On Compression of Uniform Random Intersection Graphs

Zbigniew Golebiewski Marcin Kardas Jakub Lemiesz Krzysztof Majcher

Wroclaw University of Science and Technology, Poland

AofA 2017

Structural Information

What:

- $\,\triangleright\,$ to develop compression algorithms for various structured data
- $\,\vartriangleright\,$ we focus on compression of unlabelled graphs

Why:

- many interesting combinatorial objects have structure (web graph, protein-protein interactions, collaboration networks ...)
- \triangleright they can be abstracted by (unlabelled) graphs

How:

 \triangleright using the structural entropy metric in order to measure the amount of information embodied in a graph structure

Analysis

Summary 0000

Graphs and Structural Entropy



- ZG Zbigniew Golebiewski MK1 - Marcin Kardas MK2 - Marek Klonowski KM - Krzysztof Majcher JL - Jakub Lemiesz JC - Jacek Cichon
- WS Wojtek Szpankowski
- AM Abram Magner
- OB Olivier Bodini
- DG Daniele Gardy
- BG Bernhard Gittenberger

Analysis 00000000000 Summary 0000

Graphs and Structural Entropy



ZG - Zbigniew Golebiewski MK1 - Marcin Kardas MK2 - Marek Klonowski KM - Krzysztof Majcher JL - Jakub Lemiesz JC - Jacek Cichon WS - Wojtek Szpankowski AM - Abram Magner OB - Olivier Bodini DG - Daniele Gardy

BG - Bernhard Gittenberger

{ZG, MK1, MK2, KM, JL, JC, WS, AM, OB, DG, BG}

Q: How many bits are required to describe the structure of a graph?

Summary 0000

Graphs and Structural Entropy

- $\vartriangleright\,$ let ${\mathcal G}$ be a memory-less source producing graph according to some random graph model
- $\,\triangleright\,$ the classic entropy of ${\cal G}$ is defined as

$$H(\mathcal{G}) = -\sum_{G \in \mathcal{G}} \mathbb{P}_{\mathcal{G}}(G) \lg \mathbb{P}_{\mathcal{G}}(G)$$

- $\,\vartriangleright\,$ let ${\mathcal S}$ be a random structure model for the random graph model ${\mathcal G}$
- Dash the probability of generating a given structure $S\in\mathcal{S}$ is

$$\mathbb{P}_{\mathcal{S}}(S) = \sum_{G \cong S, G \in \mathcal{G}} \mathbb{P}_{\mathcal{G}}(G)$$

 \triangleright the structural entropy for model ${\mathcal G}$ is defined as

$$H_{\mathcal{S}}(\mathcal{G}) = -\sum_{S \in \mathcal{S}} \mathbb{P}_{\mathcal{S}}(S) \lg \mathbb{P}_{\mathcal{S}}(S)$$

Summary 0000

Graphs and Structural Entropy

- \vartriangleright let ${\mathcal G}$ be a memory-less source producing graph according to some random graph model
- $\,\triangleright\,$ the classic entropy of ${\cal G}$ is defined as

$$H(\mathcal{G}) = -\sum_{G \in \mathcal{G}} \mathbb{P}_{\mathcal{G}}(G) \lg \mathbb{P}_{\mathcal{G}}(G)$$

- $\,\vartriangleright\,$ let $\,\mathcal{S}\,$ be a random *structure model* for the random graph model $\,\mathcal{G}\,$
- Dash the probability of generating a given structure $S \in \mathcal{S}$ is

$$\mathbb{P}_{\mathcal{S}}(\mathcal{S}) = \sum_{\mathcal{G}\cong\mathcal{S},\mathcal{G}\in\mathcal{G}}\mathbb{P}_{\mathcal{G}}(\mathcal{G})$$

 $\,\vartriangleright\,$ the structural entropy for model ${\mathcal G}$ is defined as

$$H_{\mathcal{S}}(\mathcal{G}) = -\sum_{S \in \mathcal{S}} \mathbb{P}_{\mathcal{S}}(S) \lg \mathbb{P}_{\mathcal{S}}(S)$$

Introduction	
000000	

Analysis

Summary 0000

Graphs and Structural Entropy



Uniform Random Intersection Graph (URIG) Model

- $Dash n \in \mathbb{N}_+$ a number of nodes
- $\, \triangleright \, m \in \mathbb{N}_+$ a number of colors
- $dash \ k \in \{1,\ldots,m\}$ a number of colors sampled (without replacement) independently by each node

By $U_{n,m,k}$ we understand a random memory-less source producing undirected graphs with *n* vertices according to the following process:

- \triangleright each vertex chooses uniformly at random a set of k colors out of m possible
- \triangleright two vertices *u* and *v* are connected if and only if both sampled at least one common color.

Uniform Random Intersection Graph (URIG) Model

- $\triangleright \ n \in \mathbb{N}_+$ a number of nodes
- $\, \triangleright \, m \in \mathbb{N}_+$ a number of colors
- $\triangleright \ k \in \{1, \dots, m\}$ a number of colors sampled (without replacement) independently by each node

By $U_{n,m,k}$ we understand a random memory-less source producing undirected graphs with *n* vertices according to the following process:

- \triangleright each vertex chooses uniformly at random a set of k colors out of m possible
- \triangleright two vertices *u* and *v* are connected if and only if both sampled at least one common color.

ntroduction	Analysis ●000000000	Summary 0000

Underlying Intersection Graph $G_{m,k}$

- ▷ has $m' = \binom{m}{k}$ vertices that correspond to all distinct *k*-element subsets of a set of colors $\{1, \ldots, m\}$
- \vartriangleright to vertices v and u are connected iff they share at least one color, i.e. $v\cap u\neq \emptyset$

emma

The number of automorphisms of the underlying intersection graph $G_{m,k}$ is

$$|Aut(G_{m,k})| = \begin{cases} \binom{m}{k}! & \text{when } m < 2k, \\ \binom{m}{k}!! & \text{when } m = 2k, \\ m! & \text{when } m > 2k. \end{cases}$$

Introduction	Analysis • • • • • • • • • • • • • • • • • • •	Summary 0000

Underlying Intersection Graph $G_{m,k}$

- ▷ has $m' = \binom{m}{k}$ vertices that correspond to all distinct *k*-element subsets of a set of colors $\{1, \ldots, m\}$
- ▷ to vertices v and u are connected iff they share at least one color, i.e. $v \cap u \neq \emptyset$

Lemma

The number of automorphisms of the underlying intersection graph $G_{m,k}$ is

$$|Aut(G_{m,k})| = \begin{cases} \binom{m}{k}! & \text{when } m < 2k, \\ \binom{m}{k}!! & \text{when } m = 2k, \\ m! & \text{when } m > 2k. \end{cases}$$

Introduction	
000000	

Underlying Intersection Graph $G_{m,k}$

proof idea:

$$|Aut(G_{m,k})| = \begin{cases} \binom{m}{k}! & \text{when } m < 2k, & \text{complete graph} \\ \binom{m}{k}!! & \text{when } m = 2k, & \begin{array}{c} \text{the complement graph} \\ \text{contains only separated} \\ \text{cliques of size } 2 \\ \\ m! & \text{when } m > 2k, & \begin{array}{c} \text{least symmetric case,} \\ \text{known result for a} \\ \text{Kneser graph, use of} \\ \text{Erdős-Ko-Rado theorem} \\ \end{cases}$$



10/22

Introduction
000000

Analysis 00000000000000

Summary 0000

Random Composition Source

▷ let $\mathcal{K}_{n,G_{m,k}}$ be a set of all *m*'-element compositions of n $(m' = \binom{m}{k})$, where the elements are indexed by vertices of $G_{m,k}$

$$\mathcal{K}_{n,G_{m,k}} = \left\{ \mathcal{K} \in \mathbb{N}^{V} : \sum_{v \in V} \mathcal{K}(v) = n \right\}$$

 \triangleright observe that

$$\mathbb{P}\left(K=(k_1,\ldots,k_{m'})\right)=\binom{n}{k_1,\ldots,k_{m'}}\frac{1}{m'^n}$$

 \triangleright therefore

$$H_{\mathcal{S}}\left(\mathcal{U}_{n,m,k}\right) = H_{\mathcal{S}}\left(\mathcal{K}_{n,G_{m,k}}\right)$$

Introduction 000000	Analysis 0000●000000				Summary	
Structural	Entropy	of	Uniform	Random	Intersection	Graph
Source						

Theorem (G., Kardas, Lemiesz, Majcher)

The structural entropy of a source generating uniform intersection graphs $U_{n,m,k}$, for $m \ge 2k$, is

$$H_{\mathcal{S}}\left(\mathcal{U}_{n,m,k}\right) = H\left(\mathcal{K}_{n,G_{m,k}}\right) - \lg |Aut\left(G_{m,k}\right)| + \mathbb{E}\left(\lg |stab(\mathcal{K})|\right) + o(1),$$

assuming that $n, \binom{m}{k} \to \infty$ in such a way that $\frac{n}{\binom{m}{k}} = \Theta(n^{\tau})$ for any $0 < \tau \leq 1$.

Introduction	Analysis	Summary
000000	00000●00000	0000
Proof idea		

 \vartriangleright for $m \geq 2k$, let $K \in \mathcal{K}_{n,G_{m,k}}$ be a positive composition, then

$$[K]_{\approx} = orb(K) \stackrel{df}{=} \{K \circ \pi : \pi \in Aut(G_{m,k})\}$$

and all compositions of $[K]_{\approx}$ are equiprobable \triangleright

$$H_{\mathcal{S}}\left(\mathcal{U}_{n,m,k}
ight) = -\mathbb{E}\left(\lg\mathbb{P}\left([\mathcal{K}]_{\approx}
ight)
ight)$$

 \triangleright by orbit-stabilizer theorem:

$$|[K]_{\approx}| = |orb(K)| = \frac{|Aut(G_{m,k})|}{|stab(K)|},$$

where $stab(K) = \{\pi \in Aut(G_{m,k}) : \pi \circ K = K\}$

lr	۱t	ro	d	u	C	ti	0	n	
0	0	00	C	0	С				

Analysis 000000000000

Structural Entropy of Uniform Random Intersection Graph Source

Theorem (G., Kardas, Lemiesz, Majcher)

The structural entropy of a source generating uniform intersection graphs $U_{n,m,k}$, for m > 2k, is

$$\begin{aligned} \mathcal{H}_{\mathcal{S}}\left(\mathcal{U}_{n,m,k}\right) = \binom{m}{k} \lg \sqrt{2\pi\alpha} - \lg \sqrt{2\pi n} + \frac{\binom{m}{k} - 1}{2\ln(2)} - \lg(m!) \\ + \frac{\binom{m}{k}}{\alpha\ln(2)} \sum_{l=0}^{\lfloor \frac{1-2\tau}{\tau} \rfloor} (-1)^l l! \mathsf{G}_{l+2} \alpha^{-l} + o(1), \end{aligned}$$

assuming that $n, \binom{m}{k} \to \infty$ in such a way that $\alpha = \frac{n}{\binom{m}{k}} = \Theta(n^{\tau})$ for any $0 < \tau \leq 1$ and G_{I} is I'th Gregory coefficient.

Introduction	Analysis	Summary
000000	00000000000	0000
Remarks		

▷ the structural entropy of the source $U_{n,m,k}$ is relatively low compared, for example, to the source $\mathcal{PA}_{n,m}$ producing preferential attachment graphs or to the source $\mathcal{G}_{n,p}$ producing Erdős-Rényi graphs

 \triangleright for the same expected number of edges in a graph we get:

$$H_{\mathcal{S}}(\mathcal{U}_{n,m,k}) = \Theta(\lg n)$$
$$H_{\mathcal{S}}(\mathcal{PA}_{n,m}) = \Theta(n \lg n)$$
$$H_{\mathcal{S}}(\mathcal{G}_{n,p}) = \Theta(n^{2})$$

 \triangleright it can be justified by the symmetries present in the graphs generated by $\mathcal{U}_{n,m,k}$

Introduction	Analysis	Summary
000000	0000000000000	0000
Proof idea		

Lemma

If
$$n$$
, $\binom{m}{k} \to \infty$ in such a way that $\alpha = \frac{n}{\binom{m}{k}} = \Theta(n^{\tau})$ for any $0 < \tau \leq 1$, then

$$H\left(\mathcal{K}_{n,G_{m,k}}\right) = \binom{m}{k} \lg \sqrt{2\pi\alpha} - \lg \sqrt{2\pi n} + \frac{\binom{m}{k} - 1}{2 \ln 2} \\ + \frac{\binom{m}{k}}{\alpha \ln 2} \sum_{l=0}^{\lfloor \frac{1-2\tau}{\tau} \rfloor} (-1)^{l} l! \mathsf{G}_{l+2} \alpha^{-l} + o(1)$$

where G_I is I'th Gregory coefficient (known also as I'th logarithmic number) defined by a Maclaurin series expansion of

$$\frac{y}{\log(1+y)} = 1 + \sum_{l=1}^{\infty} \mathsf{G}_l y^l.$$

Introduction		Analysis	Summary		
000000		0000000000	0000		
-	C + 1				

Proof idea

 \triangleright the entropy of the random composition source $\mathcal{K}_{n,G_{m,k}}$ is

$$H\left(\mathcal{K}_{n,G_{m,k}}\right) = -\sum_{K\in\mathcal{K}_{n,G_{m,k}}}\mathbb{P}(K)\lg\left(\mathbb{P}(K)\right)$$

▷ let us recall that $\mathbb{P}(K = (k_1, ..., k_{m'})) = \binom{n}{k_1, ..., k_{m'}} \frac{1}{m'^n}$ ▷ after few steps we obtain a formula with Bernoulli sum

$$H\left(\mathcal{K}_{n,G_{m,k}}\right) = n \lg(m') - \lg(n!) + m' \sum_{t=0}^{n} \lg\left(t!\right) \binom{n}{t} \left(\frac{1}{m'}\right)^{t} \left(1 - \frac{1}{m'}\right)^{n-t}$$

▷ choosing the approach due to Knessl, i.e. using:

$$\ln A = \lim_{\epsilon \to 0} \int_{\epsilon}^{\infty} \frac{e^{-x} - e^{-Ax}}{x} dx, \qquad A > 0,$$

seems to give us the biggest freedom of choosing $\alpha = n/\binom{m}{k}$

Introduction	Analysis	Summary
000000	0000000000	0000
Proof idea		

Lemma

Let
$$m > 2k$$
, $k \ge 2$ and $n/\binom{m}{k} \to \infty$, then

$$\mathbb{E}\left(\left|\mathsf{g}\left|\mathsf{stab}(\mathsf{K})\right|\right)=o(1)\quad \text{as }\mathsf{n}
ightarrow\infty.$$

proof idea:

▷ for
$$n = \omega \left({\binom{m}{k}}^5 \right)$$
 we have different bins loads whp

▷ otherwise

- let π ∈ Aut (G_{m,k}) be a non-trivial automorphism of G_{m,k} that is also a stabilizer of some random composition K
- we can show that π has to move at least $2\binom{m-2}{k-1}$ vertices of ${\cal G}_{m,k}$
- let $\pi = C_1 \circ \ldots \circ C_\ell$ then for all *i*: $K = K \circ C_i$ and therefore all elements moved by the cycle C_i has to be equal
- using poissonization technique we can bound the probability of such event

Optimal Compression Algorithm - work in progress

A lossless compression algorithm of a structure of URIG follows directions given by the analysis of the structural entropy.

Compress $(U \in U_{n,m,k})$:

- \triangleright contract vertices of U with the same colors subsets
- $\,\vartriangleright\,$ associate colors with the vertices of the underlying intersection graph

to this point we have reconstructed the underlying intersection graph $G_{m,k}$ and the composition K that corresponds to the graph U

 \triangleright use arithmetic encoding to compress the composition K

 \triangleright output: (m, k, compressed(K))

Introduction 000000	Analysis 0000000000	Summary 0000

Future work

- \triangleright quantitative result for the case m = 2k: stabilizers counts!
- \triangleright case when *n* is smaller or comparable to $\binom{m}{k}$: contracted graph is a subgraph of the underlying intersection graph $G_{m,k}$
- ▷ structural entropy of the binomial intersection graphs source
- ▷ other symmetric graphs sources

Introduction	Analysis	
000000	000000000	

Lessons learned

- \vartriangleright asymmetric graphs \rightarrow the difficulty can come from a source probability distribution
- ▷ with growth of the symmetry in the graphs generated by a source → the difficulty can come from both: a graph symmetries and a source probability distribution
- $\,\vartriangleright\,$ the most symmetric case is a clique $\,\rightarrow\,$ the analysis is trivial

Introduction	Analysis	Summary
000000	000000000	0000

Thank you!