

CHARACTERISTIC LENGTH AND CLUSTERING

OLIVER KNILL

ABSTRACT. We explore relations between various variational problems for graphs: among the functionals considered are Euler characteristic $\chi(G)$, characteristic length $\mu(G)$, mean clustering $\nu(G)$, inductive dimension $\iota(G)$, edge density $\epsilon(G)$, scale measure $\sigma(G)$, Hilbert action $\eta(G)$ and spectral complexity $\xi(G)$. A new insight in this note is that the local cluster coefficient $C(x)$ in a finite simple graph can be written as a relative characteristic length $L(x)$ of the unit sphere $S(x)$ within the unit ball $B(x)$ of a vertex. This relation $L(x) = 2 - C(x)$ will allow to study clustering in more general metric spaces like Riemannian manifolds or fractals. If η is the average of scalar curvature $s(x)$, a formula $\mu \sim 1 + \log(\epsilon)/\log(\eta)$ of Newman, Watts and Strogatz [31] relates μ with the edge density ϵ and average scalar curvature η telling that large curvature correlates with small characteristic length. Experiments show that the statistical relation $\mu \sim \log(1/\nu)$ holds for random or deterministic constructed networks, indicating that small clustering is often associated to large characteristic lengths and $\lambda = \mu/\log(\nu)$ can converge in some graph limits of networks. Mean clustering ν , edge density ϵ and curvature average η therefore can relate with characteristic length μ on a statistical level. We also discovered experimentally that inductive dimension ι and cluster-length ratio λ correlate strongly on Erdős-Renyi probability spaces.

1. INTRODUCTION

The interplay between global and local properties often appears in geometry. As an example, the Euler characteristic can by Gauss-Bonnet be written as an average of local curvature and by Poincaré-Hopf as an average of local index density i_f for Morse functions f . Given a globally defined quantity on a geometric space, one can ask to which extent the functional can be described as an average over local properties. Also of interest is the relation between the various functionals. We will look at some examples on the category of finite simple graphs and comment

Date: October 12, 2014.

1991 Mathematics Subject Classification. 05C35, 52Cxx .

Key words and phrases. Characteristic length, Clustering, Euler characteristic, Variational problems in graph theory.

on both problems. We look then primarily at characteristic length μ , which is the expectation of non-local mean distance $\mu(x)$ on a metric space equipped with a probability measure. On finite simple graphs one has a natural geodesic distance and a natural counting measure so that networks allow geometric experimentation on small geometries. Most functionals are interesting also for Riemannian manifolds, where finding explicit formulas for the characteristic length can already lead to challenging integrals. Characteristic length is sometimes also defined as the statistical median of $\mu(x)$ [37]. Since the difference is not essential, we use the averages

$$\mu(G) = \frac{1}{n^2 - n} \sum_{x \neq y \in V^2} d(x, y), \nu(G) = \frac{1}{n} \sum_{x \in V} \frac{2e(x)}{n(x)(n(x) - 1)},$$

where $n(x), e(x)$ are the number of vertices and edges in the sphere $S(x)$ of the vertex x . μ is the average of the non-local quantity $D(x) = \frac{1}{n-1} \sum_{y \neq x} d(x, y)$ and ν is the average of the local quantity $C(x) = |e(x)|/B(n(x), 2)$ giving the fraction of connections in the unit sphere in comparison to all possible pairs in the unit sphere.

This averaging convention for μ is common. [10] showed already that $\mu(G) \leq \text{diam}(G) - 1/2$. An other notion is the variance $v(G) = \max_x d(x) - \min_x d(x)$ where $d(x) = \sum_y d(x, y)$ for which Ore has shown that on graphs with n vertices has the maximum taken on trees. [11]. The definition of $\nu(G)$ is to average the edge density of the sphere relative to the case when the sphere is the complete graph with n vertices in which case the number of edges is $n(x)(n(x) - 1)/2$. Both quantities are natural functionals. The mean cluster coefficient ν is by definition an average of local quantities. While it is impossible to find a local quantity whose average captures characteristic length exactly, there are notions which come close. We look at three such relations. The first is a formula of Newman, Watts and Strogatz [31] which writes μ as a diffraction coefficient divided by scalar curvature. This is intuitive already for spheres, where the signal speed and the curvature determines the characteristic length. Empirically, we find an other quantity which also often allows to estimate characteristic length well: it is the mean cluster density ν , the average of a local cluster density $C(x)$ as defined by Watts and Strogatz. Thirdly, we will report on some experiments which correlate the length-cluster coefficient $\lambda = \mu / \log(1/\nu)$ with the inductive dimension ι of the graph.

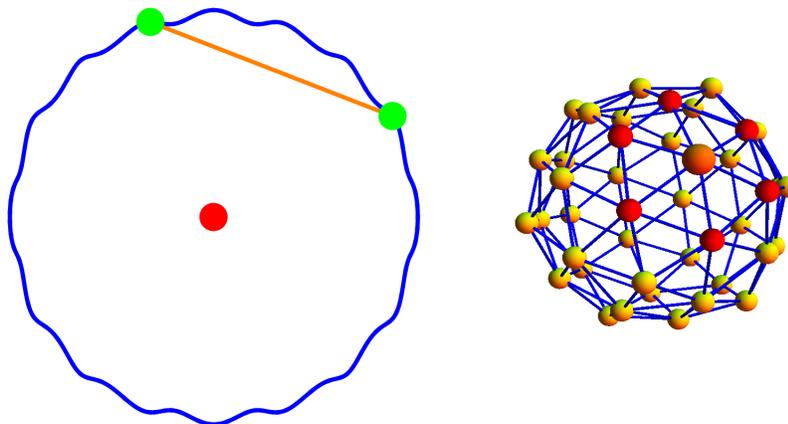


FIGURE 1. For graphs, the clustering coefficient $C(x)$ at a vertex x is related to the relative characteristic length $L(x)$ of the unit sphere $S(x)$ within the unit ball $B(x)$. In other words, $L(x)$ is the average distance between two points in $S(x)$ within $B(x)$. The relation $C(x) = 2 - L(x)$ allows so to define clustering in any metric space equipped with a measure inducing conditional probability measures on spheres. It is a local quantity which is constant in the radius r of the spheres if we are in Euclidean spaces or fractal spaces with some uniformity. The second figure shows a graph for which the local clustering coefficients $C(x)$ take the values $1/2$ or $2/5$, where the global characteristic length is $\mu = 2.9$, the mean clustering coefficient is $\nu = 3/7$ and the length-cluster coefficient $\lambda = -\mu/\log(\nu)$ is 3.4.

The starting point of this note is the observation that the local cluster property $C(x)$ if a vertex in a graph can be written in terms of relative characteristic length of the unit sphere $S(x)$ within the unit ball $B(x)$. We are not aware that this has been noted already, but it is remarkable as it allows to carry over the definition of “local cluster property” to metric spaces equipped with a probability measure as long as the measure has the property that it induces probability measures $m(\cdot, S)$ on spheres by conditional probability $m(A) = \lim_{\epsilon \rightarrow 0} m(N_\epsilon(A \cap S))/m(N_\epsilon(S))$, where $N_\epsilon(Y)$ is an ϵ neighborhood of Y . For such metric spaces, one can also define a cousin of scalar curvature $\log(2^{\iota(x)-1}\delta/\delta_2)$ by comparing the measures δ, δ_2 of spheres of radius 1 and 2. We call it “scalar curvature” because if $\iota(x)$ is the inductive dimension of space at that point, then $\log(2^{\iota(x)-1}\delta/\delta_2)$ is zero if

the volume $\delta_2(x)$ of the sphere of radius 2 is equal to $2^{\dim(S(x))} = 2^{\iota(x)-1}$ times the volume δ of a sphere of radius 1. Comparing with Riemannian manifolds, the comparison of spheres $S_{2r}(x)$ and $S_r(x)$ allows to measure scalar curvature. The dimension scaling factor $\log(2)$ does not matter much because the Hilbert action $\eta(G)$, the expectation of the scalar curvature is $\log(2)\iota(G)$ plus the expectation of $s(x) = \log(\delta/\delta_2)$. Since the later is a local function and appears also in a formula of Newman, Watts and Strogatz and because we look at the dimension $\iota(G)$ anyway, it does not matter if for simplicity $s(x)$ is called the **scalar curvature** and its average over vertex set is called the **Hilbert action** $\eta(G)$ of the graph.

Having these functionals, one can now study the relation between local cluster property, characteristic length, dimension and curvature on rather general metric spaces equipped with a natural measure. But first of all these functionals could be used to select “natural geometries”. Functionals are important in physics because most fundamental laws are of variational nature. In a geometric setup and especially in graph theory, one can consider the Euler characteristic, the characteristic length, the Hilbert action $\eta(G)$ given as an average scalar curvature or the spectral complexity $\xi(G)$, given as the product of the nonzero eigenvalues of the Laplacian L of G . An other related quantity is the number of rooted spanning forests in G which is $\theta(G) = \det(1 + L)$ by the Chebotarev-Shamis theorem [33, 32, 19, 24]. Many other extremal problems are studied in graph theory [4].

For finite simple graphs, the Euler characteristic $\chi(G) = v_0 - v_1 + v_2 - \dots$ with v_k is the number of k dimensional simplices K_{k+1} of G . In geometric situations like four dimensional geometric graphs, for which the spheres are discrete three dimensional spheres, this number can be seen as a quantized Hilbert action [21] because it is an average of the Euler characteristic over all two-dimensional subgraphs and then via Gauss-Bonnet an average over a set of sectional curvatures and all points. The characteristic length is the average length of a path between two points relating to the other variational problem in general relativity. The Hilbert action itself is an average of scalar curvature, which can be defined for a large class of metric spaces. Spectral complexity is natural because of the matrix tree theorem of Kirkhoff (see i.e. [3]) relates it with the number of trees in a graph.

While the Euler characteristic is an average of curvature by Gauss-Bonnet both in the Riemannian and graph case [16] and Hilbert action is an average of scalar curvature, both the characteristic length μ as well as the complexity ξ can not be written as an average of local properties: look at two disjoint graphs connected along a one-dimensional line graph. Cutting that line in the middle can change both quantities in a very different way depending on the two components which are obtained. If μ or ξ were local, the functional would change by a definite value, independent of the length of the “rope”. The characteristic length has been noted to be relevant for molecules: the chemist Harry Wiener found correlations between the **Wiener index** $W(G)$, the sum over all distances which is $W(G) = n(n-1)\mu(G)$ and the melting point of a hydrocarbon G . As reported in [13], the characteristic length has first been considered in [10] in graph theory. For spectral properties, see [29]. The observation that the Wiener index satisfies $W(G) \leq W(T)$ for any spanning tree T of G was made in [34].

The characteristic length $\mu(G)$ has been studied quite a bit. There are few classes of graphs, where one can compute the number explicitly: for complete graphs, we have $\mu(K_n) = 1$ for complete bipartite graphs $\mu(K_{a,b}) = (2a^2 + 2b^2 + ab)/((a+b)(a+b-1))$ for line graphs $\mu(L_n) = (n+1)/3$, for cyclic graphs $\mu(C_n) = (n+1)/4$. for odd n and $n^2/(4(n-1))$ if n is even [11]. No general relation between length and diameter exists besides the trivial $\mu(G) \leq \text{diam}(G)$. We have $1 \leq \mu(G) \leq (n+1)/3$ [10], On the class $G(n,m)$ of graphs with n vertices and m edges one has $\mu(G) \leq 2 - (2m)/(n(n-1))$ [11]. Among all graphs with n vertices the maximum $(n+1)/3$ is obtained for line graphs, the minimum 1 for complete graphs [10]. The problem to find the maximum among all graphs of given order and diameter is unknown. There are relations with the spectrum: $\mu(G) \leq b - (2(b-1)m)/(n(n-1))$ on $G(n,m)$ where b is the number of distinct Laplacian eigenvalues. There are also upper bounds in terms of the second eigenvalue. For a connected graph $\mu(G) \leq \beta(G)$ where $\beta(G)$ is the independence number [8], the maximal number of pairwise nonadjacent vertices in G . On all graphs of order n and minimal degree δ , then $\mu(G) \leq n/(\delta+1)+2$ [28]. The conjecture generating computer program Graffiti [12] suggested for constant degree δ graphs to have $\mu(G) \leq n/\delta$. There is a spectral relation $\mu \geq \text{tr}(L^+)2/(n-1)$, where L^+ is the Moore pseudo inverse of the Laplacian L and n the number of vertices([35] which is the McKay equality for trees and otherwise always a strict inequality.

The Euler characteristic is definitely one of most important functionals in geometry if not the most important one. It is a homotopy invariant and can by Poincaré-Hopf be expressed cohomologically as $\sum_{i=0}(-1)^i b_i$ using b_i the dimensions of cohomology groups $H^i(G)$. A general theme in topology is to extend the notion of Euler characteristic to larger classes of topological spaces. This essentially boils down to the construction of cohomology. The limitations are clear already in simple cases like the Cantor set which have infinite Euler characteristic as half of the space is homeomorphic to itself. Euler characteristic can be defined for a metric space if there exists a subbasis of contractible graphs. The Euler characteristic can then be defined as the Euler characteristic of the nerve graph. This illustrates already how important homotopy is in general when studying Euler characteristic. It is also historically remarkable that the first works done by Euler on Euler characteristic were of homotopy nature by deforming the graph.

Various other functionals have been considered on graphs. The **average centrality** $f(G)$ is the mean of the **local closeness centrality**

$$f(x) = \sum_y \frac{1}{\sum_{y \neq x} d(x, y)}$$

of a vertex x . An other number is the **geodetic number** $g(G)$ which is the minimum cardinality of a geodetic set in G , where a set is called **geodetic** if its geodesic closure is G [1]. The **scale measure** of a graph is defined as $\sigma(G) = s(G)/m$, where $s(G) = \sum_{e \in E} d(e)$ and $d((a, b)) = d(a)d(b)$ and $m = \max_{e \in E} d(e)$. An other important notion is the **chromatic number** $c(G)$ which can be seen as the smallest p for which a scalar function with values in Z_p exists for which the gradient field df nowhere vanishes. Related to graph coloring, we have in [25] defined **chromatic richness** $C(c)/c!$ measuring the size of the set of coloring functions modulo permutations. The **arboricity** $a(G)$ of G is the minimal number of spanning forests which are needed to cover all edges of G . One knows that $\theta(G) = \det(L + 1)$ by Chebotarev-Shamis [33, 32, 19, 24]. While the number of spanning forests is a measure of complexity, the arboricity is a measure of “denseness” of the graph. The **Nash-Williams formula** [30, 7] tells that the arboricity is the maximum of $[m_H/(n_H - 1)]$, where n_H is the number of vertices and m_H the number of edges of a subgraph H of G and where $[r]$ is the ceiling function giving the minimum of all integers larger or equal than r . For example, for $K_{4,4}$ where $m = 16, n = 8$ the arboricity must be at least $16/7 = 2.28$ and so at least 3 and one can give examples of three

forests covering all edges. For a complete graph, the arboricity is $\lceil n/2 \rceil$. The arboricity gives a bound on the chromatic number $c(G) \leq 2a(G)$ (a fact noted in [6] and follows from the fact that each forest can be colored by 2 colors). The **Laplacian ratio** $p(G) = \text{per}(L) / \prod_i d_i$, where $\text{per}(L)$ is the permanent of the Laplacian and d_i are the vertex degrees has been introduced in [5]. The **symmetry grade** of a graph is the order of the automorphism group of G . For the complete graph for example, it is $n!$ while for a cyclic graph it is $2n$, the size of the dihedral group. The **domatic number** $d(G)$ of a graph finally is the maximal size of a dominating partition of the vertex set.

Functional		Based on	Local	Spectral
Euler characteristic	χ	Curvature	yes	somehow [18]
Inductive dimension	ι	Point dimension	yes	not known
Characteristic length	μ	Distance	no	on trees [13]
Complexity	ξ	Eigenvalues	no	yes
Forest complexity	θ	Eigenvalues	no	yes
Hilbert action	η	Scalar curvature	yes	not known
Mean cluster	ν	Local cluster	yes	yes
Average degree	δ	Vertex degree	yes	yes
Graph density	ϵ	Edge number	yes	yes
Scale measure	σ	Vertex degree	yes	not known
Cluster-length-ratio	λ	Distances	no	not known
Independence number	β	Adjacency	no	not known
Variance	v	Distance	no	not known
Centrality	f	Local centrality	yes	not known
Chromatic number	c	Gradient fields	no	no [9]
Arboricity	a	Forests	no	not known
Geodetic number	g	Geodesics	no	not known
Domatic number	d	Partitions	no	not known
Symmetry grade	t	Symmetry group	no	not known
Laplacian ratio	p	Permanent	no	not known

Besides the question whether a functional is local, it would also be interesting to know more about which properties are spectral properties. Euler characteristic can be seen as a spectral property **in the wider sense**: it is the super trace of e^{-tD^2} for the Dirac operator D for every t by McKean-Singer [18]. The average degree δ can be written in terms of the adjacency matrix A as $\delta = 2\text{tr}(A^2)/\text{tr}(A^0)$ and with the Laplacian L as $\text{tr}(L)/\text{tr}(L^0)$. The **graph density** $\epsilon = \delta/(n-1) = 2v_1/(v_0(v_0-1))$ is also spectral, because both δ and $n = v_0$ are spectral.

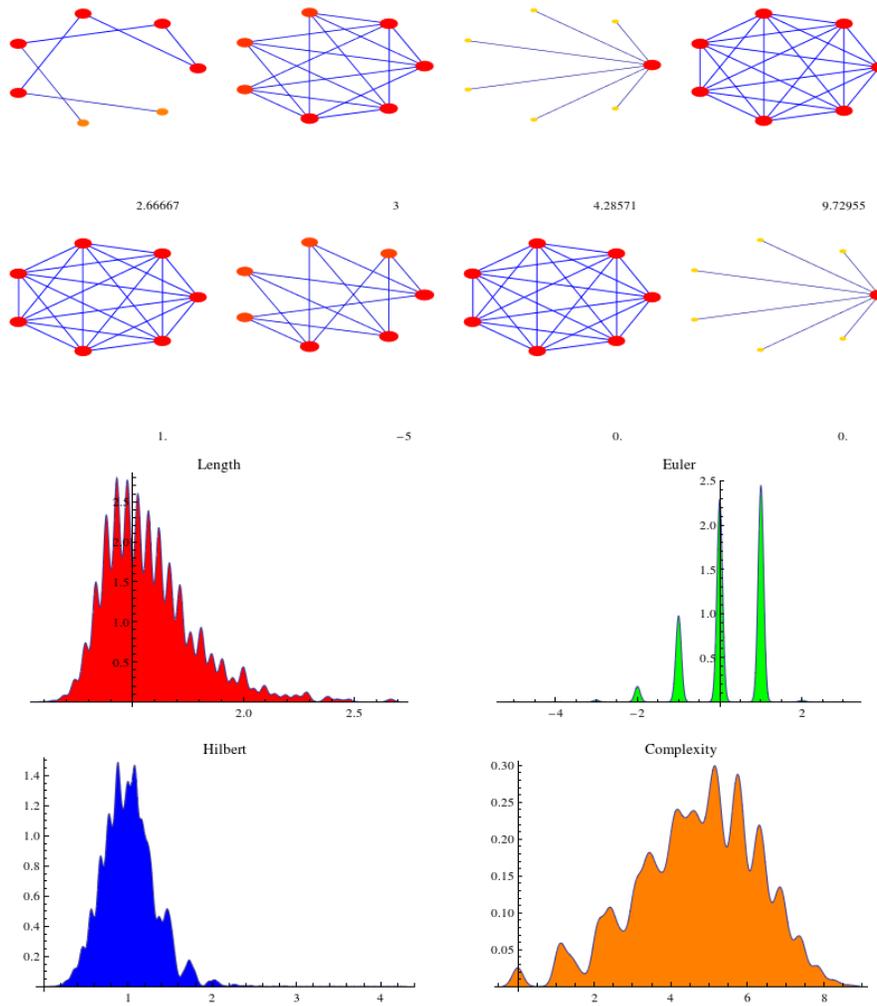


FIGURE 2. Graphs with minimal and maximal characteristic length, Euler characteristic, Hilbert action and the logarithm of the complexity among graphs all connected graphs with 7 vertices. We see then the distribution of the four functionals on this finite probability space with $1'866'256$ elements. The figure illustrates that Euler characteristic is the most interesting functional. All others appear to be extremal for line graphs, complete graphs or star graphs, all three of which do not carry interesting geometries.

2. CHARACTERISTIC LENGTH

A finite simple graph $G = (V, E)$ defines a finite metric space, where $d(x, y)$ is the geodesic distance between two vertices x, y , the length of the shortest path connecting the two points. The **characteristic length** L is the expectation of the distances between different vertices $\mu(G) = \frac{1}{n^2-n} \sum_{x \neq y} d(x, y)$, where $n = v_0 = |V|$ is the number of vertices. We have the understanding that the graph with only one vertex has zero length $\mu(G) = 0$ and that the number μ is averaged over every connected component of a graph independently. Unlike the Euler characteristic which is a homotopy invariant, the global characteristic length is a metric property as it depends on the concrete metric and is only invariant under graph isomorphisms and not under topological homeomorphisms [27] nor homotopies. Since determining μ requires to map out all the distances between any two points, it is natural to ask whether one can estimate the global length by averaging local properties. Such a relation has appeared in “formula 54” of [31],

$$\mu \sim 1 + \frac{\log(\delta/n)}{\log(\delta/\delta_2)},$$

where it was derived from generating functions. ‘Formula 54’ uses the average degree δ (the average of the degrees $\delta(x)$) and the average 2-nearest neighbors δ_2 (the average of the size $\delta_2(x)$ of the spheres of radius 2). Because $\delta = 2|E|/n$ by Euler’s handshaking lemma, we know that δ/n agrees with the **edge density** $\epsilon(G) = |E|/\binom{n}{2}$. Since $\log(2^{\dim(x)-1}\delta/\delta_2)$ measures a relation between volumes of spheres of distance r and $2r$, it can intuitively be thought of a scalar curvature, the flat case meaning the volume δ_2 of the sphere of radius 2 being larger than $2^{\iota(x)-1}$ times the volume δ of the sphere of radius 1 which has dimension $\iota(x) - 1$, where $\iota(x)$ denotes the dimension of the vertex x . The global characteristic length is according to “formula 54” an “edge density/curvature” relation which is intuitive because distances are small in spaces of positive curvature and the fact that if a material has large edge density, it allows fast travel. To summarize, “formula 54” relates the edge density ϵ , the Hilbert action η and the characteristic length μ as

$$\mu \sim 1 + \frac{\log(\epsilon)}{\log(\eta)},$$

where

$$\eta = \frac{1}{n} \sum_{x \in V} \log\left(\frac{\delta(x)}{\delta_2(x)}\right)$$

is the average of shifted scalar curvature $s(x) = \log(\delta_2/\delta)$.

Empirically, we also see a close relation between μ and $\log(1/\nu)$, where ν is the mean cluster coefficient as defined by [38]. We have mentioned this in [14] (see also [26, 23]). Clustering is a notion which is related to transitivity in sociology [36]. We see in many natural networks that these two quantities μ and $\log(\nu)$ grow linearly with $\exp(n)$ with a similar order of magnitude. Since $C(x)$ can be expressed integral geometrically using lengths, it can be pushed to natural metric spaces like compact Riemannian manifolds or metric spaces with fractal dimension.

Why is the characteristic length interesting? In physics, $d(a, b)$ is minimized by geodesics connecting a with b so that μ averages over all possible Lagrangian actions of paths between any two points. General relativity builds on two variational pillars: one is the Hilbert action η , the average scalar curvature, the other is the geodesic variational problem to minimize geodesic length between two points. The first tells how matter determines geometry, the second describes how geometry influences matter. ‘‘Formula 54’’ indicates a statistical correlation between Hilbert action, edge density and characteristic length. Since scalar curvature S appears in the expansion $|\exp_x(B_r(x))|/|B_r(x)| = 1 - Sr^2/(6(\dim + 2)) + \dots$, we can express this without referring to Euclidean space R^k as

$$\frac{|\exp_x(B_r)|}{|\exp_x(B_{2r})|} = 1 + r^2 \frac{S}{2(k+2)} + \dots$$

so that

$$S \sim \frac{2(k+2)}{r^2} \left(1 + \frac{|\exp_x(B_{2r})|}{|\exp_x(B_r)|}\right) \sim \frac{2(k+2)}{r^2} \log\left(\frac{|\exp_x(B_{2r})|}{|\exp_x(B_r)|}\right).$$

Since δ and δ_2 are averages, these are mean field approximations. We should mention that classically, scalar curvature can be described by volumes of balls not volumes of spheres. In the discrete, this does not matter. Since $|B(x)| = 1 + |S(x)|$ and $|B_2(x)| = 1 + |S(x)| + |S_2(x)|$ we have $B_2/B_1 \sim 1 + |S_2|/|S_1|$ and averaging also just changes the Hilbert action by a constant.

We would like to know more about the relation between Euler characteristic χ , graph density ϵ , dimension ι , complexity ξ , clustering ν and characteristic length μ . Also interesting are relations with corresponding notions of Riemannian manifolds. Besides average length

and dimension, complexity makes sense for Riemannian manifolds too, even so it needs zeta regularized determinants for its definition. The mean clustering ν is proportional to the volume with a proportionality factor which depends on the dimension. For a Riemannian manifold, scaling the space by a factor n scales the characteristic length in the same way and the volume by a factor n^d . For random networks, the global characteristic length typically grows like $\log(n)$ in dependence of volume n - one aspect of the “small world” phenomenon. This does not violate geometric intuition at all because dimension grows too with more nodes. We have explicit formulas [15] for the average dimension of a random Erdős-Renyi graph of size n , where edges are turned on with probability p . For other random graphs like Watts-Strogatz networks or networks generated by random permutations, we see similar growth rates. Other type of networks can show slightly different growth rates. Barabasi-Albert networks are examples, where the growth rate is slower.

Related to characteristic length is the **magnitude** $|G| = \sum_{i,j} Z_{ij}^{-1}$, where $Z_{ij} = \exp(-d(i, j))$ defined by Solow and Polasky. We numerically see that at least for small vertex cardinality n and connected graphs, the complete graph has minimal magnitude and the star graph maximal magnitude, a feature which is shared for many functionals (see Figure 2). Also the magnitude can be defined for more general metric spaces so that one can look for a general metric space at the supremum of all $|G|$ where G is a finite subset with induced metric. The **convex magnitude conjecture** of Leinster-Willington claims that for convex subsets of the plane, $|A| = \chi(A) + p(A)/4 - a(A)/(2\pi)$, where p is the perimeter and a is the area.

3. LOCAL CLUSTER COEFFICIENT

Given a subgraph H of a graph G , define the **relative characteristic length** as

$$\mu(H, G) = \frac{1}{|H|(|H| - 1)} \sum_{x,y \in H, x \neq y} d_G(x, y) .$$

The difference $\nu_H(G) = \mu(H) - \mu(H, G)$ is nonnegative. It is zero if all geodesics connecting two points in H remain in H . The notion makes sense in any metric space equipped with a probability measure. The number $\nu_H(G)$ is a measure for how far H is away from being convex within the metric space G . The notion depends on a choice of a probability measure on G . On graphs, many fractal sets, spaces on which a

Lie group acts transitively on Riemannian manifolds, there is a natural measure.

Examples:

1) For a compact surface H in three dimensional space $G = R^3$ for example, $\nu_H(G)$ is zero if and only if H is a plane. We understand that the average has been formed with respect to some absolutely continuous probability measure on the surface H which gives positive measure to every open set.

2) For a curve in a compact Riemannian manifold G , the relative characteristic length is 0 if the curve is a short enough geodesic. But the relation $\mu(H, G)/\mu(H)$ will decrease eventually to zero for aperiodic geodesic paths H as the geodesic will accumulate.

3) For a region G in Euclidean space, we have $\mu(H, G) = \mu(H)$ if and only if H is **convex** in G in the sense that for any two points $x, y \in H$ there is a geodesic in G which is also in H .

Define the **local characteristic length** as

$$\mu(x) = \mu(S(x), B(x)) ,$$

where $S(x)$ is the unit sphere and $B(x)$ is the unit ball of the vertex x . We always assume that the space is large enough so that the unit ball at every point is convex within G in the above sense. The quantity $L(x)$ is defined therefore for large enough metric spaces equipped with a probability measure m which is nice enough that it induces measures on spheres $S(x)$ by limiting conditional expectation. Examples are Riemannian manifolds or graphs.

For a finite simple graph $G = (V, E)$, the **local cluster coefficient** is defined as

$$C(x) = 2|E(x)|/(|V(x)| + 1)|V(x)| ,$$

where $E(x) = V_1(x)$ is the set of edges in the unit sphere $S(x)$ of x and $V(x) = V_0(x)$ is the set of vertices in $S(x)$. The **mean cluster coefficient** $\nu(G)$ was defined by Watts-Strogatz is the average of $S(x)$ over all $x \in V$. The local cluster coefficient gives the edge occupation rate in the sphere $S(x)$ of a vertex x . Other related quantities are the **global cluster coefficient** defined as the frequency of oriented triangles $3v_2/t_2$ within all oriented connected vertex triples. The **transitivity ratio** is the ratio v_2/s_2 of non-oriented triangles within the class of non-oriented connected vertex triples in G . We do not look at the later two notions because they are close to the mean cluster density and because intuition about them is more difficult.

We can define **higher global cluster coefficients** $C_k(G)$ of a graph as the fraction v_k/w_k , where v_k is the number of k -dimensional simplices K_{k+1} in G and w_k the number of connected k -tuples of vertices. Of course, $C_1(G)$ is the global clustering coefficient and $C_1(S(x))$ is the local clustering coefficient at a point x . One could define higher local clustering coefficient $C_k(x)$ as the global clustering coefficient $C_k(S(x))$. There are also higher dimensional characteristic lengths: define $d_k(x, y)$ as the distance between two k -dimensional simplices x, y , where the distance is the smallest l such that we can connect x, y with a sequence $x_0 = x, x_1, x_2, \dots, x_l$ of overlapping k simplices x_k . Of course $d_1(x, y)$ is the usual geodesic distance and for geometric graphs of dimension d all distances d_k are essentially the same. For a triangularization of a Riemannian manifold, where every sphere is one-dimensional cyclic graph, the distance between two triangles which overlap is 1. We mentioned these higher clustering coefficients and higher characteristic lengths because we believe they could be used to get a closer relation with inductive dimension which also uses the entire spectrum of higher dimensional simplices and not only zero dimensional vertices and one dimensional edges.

4. CLUSTER COEFFICIENT AND LOCAL LENGTH

While the following observation is almost obvious, we are not aware that it has been noticed anywhere already. The result allows to write the local cluster coefficient in terms of the relative characteristic length $L(x)$ of the unit sphere $S(x)$ within the unit ball $B(x)$. This will allow us to push the notion of local clustering coefficient to other metric spaces equipped with natural measures. Just define then $C(x) = 2 - L(x)$ there.

Lemma 1 (Cluster-Length-Lemma). $L(x) = 2 - C(x)$.

Proof. By definition, the distance function takes only two different positive values in the ball $B(x)$. The first possibility is $\text{dist}(x, y) = 1$ which is the case if one of the vertices is the center. The second possibility is $d(x, y) = 2$ if both x, y are on the sphere. If d is the degree of the vertex x then $B(x)$ has $(d + 1)$ vertices and $S(x)$ has d vertices. The distance 1 appears $C(x)d(d - 1)$ times on the sphere $S(x)$. The distance 2 appears $(1 - C(x))d(d - 1)$ times in the sphere $S(x)$. The average is

$$(1 \cdot [C(x)d(d - 1)] + 2 \cdot [(1 - C(x))d(d - 1)]) / (d(d - 1)) = 2 - C(x) .$$

□

An extremal case is a star graph S_n which has cluster coefficient 0 and dimension 0 and where the distance between any two points of the unit sphere is 2. An other extreme case is the complete graph K_{n+1} which has the cluster coefficient 1 and dimension n , where the distance between any two points of the unit sphere is 1.

The wheel graph $W_n(x)$ is an example in which each point has a 1-dimensional sphere, the local cluster coefficient of the center is $2/(n-1)$ and of the other points $2/3$. The mean cluster coefficient for W_n is $m(W_n) = (n(2/3) + 2/(n-1))/(n+1)$. The cluster-length-ratio $\lambda = -\mu/\log(\nu)$ is close to 1.

5. DIMENSION

An other important local quantity is the dimension $\dim(x)$ of a vertex and the inductive dimension of a graph, the average dimension of its vertices. It was defined in [17] as

$$\dim(\emptyset) = -1, \dim(G) = 1 + \frac{1}{|V|} \sum_{v \in V} \dim(S(v)),$$

where $S(v) = \{w \in V \mid (w, v) \in E\}$, $\{e = (a, b) \in E \mid (v, a) \in E, (v, b) \in E\}$ denotes the unit sphere of a vertex $v \in V$.

Already on a local level, there can be relations. If the dimension of a point is zero, then clearly the local length $L(x)$ is 2 and the cluster coefficient is zero. And also λ is zero.

We have shown in [15] that the expectation $E_p[\dim]$ on $G(n, p)$ satisfies the recursion

$$d_{n+1}(p) = 1 + \sum_{k=0}^n \binom{n}{k} p^k (1-p)^{n-k} d_k(p),$$

where $d_0 = -1$. Each d_n is a polynomial in p of degree $\binom{n}{2}$.

6. METRIC SPACES

The characteristic length μ of a metric space (X, d) with measure m is the expectation of length $d(x, y)$ on $X^2 \setminus D$, where D is the diagonal $\{(x, x)\}$ in $X \times X$. The local length $L(x)$ is the characteristic length of the unit sphere $S(x)$ within the unit ball $B(x)$. Motivated by the above, we call $C(x) = 2 - L(x)$ the local cluster coefficient of the pint in (X, d) and its expectation m , the mean cluster coefficient. Lets look at the quantity $-\mu r/(\log(\nu))$, where r is the radius of the small ball and where the volume of the manifold is 1. These are integral

geometric questions. We in general assume that r is scaled in such a way that the radius of injectivity is larger than 1 and the unit ball is contractible and convex. For a Riemannian manifold, the number C is a dimension-dependent constant, the average distance between two points in the unit sphere. For manifolds, studying λ is therefore equivalent than to study the characteristic length. The question is however interesting for fractals. One can ask for example, what the cluster-length ration λ is for the Sierpinsky carpet.

Example 1. For a flat torus $X = R^2/(rZ)^2$ with flat Riemannian geodesic distance and area measure, we have

$$\nu = \frac{2}{\pi} \int_0^{\pi/2} \left(2 - \frac{2 \tan(\phi)}{\sqrt{1 + \tan^2(\phi)}} \right) d\phi = \left(2 - \frac{4}{\pi} \right) = 0.72676\dots$$

and $\mu = r \frac{1}{6} (\sqrt{2} + \sinh^{-1}(1)) = r(\sqrt{2} + \operatorname{arcsinh}(1))/6 = r \cdot 0.382598\dots$
Therefore,

$$\lambda = \mu/(r \log(1/\nu)) = 0.382598/\log(0.72676) = 1.17837\dots$$

Example 2. Also for a three-dimensional flat torus, the clustering coefficient $C(x)$ is constant and given as the average distance of two points on a sphere, which is $4/3$. The characteristic length on the other hand is 0.480296, a numerical integral which we were not able to evaluate analytically. The quotient is $\lambda = 1.18454$.

Example 3. For a two-dimensional sphere of surface area 1, the characteristic length μ is $1/2\sqrt{\pi}$ times the characteristic length $L(S(x))$ of the unit sphere which is intrinsic and uses the geodesic length within the surface and not from an embedding. We measure it in R^3 to be 1.57032... so that $\mu = 0.442979\dots$ can be computed by using that the geodesic distance between two points given in spherical coordinates as $(\phi_1, \theta_1), (\phi_2, \theta_2)$ is given by the Haverside formula $H(\phi_1 - \phi_2) + \sin(\phi_1) \sin(\phi_2) H(\theta_1 - \theta_2)$ where $H(x) = \sin^2(x/2)$ is the Haverside function. Random points on the sphere can be computed with the uniform distribution in θ and the arccos distribution in ϕ . Assuming the same $C(x)$ value as in the plane (which uses that near a specific point we can replace the sphere with its tangent space) and get $\lambda = 1.36\dots$

7. CLUSTER-LENGTH RATIO

Dimension plays a role for characteristic length. We measure experimentally that the global length-cluster ratio quantity

$$\lambda = -\mu/\log(\nu)$$

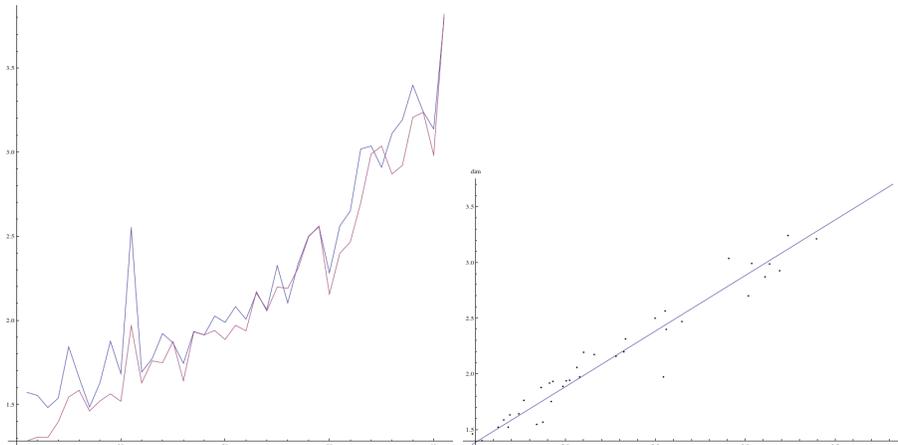


FIGURE 3. The plot of the $\lambda(p)$ and $\nu(p)$ in the Erdős-Rényi case for probabilities p between 0.3 and 0.7, where $n = 15$. The graphs indicate a clear correlation between cluster-length ratio and dimension. The second plot shows the two quantities together with a regression line.

is correlated to dimension for Erdős-Rényi graphs:

Here is some intuition, why the limit should exist: the cluster coefficient is related to the existence of triangles in the graph. For orbital networks [20, 22, 23] defined by polynomial maps the number of triangles is bounded by a constant C . This implies that $\nu(G) \leq C/n$ and $\log(\nu(G)) \leq \log(C) - \log(n)$ for orbital networks, we should get that $\log(\nu(G))/\log(n)$ has a finite interval as accumulation points. To show that $\mu(G)$ grows like $\log(n)$, we don't want too many relations $T^w = T^v$ with different words w, v . We call this a collision. If d is the number of generators, then, if there were no collisions, the relation $d^\mu = n$ holds. With C_2 double collisions, assuming no triple collisions, we have the relation $d^\mu - C = n$ and so $\mu = \log(n + C)/\log(d)$. Together with $-\log(\nu) = \log(n) - \log(C)$ we have $\gamma = \log(n + C)/(\log(d)(\log(n) - \log(C_2(n))))$. So, if we can show C_2 to be of the order $\log(n)$ and the number of triangles to be of the order $o(n)$, and triple collisions are rare, then we should be able to prove that the limit exists.

8. CLASSES OF NETWORKS

The space $E(n, p)$ of all graphs on a vertex set with n nodes, where every node is turned on with probability p is a probability space. The

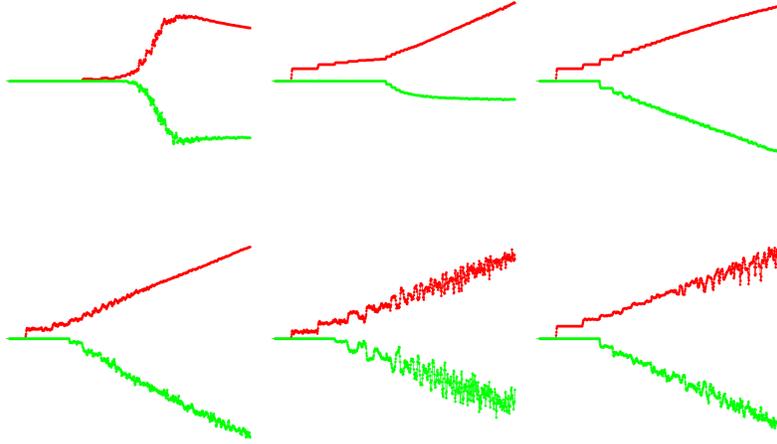


FIGURE 4. The characteristic length μ (red), the log of the cluster coefficient $\log(\nu)$ (green) and the average vertex degree d (blue) shown together as a function of the number of vertices in the network. n . We plot with a logarithmic scale in n so that logarithmic growth is seen linearly. We first use the Erdős-Rényi probability space (where each edge is turned on with probability $p = 0.1$), then for Watts-Strogatz for $k = 4, p = 0.1$ and then for Barabasi-Albert [2] networks. In the second row, we see first the case of two quadratic maps [23] then two random permutations and finally a case with correlated generators, where the clustering is extremely small.

limit

$$\lambda = -\mu / \log(\nu)$$

for $n \rightarrow \infty$ exists almost surely. We see that the value is close to r/dim , where r is the radius of the graph and dim is the dimension of the network.

The quotient (1) is interesting because μ is a global property and ν is the average of a local property. Intuitively, such a relation is to be expected because a larger $C(x)$ allows to tunnel faster through a ball $B(x)$ and allows for shorter paths. If the limit exists, then $\mu = \lambda \log(\nu)$. Knowing λ is important because the characteristic length is more costly to compute while the clustering coefficient C is easier to determine as

a simple average of local quantities. To allow an analogy from differential geometry, we could compare $C(x)$ with curvature, because a metric space with larger curvature has a smaller average distance between two points on the unit sphere.

We could look at graphs with a given dimension and volume and minimize the average path length between two points among all graphs. It is a long shot but one can ask whether there is the relation between graphs minimizing μ and graphs minimizing Euler characteristic χ . We can only explore this so far for very small graphs. The reason for asking this is that Euler characteristic can also be seen as an average of scalar curvature and therefore a quantized Hilbert action [21].

9. RELATED QUESTIONS

Variational problems on graphs usually need some constraints because the functionals are often trivial without restrictions. We can restrict the number of vertices or edges and look at the maximum or minimum on that space. More generally, we can use a Lagrange type problem and look at all the graphs for which one functional is constant and extremize the other on that class. This leads to more questions and most of them seem not have been studied. Instead of restricting to a “level surface” we can also look at the functionals on an equivalence class of graphs. One interesting example is to look at homotopy as an equivalence relation. A homotopy step $G \rightarrow G'$ is given by choosing a contractible subgraph H of G and connect each vertex of H with a new vertex v . An other homotopy step is the reverse operation: remove a vertex for which the unit sphere is contractible. The notion of contractible if a sequence of homotopy steps transforms it to a one point graph.

Lets look at the example of minimizing the dimension $\iota(G)$ in a homotopy class. The homotopy class of a circle contains graphs of arbitrary large dimension; it contains for example discretization of a solid torus (dimension 3) or an annulus (dimension 2).

We can also find one-dimensional graphs homotopic to the circle which are not C_n . We can for example attach one dimensional hairs to the circle without changing dimension, nor homotopy. We have now a new functional $\iota'(G)$ which is the minimal dimension among all graphs H homotopic to G . For a contractible graph, the minimal dimension is 0. On the class of graphs homotopic to the circle the minimal dimension

is 1 and for all graphs homotopic to an icosahedron it is 2. A similar modified dimension ι' can be defined in the continuum: define the homotopy dimension of a space M as the minimum of the Hausdorff dimensions of all compact metric spaces (X, d) homotopic to M .

The question is whether the minimum is always attained by a geometric graph or smooth manifolds.

REFERENCES

- [1] M. Kovse B. Bresar and A. Tepeh. Geodetic sets in graphs. In M. Dehmer, editor, *Structural Analysis of Networks*. Birkhäuser, 2011.
- [2] A-L. Barabási and R. Albert. Emergence of scaling in random networks. *Science*, 286(5439):509–512, 1999.
- [3] N. Biggs. *Algebraic Graph Theory*. Cambridge University Press, 1974.
- [4] B. Bollobas. *Extremal Graph Theory*. Dover Courier Publications, 1978.
- [5] R.A. Brualdi and J.L. Goldwasser. Permanent of the laplacian matrix of trees and bipartite graphs. *Discrete Mathematics*, 48:1–2, 1984.
- [6] S. Butler. Relating the arboricity with the chromatic number of a graph. <http://www.math.iastate.edu/butler/PDF/arboricity.pdf>, accessed, July 20, 2014.
- [7] B. Chen, M. Matsumoto, J. Wang, Z. Zhang, and J. Zhang. A short proof of nash-williams' theorem for the arboricity of a graph. *Graphs and Combinatorics*, 10:27–28, 1994.
- [8] F.R.G. Chung. The average distance and independence number. *J. Graph Theory*, 12:229–235, 1988.
- [9] D.M. Cvetkovic. Chromatic number and the spectrum of a graph. *Publications de l'insitute Mathematique*, 14(28):25–38, 1972.
- [10] J. Doyle and J. Graver. Mean distance in a graph. *Discrete Math*, 17:147–154, 1977.
- [11] R.C. Entringer, E.E. Jackson, and D.A. Snyder. Distance in graphs. *Czechoslovak Mathematical Journal*, 26:283–296, 1976.
- [12] S. Fajtlowicz. Toward fully automated fragments of graph theory. *Graph Theory Notes N. Y.*, 42:18–25, 2002.
- [13] W. Goddard and O.R. Oellermann. Distance in graphs. In M. Dehmer, editor, *Structural Analysis of Networks*. Birkhäuser, 2011.
- [14] O. Knill. Natural orbital networks. <http://arxiv.org/abs/1311.6554>.
- [15] O. Knill. The dimension and Euler characteristic of random graphs. <http://arxiv.org/abs/1112.5749>, 2011.
- [16] O. Knill. A graph theoretical Gauss-Bonnet-Chern theorem. <http://arxiv.org/abs/1111.5395>, 2011.
- [17] O. Knill. A discrete Gauss-Bonnet type theorem. *Elemente der Mathematik*, 67:1–17, 2012.
- [18] O. Knill. The McKean-Singer Formula in Graph Theory. <http://arxiv.org/abs/1301.1408>, 2012.
- [19] O. Knill. Counting rooted forests in a network. <http://arxiv.org/abs/1307.3810>, 2013.

- [20] O. Knill. Dynamically generated networks. <http://arxiv.org/abs/1311.4261>, 2013.
- [21] O. Knill. The Euler characteristic of an even-dimensional graph. <http://arxiv.org/abs/1307.3809>, 2013.
- [22] O. Knill. Natural orbital networks. <http://arxiv.org/abs/1311.6554>, 2013.
- [23] O. Knill. On quadratic orbital networks. <http://arxiv.org/abs/1312.0298>, 2013.
- [24] O. Knill. A Cauchy-Binet theorem for Pseudo determinants. *Linear Algebra and its Applications*, 459:522–547, 2014.
- [25] O. Knill. Curvature from graph colorings. <http://arxiv.org/abs/1410.1217>, 2014.
- [26] O. Knill. Dynamically generated networks. 2014. <http://arxiv.org/abs/1311.4261>.
- [27] O. Knill. A notion of graph homeomorphism. <http://arxiv.org/abs/1401.2819>, 2014.
- [28] M. Kouider and P.Winkler. Mean distance and minimum degree. *J. Graph Theory*, 25(1):95–99, 1997.
- [29] P. Van Miegham. Graph spectra for complex networks. 2011.
- [30] C. Nash-Williams. Edge-disjoint spanning trees of finite graphs. *J. London Math. Soc.*, 36:455–450, 1961.
- [31] M.E.J. Newman, S.H. Strogatz, and D.J. Watts. Random graphs with arbitrary degree distributions and their applications. *Physical Review E*, 64, 2001.
- [32] P.Chebotarev and E. Shamis. Matrix forest theorems. arXiv:0602575, 2006.
- [33] E.V. Shamis P.Yu, Chebotarev. A matrix forest theorem and the measurement of relations in small social groups. *Avtomat. i Telemekh.*, 9:125–137, 1997.
- [34] N.S. Schmuck. The wiener index of a graph. Diploma thesis, TU Graz, 2010.
- [35] S. Sivasubramanian. Average distance in graphs and eigenvalues. *Discrete Mathematics*, 309:3458–3462, 2009.
- [36] S. Wasserman and K. Faust. *Social Network analysis: Methods and applications*. Cambridge University Press, 1994.
- [37] D. J. Watts. *Small Worlds*. Princeton University Press, 1999.
- [38] D. J. Watts and S. H. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 393:440–442, 1998.

DEPARTMENT OF MATHEMATICS, HARVARD UNIVERSITY, CAMBRIDGE, MA, 02138