMs which exhibit clustering; GIRGs do at least whp produce a unique giant component - one with a linear number of nodes. Hence for a fair comparison $d(G, G')$, we also need to post-process the GGM generated $G' \leftarrow \text{shrinkToGCC}(G')$.

Bläsius’s only geometric model is Hyperbolic Random Graphs, for which they use a fitting algorithm that actually estimates a higher number of nodes $n > |V|$ in order to approximately have the largest connected component of $G' \sim \mathcal{G}$ be of size $|V|$. We don’t do this for our GIRGs however, as we find this algorithm prone to error, and unnecessary, at least for the socfb Facebook graphs - see table 3.1. We accept instead that our GGMs may end up with slightly fewer nodes than the real graphs, and hope that this doesn’t affect the SVM classification too much.

Restricting graphs to a connected component also has the added benefit of making some graph statistics more meaningful/sensical - for instance diameter and path lengths.

Note that GIRGs might still do the best job, compared to other non geometric GGMs, of fitting multiple large components, as they have the property of containing sublinear separators. Essentially if you divide the torus with a hyperplane (or divide out a sub-cube), $\chi = \chi_A \sqcup \chi_B$, to produce disjoint node subsets $V_A \sqcup V_B = V$ and their induced subgraphs G_A, G_B , then whp G_A has sublinear of $|V_A|$ edges to G_B (despite having linear number of edges within itself). Also interestingly G_A is stochastically still a (smaller) GIRG, just one with a non-toroidal geometry.

3.2.2 Fitting power law distribution exponent τ for node weight sampling

We fit τ to the tail of the *degree distribution* of G , using the python package `powerlaw`. Using this fit parameter we are enabled to generate new weight sequences $w_u \sim \text{powerlaw}(\tau)$ for sampling new GIRGs $G' \sim \mathcal{G}_{\text{GIRG}}(\tau, \dots)$.

The degree distribution of the graph G is defined as $dd(x) := \frac{|\{v \in V: d(v)=x\}|}{|V|}$, i.e. the fraction of nodes with degree x . For graphs generated by the GIRG model, we saw that the geometry marginalisation property means that $E[d_u] \propto w_u$. Hence if

graph name	GGM	nodes
socfb-American75	real-world	6370
socfb-American75	1d-girg	6370
socfb-American75	2d-girg	6370
socfb-American75	3d-girg	6370
socfb-American75	ER	6370
socfb-American75	chung-lu	6279
socfb-American75	hyperbolic	6583
socfb-Amherst41	real-world	2235
socfb-Amherst41	1d-girg	2235
socfb-Amherst41	2d-girg	2235
socfb-Amherst41	3d-girg	2235
socfb-Amherst41	ER	2235
socfb-Amherst41	chung-lu	2221
socfb-Amherst41	hyperbolic	2282
...
bio-diseasome	real-world	516
bio-diseasome	1d-girg	258
bio-diseasome	2d-girg	491
bio-diseasome	3d-girg	496
bio-diseasome	ER	512
bio-diseasome	chung-lu	459
bio-diseasome	hyperbolic	125

Table 3.1: For the socfb Facebook graphs, the hyperbolic model consistently has a few more nodes than the input real graph due to its fitting algorithm not quite working perfectly (and stochasticity). On other real graphs there can be larger discrepancies, especially for smaller extra sparse graphs. Numbers of nodes in the output (shrunk to GCC) graph are colored **cyan** if less than the real-world graph, and **orange** if more (only possible in hyperbolic case).

node weights are distributed following $\text{powerlaw}(\tau)$, we expect to see the tail of the degree distribution, $dd(x)$ as $x \rightarrow n$ to look like a discrete power law distribution $dd(x) \propto x^{-\tau}$.

Sketch of why GIRG power law weights \rightarrow power law degree distribution By treating node degrees as roughly iid, $dd(x) \xrightarrow{n \rightarrow \infty} P(d_u = x)$. For a fixed weight w_u , node degree d_u closely follows a binomial distribution $\text{binomial}(n-1, \Theta(w_u/n))$, so we'd like to calculate

$$P(d_u = x) = \int P(d_u = x | w_u = w) \Theta(w^{-\tau}) dw \quad (3.1)$$

For sufficiently large x (in the tail of the decaying distribution of $dd(x)$), the binomial distribution of mean and variance w is sharply peaked around its mean: if $(x - w) > w^{1/2+\epsilon}$ then $P(d_u = x | w_u = w) \approx 0$. So essentially all the probability mass is on the region $[x - x^{1/2+\epsilon}, x + x^{1/2+\epsilon}]$, on which $\Theta(w^{-\tau}) \approx \Theta(x^{-\tau})$ doesn't change much at all as e.g. $(x - x^{1/2+\epsilon})^{-\tau} = (x[1 - x^{-1/2+\epsilon}])^{-\tau} \cong x^{-\tau}$ for large x . Hence $dd(x) \xrightarrow{n \rightarrow \infty} \Theta(x^{-\tau})$ for large x .

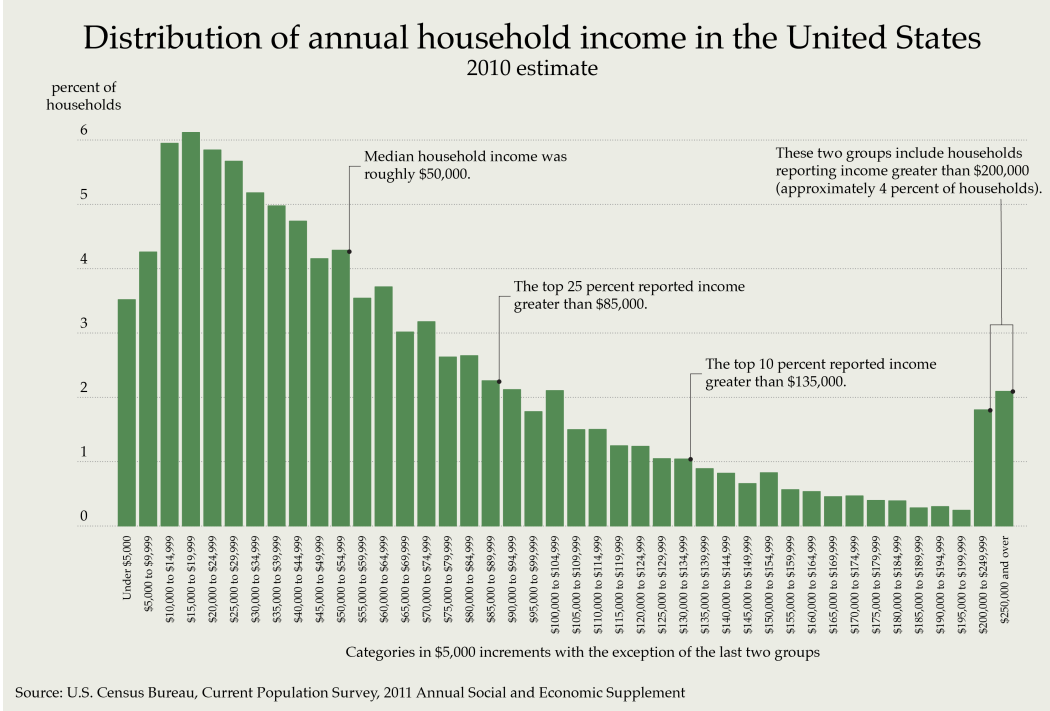


Figure 3.1: Example of a more natural weight distribution with a power law tail but not for smaller weights.

powerlaw package fitting Therefore if G is a τ exponent power law weighted GIRG, we expect its degree distribution $dd(x)$ for large degrees x to look like $dd(x) \propto x^{-\tau}$. Hence we can fit τ on this tail. `powerlaw` does this by simultaneously finding a lower bound x_{\min} for where in the tail the power law behaviour starts, and fitting τ to the region $x > x_{\min}$ using maximum likelihood estimation. The optimal x_{\min} is chosen to minimise the Kolmogorov-Smirnov distance between the power law fit and the region's empirical degree distribution.

Having fit τ , we can then sample weights $w_u \sim \text{powerlaw}(\tau)$.

3.2.3 Power Law weight alternative: weight copying

Generating weights $w_u \stackrel{iid}{\sim} \text{powerlaw}(\tau)$ is fine as a model prior, however it's not a perfect fit to real-world data. This is particularly true as real graph degrees generally only follow a power law for large degrees; they might be better modelled by a GIRG with node weights with a power law tail (e.g. upper quartile of weights) and a different distribution for the rest of the weights. Most obviously, instead of having peak number of nodes with weight $w_u = x_{\min}$, the weight distribution might follow a more natural curve like in fig. 3.1.

An easy improvement to better fit a specific real graph is to take the sequence of node degrees as weights¹. For the classification comparison framework, Blasius actually uses weight copying for fitting the Chung-Lu GGM, but not for the hyperbolic GGM (which is odd). This doesn't make a fair / like for like comparison, so we cover all bases by having both weight copied and power law fit GIRGs, as

¹these are now all positive integers ≥ 1 as opposed to real numbers $\geq x_{\min}$, since in the largest connected component, minimum degree is 1

well as adding in power law fit Chung-Lu. ER and BA GGMs are still just fit on the courser metric of overall graph average degree.

Comparing weight copied Chung-Lu with weight copied GIRGs was additionally interesting for harder feature sets where other GGMs struggled to get below 98% classification accuracy - having such high constrained numbers makes it harder to compare meaningfully between GGMs (drawing conclusions from 98% vs 99% is tough)

3.2.4 Fitting c and α

These are our last two GIRG parameters to fit. The parameter c is most directly linked to the generated graph's average degree, whereas the parameter α is more linked to the power of the geometry - larger α decreases the probability of longer distance edges which have $\rho_{uv} < 1$, leaving edges more dominated by shorter (and hence more geometric and clustered) edges. We therefore follow [Bläsius et al., 2018] by fitting α with the clustering coefficient. Regrettably we used in our experiments the mean local clustering coefficient (LCC) to fit our GIRGs and Hyperbolic Random Graphs,

$$\text{LCC}(G) := \frac{1}{|V|} \sum_{u \in V} \frac{|\{v, v' \in \Gamma(u) : v \sim v'\}|}{\binom{|\Gamma(u)|}{2}} \quad (3.2)$$

instead of the global clustering coefficient as in [Bläsius et al., 2018]. The global clustering coefficient is the fraction of total "V" shapes that indeed are completed as full triangles and is likely a better metric than the LCC. Luckily in practice they generally differ by less than 1%.

Unfortunately c and α are not wholly independent, so we must fit them in tandem, fitting c for a given α so as to match the average degree of G , and then α for that given c to match the LCC of G . This then looks like coordinate ascent in 2D: we alternately set $c \leftarrow \hat{c}_1; \alpha \leftarrow \hat{\alpha}_1; c \leftarrow \hat{c}_2; \alpha \leftarrow \hat{\alpha}_2; \dots$

Fitting $\alpha \leftarrow \hat{\alpha}_i$ based on LCC is precisely ABC like - we propose potential α_i values using binary search, for which we generate a graph $G' \sim \mathcal{G}_{\text{GIRG}}(n, d, \hat{c}_i, \alpha_i, \tau)$, and use the distance metric $\delta(G, G') = |\text{LCC}(G) - \text{LCC}(G')|$ to both propose the next candidate α_i and eventually accept the final best $\hat{\alpha}_i$. Expected LCC is monotone in α so we know which direction to make a new proposal - we follow [Bläsius et al., 2018] in doing binary search on $\alpha^{-1} \in (0, 1)$.

[Bläsius et al., 2022] luckily gives a more efficient method to fit c given α, d , and a pre-sampled set of weights $(w_u)_{u \in V}$, that doesn't involve fully sampling a new $G' \sim \mathcal{G}_{\text{GIRG}}$. They derive a formula for the expected average degree $\mathbb{E}[\text{deg}(G')]$ of the GIRG which looks like, emphasizing the dependence on c ,

$$\mathbb{E}[\overline{\text{deg}}(G')|c] = cA + c^{1/\alpha}B \quad (3.3)$$

Hence we just have to numerically solve eq. (3.3) for $\hat{c} : \mathbb{E}[\overline{\text{deg}}(G')|\hat{c}] = \overline{\text{deg}}(G) = 2|E|/|V|$ in order to obtain our desired average degree.

The expected average degree formula of eq. (3.3) miraculously holds for all volume based toroidal GIRGs (regardless of exact distance function $r(x_u, x_v) = r_{uv}$), and we can even adapt the formula to be independent of dimension d .

3.2.5 Volume formulation torus GIRG expected average degree formula

We'll calculate $E_r[p_{uv}(r)|c]$ for a general volume function $Vol(r)$, and fixed weight sequence $(w_u)_{u \in V}$ and use this to recover eq. (3.3) via

$$E[\overline{deg}(G')|c] = \frac{1}{n} \sum_{u \in V} \sum_{v \in V} E_r[p_{uv}(r)] \quad (3.4)$$

Now $E_r[p_{uv}(r)] = \int_r p_{uv}(r) p(r) dr$. We'll break this down into $\int_{r:\rho_{uv} \geq 1} + \int_{r:\rho_{uv} < 1}$ and write $\hat{r} : Vol(\hat{r}) = c^{1/\alpha} (\frac{w_u w_v}{W})$ as the boundary at which $\rho_{uv} = 1$.

Hence we get $E_r[p_{uv}(r)] = \int_0^{\hat{r}} p(r) dr + \int_{\hat{r}}^{r_{\max}} p_{uv}(r) p(r) dr$.

Substituting $Vol = Vol(r)$; $dr = dVol \frac{dr}{dVol} = \frac{dVol}{p(r)}$, we get:

$$E_r[p_{uv}(r)] = \int_0^{Vol(\hat{r})} dVol + \int_{Vol(\hat{r})}^{Vol(\mathbb{T})} p_{uv}(Vol) dVol \quad (3.5)$$

$$= Vol(\hat{r}) + \int_{Vol(\hat{r})}^{Vol(\mathbb{T})} c \left(\frac{w_u w_v}{W} \right)^\alpha Vol^{-\alpha} dVol \quad (3.6)$$

$$= Vol(\hat{r}) + c \left(\frac{w_u w_v}{W} \right)^\alpha \frac{1}{\alpha-1} \left[\left(\frac{w_u w_v}{W} \right)^{1-\alpha} c^{\frac{1-\alpha}{\alpha}} - 1 \right] \quad (3.7)$$

$$= c^{1/\alpha} \left(\frac{w_u w_v}{W} \right) \left[1 + \frac{1}{\alpha-1} \right] - \frac{c}{\alpha-1} \left(\frac{w_u w_v}{W} \right)^\alpha \quad (3.8)$$

Where in the last two lines we sub in $Vol(\hat{r}) = c^{1/\alpha} (\frac{w_u w_v}{W})$, and $Vol(\mathbb{T}) = 1$. I.e. across all GIRGs of different distance functions and dimensions d , we get the same resultant edge probabilities when marginalising out node locations. Note this is also $\Theta(\frac{w_u w_v}{W})$ since $\alpha > 1$, which shows the similarity between GIRGs and Chung-Lu!

Following [Bläsus et al., 2022] Appendix, the formula just needs a correction for pairs u, v such that $Vol(\hat{r}) = c^{1/\alpha} (\frac{w_u w_v}{W}) > 1$, for which the second integral is unnecessary, and the first's upper bound is capped lower at 1, the volume of the torus. These are pairs u, v which have such giant weights w_u, w_v that $p_{uv}(r) = 1$ identically, no matter how far apart x_u, x_v are placed in the torus.

3.2.6 Estimating const c in Cube GIRGs

To fit c for a specific real graph G , for torus GIRGs we solved for the equation $f(c) = \overline{deg}$, given fixed α . Our distance based Cube GIRGs have fewer edges in expectation than their toroidal counterparts, and no nice formula for the expected average degree. Instead we estimate \hat{c} by starting with an initial guess $c_0 : f(c_0) = \overline{deg}$, and iteratively updating by each time generating a graph $G_i \sim \text{GIRG}(c_i)$ and setting $c_{i+1} \leftarrow c_i \frac{n \overline{deg}}{2|E(G_i)|}$ until convergence.

3.3 Realism Framework Results

As introduced in section 3.1.3, we fit a selection of GGMs, including various different kinds of GIRG on around 100 socfb Facebook graphs, whose sizes range from 1000 to 100,000 nodes. For each GGM we produce a mirrored dataset of fake graphs, with which we train a sequence of SVM classifiers to distinguish between the real and fake datasets, differing on which high level graph features

they use as input. Classification accuracy using a particular dataset for a GGM could be e.g. 100% if the fake graphs are highly unrealistic, or as good as 50% if they're indistinguishable from the real graphs. In fig. 3.5 and fig. 3.7 the numbers displayed are actually the *Misclassification rate*, i.e. $1 - \text{accuracy}$, so that higher is more distinguishable.

We present in fig. 3.5 the harder classification test of using a whole distribution of node level features - i.e. for one node level property "deg" would actually include 5 features of the distribution of node degrees: mean, median, lower quartile, upper quartile and standard deviation. In contrast the classification results in fig. 3.7 use just the mean node level features, for which it is easier to fool the classifier. Having more significantly higher than zero misclassification rate numbers in fig. 3.7 allows more comparison between different GGMs and across a wider range of features.

Types of GIRG tested in classification framework

- The "base" GIRG is max norm torus GIRGs for dimension $d = 1, 2, \dots, 7$
- MCD torus GIRGs for $d = 2, 3, 4, 5$ ($d = 1$ is equivalent to max norm)
- Min/max mixed torus GIRGs with groupings 1-23; 1-234; 12-34; 1-2-34. E.g. 1-23 denotes $\|x - y\| = \min\{|x_1 - y_1|_C, \max(|x_2 - y_2|_C, |x_3 - y_3|_C)\}$
- Max norm cube GIRGs - though we were held back due to the increased computation when fitting the average degree with constant c . Hence we just had $d = 1, 2, 3$ cube GIRGs, and $d = 1, 2, 3, 4, 5$ copy-weight cube GIRGs (an attempt for maximal realism)
- The host of simpler models CL, BA, ER as in [Bläsius et al., 2018], as well as Hyperbolic Random Graphs which are very similar to 1d torus GIRGs

New features We additionally added features that weren't assessed in [Bläsius et al., 2018]. Shown in fig. 3.5 are the non node-level "k-cores, comms". These denote k-core sizes and community sizes (communities found using the Louvain method). However they can be rather high variance numbers, for example a 3000 node 2d max torus GIRG with average degree 60 can end up with 40 k-cores and just 8 communities. This means that the distribution features of community size would be the average size of these 8 communities, as well as lower, middle and upper quartile size and standard deviation of these sizes.

Unfortunately for community sizes, even Erdos-Renyi graphs have a misclassification rate comparable to the rest, so this feature seems useless.

3.3.1 Degree distribution features

As noted in [Bläsius et al., 2018] features like Katz centrality, PageRank, k-core number / k-core sizes, node degrees and node degree power law are more related to the degree distribution. This means they are less useful for distinguishing between different types of GIRG geometries.

For this group we observe much better performance for copy-weight GIRGs / copy-weight Chung-Lu. For all these groups except "n, m, Katz" in fig. 3.5 (distribution features), non copy-weight GGMs have all zero misclassification rates.

They also feature monotonically increasing rates with increasing GIRG dimension, as this moves the graphs to be closer to the Chung-Lu degree distribution, see [Friedrich et al., 2023].

For just mean features there are some non-zero rates for non-copy-weight GGMs, but generally worse numbers than the copy-weight GGMs. For whatever reason, non-copy-weight 3d cube GIRGs are much better on k-core (16%) and decently better on Katz centrality (17%) than the other GIRG types. This could be due to a more realism of extreme nodes near the edge of the cube that both have a lowered degree, and whose neighbours are also similarly low degree and extreme. The copy weight GGMs do get up to 47% on Katz, however oddly copy weight cube GIRGs have just 7 – 9% misclassification rate.

3.3.2 Geometric Features

The features betweenness centrality, closeness centrality, local clustering coefficient, diameter and effective diameter are more related to the geometry. GIRGs show on these features increased reality in comparison to the non geometric GGMs, and we can also compare between the different max / MCD / mixed; Torus / Cube GIRG types. We focus solely on the non-copy-weight GIRGs, with numbers from the mean based classifiers in fig. 3.7 rather than the distribution based classifiers in fig. 3.5 which have too low numbers for a good comparison. In the former, the copy-weight cube GIRGs have similar performance on these features - their more realistic degree distribution only help in the latter harder simulation task.

Chung-Lu graphs are a good comparison point against GIRGs by being a similarly inhomogenous graph but without any geometry.²

Comparing between non-copy-weight GIRGs, 1d max torus GIRGs are surprisingly best or close to best on many feature sets. This is evidence that the geometric component of real graphs is most strong in just one dimension, which we also see in chapter 5. It would be interesting to see if 2d or 3d distorted GIRGs could do better on these features, as they would precisely exhibit geometry with one dominating dimension and a few auxiliary. On betweenness, 1d max torus GIRGs perform the best, with a slight decay for 2d-4d. Mixed GIRGs hold up well, but MCD GIRGs do much worse for some reason; cube GIRGs are a little worse. However MCD and mixed GIRGs do better on closeness than max GIRGs, especially 1-234 mixed GIRGs. This is actually evidence for the “one strong dimension” hypothesis, as 1-234 mixed GIRGs give more significance to their dim 1 over dims 2,3,4.

Effective diameter and cube GIRGs Cube GIRGs generally underperform, except on effective diameter, where they shine through as the best, with $1cu < 2cu < 3cu$. This makes sense as cube geometry doesn’t impact many shortest paths, but rather does so on the longest shortest paths, between nodes at opposite extremes of the cube which would otherwise have torus shortcutted. We can dive into see that our intuition is correct: real effective diameters of FB graphs range from 3.0 to 4.4; 2d

²Another potential GGM geometry null hypothesis would have been the configuration model which tries to mimick a degree sequence, like Chung-Lu, but more strictly. Each node is given a set number of “half edges” which gives it a set degree (the degree sequence can be random e.g. power law generated, or given e.g. copied directly from a real graph. Half edges are then joined together at random.)

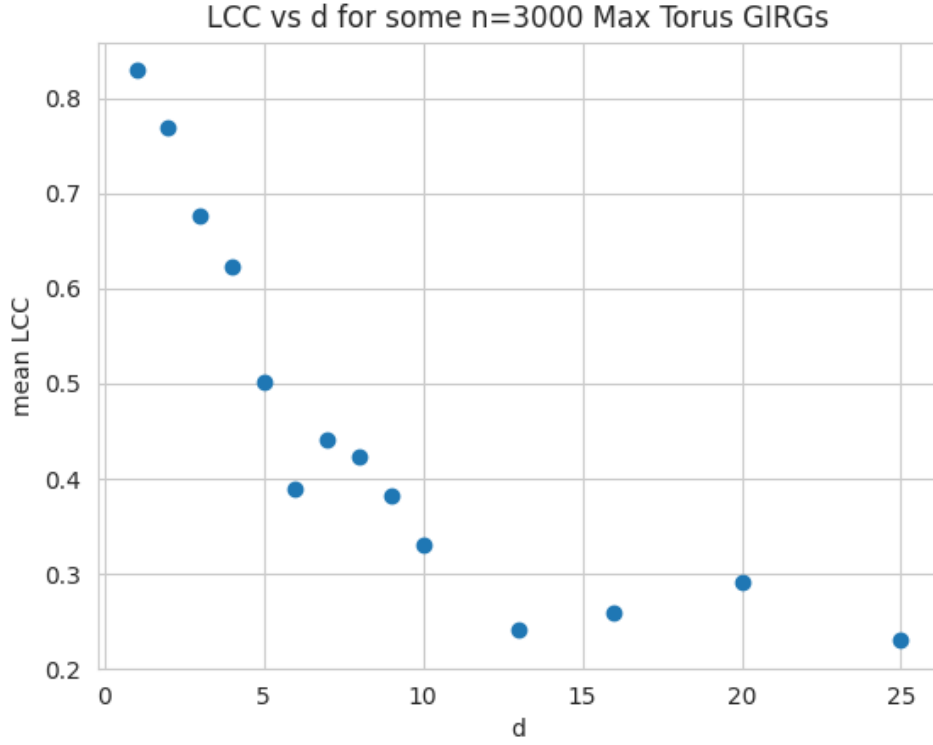


Figure 3.2: Increasing the dimension d of a GIRG will decrease the mean LCC, but only up to a $\Theta(1)$ threshold due to the triangle inequality. Each point is one generated GIRG, all other parameters remaining fixed.

torus GIRGs are too small, on average 0.17 less than their real counterparts, while cube GIRGs are larger and on average just 0.04 less.

Local Clustering Coefficient (LCC) due to their geometry, GIRGs exhibit significant clustering just like real graphs. We explicitly match for the mean LCC in our GIRG fitting procedure, so it’s not surprising that the “n, m, LCC” row in fig. 3.7 has near perfect performance for most GIRGs, whereas all other GGMs have 0% misclassification rate as they have tiny mean LCC. This drops to 0% for all GIRGs except copy-weight cube GIRGs when classifying on LCC distribution. [Bläsius et al., 2018] explained this by finding their fit hyperbolic graphs to have too low variance, which is likely the case here too.

For mean LCC only 4d, 5d MCD GIRGs and 6d, 7d max GIRGs have low misclassification rate, due to their inability to match the real graphs’ high mean LCC. E.g. on socfb-Caltech36, the real graph has mean LCC of 0.41 , however capping at $\alpha = 100.0$ ($\alpha = \infty$ wouldn’t improve the situation much either), the max norm 4d GIRG can only reach 0.33 , and the 5d GIRG just 0.24 . LCC is guaranteed to be non-tiny due to the triangle inequality: if $u \sim v_1$ and $u \sim v_2$ then likely $r(u, v_1), r(u, v_2)$ are small, which makes $r(v_1, v_2) \leq r(u, v_1) + r(u, v_2)$ also small. However for higher dimension GIRGs, $r(v_1, v_2)$ is more likely to be closer to this upper bound. LCC decays already for 4d MCD GIRGs as the MCD distance function isn’t actually a metric so only obeys the triangle inequality stochastically (some percentage of the time. E.g. if u is close to v_1 in dimension 1, and close to v_2 in dimension 2 then v_1, v_2 are generally not close to each other).

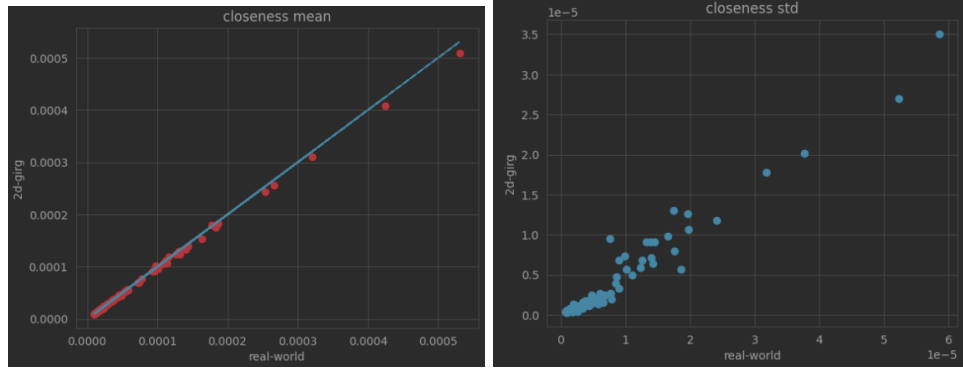
3.3.3 Case study of distribution vs mean accuracy: closeness centrality

In fig. 3.5 which provides feature distributions as input to the SVM classifier, it's disappointing that non copy-weight GIRGs have such poor performance, even on "geometric" features like closeness, betweenness and LCC (at least on closeness and betweenness they still recognizably outperform CL, BA and ER). E.g. 2d max Torus GIRGs have a misclassification rate of just 1% and 2% on "n, m, betw" and "n, m, close", compared to much better 38% and 28% on the mean based classifier equivalents. To explain this large disparity we do a case study on closeness, which happily also assuages some of our doubts.

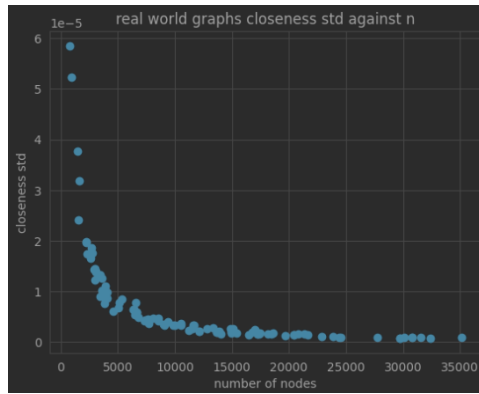
fig. 3.3a shows that mean closeness fits very accurately to the real graphs. Nonetheless, due to a consistently lower standard deviation of closeness in the GIRGs and the fact that the standard deviation of closeness in real graphs fits well as a function of number of nodes, the 2d GIRGs are distinguishable from real graphs on this feature set. We hence find the 2% misclassification rate misleadingly low, and the GIRG's realism in closeness somewhat vindicated.

Unfortunately classification accuracy is a crude metric. One drawback of this framework is that the attainable accuracy is also affected by the level of variance of a feature set across the real graph dataset relative to the fake dataset. GGMs generally end up with their features concentrated in a small subsection of the feature space, so if the real graphs are more spread out, it may be easier to overlap with a portion of their feature data points and attain a low but appreciable misclassification rate. However if the real graphs are less spread out they can miss altogether and get a tiny misclassification rate. See fig. 3.4 for an illustration of different scenarios. In this case of closeness, the real graph feature vectors fit quite a precise pattern and hence the GIRG feature vector clump, though just slightly off, has little overlap.

3.3.4 SVM Classifier Misclassification Rates Tables



(a) 2d GIRGs have almost identical mean node closeness to real-world FB graphs; $y = x$ line shown in blue. (b) 2d GIRGs have consistently lower standard deviation of node closeness than real-world FB graphs



(c) Standard deviation of node closeness in real-world FB graphs fits well as a function of number of nodes

Figure 3.3: Closeness centrality case study on fit 2d max torus GIRGs against their real counterparts. In plots (a) and (b), each dot is a real / fake graph pair, with the x/y value being the specifid feature's vale for the real/fake graph respectively. In plot (c), each dot is a real graph.

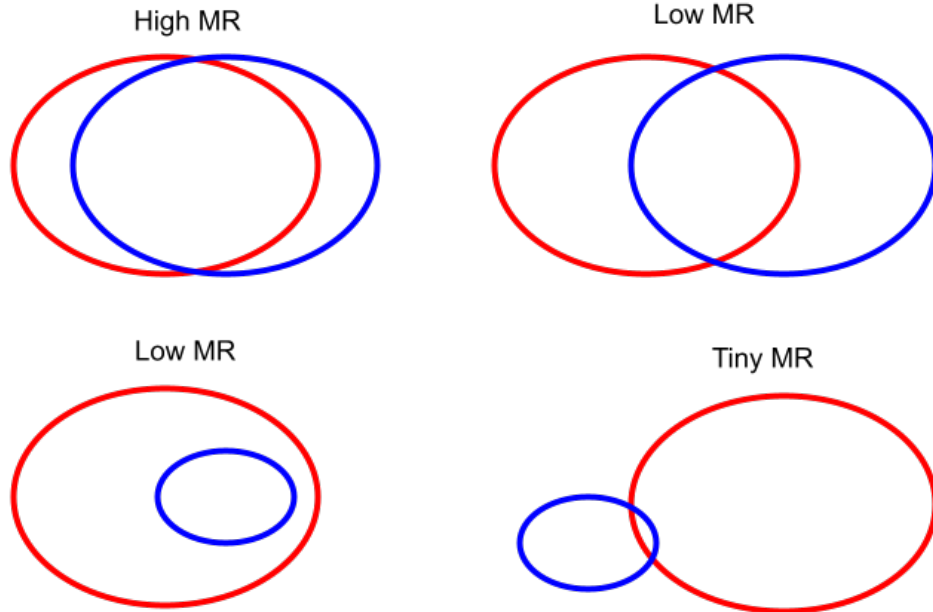


Figure 3.4: Illustration of SVM misclassification rates (MR) given different positions and variance of binary clusters. Red and blue ovals represent roughly gaussian distributed clusters of real and fake data points in a 2 dimensional feature space. An SVM would divide the space so as to maximise classification accuracy. How high accuracy it achieves is directly related to the area/density of overlap between the two clusters.

Feature Set	1ccu	2ccu	5ccu	CL-c	1-23	1-234	2m	5m	1cu	2cu	3cu	1d	2d	7d	CL	BA	ER	hyper
n, m, betw	18	16	5	5	1	1	2	0	1	0	0	1	1	1	1	0	0	1
n, m, close	8	15	16	10	1	0	0	3	2	3	5	2	2	2	1	1	0	3
n, m, LCC	4	7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n, m, diam	4	3	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n, m, eff-diam	22	25	22	9	8	17	2	0	17	23	24	13	16	13	0	0	0	17
n, m, k-core	6	9	10	14	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n, m, deg	14	23	34	46	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n, m, Katz	25	25	36	45	6	6	6	9	7	3	1	6	7	5	3	0	1	6
n, m, PR	20	25	23	23	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n, m, comms	3	4	4	9	2	2	1	1	6	6	6	5	3	11	10	5	8	1
n, m, k-cores	6	8	14	17	0	0	0	0	0	0	0	0	0	0	0	0	0	0

Figure 3.5: [Bläsius et al., 2018] GGM comparison framework on Facebook graphs, extended to different GIRG models. Some models are missing from the table e.g. 3d-6d GIRGs as their results follow a trend from low to high dimension.

1-ccu	1d copy-weight cube GIRG	betw	betweenness centrality
1-23	$1 \vee (2 \wedge 3)$ mixed min/max GIRG	k-core	node k-core number
2-min	$1 \vee 2$ 2d min GIRG	close	closeness centrality
3-cu	3d cube GIRG	LCC	node local clustering coefficient
7d	7d GIRG	deg	node degree
CL	Chung-Lu	Katz	node Katz centrality
CL-c	copy-weight Chung-Lu	PR	node PageRank
BA	Barabasi-Albert	comms	community sizes
ER	Erdos-Renyi	k-cores	k-core sizes
hyper	Hyperbolic Random Graph	diam	graph diameter
		eff-diam	graph effective diameter

Figure 3.6: Graph Generative Model abbreviations; feature name abbreviations

Features	1ccu	2ccu	3ccu	4ccu	5ccu	CL-c	1-23	1-234	12-34	1-2-34	2m	3m	4m	5m	1cu	2cu	3cu	1d	2d	3d	4d	5d	6d	7d	hyper	CL	BA	ER
betw, close, diam	7	5	3	5	6	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
betw, close, eff-diam	28	29	30	23	28	15	3	11	17	3	3	3	3	1	0	16	19	23	17	15	14	13	12	10	13	14	1	1
LCC, tau, diam	6	8	5	4	6	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0
LCC, tau, eff-diam	27	26	31	31	31	0	6	17	19	5	5	3	1	0	22	27	26	19	20	20	18	19	1	1	0	19	0	0
n, m	49	48	48	48	48	48	49	50	49	50	49	48	47	47	48	50	49	50	50	50	50	49	49	49	49	46	50	50
n, m, betw	39	38	38	36	32	9	16	34	37	15	14	11	9	7	37	32	23	39	38	38	35	28	27	26	34	2	2	3
n, m, betw, close, diam	5	5	3	3	5	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n, m, betw, close, eff-diam	21	25	22	21	21	9	2	8	16	2	2	1	1	1	14	18	18	16	16	14	13	12	9	14	16	1	0	2
n, m, close	33	33	32	34	27	10	29	42	35	35	35	20	16	11	24	20	17	31	28	29	21	16	24	42	17	8	2	12
n, m, close, LCC	38	39	37	36	29	0	39	47	36	39	40	27	5	0	26	21	22	36	29	26	22	16	1	0	24	0	0	0
n, m, close, deg	40	40	36	35	28	14	37	47	37	37	38	28	20	17	30	21	21	36	32	32	24	18	26	47	23	14	7	16
n, m, LCC	48	47	46	47	47	0	47	50	49	49	49	37	6	0	49	50	48	49	50	49	50	34	1	0	45	0	0	0
n, m, LCC, betw	35	36	36	36	35	0	16	38	39	18	16	11	3	0	41	32	24	43	39	38	35	17	1	0	37	0	0	0
n, m, LCC, tau	27	28	34	36	39	0	44	37	36	44	43	39	7	0	35	31	29	35	38	41	40	33	1	0	35	0	0	0
n, m, LCC, tau, diam	4	2	5	2	5	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0
n, m, LCC, tau, eff-diam	21	18	22	19	19	0	2	11	17	2	2	1	0	0	22	26	25	15	14	17	19	18	1	0	18	0	0	0
n, m, LCC, diam	5	5	3	2	5	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0
n, m, LCC, eff-diam	22	24	23	22	23	0	2	9	16	1	3	0	0	0	21	26	29	16	16	14	18	16	1	0	20	0	0	0
n, m, diam	4	3	2	1	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
n, m, eff-diam	22	25	21	22	22	9	2	8	17	2	2	1	0	0	17	23	24	13	16	15	17	19	11	13	17	0	0	0
n, m, k-core	9	9	7	9	10	10	1	2	2	2	1	2	2	1	1	5	16	2	2	1	2	2	1	2	2	1	7	0
n, m, k-core, LCC	11	11	12	12	12	0	4	4	4	2	3	1	0	0	5	10	18	4	5	4	4	2	0	0	5	0	0	0
n, m, k-core, deg	5	5	4	5	5	9	0	0	0	0	0	0	0	0	0	1	14	0	0	0	0	0	0	0	0	0	5	0
n, m, deg	47	49	48	48	48	48	49	49	50	49	48	46	46	48	47	49	48	50	50	50	50	49	48	47	43	50	50	50
n, m, deg, tau	26	28	33	35	38	47	42	37	38	41	44	45	38	37	35	32	29	35	38	40	40	42	43	43	36	42	0	0
n, m, deg, betw	39	39	40	37	32	12	18	36	41	16	16	10	9	9	39	34	24	41	40	42	36	32	29	29	33	3	2	4
n, m, deg, LCC	48	45	47	47	46	0	49	49	47	50	49	42	6	0	48	49	48	49	49	49	49	36	2	0	43	0	0	0
n, m, Katz	47	45	45	46	47	45	13	12	11	12	12	14	16	14	11	13	17	12	12	12	11	11	13	12	13	12	21	4
n, m, Katz, LCC	47	45	46	45	45	0	14	13	12	12	14	15	6	0	13	14	20	12	13	14	12	14	1	0	15	0	0	0
n, m, Katz, deg	46	45	48	45	47	47	14	14	13	13	12	16	17	14	14	15	18	13	13	14	15	13	14	14	14	13	22	6
n, m, PR	48	49	49	49	49	48	49	50	49	50	49	47	47	47	48	50	50	50	50	50	50	49	49	48	46	49	49	50
n, m, PR, LCC	48	47	47	46	46	0	49	49	50	48	49	39	5	0	49	49	48	49	48	48	50	34	1	0	45	0	0	0
n, m, PR, deg	47	49	47	48	48	47	48	49	49	49	48	46	44	47	47	49	49	50	50	50	50	48	47	47	44	50	48	50

Figure 3.7: Blasius Framework on just mean values

Diffusion Maps

Contents

4.1	Introduction	29
4.2	Diffusion Maps Theory	30
4.2.1	Random Walk Formulation	30
4.2.2	Improving Diffusion Map Geometry	32
4.2.2.1	Expected distance approach	32
4.2.2.2	Homogeneous random walk approach	33
4.2.2.3	Bounded diffusion hop distances	34
4.3	Empirical results and post processing	35
4.3.1	Inferring dimension d	35
4.3.2	Rescaling/Shifting Diffusion Maps	37

4.1 Introduction

The GGM comparison framework detailed in chapter 3 was designed to assess the similarity of two distributions of graphs, real-world and fake generated, based on higher level features. The fake GGM generates are only loosely fit on their real counterparts, such that the graphs have no node-node or edge-edge correspondences. Hence the comparison features like mean (or median / standard deviation) node degree, node closeness centrality, or effective diameter etc. are node permutation invariant.

An alternative framework is to try to fit a GGM model to a graph so as to compare on an edgewise and likelihood basis. There is little hope for Chung-Lu, Erdos-Renyi and Barabasi-Albert to achieve this. However the GIRG GGM has more parameters, so could fit not just $\hat{\alpha}$ to match the local clustering coefficient, \hat{c} to match the number of edges, and $\hat{\tau}$ to match the degree distribution tail, but to further actually try and infer individual node weights $w_u \in \mathbb{R}^+$, and positions $x_u \in \chi$.

We saw this already to some extent with the copy-weight GIRGs - where $\hat{w}_u = d_u$ is fit, as the observed real graph node degrees. The x_u are much harder to fit - for a start any maximum likelihood fit \hat{x}_u would be rotation, reflection, and translation

invariant (isometries). Finding any one maximum likelihood fit is impractical for all but the smallest graphs.

In this chapter we explore the method of *diffusion maps* for fitting an initial good guess of the $(x_u)_{u \in V}$ positions.

4.2 Diffusion Maps Theory

Diffusion Maps [Coifman and Lafon, 2006] are a technique originally intended to find a lower dimensional representation of some vector points in a high dimensional space, which actually conform to a lower dimensional manifold \mathcal{M} . E.g. points on a 3d sphere $x_1^2 + x_2^2 + x_3^2 = 1$, but embedded in \mathbb{R}^5 . Or a “swiss roll”, an inwardly spiralling 2d manifold in 3d space. The method involves first converting points to a (weighted) adjacency graph W via a distance kernel, and then decomposing a diffusion process on this graph to produce a lower dimensional vector representation of the original points.

This method is hence suitable for our purposes of x_u inference, as though we don’t have a set of high dimensional points, we do have a simple graph $G = (V, E)$ that, like the distance kernel produced graph of diffusion maps, is assumed to have been produced by a geometry influenced process (e.g. generated from a GIRG).

Relation to node location likelihood maximisation [Belkin and Niyogi, 2001] presents the method of Laplacian eigenmaps, which is an interpretation of diffusion maps and shows that it’s solving for the objective

$$\min_X \sum_{u,v} W_{uv} \|x_u - x_v\|^2 \quad \text{subject to } X^T D X = I, \quad X^T D \mathbf{1} = \mathbf{0}$$

where $X \in \mathbb{R}^{n \times d}$ is a matrix whose rows $x_u \in \mathbb{R}^d$ are the lower dimensional node representations, and D is the diagonal degrees matrix. I.e. if nodes u, v are connected directly by $W_{uv} = 1$, we want that their geometric locations x_u, x_v are close by. For fitting GIRG locations we would actually like to solve the objective of likelihood maximisation

$$\begin{aligned} & \max_{(x_u \in \mathbb{T}^d)_{u \in V}} \sum_{u,v} W_{uv} \log p_{uv} + \bar{W}_{uv} \log \bar{p}_{uv} \\ & \text{where } p_{uv} = \min \left\{ 1, c \left(\frac{w_u w_v / \sum_j w_j}{\|x_u - x_v\|^d} \right)^\alpha \right\} \end{aligned}$$

Although the two objectives are different, we hope that the efficiently solvable diffusion map provides a good first approximation of the maximum likelihood locations.

4.2.1 Random Walk Formulation

The random walk characterisation of diffusion maps helps to give us intuition for how to effectively use the method with GIRG generated graphs.

The idea is to analyse the diffusion process of randomly walking on the edges of the graph, characterising the probability cloud starting from one node in the

graph as a sum of decreasingly important orthogonal contributions - coming from the eigenvectors of decreasing eigenvalues of a diagonalisable matrix. The top d contributions can then be used as a coordinate system to describe each point in the graph. By taking a large timestep diffusion cloud, the general relative location of the initial point is the main signal. The hope is that if connections (probabilistically) follow a geometry of d -dimensions, then the diffusion map coordinate system will capture / align with this real geometry.

The relevant quantities of the markovian random walk $u_t : t \in \mathbb{N}$ are defined

$$\begin{aligned}
 M_{uv} &:= P(u_{t+1} = v | u_t = u) = \frac{W_{uv}}{\deg(u)} && u \rightarrow v \text{ transition probability} \\
 M &= D^{-1}W && \text{transition matrix} \\
 D_{uu} &= \sum_{v \in V} W_{uv} && \text{diagonal degree matrix} \\
 W_{uv} &= \begin{cases} 1 & u \sim v \\ 0 & u \not\sim v \end{cases} && \text{node adjacency matrix} \\
 S &= D^{-1/2}WD^{-1/2} && \text{symmetric matrix} \\
 &= V\Lambda V^T && \text{diagonalisation into orthonormal e-vectors} \\
 \Phi &= D^{-1/2}V = [\phi_1, \phi_2, \dots, \phi_n] && \Psi = D^{1/2}V = [\psi_1, \psi_2, \dots, \psi_n] \\
 \Phi^T \Psi &= \Psi^T \Phi = I_{n \times n} && \text{due to orthonormality of } V
 \end{aligned}$$

We use the diagonalisation of S to write M as

$$\begin{aligned}
 M &= D^{-1/2}SD^{1/2} = \Phi\Lambda\Psi^T && \text{diffusion map representation} \\
 &= \sum_{k=1}^n \lambda_k \phi_k \psi_k^T
 \end{aligned}$$

The diagonalisation of sparse¹ symmetric matrix S can be done quite efficiently: we use `scipy.sparse.linalg.eigsh` to find the top $d + 1$ eigenvalues.

The biorthonormality $\langle \phi_i, \psi_j \rangle = \delta_{ij}$ means that $M\phi_i = \lambda_i \phi_i$ and $M^T \psi_i = \lambda_i \psi_i$. For a diffusion map representation of nodes, we order eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots$. Notably the transition matrix satisfies $M\mathbf{1} = \mathbf{1}$ since $\sum_j \frac{W_{ij}}{\deg(i)} = 1$, which can be shown to be the largest eigenvalue $\lambda_1 = 1$; if the graph is connected then also $\lambda_2 < 1$. In particular then $\phi_1 = l\mathbf{1}$, some scalar multiple of the all ones vector.

The diffusion map representation of a node is then

$$\begin{aligned}
 e_u^T M &= \sum_{k=1}^n \phi_u(k) \lambda_k^t \psi_k && \text{diff map of node } i \text{ after } t \text{ steps} \\
 \phi_t(u) &:= (\phi_2(u) \lambda_2^t, \dots, \phi_{d+1}(u) \lambda_{d+1}^t) && \text{d-trunc diff map representation} \\
 D_t(v_i, v_j)^2 &:= \sum_k \frac{1}{d_k} (M_{ik}^t - M_{jk}^t)^2 && \text{diffusion distance} \\
 D_t(v_i, v_j)^2 &= \|\phi_t(v_i) - \phi_t(v_j)\|^2 && \text{provable equality}
 \end{aligned}$$

The truncated diffusion map summarises the distinguishing features of a node's random walk cloud after t steps, removing the tail of decaying contributions which are scaled by smaller $\lambda_{i>d+1}$. Since $\phi_1 = l\mathbf{1}$, the first coordinate is dropped as it's useless for distinguishing nodes. This corresponds to the fact that the diffusion cloud starting at any node converges as $t \rightarrow \infty$ to the same *stationary distribution* $\pi_i = \frac{d_i}{\sum_j d_j}$, the unique normalised solution to $\pi^T M = \pi^T$.

¹nodes in GIRGs have $O(1)$ average degree; nodes in socfb Facebook graphs generally have average degree less than 90 - hence they are relatively sparse.

4.2.2 Improving Diffusion Map Geometry

We ideally want the diffusion map embedding of a GIRG generated graph to have similar distances to the original geometry: $\|\varphi(u) - \varphi(v)\| \approx \|x_u - x_v\|$.

Long story short, in the conventional diffusion map setting of manifold \mathcal{M} in high dim space \rightarrow lower dim embedding, the diffusion map is achieving a discrete approximation of the Laplace-Beltrami operator on the whole continuous manifold - see [Singer, 2006]. This operator defines a continuous diffusion distance $D_t(x, y) = \|\varphi_t(x) - \varphi_t(y)\|$ between points $x, y \in \mathcal{M}$, which is approximately equal to the geodesic distance $d_g(x, y)$ on the manifold, when $d_g(x, y)^2 \ll t \ll 1$, as shown in the background section of [Berry and Harlim, 2018]. Hence if we can achieve a good discrete approximation to the continuous diffusion process, we can hope to get good geometric embeddings.

With GIRGs, our “manifold” is the original geometric space χ , e.g. a torus or a cube. Each node $u \in V$ has a true location $x_u \in \chi$, and instead of having distances $d(u, v)^2$ in some higher dimensional space we just observe the presence of edges/non-edges, $u \sim v$ or $u \not\sim v$, which give a signal as to whether $\|x_u - x_v\|$ is large (if $u \not\sim v$) or small (if $u \sim v$).

We highlight three different approaches in implementing a diffusion map embedding of an input simple unweighted graph (suspected to be a GIRG) $G = (V, E)$.

G has adjacency matrix $A_{ij} = \begin{cases} 1 & \text{if } u \sim v \\ 0 & \text{if } u \not\sim v \end{cases}$, and the differences between the approaches can be encapsulated in the choice of the weighted adjacency matrix W_{ij} seen in section 4.2.1 used to define the diffusion process.

Naive approach The naive approach is just to use $W = A$. This is a valid methodology akin to taking edge weights 1 or 0 between points z_i, z_j in high dim space if they fall within some $\varepsilon > 0$ distance of each other; this is the “Simple-Minded” variant proposed in [Belkin and Niyogi, 2001]. Furthermore if we wanted to make minimal assumptions on the generative process that produced G , beyond that edges formed influenced by some geometric principle, in some finite d -dim space, then this is probably the best approach.

4.2.2.1 Expected distance approach

The approach used by [García-Pérez et al., 2019] to estimate node locations for the Hyperbolic Random Graph model uses the second “Heat Kernel” edge weighting variant of [Belkin and Niyogi, 2001]. This is specified as taking $W_{ij} = e^{-\|z_i - z_j\|^2/T}$, for some parameter $T > 0$, as adjacency matrix weights between vectors z_i in the high dim space. This is designed to optimise the discrete approximation of the continuous heat diffusion on the manifold \mathcal{M} .

The idea of [García-Pérez et al., 2019] is to substitute for $\|z_i - z_j\|^2$ the expected distance between nodes u, v in the HRG - we can do the same for a GIRG: $\mathbb{E}[\|x_u - x_v\|^2]$. This is our best distance guess given our assumption that G is GIRG generated.

If $u \sim v$, then $E[r_{uv}|w_u, w_v, u \sim v] = \Theta \left[\left(\frac{w_u w_v}{W} \right)^{1/d} \right]$, assuming that $\alpha - 1 > \frac{1}{d}$, and is otherwise $\Theta \left[\left(\frac{w_u w_v}{W} \right)^{\alpha-1} \right]$. We can use the node weight estimates by degree of $\hat{w}_u = d_u$ to get distance estimates $\hat{r}_{uv} = \left(\frac{w_u w_v}{W} \right)^{\min(1/d, \alpha-1)}$. Finally edge weights are $W_{uv} = e^{-\hat{r}_{uv}^2/T}$, i.e. decaying with predicted distance.

$$\begin{aligned}
 p(r_{uv}|w_u, w_v, u \sim v) &= \frac{p(u \sim v|w_u, w_v, r_{uv}) p(r_{uv}|w_u, w_v)}{p(u \sim v|w_u, w_v)} \\
 &= \frac{\min\{1, c(\dots)^\alpha r^{-d\alpha}\} dr^{d-1}}{\Theta(\dots)} \\
 (\hat{r}^d = (\dots)) &= \Theta(\dots)^{-1} \left[\int_0^{\hat{r}} dr^d + \int_{\hat{r}}^{0.5} c(\dots)^\alpha dr^{d(1-\alpha)} \right] \\
 &= \Theta(\dots)^{-1} \left[\frac{d}{d+1} (\dots)^{1+1/d} + \frac{dc}{1+d(1-\alpha)} (-(\dots)^{1/d+1} + (\dots)^\alpha) \right] \\
 &= \Theta \left((\dots)^{1/d} + \text{sgn}(\alpha < 1 + 1/d) [-(\dots)^{1/d+1} + (\dots)^\alpha] \right)
 \end{aligned}$$

4.2.2.2 Homogeneous random walk approach

The final approach we consider acts similarly to the expected distance approach, downweighting edges between nodes u, v which have high degree - intuitively this edge is likely to be longer in the χ -space and hence less important to the geometry. We aim to make the random walk on nodes of the graph G directly similar to a (geometric) random walk on the χ -space - i.e. is heat diffusion / brownian motion like.

Again we take d_u as the estimated weight \hat{w}_u of node u . We can modify the naive random walk M to prioritize neighbours v with lower degree d_v via

$$\begin{aligned}
 M_{uv} &\leftarrow \frac{M_{uv} d_v^{-\gamma}}{\sum_{v'} M_{uv'} d_{v'}^{-\gamma}} \\
 M &\leftarrow D' M D^{-\gamma} \quad (\text{diagonal } D' \text{ normalises probabilities})
 \end{aligned}$$

for some parameter $\gamma \geq 0$. This modification of transition probabilities is viewed more properly as modifying original edge weights such that $W_{uv} = k(d_u)k(d_v)$, where $k(d_u) = d_u^{-\gamma}$.

Picking $\gamma = 1$ seems like a good start - in powerlaw weighted GIRGs, it counteracts the tendency for a node's neighbours to have a shifted weight distribution. I.e. in the base $W_{uv} = 1$ for $u \sim v$ model, the next state of the random walk starting at u uniformly randomly samples a neighbour $v \sim u$ which then has weight $w_v \sim \text{powerlaw}(\tau - 1)$; our modified random walk has $w_v \sim \text{powerlaw}(\tau)$, matching the original node weight distribution.

Furthermore with no assumptions on the weight distribution, $\gamma = 1$ also means that the random walk's stationary distribution is more like a constant 1 distribution rather than the known $\pi_u = \frac{d_u}{\sum_v d_v} \propto d_u$ distribution. We can see this by using the fact that the modified node degrees is now $\tilde{d}_u = \sum_{v \sim u} d_u^{-1} d_v^{-1} = d_u^{-1} \sum_{v \sim u} d_v^{-1}$. For graphs where most nodes have a decent number of neighbours spread across lower and higher degree nodes, then $\sum_{v \sim u} d_v^{-1} = \Theta(d_u)$, and so $\tilde{\pi}_u \propto d_u^{-1} \Theta(d_u) = \Theta(1)$.

A homogeneous final stationary distribution is the outcome of heat diffusion on a manifold, so is perhaps desirable.

4.2.2.3 Bounded diffusion hop distances

To achieve distance preserving diffusion map point inference, we want the diffusion process to closely match brownian motion on the original manifold, in our case the geometric space $\chi = [0, 1]^d$. Consider one random walk step starting at node u and choosing a neighbour $v \sim u$ to transition to. While the direction in which we go is roughly random (e.g. for $\chi = [0, 1]^2$ a 2d square, any neighbour $v \sim u$ is relatively positioned with $x_v - x_u$ having an angle $\theta \in [0, 2\pi)$ uniformly randomly), the distance we travel is not uniformly distributed. For our random walk to approximate brownian motion we want the hop distance $r_{u_t u_{t+1}}$ to have *finite variance*.

Rescaled geometric space For $\chi = [0, 1]^d$ hops are of course bounded to be finite; we actually have to view the GIRG in a new lens for this argument to make sense. Without going too much into detail, a GIRG can be equivalently formulated by replacing $\chi \leftarrow [0, n^{1/d}]^d$ and redefining $p_{uv} = \min \left\{ 1, c \left(\frac{w_u w_v}{r_{uv}^d} \right)^\alpha \right\}$ (i.e, removing the $1/W$ scaling on the node weight product ²). This means that as we increase $n \rightarrow \infty$, the uniformly distributed node points don't get more tightly packed into the $[0, 1]^d$ cube, instead they're always spaced with mean density of occupying 1 unit of volume per point.

We can also quantify the expected random step distance in χ -space. If the starting state is node u at timestep 0, and we transition to node v in the first hop (necessarily $v \sim u$ is a neighbour), then $E[r_{uv} | w_u, u \rightarrow v]$ is the expected distance travelled in χ -space due to our transition.

$$\begin{aligned} E[r_{uv}^2 | w_u, w_v, u \sim v] &= \Theta(w_v^\eta) \quad \text{where } \eta = \min(\frac{2}{d}, \alpha - 1), \text{ dropping factors of } w_u \\ p(w_v | w_u, u \sim v) &\propto w_v^{1-\tau} \quad (\text{powerlaw weights}) \\ u \rightarrow v = u \rightarrow v \cap u \sim v &\quad \text{equivalent events} \end{aligned}$$

We also need

$$\begin{aligned} p(w_v = w | w_u, u \rightarrow v, u \sim v) &\propto p(u \rightarrow v | w_u, w_v = w, u \sim v) p(w_v = w | w_u, u \sim v) \\ &\propto p(u \rightarrow v | w_u, w_v = w, u \sim v) \Theta(w^{1-\tau}) \end{aligned}$$

And finally, for each of the different adjacency matrix weighting schemes,

$$\begin{aligned} p(u \rightarrow v | w_u, w_v = w, u \sim v) &\approx \propto \quad (\text{assuming degrees } d_i \approx w_i) \\ \begin{cases} \frac{1}{w_u} & \text{Naive approach} \\ e^{-r_{uv}^2/T} = \exp(-\frac{1}{T}(\frac{w_u w_v}{W})^{2/d}) = \exp(-Kw^{2/d}) & \text{Expected distance} \\ w^{-\gamma} & \text{Homogeneous} \end{cases} \end{aligned}$$

²note that whp $W = \Theta(n)$ for $\tau > 2$ powerlaw weights, so we've swapped the downscaling of edge probabilities by $1/n$ from the weights to the distances

Putting this together we have a formula for the expected hop distance in one discrete timestep of the diffusion process

$$\begin{aligned}
& E[r_{uv}^2 | w_u, u \rightarrow v] \\
&= \int_1^\infty E[r_{uv}^2 | w_u, u \rightarrow v, w_v = w] p(w_v = w | w_u, u \rightarrow v) dw \\
&= \int_1^\infty E[r_{uv}^2 | w_u, w_v = w, u \sim v] p(w_v = w | w_u, u \rightarrow v) dw \\
&= \int_1^\infty \Theta(w^\eta) p(w_v = w | w_u, u \rightarrow v) dw \\
&= \int_1^\infty \Theta(w^{1+\eta-\tau}) p(u \rightarrow v | w_u, w_v = w, u \sim v) dw
\end{aligned}$$

In the naive approach, similarly to how $\int_1^\infty x^\beta dx$ is bounded if $\beta < -1$, if $1 + \eta - \tau > -1$ then the hop distance is unbounded, i.e. if $\tau < 2 + \min(\frac{2}{d}, \alpha - 1)$. This is eminently possible for $\tau \in (2, 3)$. The expected distance approach has no problems as $\exp(-Kw^{2/d})$ decays swiftly, however in the homogeneous approach taking $\gamma = \min(\frac{2}{d}, \alpha - 1)$ would guarantee bounded hop distances for all $\tau \in (2, 3)$, as $2 - \tau < 1$.

In fig. 4.1 diffusion clouds out of a point are shown for the three different edge weighting approaches. Both non naive versions look more like brownian motion clouds as desired. We additionally benchmarked the approaches by comparing the diffusion map embedding of a fake graph generated from a GIRG with known χ -space embeddings. The homogeneous random walk approach has the best correlation between embeddings and true locations, although with $\gamma = 0.9$ not $\gamma = 1.0$.

4.3 Empirical results and post processing

We see in fig. 4.2 some 2d truncated diffusion maps. When the underlying graph was a 2d square GIRG, the inferred points indeed look square like. From the 1d GIRG we get a curved line shape - i.e. a 1d manifold embedded in 2d, as expected.

4.3.1 Inferring dimension d

If we have a graph G which we know is generated from a d -dimensional GIRG, we can simply extract the d -truncated diffusion map coordinates of each node as an initial estimate for the original geometric location of the node.

If the true geometric dimension d is unknown, an important first question is how to choose the output dimension d (truncation for the embedding) of the diffusion map.

Diffusion maps actually present one way to infer the dimensionality by analysing the ordered sequence of eigenvalues $\lambda_2 < \lambda_3 < \dots$. For example in fig. 4.3 we see that for graphs generated from cube GIRGs with dimension $d = 1, 2, 3, 4$, the diffusion map eigenvalue have a clear cutoff point after the first $d + 1$ eigenvalues, with $\lambda_2, \dots, \lambda_{d+1}$ being all approximately equal. Hence an eigenvalue cutoff can be used to infer geometric dimensionality of a graph. This is inherently manual/empirical, and so more difficult for real world graphs.

Toroidal GIRGs Interestingly Toroidal GIRGs are differentiated from Cube GIRGs in that the diffusion map requires $2d$ coordinates to capture the Torus geometry

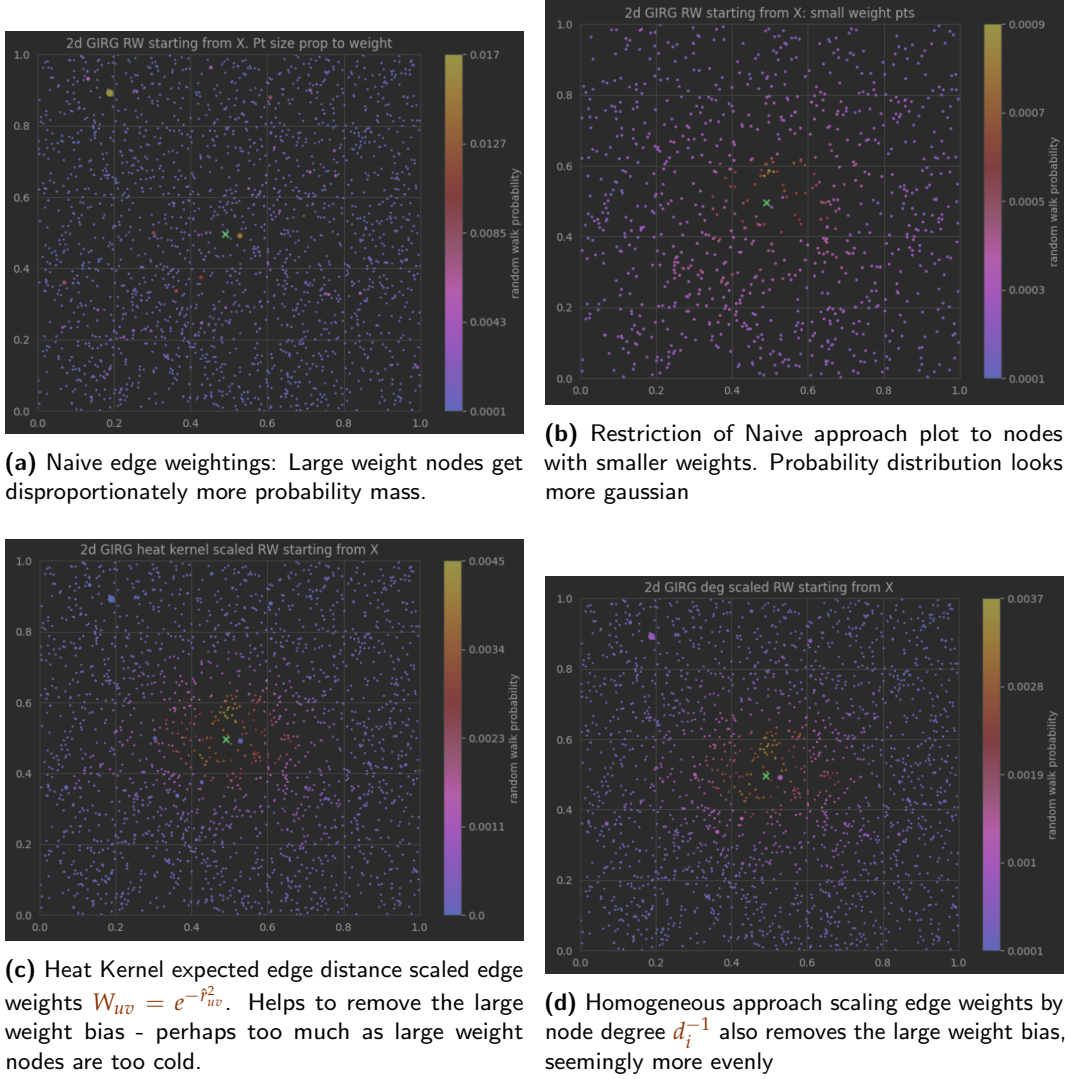


Figure 4.1: $t = 6$ step random walk diffusion cloud on a 2d torus GIRG. Each point is a node shown at its location $x_u \in [0, 1]^2$. $u_{t=0}$ is the node marked with the green x. Point colour denotes probability mass in the diffusion cloud. Point radius is proportional to node weight w_u . Three different random walks are shown corresponding to our three different graph edge weighting (transition probability) schemes.

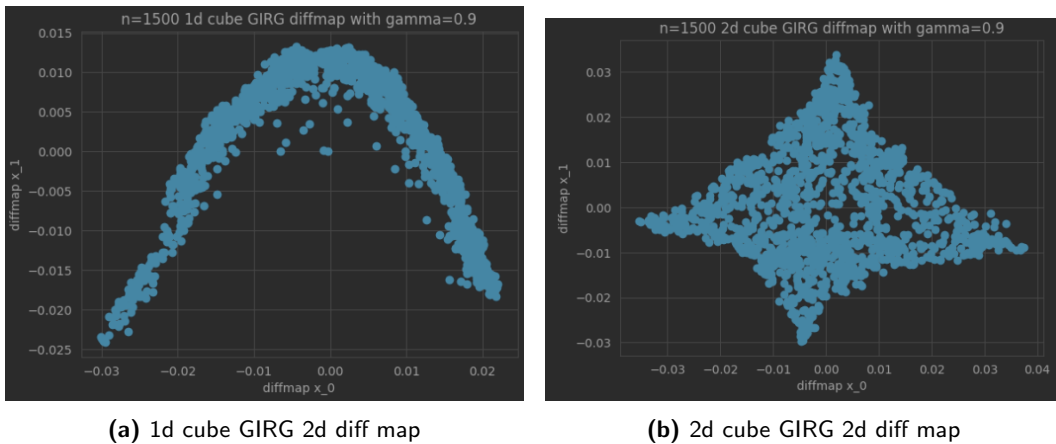


Figure 4.2: Diffusion map scatter plot of the 2d truncated diffusion maps of (a) graph from a 1d cube GIRG, (b) graph from a 2d cube GIRG. Each plotted point is the inferred location of a single node.

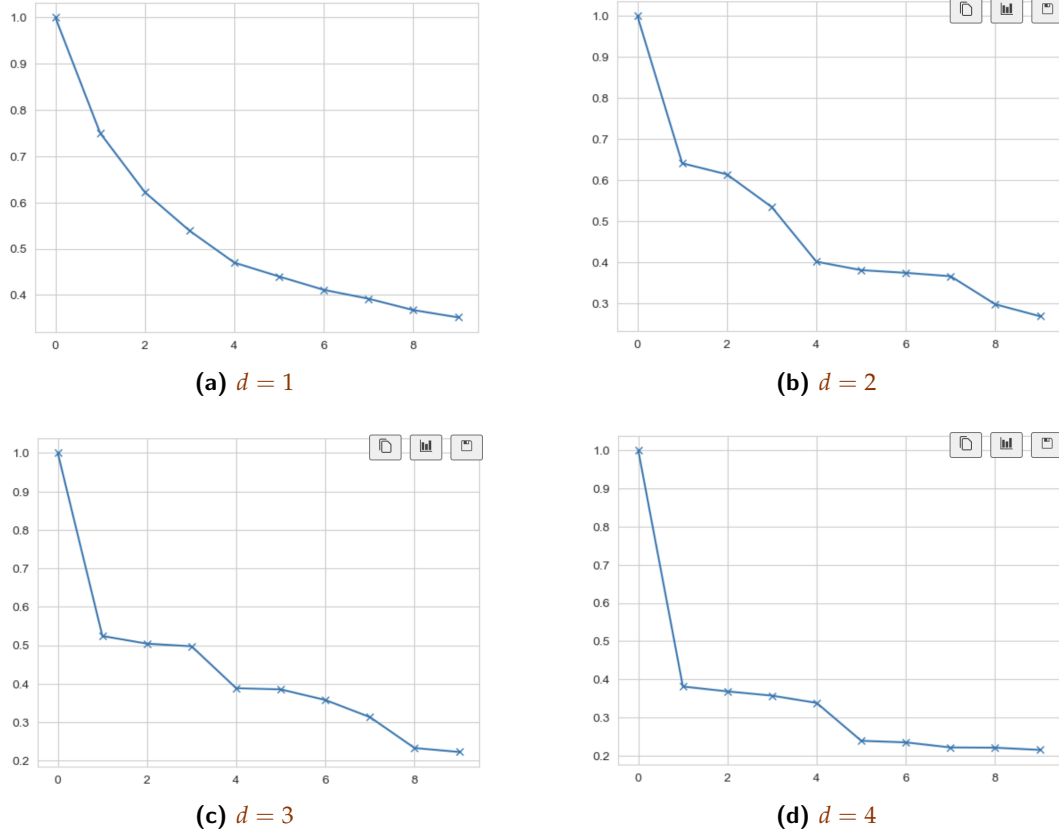


Figure 4.3: Eigenvalues of the diffusion map laplacian matrix (including $\lambda_1 = 1$) for 4 different $n = 2000, \tau = 2.5, \alpha = 1.3$ cube GIRGs: with dimensions $d = 1, 2, 3, 4$. X-axis is the i th eigenvalue, y-axis is the eigenvalue value.

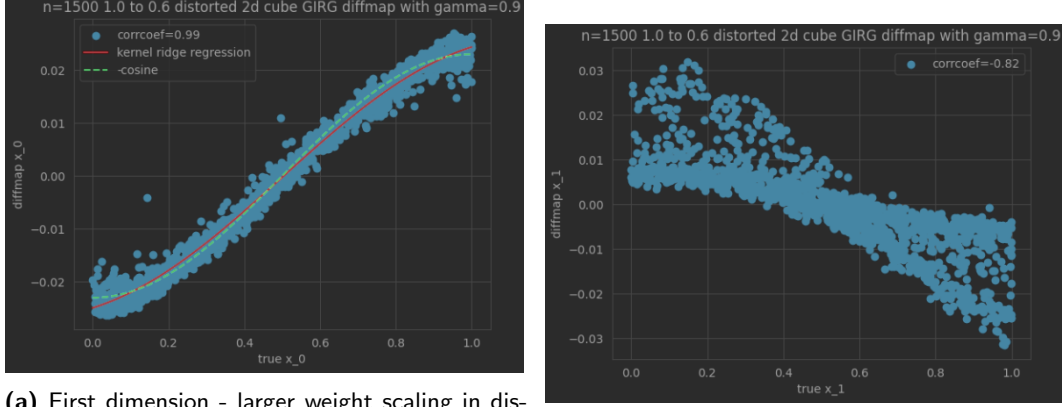
instead of just d for the cube. The natural diffusion map of a 1D Torus GIRG ends up being a 2D circle about the origin; 2 large eigenvalues λ_2, λ_3 are necessary.

4.3.2 Rescaling/Shifting Diffusion Maps

From section 4.2.2 we know the embedding should be roughly isometric up to a scale factor - so rescaling and shifting into the unit cube/torus is important to make inter-point distances meaningful/consistent across multiple graphs (even though this could be absorbed into the probability scaling constant c). We detail some post processing steps that can be done to improve the diffusion map embedding for GIRGs and real graphs.

The diffusion map embeddings will be centered around the origin, as all $\lambda_2, \lambda_3, \dots$ eigenvalue eigenvectors are orthogonal to $\phi_1 = k\mathbf{1}$: they represent a deviation from the stationary distribution - e.g. in a 1d cube GIRG for a node with $x_u \approx 0$ at the extreme geometric left end, its diffusion cloud will need to put more diffusion probability on the left side nodes and less on the right side nodes than the stationary distribution. The y-axis scale is small as $\lambda_2 < 1$, and decreases for larger t .

In order to map to the $[0, 1]^d$ torus/cube, the 0 centering can easily be fixed by shifting the points $(x_u)_i \leftarrow (x_u)_i + \min_v (x_v)_i$. If we are certain that the points should be distributed within the unit cube, then we can simply rescale separately along each dimension: $x \leftarrow \frac{x - x_{\min}}{x_{\max} - x_{\min}}$. If furthermore we're certain that the distribution within the unit cube should be relatively uniform, we can perform a



(a) First dimension - larger weight scaling in distance function. As the more geometry influencing dimension, it is easier for the diffusion map to fit well, and is picked up as the largest non 1 eigenvalue. The fit is quite close to a cosine wave as shown.

(b) Second (minor) dimension - diffusion map has more trouble fitting this information as it is harder to distinguish from noise. Note a linear relationship means success - negative gradient is irrelevant.

Figure 4.4: 2d distorted cube GIRG. 2 dim truncated diffusion map embeddings shown with one plot for each dimension. Each plotted point is one node from the graph, with X-axis being its true location in the 2d cube for that dimension, and Y-axis being the diffusion map embedding.

coordinatewise **uniformify** procedure that replaces $(x_u)_i$ with its percentile value compared with other $(x_v)_i$. fig. 4.5 shows the **uniformify** procedure in action. This helps to counteract the phenomenon of terribly different diffusion map scaling, whereby the majority of the graph which is well connected ends up highly bunched up, and a small number of outlying poorly connected nodes cover most of the embedding space.

Rotated Points One issue with diffusion maps as seen in fig. 4.2b is rotations - the 2d GIRG's points are rotated into a diamond shape. As we only hope for an isometric mapping of the original point locations to the diffusion map embedding, we can't expect the diffusion map embedding to be in the same basis as the original point locations. If we want to fit it back into $[0, 1]^d$, we don't care about reflections, and we can translate and rescale, however rotating is more difficult.

That fig. 4.2b looks almost perfectly 45 degree rotated could potentially be explained as the diffusion map trying to maximise diffusion explainability - the long square diagonal has overall more diffusion along it so might be picked up as a major eigenvalue - maybe a small imbalance of more nodes in one corner can cause this. Even assuming no overall rotation bias, you're going to get a square rotated somewhere between 0 and 45 degrees.

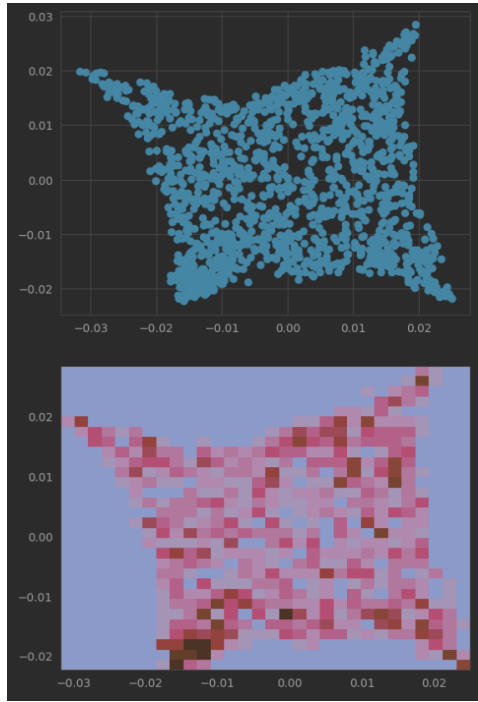
Distorted GIRGs make cuboidal diffusion map embeddings In practice, real graphs never have a perfectly equal balance in geometric dimension importance. We introduced distorted GIRGs in section 2.5 as a possible model for this, where for example the distance function could be the weighted Euclidean norm in 2d: $\|x - y\|^2 = a(x_1 - y_1)^2 + b(x_2 - y_2)^2$. If $a > b$ then the first diffusion map coordinate $\varphi_1(u)$ is likely to be very similar to $(x_u)_1$. Only if $a = b$ is there no preference between $(x_u)_1$ and $(x_u)_2$, making a rotation possible. Hence the points $(\varphi_1(u), \varphi_2(u))_{u \in V}$ are likely to end up as a rectangle, not a rotated square; the variation in first coordinate will be greater than in the second.

fig. 4.4a shows the 2d diff map embedding of such a distorted 2d cube GIRG. Note that the major coordinate is much easier to fit well, whereas the minor one has higher error (NB the diff map embedding of x_2 is negatively scaling with the true value which is fine - we can never guarantee same signs). We also see that the x_1 diffusion map embedding doesn't have quite a linear relationship with the real locations - it looks more like a cosin wave like relationship - this actually makes sense as eigenfunctions of the laplacian operator on a cube look like sine waves.

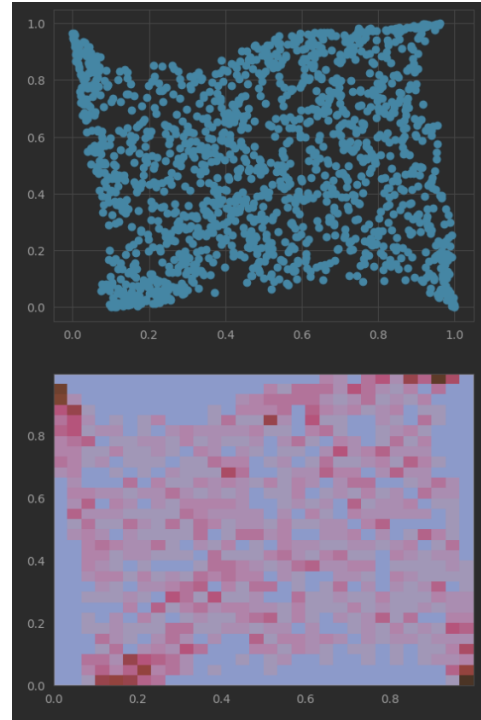
In the light of potentially non equal coordinates, the relative $\lambda_2 > \lambda_3 > \dots$ scaling can be seen as a feature not a bug, if the hypothesis space of generative graph models is to be expanded to cuboid (non-cube) GIRGs. In this case all coordinates of the diffusion map don't have to be all (non-homogeneously) rescaled to the range $[0, 1]$ and can rather be homogeneously rescaled to retain their relative size ratios. This sheds new light on the horizontalness of the $\lambda_2, \lambda_3, \lambda_4, \lambda_5$ line segment in fig. 4.3d as a testament to the underlying cubeness (non cuboid) of the GIRG that generated the graph.

Restricted Rescaling This simple post processing method aims to map points into a reasonable geometric space without strongly biasing towards a uniformly distributed prior like the [uniformify](#) procedure. Empirically, the diffusion map coordinates of real-world facebook graphs often have $\geq 90\%$ of the nodes concentrated in a small parcel, with only a few nodes having extremely far out locations. This defeats the simple rescaling method of $x \leftarrow \frac{x - x_{\min}}{x_{\max} - x_{\min}}$ as most nodes will end up very tightly packed - however if we hit all the points with the [uniformify](#) procedure hammer, we might lose some of the subtleties of the geometry picked up by the diffusion map within the highly connected central parcel.

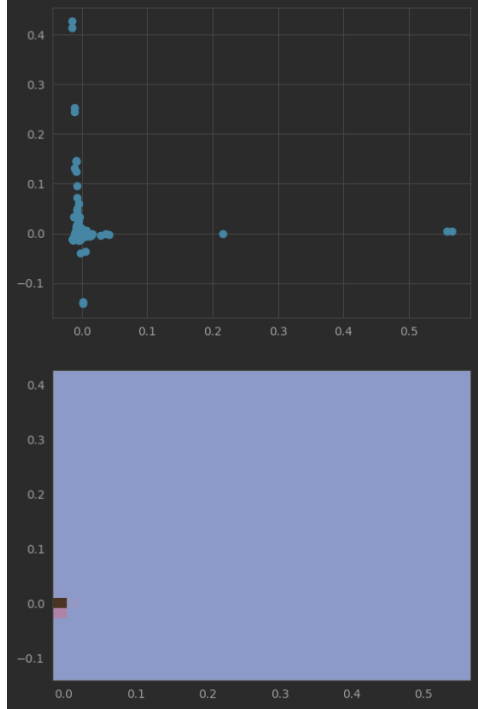
Instead we only rescale the central nodes: those whose joint coordinate-wise percentiles lie in $[5\%, 95\%]^d$; they are linearly scaled to the $[0.05, 0.95]^d$ cube. Finally the outlying nodes are percentile rescaled (non-linearly) to the upper/lower cube margins. This method is shown in comparison for a few graphs in fig. 5.6



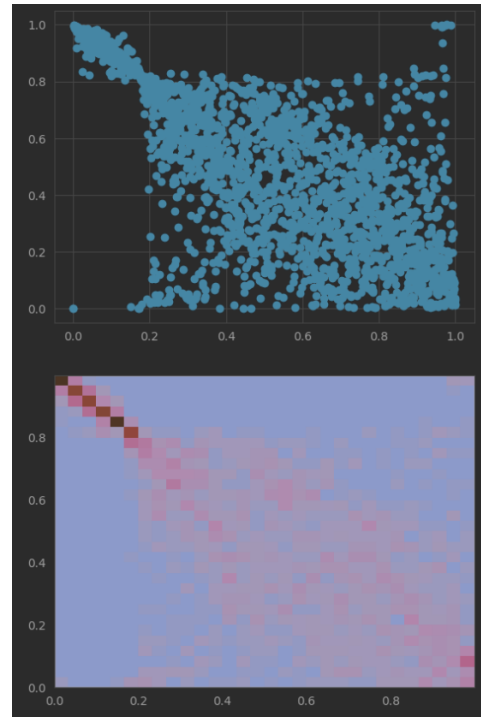
(a) 2d cube GIRG diff map



(b) 2d cube GIRG diff map with [uniformify](#) rescaling



(c) socfb-Amherst41 diff map



(d) socfb-Amherst41 diff map with [uniformify](#) rescaling

Figure 4.5: 2d truncated diffusion maps plotted - each point is the embedding of a node in the graph. (a) and (b) are on a 2d GIRG, (c) and (d) on a real-world social network graph. The effect of the [uniformify](#) point rescaling procedure is shown, in (b) and (d).

Chapter 5

Likelihood Point Estimation

Contents

5.1	Introduction	41
5.2	MCMC Formulation	41
5.2.1	MCMC Likelihood comparison of GIRG vs CL fit to real graph	43
5.2.2	Percent Edges Captured Metric Comparison	45
5.3	Direct Ordered Likelihood Maximisation	45
5.3.1	Analysis on GIRG to real graphs quality of fit	46

5.1 Introduction

Diffusion maps in chapter 4 provide a good initial point estimate for fitting GIRG node locations to a real graph. However what we'd truly like to maximise is the actual data likelihood $p(G|\theta)$, where G is the real graph and θ are all the GIRG parameters like c, α , node weights and node locations. As in [García-Pérez et al., 2019] we can use the likelihood to do a further refinement of point locations, although finding anything like a global optimum is infeasible for all but tiny graphs.

We first tried a *Markov Chain Monte Carlo (MCMC)* approach to not just improve $p(G|\theta)$ but also, in theory provide an estimate of the posterior distribution of $p(\theta|G)$. While effective, it was quite slow, so we later shifted to a *direct likelihood optimisation* with rounds of point location updates in order from highest to lowest degree node, inspired by [García-Pérez et al., 2019].

5.2 MCMC Formulation

MCMC is a method in the bayesian framework of a parametric generative model. Datapoints z are generated by first sampling $\theta \sim p(\theta)$ from a prior, and then generating from the model $z \sim p(z|\theta)$. In our case of fitting GIRG GGM, z is a real graph instance $G, \theta = (\alpha, c, \{w_u\}_{u \in V}, \{x_u\}_{u \in V})$, and we focus in particular on the node locations x_u .

The posterior likelihood $p(\theta|z) = \frac{p(z|\theta)p(\theta)}{p(z)}$ is infeasible to compute due to the normalising factor $p(z)$. MCMC instead uses the non normalised $Q(\theta) = p(z|\theta)p(\theta)$ which can be evaluated, to set up a *Markov Chain (MC)* with states $\theta \in \Theta$, and transition probabilities derived from $Q(\theta)$. In particular we will use a *Metropolis-Hastings* style MC. The two basic components are a *proposal distribution* to propose a new MC state $\theta' \sim p_{\text{prop}}(\theta'|\theta)$, and an *acceptance probability* $p_{\text{acc}}(\theta'|\theta)$ to decide whether to accept θ' . It might take a few rounds of proposal and rejection to produce one actual MC transition. This defines a random walk on the MC state space that converges in the limit to the posterior distribution $\theta \sim p(\theta|z)$. If the j th step of the random walk yields θ^j , given sufficient burn in and spacing, we can sample $\theta^1, \theta^2, \dots$ from the posterior. We will be content with just one posterior sampled $\hat{\theta}$ (NB not a maximum likelihood estimator), and use this to evaluate our overall GIRG model fit to the real graph.

We transition $\theta \rightarrow \theta'$ with a *Gibbs sampling* approach. This means breaking down θ into subcomponents θ_i , randomly choosing one to propose a new state for, i.e. $\theta' = (\theta'_i, \theta_{j:j \neq i})$, and using the marginal $Q_i(\theta_i)$ instead of $Q(\theta)$ in order to compute the transition probabilities. In our case $Q(\theta) = p(G|(\mathbf{x}_u)_{u \in V})p((\mathbf{x}_u)_{u \in V})$, so we can ignore the uniform prior $p((\mathbf{x}_u)_{u \in V})$ as it's constant for all node locations. Then we have

$$Q(\theta) \propto \prod_{(u,v) \in \binom{V}{2}} p_{uv}|w_u, w_v, x_u, x_v, \alpha, c, G \quad (5.1)$$

which is easily marginalised as $Q_u(\theta_u) = \prod_{v \neq u} p_{uv} | \dots$

Burn in time of the MC can be quite long. With the GIRG model, the natural initialisation for \mathbf{x}_u would be to follow the prior $\mathbf{x}_u \sim U([0, 1]^d)$. By using a diffusion map embedding initialisation this can be greatly reduced.

The proposal distribution should be designed to maximise chances of acceptance. It seemed reasonable to randomly propose one of a small local perturbation, or a random jump to anywhere in the cube. A random uniform jump is useful to try and find a completely new location for \mathbf{x}_u that could suit (hopefully near to its neighbours) - in fact another good proposal would be to randomly choose a neighbour $v \sim u$, and move \mathbf{x}_u to a random offset of \mathbf{x}_v , although this would unfortunately make p_{prop} no longer symmetric. A small perturbation $\mathbf{x}'_u = \mathbf{x}_u + \epsilon$ is good as assuming \mathbf{x}_u has high likelihood, somewhere nearby might have even higher.

Acceptance probability as defined by the Metropolis-Hastings's algorithm is

$$p_{\text{acc}}(\mathbf{x}'_u | \mathbf{x}_u) = \min \left(1, \frac{p_{\text{prop}}(\mathbf{x}_u | \mathbf{x}'_u) p(G | \mathbf{x}'_u) p(\mathbf{x}'_u)}{p_{\text{prop}}(\mathbf{x}'_u | \mathbf{x}_u) p(G | \mathbf{x}_u) p(\mathbf{x}_u)} \right) \quad (5.2)$$

$$= \min \left(1, \frac{p(G | \mathbf{x}'_u)}{p(G | \mathbf{x}_u)} \right) \quad (5.3)$$

$$\text{as } p_{\text{prop}} \text{ is symmetric and } p(\mathbf{x}_u) = 1 \ \forall \mathbf{x}_u \quad (5.4)$$

5.2.1 MCMC Likelihood comparison of GIRG vs CL fit to real graph

How well do the converged MCMC points do at replicating the real graph? As suggested in the classification comparison framework, Chung-Lu model is a good non-geometric null hypothesis to compare against - i.e. we hope the MCMC fit GIRG model will do better than Chung-Lu. It's also interesting to see how large dimension d GIRG is necessary for a good fit.

As a first sanity check, we compare the MCMC fit $\hat{\theta}_{\text{GIRG}}$ to the more simply fit copy-weight Chung-Lu parameters $\hat{\theta}_{\text{CL}}$, by comparing $p(G|\hat{\theta}_{\text{GIRG}}, \mathcal{G}_{\text{GIRG}})$ vs $p(G|\hat{\theta}_{\text{CL}}, \mathcal{G}_{\text{CL}})$.

For most of the facebook graphs however, $p(G|\hat{\theta}_{\text{GIRG}}, \mathcal{G}_{\text{GIRG}}) < p(G|\hat{\theta}_{\text{CL}}, \mathcal{G}_{\text{CL}})$, despite GIRGs having more parameters / flexibility to fit G !

The GIRG model is far too confident about edges, giving $p_{uv} \approx 1$ for small $\|x_u - x_v\|$ and $p_{uv} \approx 0$ for large $\|x_u - x_v\|$. Hence a mistake on a few edges can lead to a large penalisation in likelihood. The Chung-Lu model is much more forgiving, with all probabilities more medium sized.

One reasonable tweak is to introduce a *failure rate* $0 \leq f \leq 1$ to the GIRG model:

$$p_{uv} = (1 - f) \min \left\{ 1, c \left(\frac{w_u w_v / W}{\|x_u - x_v\|^d} \right)^\alpha \right\} \quad (5.5)$$

For a social network this means that two highly similar people are not guaranteed/forced to be friends. Indeed there may be a very like-minded person who lives next door to you that you've never met, or had no opportunity / free time to properly get to know.

A failure rate will lower the impact of short non-edges (mistakenly predicted high probability of edge), but for long edges (mistakenly predicted low probability of edge) we need a baseline edge probability to prevent p_{uv} being too small. A simple fix is a *mixture probability* of deriving edges instead from the Chung-Lu model - i.e. letting

$$p_{uv} = \eta p_{uv}^{\text{CL}} + (1 - \eta) p_{uv}^{\text{GIRG}} \quad (5.6)$$

for mixture probability $0 \leq \eta \leq 1$. This should not be seen strictly as an ensemble of models or gross addition to the number of parameters of the GIRG model, as the GIRG model already contains fit node weights $(\hat{w}_u)_{u \in V}$. The mixture probability η does also help a lot on short non-edges similarly to failure rate f , but it could still be good to have both as even the Chung-Lu model can demand an edge $u \sim v$ with high probability if w_u, w_v are both very large - though this is quite rare.

The intuitive social network interpretation of a GIRG Chung-Lu mixture model is that it allows for some "random" friendships - ones that wouldn't normally have been predicted by the GIRG model. For instance if you happen to sit next to a stranger on a plane, as long as they don't fall into a sworn enemy category you may still add them as a facebook friend with a vague promise of staying in touch. Of course you and the stranger must bear at least some small similarity to have been on the same plane, but this would likely need a very high dimensional GIRG model to fully capture this.

With these two new parameters f and η , the augmented GIRG model does have a higher specific fit likelihood than the Chung-Lu model: $p(G|\hat{\theta}_{\text{GIRG}}, \mathcal{G}_{\text{GIRG}}) >$

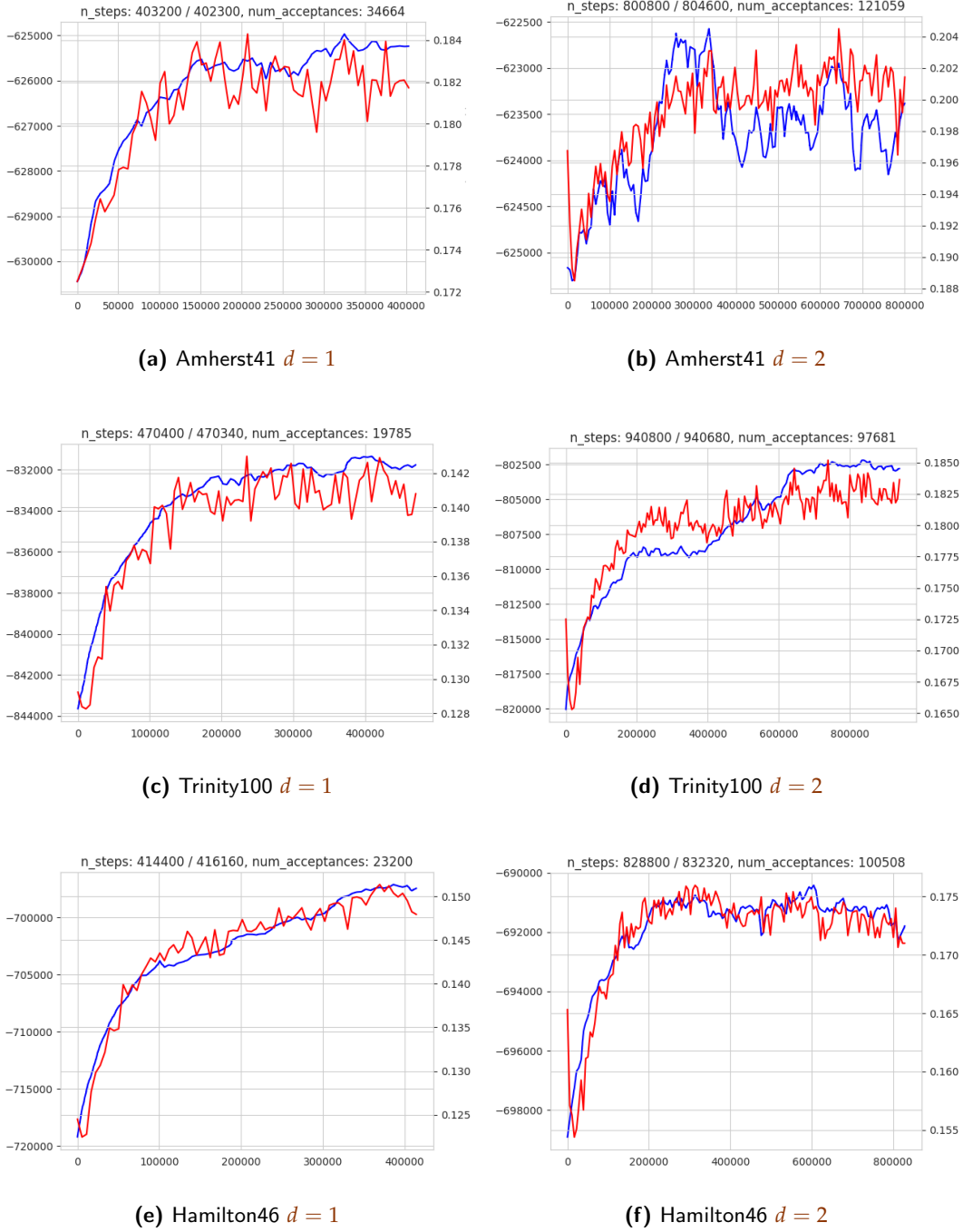


Figure 5.1: MCMC runs for socfb-Amherst41, socfb-Trinity100, socfb-Hamilton46 with failure rate 0.3 and Chung-Lu Mixin rate 0.5, with a cube GIRG model. Log likelihood wise, Amherst 1D GIRG does best; for Bowdoin it's 2D GIRG. x-axis is number of steps; left y-axis (blue curve) is log likelihood; right y-axis (red curve) is PEC metric

$p(G|\hat{\theta}_{CL}, \mathcal{G}_{CL})$. In a holistic MCMC setup these two parameters could also have a prior and be sampled from the posterior, but for simplicity we set them to $f = 0.3, \eta = 0.5$. See fig. 5.1 for some likelihood convergence curves on some example graphs.

5.2.2 Percent Edges Captured Metric Comparison

Another simple framework to compare quality of fit without taking into account increased parametrisation is to analyse the “accuracy” on successfully producing edges / non-edges.

Our MCMC posterior is constrained by directly fitting \hat{c} to produce a similar number of edges $|E|$ as in the real graph. Hence in the standard classification confusion matrix

$$\begin{array}{|c|c|} \hline TP & FN \\ \hline FP & TN \\ \hline \end{array} \quad (5.7)$$

we can equivalently count $\frac{TP}{TP+FN}$ (Recall), the fraction of edges in the real graph that are successfully predicted, or $\frac{TP}{TP+FP}$ (Precision), the fraction of edges in the predicted graph that are also in the real graph, these numbers will be very similar. We call this metric *Percent Edges Captured (PEC)*, and focus on this instead of the alternative Percent Non-Edges Captured (PNEC). PEC is preferable as our graphs are all relatively sparse - hence the PNEC is always high as it is dominated by the large number of non-edges.

5.3 Direct Ordered Likelihood Maximisation

MCMC proved too slow and didn’t achieve such high PECs. Hence we shifted to a simpler approach inspired by [García-Pérez et al., 2019].

Instead of picking points to update randomly, we sequentially update the positions of every single node, ordered from highest to lowest degree. Updating point x_u is done by proposing 100 new locations - near to current location, near to a neighbour, or random uniform in the cube. Instead of using an acceptance probability as in MCMC, we simply pick the best of all the proposals, by marginal edge likelihood $\prod_{v \neq u} p(e(u, v) | \dots)$.

Updating points in order of highest degree makes sense as in essence a high degree node “drags” its local lower degree neighbourhood of nodes along with it, and so hence should be updated in that order - otherwise if its lower degree neighbours are moved first, they might spread out more nonsensically and leave the higher degree node no good place to move to. We think this happened to some extent with the MCMC method, evidenced by initial drops in the PEC in fig. 5.1 for 2d GIRG fits, and also observations of the final point locations being too spread out.

After each round of point updates, α is fit to maximise the likelihood, and c is refit to match the number of edges in the real graph. Multiple rounds are repeated until convergence. One optimisation that [García-Pérez et al., 2019] implements but we missed was to also after each round update node weight estimates: comparing two nodes with similar degree, if the first one has more geometrically close nodes than the second, the second should have a relatively higher estimated weight.

Gradients to speed up convergence The relatively slow speed of convergence to good point locations hindered us from experimenting on graphs of much more than $n = 10000$ nodes. When randomly proposing 100 new locations for each node, the resultant likelihood contribution of each location has to be computed. This might

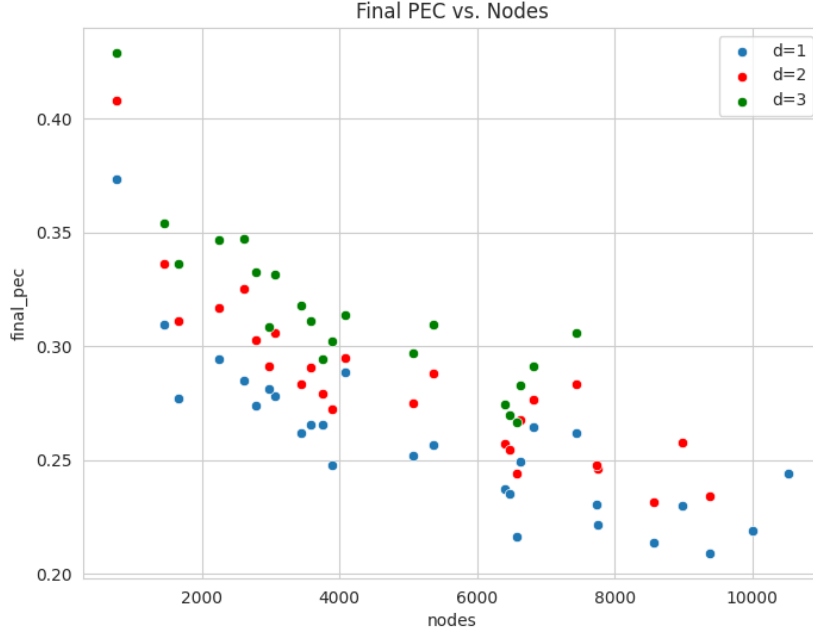


Figure 5.2: Final converged PECs for fitting 1d, 2d and 3d cube GIRGs to real-world Facebook graphs using direct ordered likelihood maximisation. Missing some 3d numbers as batch job exceeded allotted time.

even need to be increased exponentially for larger dimensions d . We can speed this up by instead calculating the point location gradient of the likelihood contribution of each node, and taking a small step in the proposed direction. While indeed a lot faster, we had mixed final PEC scores, probably due to the learning rate being improperly configured across different sized graphs and different dimensions d (e.g. for $d = 1$, a node has neighbours and non-neighbours to the left and right, and will end up with quite a strong gradient pointing in one of the two directions. For $d \geq 2$ it's much more likely that there will be some more balancing of forces on the node such that the resultant gradient will be much smaller). Gradient clipping could potentially help. Optimizing a gradient enhanced method is an area for future work.

5.3.1 Analysis on GIRG to real graphs quality of fit

We see in fig. 5.2 the final achieved PECs of different dimensional cube GIRGs fit to real-world graphs. At first blush, PECs seem worryingly low at just 22% to 32%. You wouldn't be happy with a cat image classifier with only 30% accuracy! Furthermore, PEC decreases with increasing node count which seems to imply that our model fitting could be a fluke that only works on smaller graphs, and that might scale to irrelevance for even larger graphs. Luckily all these signs of disaster have good explanations, and we can restore faith in our point fitting methods.

Inherently low GIRG PECs A first thing to note is that edge probabilities p_{uv} for a GIRG generally are not all $p_{uv} \in \{0, 1\}$ (at least for $\alpha < \infty$). Hence we cannot expect too high a PEC, even if we fit every single GIRG parameter perfectly to identically reproduce all the p_{uv} . We observe this in fig. 5.4, where PEC starts

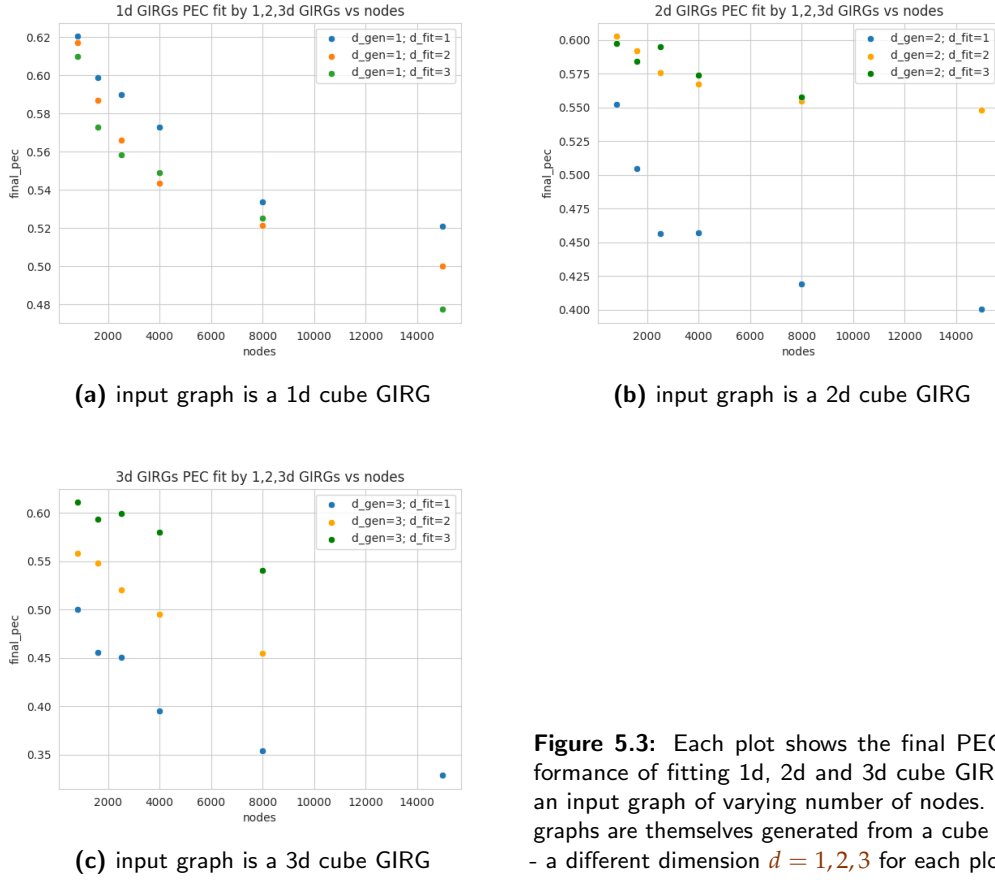


Figure 5.3: Each plot shows the final PEC performance of fitting 1d, 2d and 3d cube GIRGs to an input graph of varying number of nodes. Input graphs are themselves generated from a cube GIRG - a different dimension $d = 1, 2, 3$ for each plot.

for the smallest graphs of $n = 800$ nodes at 64%. Note the plot is of cube GIRGs generated with $\alpha = 1.3$ which is quite typical an LCC fit on the real-world Facebook graphs, so indeed their edge probabilities range broadly in $[0.0, 1.0]$.

Fixed average degree \Rightarrow GIRG lower $p_{uv} \Rightarrow$ lower PECs for higher n Also striking in fig. 5.4 is a similar curve to fig. 5.2 of decreasing PEC with increasing node count. We think this is an artifact of the fixed average degree of $\overline{deg} = 60$ used across these varying node count GIRGs. We fixed \overline{deg} because the socfb graphs have basically stable average degrees - correlation coefficient of just 0.1 with node count, and their mean average degree across graphs is 77. We think that essentially as n grows for a fixed \overline{deg} , since long distance edges in a GIRG contribute a constant fraction of total edges, in order to keep average degree constant the edge probability constant c is lowered. Hence for small n all edge probabilities p_{uv} are generally larger and less numerous ($O(n^2)$ potential edges); for large n there are many more potential edges, and generally smaller edge probabilities p_{uv} . This naturally brings about decreased PEC for larger n .

Fixed average degree \Rightarrow lower PECs in general Another issue with fixed average degree and increasing n is that even in the dumbest case of fitting an Erdos-Renyi graph to an input graph, we would observe a $1/n$ shaped PEC curve with increasing n . For each present edge in the input graph, there is a $p = \frac{\overline{deg}}{n}$

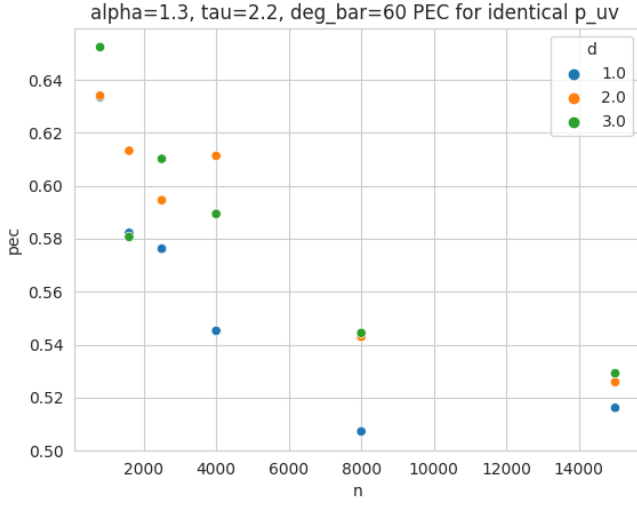


Figure 5.4: Exact p_{uv} replication has $\text{PEC} \in [0.52, 0.64]$, and decreasing with increasing node count. Done for $\alpha = 1.3, \tau = 2.2$ cube GIRGs of dimension $d = 1, 2, 3$ and average degree 60.

chance that it is also present in the Erdos-Renyi graph. Hence PEC would be $\frac{\overline{\text{deg}}}{n}$ in expectation.

This means that if we assume real-world graphs are only partially GIRG explainable, and that their residual non geometric random component is hence hard to fit with our low dimensional cube GIRGs, then the observed PEC against node count n plot would still look like a $1/n$ curve, but with a higher non-zero asymptote.

Chung-Lu control shows GIRGs not overfitting on real-world As a control, we fit GIRGs using direct ordered likelihood maximisation to Chung-Lu graphs, and observed PEC's is not much higher than between two resampled Chung-Lu graphs, at 10% for graphs of $n = 1000$ nodes, decaying to 3% for $n = 5000$. While there is roughly a 3% PEC boost from increasing dimension d which is clearly overfitting, this vindicates that the relatively higher achieved PECs of 22% – 33% in fig. 5.2 represent a true geometric quantity of real-graphs that we are successfully fitting.

Overall we find that cube GIRGs get quite a decent fit on the real-world facebook graphs, and that again 1 dimension seems already to bring a high amount of accuracy - increasing d only brings small gains, e.g. looking at graphs of size about $n = 8000$ nodes, the PEC increase from 1d to 2d cube GIRG fitting is just 22% \rightarrow 25%. In contrast for fitting a 2d cube GIRG with a 1d cube GIRG, PEC jumps from 42% \rightarrow 55%; fitting a 3d with 1d PEC jumps from 0.35 \rightarrow 0.55. We can observe in fig. 5.6 the final fit point locations of 2d GIRGs to a few graphs (2d because this is easiest for visualisation in a plot) - we see that points for socfb-Bowdoin47 and especially socfb-Haverford76 look like they could be pretty well approximated by 1d fit locations, whereas socfb-Rice31 has more complexity going on, and an extra dimension seems important to fully capture the graph structure.

TODO add in GIRG generated graphs final PEC vs Nodes plot.

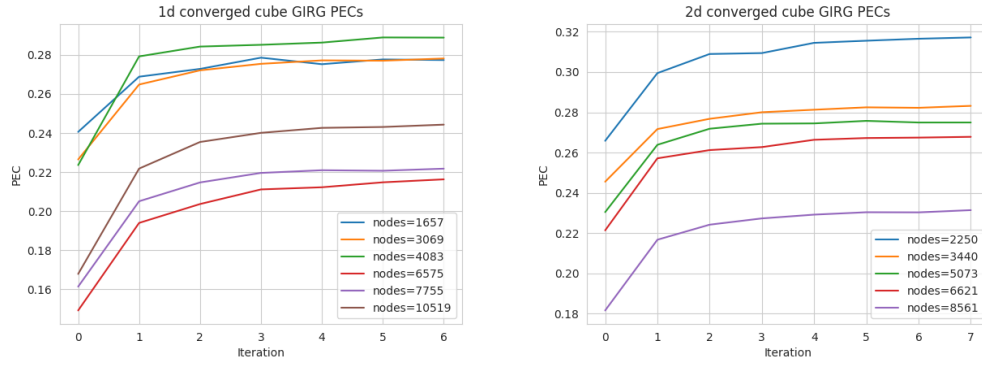
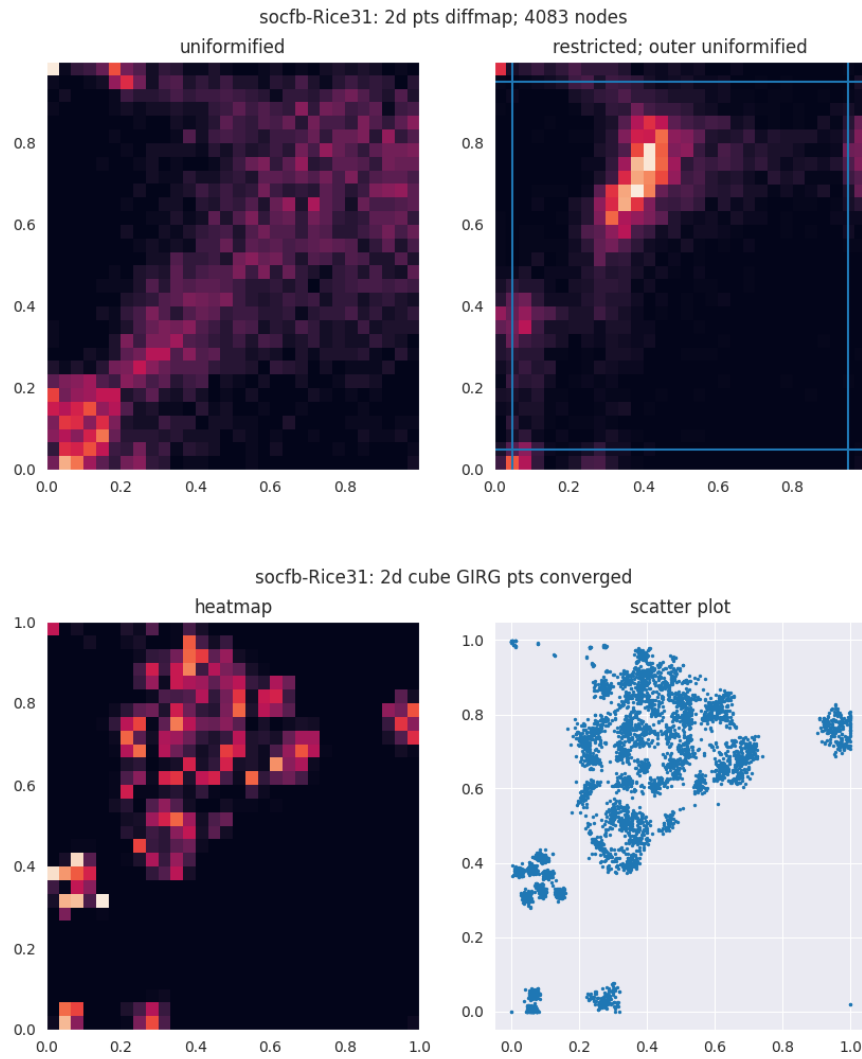
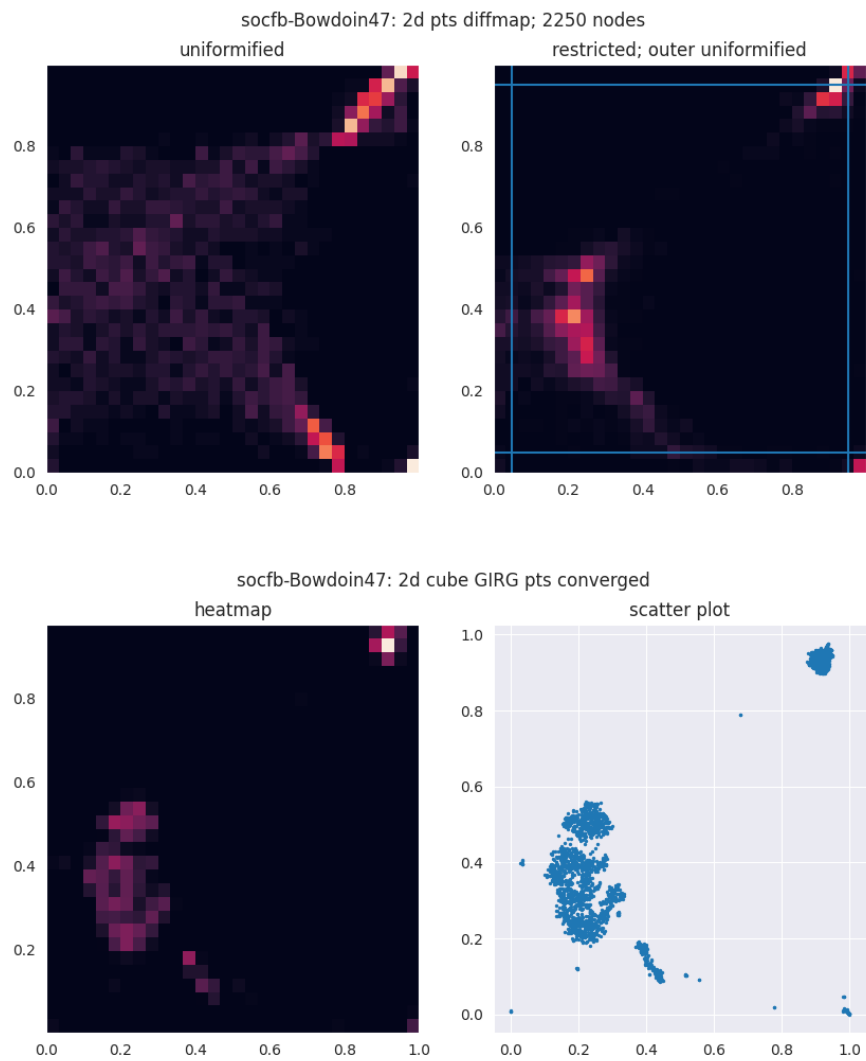
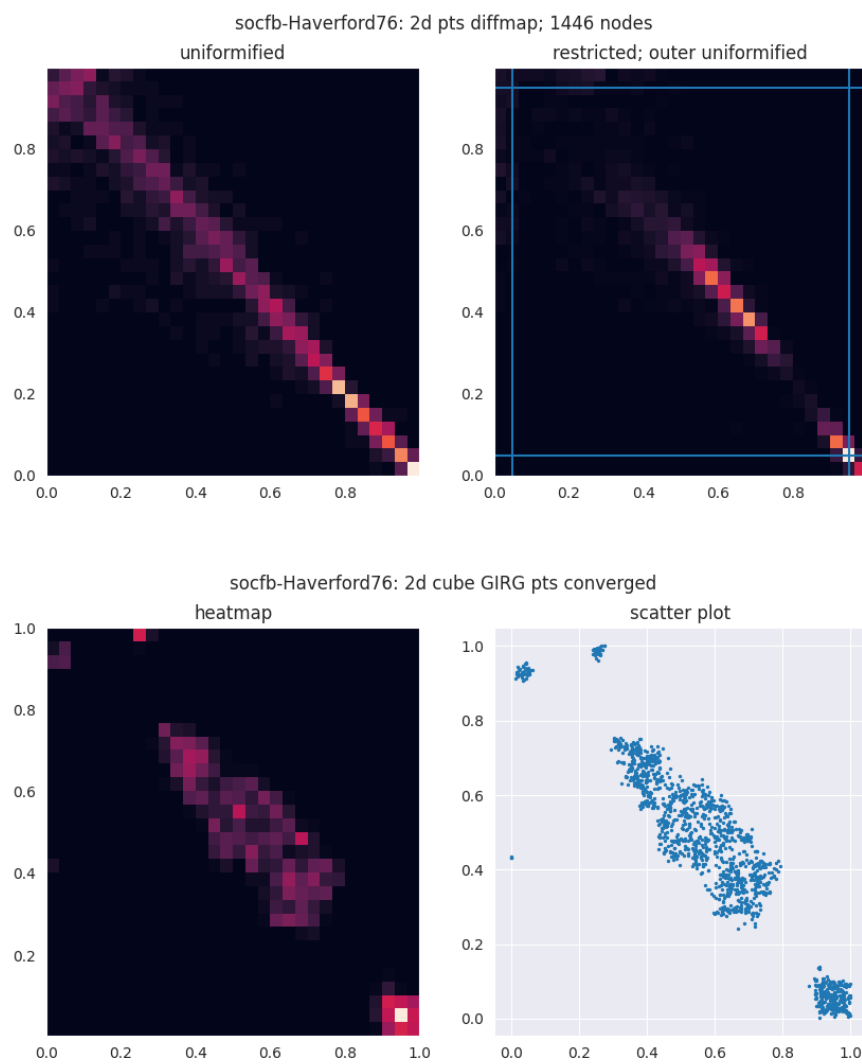


Figure 5.5: Direct ordered likelihood maximisation PEC convergence over iterations for a range of example socfb graphs with different number of nodes. Each curve is a single fit graph. X-axis is round number, where each round every node is moved once to a new location. Y-axis is PEC after that round. We claim convergence due to plateaued PEC, and note that in general the final achieved PEC decreases with graph node count.

Figure 5.6: 2d cube GIRG fitting examples, to three different Facebook graphs. Two initial diffusion map embeddings with different post-processing methods. The blue lines for the restricted version show the border at which points are outer uniformified. Then heat map / scatter plot for further refined point locations using direct ordered likelihood maximisation.







Graph Kernels

Contents

6.1	Bayes Factor Theory	53
6.1.1	Graph isomorphism problem in computing model evidence	54
6.2	Graph Kernel Introduction	55
6.2.1	Random walk kernel	55
6.2.2	Weisfeiler-Lehman kernel	55
6.2.3	Experiments	56

In this chapter we try to compare the bayesian evidence based plausibility of different GGMs in generating Facebook graphs by using the method of *graph kernels*. Unfortunately the whole endeavour was unsuccessful, but we nonetheless record here our experiments and the framework.

6.1 Bayes Factor Theory

A bayesian approach to model selection is to compare the posterior probability of different models $p(M_1|D)$ vs $p(M_2|D)$, given the model generative processes $P(D|M)$. M_1, M_2 are possible models of the data, e.g. $M_1 = \mathcal{G}_{1D} \text{ GIRG}$ and $M_2 = \mathcal{G}_{CL}$. D is a single graph instance G , or alternatively a whole dataset of our 100 Facebook graphs, assuming that they all come from the same generative model family. We have some prior of $p(M_1)$ vs $p(M_2)$, e.g. $50 : 50$. We could even possibly have a set of models $M_{d=1}, M_{d=2}, M_{d=3}, \dots$ for different dimensional GIRGs with some kind of decaying $p(1d \text{ GIRG}) > p(2d \text{ GIRG}) > p(3d \text{ GIRG}) > \dots$

For now focussing on just M_1 vs M_2 ,

$$p(M_1|D) = \frac{p(D|M_1)p(M_1)}{p(D)} \implies \frac{p(M_1|D)}{p(M_2|D)} = \frac{p(D|M_1)p(M_1)}{p(D|M_2)p(M_2)} \quad (6.1)$$

Model selection is done with the ratio $\frac{p(M_1|D)}{p(M_2|D)}$ which is called the *Bayes Factor*. In order to compute the bayes factor, we will for now ignore the priors $p(M_1), p(M_2)$, and hence focus just on the likelihoods $p(D|M_1), p(D|M_2)$.

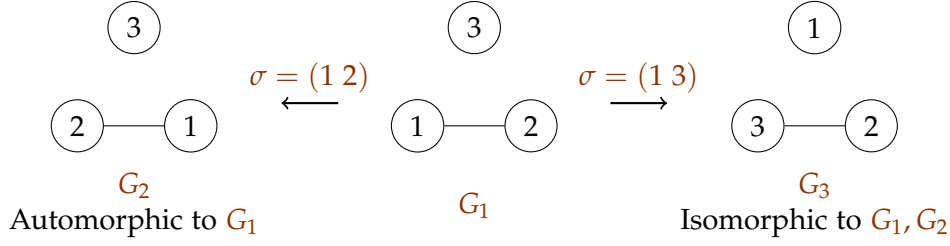


Figure 6.1: Example automorphic/isomorphic graphs

6.1.1 Graph isomorphism problem in computing model evidence

In our case if M_1 is a 1D GIRG, and we've decided to fix $\alpha, \{w_u\}_{u \in V}$ and just vary $\theta = c, \{x_u\}_{u \in V}$, then we could Monte Carlo sample $\theta \sim p(\theta)$ from the prior to get an estimate for

$$p(D|M_1) = p(G|\mathcal{G}_{1d-GIRG}) = \int_{\theta} p(G|\theta, \mathcal{G}_{1d-GIRG}) p(\theta|\mathcal{G}_{1d-GIRG}) d\theta \quad (6.2)$$

The simplified notation is dropping the $|M_1$ as we fix into some dimensional GIRG universe, $\int_{\theta} p(G|\theta) p(\theta) d\theta$.

There's a problem here. Since our data is a graph $D = G$, in the computation we actually need to inspect $p(G|\theta) = \sum_{\sigma} p(G' \stackrel{\sigma}{\cong} G|\theta) p(\sigma)$. Here σ is a permutation of the node ID, and $\stackrel{\sigma}{\cong}$ denotes *graph isomorphism* under this permutation. Note here our abuse of notation - normally $p(G|\theta)$ is not a sum over permutations, but rather the probability of generating this exact graph G ; this will have to be contextually apparent. So instead of evaluating $p(G|\theta)$, we rather want to evaluate "the probability of producing G' , given parameters θ , which is isomorphic to G ".

Small concrete example of a labelled graph $G_1 = (V, E) = (\{1, 2, 3\}, \{(1, 2)\})$ which is shown in fig. 6.1. G_1 is isomorphic to any other labelled graph H if they look the same if you anonymise the nodes. More formally, if $\sigma : V(G) \rightarrow V(H)$ is a bijection mapping nodes to nodes, such that $u \sim v$ in $G \iff \sigma(u) \sim \sigma(v)$ in H , then G is isomorphic to H with permutation σ , written $G \stackrel{\sigma}{\cong} H$. So G_1 is isomorphic to G_3 with edge set $\{(3, 2)\}$. You can even say that it is automorphic to the graph H with edge set $\{(2, 1)\}$, i.e. just switching nodes 1, 2. On this graph G_1 of 3 nodes, there are $3! = 6$ node ID permutations, each of which produces an isomorphic graph G' . For our particular G_1 , we could group together these 6 isomorphic graphs G' into 3 pairs of automorphic graphs.

Finally repeat for further θ samples and take an outer average.

Consider a simplified Chung-Lu / GIRG model of three nodes where of $\theta = (x_1, x_2, x_3)$ and all nodes have the same weight 1. The Chung-Lu has identical $p(G')$ for each of the 6 isomorphisms. The GIRG model does not - if x_1, x_2 are close, and x_3 is far from both, then it awards higher probability to $G' = (V, \{(1, 2)\})$ and $G' = (V, \{(2, 1)\})$. Furthermore both models always award the same probability to any equivalent class of automorphic graphs - this is because the adjacency matrix is the same for automorphic graphs, which is all that the models care about.

Unfortunately computing the $n!$ length sum $\sum_{\sigma} p(G' \stackrel{\sigma}{\cong} G|\theta) p(\sigma)$ is infeasible for all but tiny graphs. As an alternative, we hoped to use a graph similarity ker-

nel k which compares two graphs G, H , giving some measure $k(G, H)$ of how similar they are up to isomorphism (in a node permutation invariant way). Therefore the idea would be to replace the incorrect $p(G|\mathcal{G}_{\text{GIRG}}) = \int_{\theta} p(G|\theta)p(\theta)d\theta$ with $E_{G' \sim \mathcal{G}_{\text{GIRG}}} [k(G, G')] \approx \frac{1}{m} \sum_{i=1}^m k(G, G'_i \sim \mathcal{G}_{\text{GIRG}})$, computable by Monte Carlo sampling.

6.2 Graph Kernel Introduction

Graph Kernels provide a means for a similarity metric between graphs. They're ideally a positive semidefinite function $k : \Gamma \times \Gamma \rightarrow \mathbb{R}$, where Γ is the set of all graphs. Such a function exists if and only if there is a corresponding feature map resenation of $\phi : \Gamma \rightarrow \mathcal{H}_{\phi}$ where \mathcal{H}_{ϕ} is a Hilbert space, and $k(G, G') = \langle \phi(G), \phi(G') \rangle_{\mathcal{H}_{\phi}}$ is just an inner product.

Graph kernels can be contrasted with the Blasius' GGM realism framework . Blasius compares multiple different graph feature combinations on which to train an SVM for distinguishing two graph datasets. Viewing the graph kernel as an inner product of graph feature maps, these hopefully encapsulate all the relevant information about the graph including the specifically chosen features in Blasius' framework like graph diameter, average degree, etc. The question of which chosen graph features to consider then shifts to which graph kernel to use!

There are a range of graph kernels to choose from, however given the relatively large size of our graphs, runtime can be an issue. We ended up testing two kernels; alas both proved unsatisfactory. Without further expertise on graph kernels, we were forced to abandon this line of attack. We think that graph kernels might be more effective on smaller graphs with more unique structure (our random graph models and the real graphs themselves have a lot of homogenous structure) and particularly graphs with meaningful node labels.

6.2.1 Random walk kernel

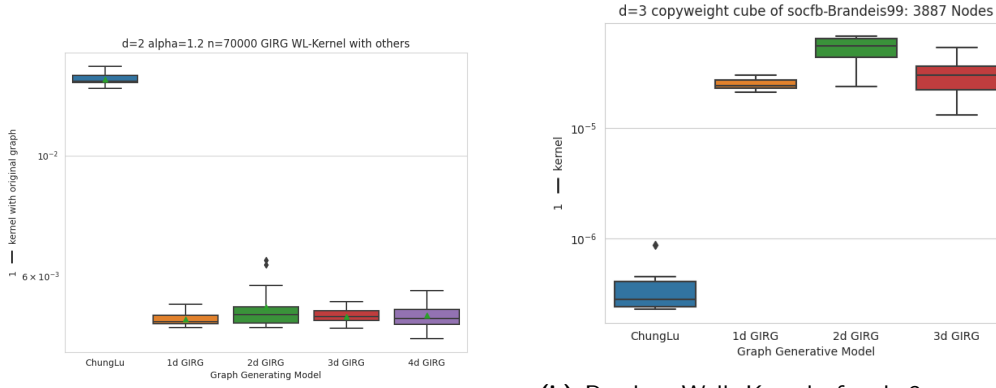
The *random walk kernel* between two graphs is just a count of the number of common walks in the pair, of any length l . This is generally an infinite sum, so geometric weighting with the factor λ^l is used to decay the contribution of longer walks, where $0 < \lambda < 1$ (we tested with $\lambda = 10^{-5}$).

The product graph of G and G' is defined as $G_{\times} = (V_{\times}, E_{\times})$ where $V_{\times} = \{(v, v') | v \in V, v' \in V'\}$ and $E_{\times} = \{((v_1, v'_1), (v_2, v'_2)) | v_1 \sim v_2 \in E_1, v'_1 \sim v'_2 \in E_2\}$. This definition can be extended to node-labelled graphs, where we only take product vertices $(v, v') \in V_{\times}$ if $v \in V, v' \in V'$ have the same label. Common walks in the pair of graphs G, G' are just walks in the product graph.

The naive implementation of the random walk kernel has complexity $O(n^6)$, however this can be sped up to $O(n^3)$. This is still too slow for our purposes. We only made limited tests with small graphs of $n \leq 3000$ nodes - see e.g. fig. 6.2b.

6.2.2 Weisfeiler-Lehman kernel

The *Weisfeiler-Lehman kernel* runs much faster in $O(h|E|)$ time, where $|E|$ is the number of edges in the graph and h is a chosen parameter, the number of rounds of



(a) Weisfeiler-Lehman Kernel of a $d=2$, $\alpha=1.2$, $n=70000$ Torus GIRG with other generated graphs (13 per model). All the GIRGs are more similar to the original than Chung-Lu, which is good, but we cannot differentiate between GIRGs of different dimensions.

(b) Random Walk Kernel of a $d=3$ copy weight cube GIRG fit to socfb-Brandeis99 (matching number of edges and mean LCC), with other graphs generated from different GGMs (6 per model type). We hoped for 3d GIRGs to have the highest similarity to the input graph, however instead Chung-Lu graphs matched closest.

relabelling within the algorithm. The algorithm requires some starting node labels $(l(v))_{v \in V}$ - we colour nodes into a small discrete set of colours by grouping together nodes with similar sized degrees. It also uses a base kernel k_{base} which is computed at each relabelling round - we used the simplest/speediest node label histogram dot product kernel: $k(G, G') = \langle f, f' \rangle$ where $f = (f_1, \dots, f_k)$, $f_i = |\{v \in V : l(v) = i\}|$ is the number of nodes with label i .

The Weisfeiler-Lehman kernel is then defined as $k_{WLK}(G, G') = k_{\text{base}}(G_1, G'_1) + \dots + k_{\text{base}}(G_h, G'_h)$. $G = G_1$, $G' = G'_1$, and G_{i+1} is computed iteratively as relabelling each node in G_i with $l(v) \leftarrow (l(v), (l(u))_{u \sim v})$. This looks odd, but in practice each label is just rehased as a new integer rather than becoming a highly nested tuple.

6.2.3 Experiments

Unfortunately both kernels failed to pass the basic test of showing higher similarity between two graphs generated from the same GGM than from a different GGM. We tried a variety of combinations, but didn't get very reliable results: see fig. 6.2a and fig. 6.2b.

Conclusion

Contents

7.1 Recap	57
---------------------	----

7.1 Recap

The goal of this thesis was to explore the ability of different types of GIRG model to match real social network graphs with some suspected inherent geometry and power law degree distribution. In documenting our various experiments, we also detail our methods and ideas for working with GIRGs and graph data which may prove useful to other practitioners.

Our first effort was to follow the realism framework of [Bläsius et al., 2018], to compare a range of generative graph models, including many GIRG subtypes, in their ability to replicate global statistics of a set of ~ 100 social network graphs. With this framework we did find GIRGs able to outperform on a few statistics: closeness centrality, betweenness centrality, effective diameter, and mean local clustering coefficient. We identified differences between different GIRG variants, with max cube GIRGs excelling at diameter replication, MCD GIRGs on closeness centrality, and max torus GIRGs on betweenness centrality. Not much difference was found by increasing the geometric dimension, suggesting that 1-3 dimensions may be sufficient for high level similarity to real graphs, and indeed 1d is often the best.

We then explored fitting GIRG node geometric location parameters to the real graphs, to see if GIRGs are capable of exact replication on a node and edge level. This proved quite successful, showing that GIRGs can capture a high percentage of edges, even with just 1 dimension. Adding further dimensions is computationally intensive and doesn't increase the fit substantially. Our range of fitting methods might be further refined in practice according to the use case.

Future work might investigate more into how 1 dimensional geometry already fits the data quite well, in particular seeing if the variants of distorted GIRGs or min/max mixed GIRGs can explain a relative ordering of dimension importance. Intuitively 1 dimension may be the most influential for our social connections, but surely our friendships are more multifaceted.

Thanks to my supervisors Marc Kaufmann and Ulysse Schaller, as well as Johannes Lengler, for their valuable input and guidance throughout the thesis. Many thanks to Angelika Steger and her lab for their friendly and welcoming research environment. Finally much appreciation to friends, family, ETH, etc. for their support and existence. And thanks to the reader for making it this far! I hope it's been an enjoyable/interesting reading experience.

Bibliography

- [Belkin and Niyogi, 2001] Belkin, M. and Niyogi, P. (2001). Laplacian eigenmaps and spectral techniques for embedding and clustering. *Advances in neural information processing systems*, 14.
- [Berry and Harlim, 2018] Berry, T. and Harlim, J. (2018). Iterated diffusion maps for feature identification. *Applied and Computational Harmonic Analysis*, 45(1):84–119.
- [Bläsius et al., 2018] Bläsius, T., Friedrich, T., Katzmann, M., Krohmer, A., and Striebel, J. (2018). Towards a systematic evaluation of generative network models. In *Algorithms and Models for the Web Graph: 15th International Workshop, WAW 2018, Moscow, Russia, May 17-18, 2018, Proceedings 15*, pages 99–114. Springer.
- [Bläsius et al., 2022] Bläsius, T., Friedrich, T., Katzmann, M., Meyer, U., Penschuck, M., and Weyand, C. (2022). Efficiently generating geometric inhomogeneous and hyperbolic random graphs. *Network Science*, 10(4):361–380.
- [Bringmann et al., 2016] Bringmann, K., Keusch, R., and Lengler, J. (2016). Average distance in a general class of scale-free networks with underlying geometry. *arXiv preprint arXiv:1602.05712*.
- [Bringmann et al., 2019] Bringmann, K., Keusch, R., and Lengler, J. (2019). Geometric inhomogeneous random graphs. *Theoretical Computer Science*, 760:35–54.
- [Coifman and Lafon, 2006] Coifman, R. R. and Lafon, S. (2006). Diffusion maps. *Applied and computational harmonic analysis*, 21(1):5–30.
- [Friedrich et al., 2023] Friedrich, T., Göbel, A., Katzmann, M., and Schiller, L. (2023). Cliques in high-dimensional geometric inhomogeneous random graphs. *arXiv preprint arXiv:2302.04113*.
- [García-Pérez et al., 2019] García-Pérez, G., Allard, A., Serrano, M. Á., and Boguñá, M. (2019). Mercator: uncovering faithful hyperbolic embeddings of complex networks. *New Journal of Physics*, 21(12):123033.
- [Rossi and Ahmed, 2015] Rossi, R. and Ahmed, N. (2015). The network data repository with interactive graph analytics and visualization. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29.

[Singer, 2006] Singer, A. (2006). From graph to manifold laplacian: The convergence rate. *Applied and Computational Harmonic Analysis*, 21(1):128–134.



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Declaration of originality

The signed declaration of originality is a component of every semester paper, Bachelor's thesis, Master's thesis and any other degree paper undertaken during the course of studies, including the respective electronic versions.

Lecturers may also require a declaration of originality for other written papers compiled for their courses.

I hereby confirm that I am the sole author of the written work here enclosed and that I have compiled it in my own words. Parts excepted are corrections of form and content by the supervisor.

Title of work (in block letters):

Authored by (in block letters):

For papers written by groups the names of all authors are required.

Name(s):

First name(s):

With my signature I confirm that

- I have committed none of the forms of plagiarism described in the '[Citation etiquette](#)' information sheet.
- I have documented all methods, data and processes truthfully.
- I have not manipulated any data.
- I have mentioned all persons who were significant facilitators of the work.

I am aware that the work may be screened electronically for plagiarism.

Place, date

Signature(s)

For papers written by groups the names of all authors are required. Their signatures collectively guarantee the entire content of the written paper.

