

Conditional graph entropy as an alternating minimization problem

Viktor Harangi*, Xueyan Niu[†], Bo Bai[†]

*Alfréd Rényi Institute of Mathematics, Budapest, Hungary
harangi@renyi.hu

[†]Theory Lab, Central Research Institute, 2012 Labs, Huawei Technologies Co. Ltd., Hong Kong SAR, China
{niuxueyan3,baibo8}@huawei.com

Abstract

Conditional graph entropy is known to be the minimal rate for a natural functional compression problem with side information at the receiver. In this paper we show that it can be formulated as an alternating minimization problem, which gives rise to a simple iterative algorithm for numerically computing (conditional) graph entropy. This also leads to a new formula which shows that conditional graph entropy is part of a more general framework: the solution of an optimization problem over a convex corner. In the special case of graph entropy (i.e., unconditioned version) this was known due to Csiszár, Körner, Lovász, Marton, and Simonyi. In that case the role of the convex corner was played by the so-called vertex packing polytope. In the conditional version it is a more intricate convex body but the function to minimize is the same. Furthermore, we describe a dual problem that leads to an optimality check and an error bound for the iterative algorithm.

I. INTRODUCTION

We consider the problem of computing *conditional graph entropy*. Orłitsky and Roche [14] used this entropy notion to characterize the optimal rate of lossless functional compression with side information at the decoder. Despite playing an inherent role in data compression, little can be found in the literature about conditional graph entropy.

A. Background and related work

Conditional graph entropy: Suppose that the random variable X takes values in a finite alphabet \mathcal{X} , where not every pair of letters can be distinguished. Let G be a graph with vertex set $V(G) = \mathcal{X}$ describing which pairs are distinguishable: $x, x' \in \mathcal{X}$ can be distinguished if and only if xx' is an edge of G . Furthermore, we say that the sequences x_1, \dots, x_ℓ and x'_1, \dots, x'_ℓ are distinguishable if x_i and x'_i are distinguishable for at least one index i . We wish to encode an i.i.d. sequence X_1, \dots, X_ℓ with high probability in a way that distinguishable sequences are mapped to different codewords. An *independent set* of G contains no edges, and hence any two letters in the set are indistinguishable. Therefore one possible strategy is to replace each X_i with an independent set $J_i \subseteq \mathcal{X}$ containing X_i , and encode the sequence J_1, \dots, J_ℓ instead. If we do this randomly in a way that (X_i, J_i) are i.i.d. samples of some (X, J) where J is a random independent set containing X , then we can encode the J_i sequence with rate $H(J)$ (asymptotically as $\ell \rightarrow \infty$). Note that the number of times any given typical X_i sequence is covered has exponential rate $H(J|X)$. Based on this, one can design an encoding with rate $H(J) - H(J|X) = I(X; J)$. Then, for a given X , one needs to choose (X, J) in a way that the mutual information $I(X; J)$ is as small as possible. Körner showed that this is the best achievable code rate and introduced the corresponding notion of *graph entropy* [12]:

$$H_G(X) = \min_J I(X; J), \quad \text{where } J \text{ is a random independent set of } G \text{ such that } X \in J. \quad (1)$$

The analogous problem with side information Y_i at the receiver leads to the notion of *conditional graph entropy* $H_G(X|Y)$. Let (X, Y) be discrete random variables of some given joint distribution and let (X_i, Y_i) be i.i.d. samples. We assume that the decoder knows the sequence Y_1, Y_2, \dots . If we want to use the same approach (i.e., choosing a random J), then J and Y should be independent conditioned on X (because the sender does not know Y_i when choosing J_i). This can be made rigorous, leading to the following formula:

$$H_G(X|Y) = \min_J I(X; J|Y) = \min_J \left(H(J|Y) - \underbrace{H(J|X, Y)}_{H(J|X)} \right), \quad (2)$$

where J is a random independent set of G such that $X \in J$, and J and Y are conditionally independent conditioned on X (which is equivalent to saying that $Y - X - J$ is a Markov chain).

A lot of work has been done regarding graph entropy since Körner [12] introduced the notion in 1973; see the surveys [16], [17]. In particular, Csiszár, Körner, Lovász, Marton, and Simonyi [5] found a new way to express graph entropy based on a

During this project the first author received partial support from NRD (grant KKP 138270) and from the Hungarian Academy of Sciences (János Bolyai Scholarship).

beautiful connection to the so-called *vertex packing polytope* $\text{VP}(G)$, leading to, among other things, an elegant information theoretic characterization of *perfect graphs*. This connection motivated the study of a more general framework, namely, *entropy functions* corresponding to *convex corners*. (See [19, Proposition 5.4] for a recent characterization of such functions.) Besides $\text{VP}(G)$, another notable convex corner associated to graphs is the *theta body* $\text{TH}(G)$ defined by Grötschel, Lovász, and Schrijver [10]. It is closely related to the *Lovász number* (or ϑ function), originally introduced in [13] for bounding the Shannon capacity of a graph.

Much less is known about conditional graph entropy, however. Let us first describe its connection to functional compression.

Compression with side information: Suppose now that the receiver wishes to recover the values $f(X_i, Y_i)$ of some given function $f: \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{Z}$ (with high probability, over long blocks) as depicted in Figure 1. Orłitsky and Roche [14] showed that the minimal rate of information that needs to be transmitted is precisely the conditional graph entropy of the so-called *characteristic graph*, which is defined on the vertex set \mathcal{X} as follows: vertices $x_1, x_2 \in \mathcal{X}$ are connected with an edge if and only if

$$\exists y \in \mathcal{Y} \text{ s.t. } (f(x_1, y) \neq f(x_2, y) \ \& \ \mathbb{P}(X = x_1, Y = y) > 0 \ \& \ \mathbb{P}(X = x_2, Y = y) > 0).$$

(This definition goes back to Witsenhausen [20].) We mention that in the special case $f(x, y) = x$, which was already studied in Shannon's classical work [15], the optimal rate is given by the conditional entropy $H(X|Y) = H(X, Y) - H(Y)$.

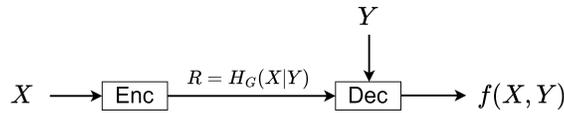


Fig. 1. The problem of functional compression with side information at the decoder

Also note that these problems are naturally connected to graph coloring. Doshi et al. [7] extended the notion of *chromatic entropy* [1] and defined *conditional chromatic entropy*. They showed that first coloring a sufficiently large power graph and then encoding the colors achieves conditional graph entropy.

Alternating optimization: As we have mentioned, it turns out that conditional graph entropy can be obtained by alternating optimization. This family includes a large number of problems from a variety of fields. A usual feature is that although the optimum has no closed-form expression, there is an efficient way to optimize in each variable. Thus, optimizing in the different variables in turns can lead to a good numerical approximation of the optimum. A well-known example is the expectation—maximization (EM) algorithm. Another prominent example is the Blahut–Arimoto (BA) algorithm [2], [3], which deals with the capacity of discrete memoryless channels. To find the capacity achieving input, the algorithm turns the objective into a double supremum and alternately optimizes over the distribution parameters from random initialization. Also, [8] provides a numerical method for the Gel'fand-Pinsker problem [9], where noncausal state information is known at the encoder. Generalization of the BA algorithm for finite-state channels, proposed in [18], also takes advantage of the iterative nature of alternating optimization. To the best of our knowledge, no similar procedure was proposed for the Orłitsky–Roche problem [14] beforehand.

It is important to point out that Csiszár and Tusnády [6] initiated the systematic study of such problems in the 1980s already. The cornerstone of their theory is a collection of inequalities called 3-point, 4-point, and 5-point properties. They will play a key role in our problem as well.

B. Notations

Random variables are denoted by uppercase letters (X, Y, J), while their realizations are denoted by lowercase letters (x, y, j). The (discrete) alphabet of a random variable is denoted by the corresponding script letter ($\mathcal{X}, \mathcal{Y}, \mathcal{J}$). For brevity, we write \sum_x for $\sum_{x \in \mathcal{X}}$, and \sum_y for $\sum_{y \in \mathcal{Y}}$. We use $\mathbb{P}(\cdot)$ to denote the probability of an event, and the following shorthand notations will be used as well:

$$\begin{aligned} p_{x,y} &:= \mathbb{P}(X = x, Y = y); \\ p_x &:= \mathbb{P}(X = x) = \sum_y p_{x,y}; \\ p^y &:= \mathbb{P}(Y = y) = \sum_x p_{x,y}; \\ p_{x|y} &:= \mathbb{P}(X = x|Y = y) = p_{x,y}/p^y; \\ p^{y|x} &:= \mathbb{P}(Y = y|X = x) = p_{x,y}/p_x. \end{aligned}$$

In most settings j denotes a subset of \mathcal{X} . When a graph G is given on the vertex set \mathcal{X} , then j always denotes an independent set: $j \subseteq \mathcal{X}$ is a set of vertices such that the induced subgraph $G[j]$ contains no edge. In this setting J stands for a random independent set. Then \mathcal{J} denotes the set of all independent sets, while \mathcal{J}_x is the set of independent sets containing x .

C. Contributions

Alternating minimization: Let us consider the following optimization problem.

Problem. Suppose that we have two finite families of probability measures on a given finite set \mathcal{J} : $\mu_x, x \in \mathcal{X}$ and $\nu_y, y \in \mathcal{Y}$. In the first family for each $x \in \mathcal{X}$ we have a constraint: the support $\text{supp } \mu_x$ must be contained in a given subset \mathcal{J}_x of \mathcal{J} . Find the measures μ_x, ν_y that minimize the weighted sum of the Kullback–Leibler divergences:

$$\sum_{x,y} p_{x,y} D_{\text{KL}}(\mu_x \parallel \nu_y) \text{ for some given weights } p_{x,y} \geq 0.$$

That is, given $\mathcal{J}_x \subseteq \mathcal{J}, x \in \mathcal{X}$ and $p_{x,y} \geq 0, x \in \mathcal{X}, y \in \mathcal{Y}$, find the minimum of the above sum under the constraint $\text{supp } \mu_x \subseteq \mathcal{J}_x$.

In our setting we have random variables X and Y taking values in the finite sets \mathcal{X} and \mathcal{Y} , respectively, and G is a graph on the vertex set \mathcal{X} . Then by j we denote an independent set of G , hence each j is a subset of \mathcal{X} . We choose \mathcal{J} to be the set of all j , while

$$\mathcal{J}_x := \{j : x \in j\}$$

consists of the independent sets containing a fixed x . With this setup and with $p_{x,y} := \mathbb{P}(X = x, Y = y)$, the minimum of the problem above turns out to be precisely $H_G(X|Y)$.

To get concrete formulas, let us represent the distributions μ_x and ν_y by the following vectors:

$$\begin{aligned} \mathbf{q} &= (q_{j|x})_{(j,x) \in \mathcal{J} \times \mathcal{X}} \in \mathbb{R}^{\mathcal{J} \times \mathcal{X}}; \\ \mathbf{r} &= (r_{j|y})_{(j,y) \in \mathcal{J} \times \mathcal{Y}} \in \mathbb{R}^{\mathcal{J} \times \mathcal{Y}}, \end{aligned}$$

where $q_{j|x}$ and $r_{j|y}$ stand for $\mu_x(\{j\})$ and $\nu_y(\{j\})$, respectively.¹ The constraints for \mathbf{q} and \mathbf{r} lead to the following definition.

Definition 1. We define the convex polytopes $K_{\mathbf{q}} \subset \mathbb{R}^{\mathcal{J} \times \mathcal{X}}$ and $K_{\mathbf{r}} \subset \mathbb{R}^{\mathcal{J} \times \mathcal{Y}}$ as

$$K_{\mathbf{q}} := \left\{ \mathbf{q} = (q_{j|x}) : q_{j|x} \geq 0; \sum_{j \ni x} q_{j|x} = 1 (\forall x \in \mathcal{X}); q_{j|x} = 0 \text{ if } x \notin j \right\}$$

and

$$K_{\mathbf{r}} := \left\{ \mathbf{r} = (r_{j|y}) : r_{j|y} \geq 0; \sum_j r_{j|y} = 1 (\forall y \in \mathcal{Y}) \right\}.$$

By $\text{int}(K_{\mathbf{q}})$ and $\text{int}(K_{\mathbf{r}})$ we denote the relative interiors of the polytopes (within their affine hull).

In the sequel we will always assume that $\mathbf{q} \in K_{\mathbf{q}}$ and $\mathbf{r} \in K_{\mathbf{r}}$. Then

$$D_{\text{KL}}(\mu_x \parallel \nu_y) = \sum_j q_{j|x} \log \frac{q_{j|x}}{r_{j|y}}.$$

Therefore we need to minimize the function

$$\varphi(\mathbf{q}, \mathbf{r}) := \sum_{x,y,j} p_{x,y} q_{j|x} \log \frac{q_{j|x}}{r_{j|y}} \quad (3)$$

over $\mathbf{q} \in K_{\mathbf{q}}$ and $\mathbf{r} \in K_{\mathbf{r}}$.

As we will see, this is an alternating minimization problem. The point is that if we fix one of the two variables \mathbf{q} and \mathbf{r} , then there are explicit formulas for the optimal choice of the other variable: we will define maps

$$Q: K_{\mathbf{r}} \rightarrow K_{\mathbf{q}} \text{ and } R: K_{\mathbf{q}} \rightarrow K_{\mathbf{r}}$$

such that $\mathbf{r} = R(\mathbf{q})$ is the optimal choice for a fixed \mathbf{q} , and similarly $\mathbf{q} = Q(\mathbf{r})$ is optimal for a fixed \mathbf{r} ; that is, for any \mathbf{q} and \mathbf{r} we have

$$\varphi(\mathbf{q}, \mathbf{r}) \geq \varphi(\mathbf{q}, R(\mathbf{q})) \text{ and } \varphi(\mathbf{q}, \mathbf{r}) \geq \varphi(Q(\mathbf{r}), \mathbf{r}).$$

Using Q and R we can explicitly define the following functions:

$$\begin{aligned} \varphi_{\mathbf{q}}(\mathbf{q}) &:= \varphi(\mathbf{q}, R(\mathbf{q})) = \min_{\mathbf{r}} \varphi(\mathbf{q}, \mathbf{r}) \text{ and} \\ \varphi_{\mathbf{r}}(\mathbf{r}) &:= \varphi(Q(\mathbf{r}), \mathbf{r}) = \min_{\mathbf{q}} \varphi(\mathbf{q}, \mathbf{r}), \end{aligned}$$

¹We index the coordinates/variables by $j|x$ and $j|y$ to emphasize the fact that they express certain conditional probabilities, see the proof of Proposition 8 for details. This notation may also serve as a reminder that $q_{j|x}$ and $r_{j|y}$ have to sum up to 1 for any fixed x and y , respectively.

They clearly have the same minimum as φ . When we work out the details in Section II, we will see that the q-problem $\min \varphi_q$ is actually equivalent to the original formula (2) for conditional graph entropy (see Proposition 8).

Theorem 2. *We have the following formulas for conditional graph entropy:*

$$H_G(X|Y) = \min_{K_q \times K_r} \varphi = \min_{K_q} \varphi_q = \min_{K_r} \varphi_r.$$

Algorithm: When trying to find the minimum of $\varphi(\mathbf{q}, \mathbf{r})$, the fact that we can easily optimize in either variable (while the other is fixed) gives rise to the following simple iterative algorithm. Let us start from a point $\mathbf{q}^{(0)}$ and apply R and Q alternately:

$$\mathbf{q}^{(0)} \xrightarrow{R} \mathbf{r}^{(0)} \xrightarrow{Q} \mathbf{q}^{(1)} \xrightarrow{R} \mathbf{r}^{(1)} \xrightarrow{Q} \mathbf{q}^{(2)} \xrightarrow{R} \mathbf{r}^{(2)} \dots \quad (4)$$

The corresponding φ -value decreases at each step:

$$\begin{aligned} \varphi(\mathbf{q}^{(0)}, \mathbf{r}^{(0)}) &= \varphi_q(\mathbf{q}^{(0)}) \\ &\quad \Downarrow \\ \varphi(\mathbf{q}^{(1)}, \mathbf{r}^{(0)}) &= \varphi_r(\mathbf{r}^{(0)}) \\ &\quad \Downarrow \\ \varphi(\mathbf{q}^{(1)}, \mathbf{r}^{(1)}) &= \varphi_q(\mathbf{q}^{(1)}) \\ &\quad \Downarrow \\ \varphi(\mathbf{q}^{(2)}, \mathbf{r}^{(1)}) &= \varphi_r(\mathbf{r}^{(1)}) \\ &\quad \Downarrow \\ \varphi(\mathbf{q}^{(2)}, \mathbf{r}^{(2)}) &= \varphi_q(\mathbf{q}^{(2)}) \\ &\quad \vdots \end{aligned}$$

One can also think of this alternating optimization as “jumping” between the q-problem $\min_{K_q} \varphi_q$ and the r-problem $\min_{K_r} \varphi_r$ using the maps $Q: K_r \rightarrow K_q$ and $R: K_q \rightarrow K_r$. The value to minimize (i.e., the φ_q -value and the φ_r -value, respectively) always decreases, so with each step we get closer to the optimum.

Following the footsteps of the general theory of Csiszár and Tusnády [6], we will show that, for an arbitrary starting point $\mathbf{q}^{(0)}$ in the relative interior $\text{int}(K_q)$, the iterative process converges to the minimum.

Theorem 3. *For an arbitrary starting point $\mathbf{q}^{(0)} \in \text{int}(K_q)$ consider the sequence (4) obtained by alternating optimization. Then $\varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n)})$ is a decreasing sequence that converges to $\min_{K_q \times K_r} \varphi$ as $n \rightarrow \infty$.*

We implemented the algorithm in Python and made the codes publicly available in a GitHub repository [11].

Convex corners: As we have mentioned, the q-problem $\min \varphi_q$ gives back the original formula (2). On the other hand, the r-problem $\min \varphi_r$ gives us a new formula. We make this new formula explicit in the next theorem because it shows how conditional graph entropy is related to convex corners (a known phenomenon in the unconditioned case).

Theorem 4. *For any (maximal) independent set j of G and any possible value y of Y we have a variable $r_{j|y}$. Then conditional graph entropy can be expressed as the solution of the following optimization problem:*

$$H_G(X|Y) = \min - \sum_x p_x \log \left(\sum_{j: x \in j} \prod_y (r_{j|y})^{p^{y|x}} \right), \quad (5)$$

where the minimum is taken over all choices of $r_{j|y} \geq 0$ satisfying $\sum_j r_{j|y} = 1$ for each fixed y .

We can easily turn this new r-problem into another one (that we will call the a-problem) which attests that conditional graph entropy is a special case of a more general entropy notion defined for convex corners.²

To see this connection, note that (5) is in the form

$$H_G(X|Y) = \min_{K_a} - \sum_x p_x \log A_x,$$

where A_x is a function of the variables $r_{j|y}$. The key property is that $A_x: K_r \rightarrow [0, 1]$ is a concave function for each x . It means that the set of image points $\mathbf{a} = (a_x)_{x \in \mathcal{X}}$ with $a_x = A_x(\mathbf{r})$, as \mathbf{r} ranges over K_r , (essentially) defines a convex corner K_a in $\mathbb{R}^{\mathcal{X}}$. Then we have

$$H_G(X|Y) = \min_{K_a} \varphi_a, \text{ where } \varphi_a(\mathbf{a}) = - \sum_x p_x \log a_x.$$

A nice feature of this a-problem is that the minimum is attained at a single point $\mathbf{a} \in K_a$ because φ_a is strictly convex (provided that $p_x = \mathbb{P}(X = x) > 0$ for each x). Also note that φ_a depends only on the distribution of X , while the convex

²A convex corner of $\mathbb{R}^{\mathcal{X}}$ is a convex compact set in the positive orthant $[0, \infty)^{\mathcal{X}}$ that is *downward closed*, i.e., if we take any point in the set and decrease some of its coordinates, then the new point still lies in the set (see Definition 15).

corner K_a depends only on the graph G and the conditional distributions $Y|X=x$ for any given x . Thus, the parameters of the problem are, so to say, split between φ_a and K_a .

Moreover, we will define another convex corner, denoted by L , that can be regarded as the *dual problem*. To keep the introduction concise, we will postpone the actual definition of L and the precise statements until Section IV-C. In short, we will show that

$$H_G(X|Y) = \min_{K_a} \varphi_a = -\min_L \varphi_a,$$

and a vector $\mathbf{a} = (a_x)_{x \in \mathcal{X}} \in K_a$ is optimal (i.e., the minimum point of φ_a) if and only if $\mathbf{a}^{-1} := (a_x^{-1})_{x \in \mathcal{X}} \in L$. This provides a fairly simple way to check optimality, and even leading to an error bound for our iterative algorithm as we will explain in Section V.

In the special case when Y is trivial, i.e., \mathcal{Y} is a one-element set, we get back Körner's original setting of graph entropy, and things simplify considerably. For example, K_a is simply a polytope: the aforementioned *vertex packing polytope* (or *independent/stable set polytope*) $VP(G)$. We did not find any mention in the literature of the fact that graph entropy can be considered as an alternating minimization problem. In particular, to the best of our knowledge, the corresponding iterative algorithm has not been used or proposed before even in this unconditioned setting.

Outline of the paper

We give more details of our alternating minimization problem in Section II, and collect its key properties in Section III, proving, in particular, the convergence of the iterative process. In Section IV we discuss the connection to convex corners and introduce the dual problem. In Section V we discuss some details of the iterative algorithm; in particular, an error bound based on the dual problem and a tweak for speeding the convergence up.

II. THE ALTERNATING MINIMIZATION PROBLEM

In this section we rigorously introduce the optimization problem $\min \varphi(\mathbf{q}, \mathbf{r})$ described in the introduction. We will use the notations outlined in Section I-B.

A. Assumptions

For the sake of simplicity, we will work under the following three assumptions that do not actually reduce generality.

- **Each** $p_x = \sum_y p_{x,y}$ **and** $p^y = \sum_x p_{x,y}$ **is strictly positive.** (Otherwise we simply delete the corresponding elements from \mathcal{X} and \mathcal{Y} .) Note that under this assumption the conditional probabilities $p_{x|y}$ and $p^{y|x}$ all exist.
- **The sets j cover \mathcal{X} ,** that is, $\forall x \in \mathcal{X} \exists j \in \mathcal{J}$ s.t. $x \in j$. (Otherwise the minimum we will consider would be ∞ anyway.)
- **\mathcal{J} contains inclusion-wise maximal sets.** (Removing subsets of other sets from \mathcal{J} does not change the minimum.)

Also, all the results will be true under the more general setting when $\mathcal{J} \subseteq \mathcal{P}(\mathcal{X})$ is any set of subsets of \mathcal{X} . That is, the sets $j \in \mathcal{J}$ are subsets of \mathcal{X} but they do not necessary need to be independent sets of some graph G on \mathcal{X} . In conclusion, our setup essentially has the following fixed parameters: the probabilities $p_{x,y}$ and a binary relation \in on $\mathcal{X} \times \mathcal{J}$: whenever x is in the set j we write $x \in j$.³

B. The mappings

Recall the convex polytopes K_q and K_r defined in Section I-C of the introduction. Now we explicitly define the mappings Q and R between these polytopes along with an auxiliary mapping A . In fact, the formula defining Q will make sense only on the subset K_r^* where none of the coordinates of A vanishes. We define Q arbitrarily outside K_r^* .

Definition 5. We define the mappings $A: K_r \rightarrow \mathbb{R}^{\mathcal{X}}$; $Q: K_r \rightarrow K_q \subset \mathbb{R}^{\mathcal{J} \times \mathcal{X}}$; $R: K_q \rightarrow K_r \subset \mathbb{R}^{\mathcal{J} \times \mathcal{Y}}$ by the following coordinate-wise functions $Q_{j|x}$, $R_{j|y}$, A_x :

$$\begin{aligned} R_{j|y}(\mathbf{q}) &:= \sum_{x \in j} p_{x|y} q_{j|x}; \\ A_x(\mathbf{r}) &:= \sum_{j \ni x} \prod_y (r_{j|y})^{p^{y|x}}; \\ Q_{j|x}(\mathbf{r}) &:= \begin{cases} 0 & \text{if } x \notin j; \\ \prod_y (r_{j|y})^{p^{y|x}} / A_x(\mathbf{r}) & \text{if } x \in j. \end{cases} \end{aligned}$$

³Equivalently, we may write $j \ni x$. In particular, $\sum_{j \ni x}$ means that the sum runs over sets $j \in \mathcal{J}$ containing the (fixed) element x .

Since the formula for $Q_{j|x}$ involves a division by A_x , it only defines Q over the subset

$$K_r^* := K_r \setminus \bigcup_x A_x^{-1}(0). \quad (6)$$

For $\mathbf{r} \in K_r \setminus K_r^*$ let $Q(\mathbf{r})$ be an arbitrary point in $\text{int}(K_q)$.

In the formulas above we define $t^0 = 1$ even for $t = 0$. This ensures that A is continuous over the entire K_r even when $p^{y|x} = 0$ for some pairs x, y . It is also consistent with the convention $0 \cdot \log 0 = 0$ which is implicit in the definition of Shannon entropy.

It is straightforward to check that $Q(\mathbf{r}) \in K_q$ and $R(\mathbf{q}) \in K_r$ always hold. For example, in the definition of $Q_{j|x}(\mathbf{r})$, dividing by $A_x(\mathbf{r})$ ensures that their sum is 1 for any fixed x .

Remark 6. Note that R is a linear map and it actually describes how the conditional distributions $J|Y = y$ can be expressed in terms of $J|X = x$ in a Markov chain $Y - X - J$; see the proof of Proposition 8 for details.

C. The functions

Now we can turn our attention to the functions to be minimized. We already gave an explicit formula (3) for $\varphi(\mathbf{q}, \mathbf{r})$ in the introduction. However, we did not mention a few subtleties there. In particular, we need to specify the function values when some of the variables $q_{j|x}$ or $r_{j|y}$ are 0.

Definition 7. For $u, v \in [0, 1]$ let

$$f(u, v) := u \log u - u \log v$$

with the usual conventions $\log 0 = -\infty$ and $0 \cdot \infty = 0$ so that

$$f(0, v) = 0 \text{ if } v \in [0, 1] \quad \text{and} \quad f(u, 0) = \infty \text{ if } u \in (0, 1].$$

Then

$$\varphi(\mathbf{q}, \mathbf{r}) := \sum_{x,y,j} p_{x,y} f(q_{j|x}, r_{j|y}) \quad (7)$$

is well-defined for any $\mathbf{q} \in \mathbb{R}^{\mathcal{J} \times \mathcal{X}}$ and $\mathbf{r} \in \mathbb{R}^{\mathcal{J} \times \mathcal{Y}}$. Note that we may restrict the sum for $x \in j$ because otherwise $q_{j|x} = 0$, and hence the summand is 0 anyway.

Let us also define the following auxiliary function that we will need for establishing the so-called 3-point and 4-point properties.

$$\delta(\mathbf{q}, \mathbf{q}') := \sum_x p_x \sum_{j \ni x} f(q_{j|x}, q'_{j|x}) = \sum_x p_x \sum_{j \ni x} q_{j|x} \log q_{j|x} - q_{j|x} \log q'_{j|x}. \quad (8)$$

One may think of $\varphi(\mathbf{q}, \mathbf{r})$ and $\delta(\mathbf{q}, \mathbf{q}')$ as (non-symmetric) squared distances between these points. We mention that both functions are convex combinations of certain Kullback–Leibler divergences. In particular, they are nonnegative and they may be ∞ . For example, $\varphi(\mathbf{q}, \mathbf{r}) = \infty$ if and only if there exist x, y, j such that $r_{j|y} = 0$ while $p_{x,y} > 0$ and $q_{j|x} > 0$. It is easy to see that if $\mathbf{r} \in K_r \setminus K_r^*$, then $\varphi(\mathbf{q}, \mathbf{r}) = \infty$ for any choice of $\mathbf{q} \in K_q$. (The two-line proof of this fact is included in the proof of Proposition 13.)

We include here two useful equivalent formulas for φ . On the one hand, summing $q_{j|x} \log q_{j|x}$ and $q_{j|x} \log r_{j|y}$ separately gives

$$\varphi(\mathbf{q}, \mathbf{r}) = \sum_x p_x \sum_{j \ni x} q_{j|x} \log q_{j|x} - \sum_y p^y \sum_j R_{j|y}(\mathbf{q}) \log r_{j|y}. \quad (9)$$

On the other hand, for $\mathbf{r} \in K_r^*$ we can write

$$\begin{aligned} \varphi(\mathbf{q}, \mathbf{r}) &= \sum_x p_x \sum_{j \ni x} q_{j|x} \left(\log q_{j|x} - \sum_y p^{y|x} \log r_{j|y} \right) \\ &= \sum_x p_x \sum_{j \ni x} q_{j|x} \log \frac{q_{j|x}}{\prod_y (r_{j|y})^{p^{y|x}}} = \sum_x p_x \sum_{j \ni x} q_{j|x} \log \frac{q_{j|x}}{Q_{j|x}(\mathbf{r}) A_x(\mathbf{r})} \end{aligned} \quad (10)$$

with the remark that if $q_{j|x}$ and $Q_{j|x}(\mathbf{r})$ are both 0, then the fraction in the log should simply be $1/A_x(\mathbf{r})$.

Next we define φ_q and φ_r . At this point we simply express them using Q , R , and φ , but we will shortly see that they are indeed the minimum of φ with one of the variables fixed. Using (9) and (10) we get the following specific formulas: for $\mathbf{q} \in K_q$ and $\mathbf{r} \in K_r$ let

$$\varphi_q(\mathbf{q}) := \varphi(\mathbf{q}, R(\mathbf{q})) = \sum_x p_x \sum_{j \ni x} q_{j|x} \log q_{j|x} - \sum_y p^y \sum_j R_{j|y}(\mathbf{q}) \log R_{j|y}(\mathbf{q}); \quad (11)$$

$$\varphi_r(\mathbf{r}) := \varphi(Q(\mathbf{r}), \mathbf{r}) = - \sum_x p_x \log A_x(\mathbf{r}) = - \sum_x p_x \log \sum_{j \ni x} \prod_y (r_{j|y})^{p^{y|x}}. \quad (12)$$

Note that (12) works even for $\mathbf{r} \notin K_r^*$ as all the expressions are ∞ in that case.

Proposition 8. *If \mathcal{J} is the set of independent sets of some graph G on the vertex set \mathcal{X} , then*

$$H_G(X|Y) = \min_{K_q} \varphi_q.$$

Proof. Recall that the original formula (2) for $H_G(X|Y)$ involves minimization over random J containing X and independent from Y when conditioned on X (in other words, $Y - X - J$ is a Markov chain). To define such a J one needs to specify the conditional probabilities $\mathbb{P}(J = j|X = x)$ whenever $x \in j$. These conditional probabilities can be represented by a vector $\mathbf{q} \in K_q$. Due to the conditional independence, we have the expression

$$\mathbb{P}(J = j|Y = y) = \sum_{x \in j} p_{x|y} \mathbb{P}(J = j|X = x). \quad (13)$$

Note that we defined R using the same linear combinations, see Definition 5. Consequently, if the $\mathbb{P}(J = j|X = x)$'s are represented by \mathbf{q} , then the $\mathbb{P}(J = j|Y = y)$'s are represented by $\mathbf{r} = R(\mathbf{q})$. Therefore

$$\begin{aligned} H(J|X) &= - \sum_x p_x \sum_{j \ni x} q_{j|x} \log q_{j|x}; \\ H(J|Y) &= - \sum_y p^y \sum_j R_{j|y}(\mathbf{q}) \log R_{j|y}(\mathbf{q}), \end{aligned}$$

and hence the conditional mutual information $I(X; J|Y) = H(J|Y) - H(J|X)$ is precisely $\varphi_q(\mathbf{q})$ according to (11), proving $H_G(X|Y) = \min_{K_q} \varphi_q$. \square

III. CONVERGENCE

In this section we derive various properties of the the minimization problems introduced in Section II. They will culminate in the proof that alternating optimization converges to the true minimum (Theorem 3). We will also prove Theorems 2 and 4 along the way.

Proposition 9. *The functions φ and δ are nonnegative, lower semicontinuous, and convex. Moreover, $\delta(\mathbf{q}, \mathbf{q}') = 0$ if and only if $\mathbf{q} = \mathbf{q}'$.*

Proof. Recall that φ and δ were defined using the function $f: [0, 1]^2 \rightarrow (-\infty, \infty]$ in Definition 7. It is well known and easy to show that f is convex and lower semicontinuous, and hence so are φ and δ .

Using the convexity of f we get that for any fixed x, y :

$$\text{both } \sum_j f(q_{j|x}, r_{j|y}) \text{ and } \sum_j f(q_{j|x}, q'_{j|x}) \geq |\mathcal{J}| f(1/|\mathcal{J}|, 1/|\mathcal{J}|) = 0$$

showing that $\varphi, \delta \geq 0$. (This, of course, also follows from their representations as the sum of Kullback–Leibler divergences.) \square

Note that lower semicontinuity implies that φ attains its minimum over any compact set. In particular, it has a minimum over $K_q \times K_r$.

Proposition 10 ($\mathbf{r} = R(\mathbf{q})$ is optimal for fixed \mathbf{q}). *We have*

$$\varphi(\mathbf{q}, \mathbf{r}) \geq \varphi(\mathbf{q}, R(\mathbf{q})) = \varphi_q(\mathbf{q}) \text{ for any } \mathbf{q} \in K_q; \mathbf{r} \in K_r.$$

Equality holds if and only if $\mathbf{r} = R(\mathbf{q})$.

Proof. Using formula (9) for a fixed \mathbf{q} , it immediately follows from Gibbs' inequality (applied for each y in the second sum) that the unique optimal choice for \mathbf{r} is $R(\mathbf{q})$. \square

Proposition 11 ($\mathbf{q} = Q(\mathbf{r})$ is optimal for fixed \mathbf{r}). *We have*

$$\varphi(\mathbf{q}, \mathbf{r}) \geq \varphi(Q(\mathbf{r}), \mathbf{r}) = \varphi_r(\mathbf{r}) \text{ for any } \mathbf{q} \in K_q; \mathbf{r} \in K_r.$$

If $\mathbf{r} \notin K_r^$, then both sides are ∞ . Furthermore, for $\mathbf{r} \in K_r^*$ equality holds if and only if $\mathbf{q} = Q(\mathbf{r})$.*

Proof. This is an immediate consequence of the 3-point property (that we will shortly state in Proposition 13) and the fact that $\delta \geq 0$. \square

Corollary 12. *Propositions 10 and 11 clearly show that*

$$\varphi_q(\mathbf{q}) = \min_{\mathbf{r} \in K_r} \varphi(\mathbf{q}, \mathbf{r}) \text{ and } \varphi_r(\mathbf{r}) = \min_{\mathbf{q} \in K_q} \varphi(\mathbf{q}, \mathbf{r}).$$

In particular, φ , φ_q , φ_r have the same minimum over their respective convex domains:

$$\min_{K_q \times K_r} \varphi = \min_{K_q} \varphi_q = \min_{K_r} \varphi_r.$$

It also follows that both φ_q and φ_r are convex as they can be obtained as minimizing the convex $\varphi(\mathbf{q}, \mathbf{r})$ in one of the variables.

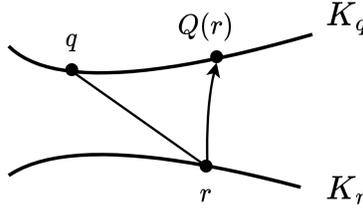
Note that, combined with Proposition 8, this completes the proof of Theorem 2. Moreover, Theorem 4 also follows as we simply need to substitute (12), which expresses φ_r , into $H_G(X|Y) = \min_{K_r} \varphi_r$.

From this point on we follow the footsteps of the general theory [6] of alternating minimization problems by proving the so-called *3-point and 4-point properties*, and show how they imply convergence to the minimum through the 5-point property.

The following identity can be thought of as a Pythagorean theorem for the “squared distances” φ and δ . Csiszár and Tusnády refer to it as the 3-point property. (In their general setting it may hold only as an inequality \geq but in our case we always have equality.)

Proposition 13 (3-point property). *For any $\mathbf{q} \in K_q$ and $\mathbf{r} \in K_r$ we have*

$$\varphi(\mathbf{q}, \mathbf{r}) = \delta(\mathbf{q}, Q(\mathbf{r})) + \underbrace{\varphi(Q(\mathbf{r}), \mathbf{r})}_{=\varphi_r(\mathbf{r})}.$$



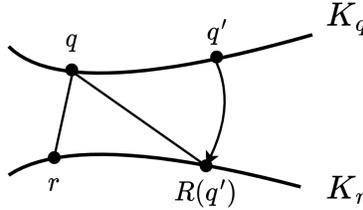
Proof. First let us consider the cases when one of the terms on the right-hand side is ∞ . In both cases we need to show that one can find x, y, j with $p_{x,y} > 0$, $q_{j|x} > 0$, $r_{j|y} = 0$ so that we can conclude that the left-hand side $\varphi(\mathbf{q}, \mathbf{r})$ is also ∞ .

- We have $\varphi_r(\mathbf{r}) = \infty$ if and only if $A_x(\mathbf{r}) = 0$ for some x . Fix such an x and take a $j \ni x$ with $q_{j|x} > 0$, which must exist as their sum is 1. Since $A_x(\mathbf{r}) = 0$, there must exist y such that $r_{j|y} = 0$ and $p_{x,y} > 0$.
- We have $\delta(\mathbf{q}, Q(\mathbf{r})) = \infty$ if and only if there exist j, x such that $q_{j|x} > 0$ but $Q_{j|x}(\mathbf{r}) = 0$, which means, by the definition of $Q_{j|x}$, that there exists y such that $p_{x,y} > 0$ and $r_{j|y} = 0$.

Otherwise we can simply combine formula (10) for $\varphi(\mathbf{q}, \mathbf{r})$, formula (12) for $\varphi_r(\mathbf{r})$, and formula (8) for $\delta(\mathbf{q}, \mathbf{q}')$ with $\mathbf{q}' = Q(\mathbf{r})$ to get the claim. \square

Proposition 14 (4-point property). *For any $\mathbf{q}, \mathbf{q}' \in K_q$ and $\mathbf{r} \in K_r$ we have*

$$\varphi(\mathbf{q}, R(\mathbf{q}')) \leq \varphi(\mathbf{q}, \mathbf{r}) + \delta(\mathbf{q}, \mathbf{q}').$$



Proof. We may assume that the right-hand side is finite, otherwise the inequality is trivial. It follows that for any triple x, y, j with $p_{x,y} > 0$, $q_{j|x} > 0$ we must have both $r_{j|y} > 0$ and $q'_{j|x} > 0$. Let $\mathbf{r}' := R(\mathbf{q}')$. Then $r'_{j|y} \geq p_{x,y} q'_{j|x} > 0$. So for any such triple all the variables are positive and we may write:

$$\varphi(\mathbf{q}, \mathbf{r}) + \delta(\mathbf{q}, \mathbf{q}') - \varphi(\mathbf{q}, \mathbf{r}') = \sum_{x,y,j} p_{x,y} q_{j|x} \log \frac{q_{j|x} r'_{j|y}}{q'_{j|x} r_{j|y}},$$

which, using that $\log t \geq 1 - 1/t$, can be bounded from below as follows:

$$\sum_{x,y,j} p_{x,y} q_{j|x} \left(1 - \frac{q'_{j|x} r_{j|y}}{q_{j|x} r'_{j|y}} \right) = 1 - \sum_y p^y \sum_j \frac{r_{j|y}}{r'_{j|y}} \underbrace{\sum_x p_{x|y} q'_{j|x}}_{=R_{j|y}(\mathbf{q}')=r'_{j|y}} = 1 - \sum_y p^y \sum_j r_{j|y} = 0,$$

and the proof is complete. \square

Now we are ready to prove that the alternating optimization process converges to the minimum.

Proof of Theorem 3. Consider the sequences $\mathbf{q}^{(n)}$ and $\mathbf{r}^{(n)}$ of alternating optimization started from some $\mathbf{q}^{(0)} \in \text{int}(K_q)$. Fix any pair $\mathbf{q} \in K_q$, $\mathbf{r} \in K_r$, and let n be a positive integer. Using Proposition 13 for the triple $\mathbf{q}, \mathbf{r}, \mathbf{r}^{(n-1)} \xrightarrow{Q} \mathbf{q}^{(n)}$ and Proposition 14 for the quadruple $\mathbf{q}, \mathbf{r}, \mathbf{q}^{(n)} \xrightarrow{R} \mathbf{r}^{(n)}$ we get that

$$\begin{aligned} \delta(\mathbf{q}, \mathbf{q}^{(n)}) + \varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n-1)}) &\stackrel{3-pt}{=} \varphi(\mathbf{q}, \mathbf{r}^{(n-1)}); \\ \varphi(\mathbf{q}, \mathbf{r}^{(n)}) &\stackrel{4-pt}{\leq} \varphi(\mathbf{q}, \mathbf{r}) + \delta(\mathbf{q}, \mathbf{q}^{(n)}). \end{aligned}$$

Since $\mathbf{q}^{(n)} \in \text{int}(K_q)$ holds for all n , we have $\delta(\mathbf{q}, \mathbf{q}^{(n)}) < \infty$. Therefore adding the two inequalities above results in

$$\varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n-1)}) + \varphi(\mathbf{q}, \mathbf{r}^{(n)}) \leq \varphi(\mathbf{q}, \mathbf{r}) + \varphi(\mathbf{q}, \mathbf{r}^{(n-1)}).$$

Since $\varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n)}) \leq \varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n-1)})$ by Proposition 10, it follows that

$$\varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n)}) + \varphi(\mathbf{q}, \mathbf{r}^{(n)}) \leq \varphi(\mathbf{q}, \mathbf{r}) + \varphi(\mathbf{q}, \mathbf{r}^{(n-1)}). \quad (14)$$

This is what Csiszár and Tusnády refer to as the *5-point property* for the points $\mathbf{q}, \mathbf{r}, \mathbf{r}^{(n-1)} \xrightarrow{Q} \mathbf{q}^{(n)} \xrightarrow{R} \mathbf{r}^{(n)}$.

Note that the second term on either side is an element of the sequence $\varphi(\mathbf{q}, \mathbf{r}^{(n)}) \geq 0$. First we assume that these elements are all finite. Then for any $\varepsilon > 0$ there must be infinitely many n such that

$$\varphi(\mathbf{q}, \mathbf{r}^{(n)}) \geq \varphi(\mathbf{q}, \mathbf{r}^{(n-1)}) - \varepsilon,$$

otherwise the sequence would converge to $-\infty$, contradicting that each element is nonnegative. For any such n we get from (14) that

$$\varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n)}) \leq \varphi(\mathbf{q}, \mathbf{r}) + \varepsilon.$$

Since $\varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n)})$ is monotone decreasing, it has a limit that must satisfy

$$\lim_{n \rightarrow \infty} \varphi(\mathbf{q}^{(n)}, \mathbf{r}^{(n)}) \leq \varphi(\mathbf{q}, \mathbf{r}) + \varepsilon$$

for any positive ε , and hence for $\varepsilon = 0$ as well.

If, on the other hand, $\varphi(\mathbf{q}, \mathbf{r}^{(n)}) = \infty$ for some n , then $\varphi(\mathbf{q}, \mathbf{r}) = \infty$ follows from the 4-point property as $\delta(\mathbf{q}, \mathbf{q}^{(n)}) < \infty$, and we have the same conclusion: the limit is at most $\varphi(\mathbf{q}, \mathbf{r})$.

Since this holds for any \mathbf{q} and \mathbf{r} , it follows that the limit must be the minimum of φ . \square

IV. CONVEX CORNERS

Convex corners are downward closed, convex subsets of $[0, \infty)^n$. It is possible to define entropy functions for convex corners, and this general theory was known to include the notion of graph entropy (via the vertex packing polytope, a convex corner associated to a graph). In this section we will show that conditional graph entropy can also be expressed as the entropy of an associated convex corner. Moreover, we will even define a dual problem in the form of another convex corner.

Besides revealing a nice theoretical connection to a general theory, this also has significant practical implications: the dual problem provides a way to check optimality in the primal problem, even yielding an error bound. The error bound comes in particularly handy when combined with alternating optimization: we can stop at any time through the iterations $\mathbf{r}^{(n)}$ and compute this error bound δ , which then ensures that we are at most δ away from the optimum:

$$\varphi_r(\mathbf{r}^{(n)}) - \delta \leq H_G(X|Y) \leq \varphi_r(\mathbf{r}^{(n)}).$$

We start by recalling the basic concepts regarding convex corners.

A. Entropy of convex corners

Definition 15. A set $K \subset [0, \infty)^{\mathcal{X}}$ is said to be *downward closed* if the following property⁴ holds:

$$\text{if } \mathbf{a} \in K, \text{ then } \mathbf{a}' \in K \text{ for all points } 0 \leq \mathbf{a}' \leq \mathbf{a}.$$

Similarly, K is *upward closed* if $\mathbf{a}' \in K$ whenever $\mathbf{a} \in K$ and $\mathbf{a}' \geq \mathbf{a}$.

We say that $K \subset [0, \infty)^{\mathcal{X}}$ is a *convex corner* if K is compact, convex, and downward closed. Usually K is also required to have nonempty interior, or equivalently, to contain a point with strictly positive coordinates.

⁴Here $0 \leq \mathbf{a}' \leq \mathbf{a}$ means that $0 \leq a'_x \leq a_x$ for each x . As before, we use the notation $\mathbf{a} = (a_x)_{x \in \mathcal{X}}$ for points in $\mathbb{R}^{\mathcal{X}}$.

Given a random variable X and the corresponding probabilities p_x , $x \in \mathcal{X}$, let φ_a denote the following $[0, \infty)^{\mathcal{X}} \rightarrow [0, \infty]$ function:

$$\varphi_a: \mathbf{a} \mapsto - \sum_x p_x \log a_x.$$

Note that φ_a depends on the distribution of X , and we write φ_a^X when we want to emphasize this dependence. The entropy is defined as the minimum of φ_a^X over K :

$$H_K(X) := \min_{\mathbf{a} \in K} \varphi_a^X(\mathbf{a}).$$

The function $H_K(\cdot)$, defined for random variables on \mathcal{X} , is sometimes referred to as the *entropy function* corresponding to the convex corner K . It can be seen that the entropy function $H_K(\cdot)$ uniquely determines K .

A related useful concept is the *antiblocker* K^* of a convex corner K :

$$K^* := \{ \mathbf{b} \geq 0 : \sum_x a_x b_x \leq 1 \text{ for all } \mathbf{a} \in K \}.$$

One can show that K^* is also a convex corner, $(K^*)^* = K$, and $H_K(X) + H_{K^*}(X) = H(X)$. For these and further properties of $H_K(\cdot)$, see [17, Sections 4.1 & 6] and [19, Section 5].

We will also use the following notations: for $\mathbf{a}, \mathbf{b} \in [0, \infty)^{\mathcal{X}}$ let \mathbf{ab} denote the vector $(a_x b_x)_{x \in \mathcal{X}}$ (coordinate-wise multiplication). Similarly, \mathbf{a}^{-1} denotes the vector with coordinates $1/a_x$ (provided that each a_x is positive). Furthermore,

$$\mathbf{a}K := \{ \mathbf{ab} : \mathbf{b} \in K \} \text{ and } K^{-1} := \{ \mathbf{b}^{-1} : \mathbf{b} \in K \}.$$

Note that $\varphi_a(\mathbf{ab}) = \varphi_a(\mathbf{a}) + \varphi_a(\mathbf{b})$ and $\varphi_a(\mathbf{a}^{-1}) = -\varphi_a(\mathbf{a})$. Finally, we denote the vector $(p_x)_{x \in \mathcal{X}}$ by \mathbf{p} . Then $\varphi_a(\mathbf{p}) = -\sum_x p_x \log p_x$ is the entropy of X .

B. Primal problem

Now we introduce the *a-problem* $\min_{K_a} \varphi_a$, which is, in fact, an equivalent formulation of the r-problem. We have already defined the function to minimize: φ_a . Next we define the convex corner K_a (associated to X, Y, \mathcal{J}) simply as the smallest downward closed set containing $A(K_r)$.

Definition 16. Let

$$K_a := \{ \mathbf{a} \in \mathbb{R}^{\mathcal{X}} : 0 \leq \mathbf{a} \leq A(\mathbf{r}) \text{ for some } \mathbf{r} \in K_r \}.$$

Proposition 17. The set K_a is a convex corner and $H_{K_a}(X) = \min_{K_a} \varphi_a = H_G(X|Y)$.

Proof. The key observation is that A_x is a concave function for each x , which follows immediately from the following claim: let $\alpha_1, \dots, \alpha_k \geq 0$ with $\alpha_1 + \dots + \alpha_k \leq 1$; then

$$f(t_1, \dots, t_k) := t_1^{\alpha_1} \dots t_k^{\alpha_k}$$

is a concave function in the positive orthant $\{(t_1, \dots, t_k) : t_1, \dots, t_k \geq 0\}$. Indeed, it is easy to see that the Hessian of f is given by

$$H_{i,j} = \begin{cases} \frac{\alpha_i \alpha_j}{t_i t_j} f(\mathbf{t}) & \text{if } i \neq j; \\ \frac{\alpha_i(\alpha_i - 1)}{t_i^2} f(\mathbf{t}) & \text{if } i = j. \end{cases}$$

Then for a vector $\mathbf{u} = (u_i)$ we have

$$\mathbf{u} H \mathbf{u}^\top / f(\mathbf{t}) = \left(\sum_i \frac{\alpha_i u_i}{t_i} \right)^2 - \sum_i \frac{\alpha_i u_i^2}{t_i^2} \leq 0$$

by the Cauchy–Schwarz inequality, proving that the Hessian is negative semidefinite.

Since each A_x is the sum of such functions, it is concave as well.

Now suppose that $\mathbf{a}, \mathbf{a}' \in K_a$. By definition, there exist $\mathbf{r}, \mathbf{r}' \in K_r$ such that $\mathbf{a} \leq A(\mathbf{r})$ and $\mathbf{a}' \leq A(\mathbf{r}')$. Then for any $t \in (0, 1)$ and for any x we have

$$t a_x + (1-t) a'_x \leq t A_x(\mathbf{r}) + (1-t) A_x(\mathbf{r}') \leq A_x(t\mathbf{r} + (1-t)\mathbf{r}'),$$

where the second inequality is due to the concavity of A_x . It follows that the convex combination

$$t\mathbf{a} + (1-t)\mathbf{a}' \leq A(\underbrace{t\mathbf{r} + (1-t)\mathbf{r}'}_{\in K_r})$$

also lies in K_a , proving the convexity of K_a .

Since K_r is compact and A is continuous, the image $A(K_r)$ is also compact, and hence so is K_a . Furthermore, if $\mathbf{r} \in K_r^* \supseteq \text{int}(K_r) \neq \emptyset$, then $A_x(\mathbf{r}) > 0$ for each x , so K_a has a nonempty interior.

Finally, to see that $H_{K_a}(X) = H_G(X|Y)$, it suffices to show that

$$\min_{K_a} \varphi_a = \min_{K_r} \varphi_r,$$

which follows immediately from $\varphi_r = \varphi_a \circ A$ and the monotonicity of φ_a : if $\mathbf{a} \leq A(\mathbf{r})$, then

$$\varphi_a(\mathbf{a}) \geq \varphi_a(A(\mathbf{r})) = \varphi_r(\mathbf{r})$$

with equality when $\mathbf{a} = A(\mathbf{r}) \in K_a$. \square

Remark 18. We make some comments regarding the a-problem.

- Note that K_a depends only on \mathcal{J} (or the graph) and the conditional distributions of $Y | X = x$ (but not on the distribution of X). In the unconditioned case K_a is the vertex packing polytope $\text{VP}(G)$ of the graph: the convex hull of the indicator functions of the independent sets. In general, K_a is not necessarily a polytope, it may be a more complicated convex set with ‘‘curvy’’ boundary. For an example, see Figure 2 in Section IV-D.
- It is easy to see that φ_a is a strictly convex function over $(0, 1]^{\mathcal{X}}$. Consequently, the a-problem always has a unique minimum point.
- Note that the dimension of the a-problem is usually much smaller than that of the q-problem or the r-problem. However, the domain is not a polytope in this case and the complexity of the a-problem is, in some sense, hidden in the definition of the domain.

C. Dual problem

Now we introduce another convex corner that will lead to a dual problem. To this end, for each j we define a function $\tau_j: [0, \infty)^{\mathcal{X}} \times [0, \infty)^{\mathcal{Y}} \rightarrow \mathbb{R}$: for $\mathbf{b} \in [0, \infty)^{\mathcal{X}}$ and $\mathbf{t} \in [0, \infty)^{\mathcal{Y}}$ we set

$$\tau_j(\mathbf{b}, \mathbf{t}) := \sum_{x \in j} p_x b_x \prod_y (t_y)^{p^{y|x}}.$$

Note that $\sum_y p^{y|x} = 1$, hence $\tau_j(\mathbf{b}, \mathbf{t})$ is homogeneous in \mathbf{t} : for any scalar $\lambda > 0$ we have $\tau_j(\mathbf{b}, \lambda \mathbf{t}) = \lambda \tau_j(\mathbf{b}, \mathbf{t})$.

Definition 19. Let

$$L_j := \left\{ \mathbf{b} \in [0, \infty)^{\mathcal{X}} : \forall \mathbf{t} \tau_j(\mathbf{b}, \mathbf{t}) \leq \sum_y p^y t_y \right\} = \left\{ \mathbf{b} \in [0, \infty)^{\mathcal{X}} : \tau_j(\mathbf{b}, \mathbf{t}) \leq 1 \text{ for all } \mathbf{t} \text{ with } \sum_y p^y t_y = 1 \right\}. \quad (15)$$

Finally, we define L as the intersection of all L_j :

$$L := \bigcap_j L_j.$$

To see that L is a convex corner, notice that for any given j and \mathbf{t} , the points \mathbf{b} for which $\tau_j(\mathbf{b}, \mathbf{t}) \leq 1$ form a downward closed polyhedron, and L is the intersection of such sets.

Remark 20. For graph entropy (i.e., the unconditioned case $|\mathcal{Y}| = 1$) it can be seen easily that

$$L = \left\{ \mathbf{b} : \forall j \sum_{x \in j} p_x b_x \leq 1 \right\} = (\mathbf{p}K_a)^*.$$

As the following lemma shows, the containment $L \subseteq (\mathbf{p}K_a)^*$ is true in general.

Lemma 21. For any $\mathbf{a} \in K_a$ and $\mathbf{b} \in L$ we have $\sum_x p_x a_x b_x \leq 1$. In other words, $L \subseteq (\mathbf{p}K_a)^*$.

Proof. We have $\mathbf{a} \leq A(\mathbf{r})$ for some $\mathbf{r} \in K_r$, thus

$$\sum_x p_x a_x b_x \leq \sum_j \sum_{x \in j} p_x b_x \prod_y (r_{j|y})^{p^{y|x}} \leq \sum_j \sum_y p^y r_{j|y} = \sum_y p^y \underbrace{\sum_j r_{j|y}}_{=1} = \sum_y p^y = 1,$$

where we used that $\mathbf{b} \in L_j$ for any given j , and hence $\tau_j(\mathbf{b}, \mathbf{t}) = \sum_{x \in j} p_x b_x \prod_y (t_y)^{p^{y|x}} \leq \sum_y p^y t_y$ holds for $t_y = r_{j|y}$. \square

Corollary 22. For any $\mathbf{a} \in K_a$ and $\mathbf{b} \in L$ we have

$$\varphi_a(\mathbf{a}) + \varphi_a(\mathbf{b}) \geq 0.$$

In other words,

$$\min_{K_a} \varphi_a + \min_L \varphi_a \geq 0. \quad (16)$$

Proof. Since $-\log$ is convex and monotone decreasing, by the above lemma we have

$$\varphi_a(\mathbf{a}) + \varphi_a(\mathbf{b}) = \varphi_a(\mathbf{ab}) = - \sum_x p_x \log(a_x b_x) \geq - \log \left(\sum_x p_x a_x b_x \right) \geq - \log(1) = 0.$$

□

We will shortly see that (16) actually holds with equality. In order to prove this, let us consider the set $L^{-1} = \{\mathbf{b}^{-1} : \mathbf{b} \in L\}$. Since L is convex and downward closed, it follows easily that L^{-1} is convex and upward closed (using the convexity of $t \mapsto 1/t$ for $t > 0$). The key observation is that K_a and L^{-1} always have a common point.

Theorem 23. *The intersection of the downward closed convex set K_a and the upward closed convex set L^{-1} is a single point \mathbf{a} , where φ_a takes its minimum over K_a and its maximum over L^{-1} . Then*

$$H_G(X|Y) = \varphi_a(\mathbf{a}) = \min_{K_a} \varphi_a = \max_{L^{-1}} \varphi_a = - \min_L \varphi_a.$$

Furthermore, K_a and L^{-1} are separated by a hyperplane with normal vector $\mathbf{pa}^{-1} = (p_x/a_x)_{x \in \mathcal{X}}$.

Before we present the proof, recall that the mappings Q and R “jump” between the q-problem and r-problem in a way that the function value decreases. In what follows we will focus on the r-problem and the corresponding stepping map

$$F_r := R \circ Q: K_r \rightarrow K_r.$$

Proposition 24. *Every minimum point of φ_r must be a fixed point of F_r .*

Proof. Combining Propositions 10 and 11 gives that for any $\mathbf{r} \in K_r^*$ we have

$$\varphi_r(F_r(\mathbf{r})) \leq \varphi_r(\mathbf{r})$$

with equality if and only if \mathbf{r} is a fixed point of F_r . In other words, if \mathbf{r} is not a fixed point, then we have strict inequality and hence $\varphi_r(\mathbf{r})$ cannot be the minimum. □

Proof of Theorem 23. Since $\varphi_a(\mathbf{b}^{-1}) = -\varphi_a(\mathbf{b})$, we have

$$\max_{L^{-1}} \varphi_a = - \min_L \varphi_a \stackrel{(16)}{\leq} \underbrace{\min_{K_a} \varphi_a}_{=H_G(X|Y)} = \varphi_a(\mathbf{a}),$$

where \mathbf{a} denotes the unique minimum point⁵ of φ_a over K_a . It remains to be shown that $\mathbf{a} \in L^{-1}$, implying the only missing inequality $\max_{L^{-1}} \varphi_a \geq \varphi_a(\mathbf{a})$ and confirming that $K_a \cap L^{-1} = \{\mathbf{a}\}$. (Note that the $\varphi_a(\mathbf{a}') > \varphi_a(\mathbf{a})$ for any $\mathbf{a}' \in K_a \setminus \{\mathbf{a}\}$, and hence $\mathbf{a}' \notin L^{-1}$.) Also, the gradient of φ_a at \mathbf{a} is $-\mathbf{pa}^{-1}$, so the hyperplane through \mathbf{a} that separates the convex sets K_a and L^{-1} must be the one with normal vector \mathbf{pa}^{-1} .

In order to prove that $\mathbf{a} \in L^{-1}$, let \mathbf{r} be such that $\mathbf{a} = A(\mathbf{r})$ so that $\varphi_r(\mathbf{r}) = \varphi_a(\mathbf{a})$, that is, \mathbf{r} minimizes φ_r over K_r . (Note that \mathbf{r} may not be unique.) By Proposition 24, \mathbf{r} is a fixed point of F_r . That is, for $\mathbf{q} := Q(\mathbf{r})$ we have $R(\mathbf{q}) = \mathbf{r}$.

For brevity, we write $\partial_{j|y}$ for the partial derivative w.r.t. the variable $r_{j|y}$, and $g_{j,x}$ for the product in the definition of A_x , that is:

$$g_{j,x} := \prod_y (r_{j|y})^{p^{y|x}} \text{ so that } A_x(\mathbf{r}) = \sum_{j \ni x} g_{j,x}.$$

If $r_{j|y} > 0$, then we have

$$\partial_{j|y} A_x(\mathbf{r}) = \begin{cases} 0 & \text{if } x \notin j; \\ \frac{p^{y|x}}{r_{j|y}} g_{j,x} & \text{if } x \in j; \end{cases}$$

and hence

$$\partial_{j|y} \varphi_r(\mathbf{r}) = - \sum_x p_x \frac{\partial_{j|y} A_x(\mathbf{r})}{A_x(\mathbf{r})} = - \sum_{x \in j} \frac{p_x p^{y|x} g_{j,x}}{r_{j|y} A_x(\mathbf{r})} = - \frac{p^y \sum_{x \in j} p_{x|y} q_{j|x}}{r_{j|y}} = - \frac{p^y R_{j|y}(\mathbf{q})}{r_{j|y}} = -p^y,$$

where we used that $R(\mathbf{q}) = \mathbf{r}$.

⁵Since K_a is compact and $\varphi_a: K_a \rightarrow [0, \infty]$ is continuous, its minimum is attained at some $\mathbf{a} \in K_a$. The minimum is finite, so we have $a_x > 0$ for each x . In that region φ_a is strictly convex, therefore \mathbf{a} is indeed unique.

Now we fix a $\hat{j} \in \mathcal{J}$ and a vector $\mathbf{t} = (t_y)$ with $t_y \geq 0$. Then we perturb \mathbf{r} in the coordinates $\hat{j}|y$ as follows: for a given $\varepsilon > 0$ we define the perturbed vector \mathbf{r}^ε as

$$r_{\hat{j}|y}^\varepsilon := \begin{cases} r_{\hat{j}|y} + \varepsilon t_y & \text{if } j = \hat{j}; \\ r_{\hat{j}|y} & \text{if } j \neq \hat{j}. \end{cases}$$

Then A_x does not change for $x \notin \hat{j}$. As for $x \in \hat{j}$, we claim that

$$A_x(\mathbf{r}^\varepsilon) - A_x(\mathbf{r}) \geq \varepsilon \prod_y (t_y)^{p^{y|x}}. \quad (17)$$

To see this, notice that the function $f: \mathbf{u} \mapsto \prod_y (u_y)^{p^{y|x}}$ is concave (as we have shown it in the proof of Proposition 17) and homogeneous (of degree 1) to conclude that

$$f(\mathbf{u} + \varepsilon \mathbf{t}) = 2f\left(\frac{\mathbf{u} + \varepsilon \mathbf{t}}{2}\right) \geq 2\frac{f(\mathbf{u}) + f(\varepsilon \mathbf{t})}{2} = f(\mathbf{u}) + \varepsilon f(\mathbf{t}),$$

which implies (17). It follows that

$$\lim_{\varepsilon \rightarrow 0^+} \frac{\varphi_r(\mathbf{r}^\varepsilon) - \varphi_r(\mathbf{r})}{\varepsilon} = \lim_{\varepsilon \rightarrow 0^+} \frac{\varphi_a(A(\mathbf{r}^\varepsilon)) - \varphi_a(A(\mathbf{r}))}{\varepsilon} \leq -\sum_{x \in \hat{j}} \frac{p_x}{a_x} \prod_y (t_y)^{p^{y|x}} = -\tau_{\hat{j}}(\mathbf{a}^{-1}, \mathbf{t}).$$

Note that $\mathbf{r}^\varepsilon \notin K_r$ anymore because moving in direction \mathbf{t} violates the linear constraints of K_r . So for each y we decrease other (positive) coordinates $r_{j|y}$ by a total of εt_y to get a point $\hat{\mathbf{r}}^\varepsilon$ in K_r . Since the partial derivative $\partial_{j|y} \varphi_r$ is $-p^y$ for such positive coordinates, it is easy to see that we get the following:

$$\lim_{\varepsilon \rightarrow 0^+} \frac{\varphi_r(\hat{\mathbf{r}}^\varepsilon) - \varphi_r(\mathbf{r})}{\varepsilon} = -\tau_{\hat{j}}(\mathbf{a}^{-1}, \mathbf{t}) + \sum_y p^y t_y. \quad (18)$$

Since $\hat{\mathbf{r}}^\varepsilon \in K_r$, we have $\varphi_r(\hat{\mathbf{r}}^\varepsilon) \geq \min_{K_r} \varphi_r = \varphi_r(\mathbf{r})$ for each ε . So the above limit must be nonnegative, that is,

$$\tau_{\hat{j}}(\mathbf{a}^{-1}, \mathbf{t}) \leq \sum_y p^y t_y.$$

This holds for any $\mathbf{t} \geq 0$, meaning that $\mathbf{a}^{-1} \in L_{\hat{j}}$. This can be done for any $\hat{j} \in \mathcal{J}$, implying $\mathbf{a}^{-1} \in L$. \square

D. The Orlitsky–Roche example

Orlitsky and Roche considered the following simple example, see [14, Examples 2&5]. Let $\mathcal{X} = \mathcal{Y} = \{1, 2, 3\}$ with the distribution

$$p_{x,y} = \begin{cases} 1/6 & \text{if } x \neq y; \\ 0 & \text{if } x = y. \end{cases}$$

Furthermore, let G be the graph on the vertex set \mathcal{X} containing a single edge $(1, 3)$ so that G has two maximal independent sets: $\{1, 2\}$ and $\{2, 3\}$. They showed that

$$H_G(X|Y) = -\frac{2}{3} \left(\frac{1}{4} \log \frac{1}{4} + \frac{3}{4} \log \frac{3}{4} \right) \approx 0.37489. \quad (19)$$

We will use this example to illustrate our results. We have

$$p_x = p^y = 1/3 (\forall x, y); \quad p_{x|y} = p^{y|x} = \begin{cases} 1/2 & \text{if } x \neq y; \\ 0 & \text{if } x = y. \end{cases}$$

We will use the notations $\alpha := \{1, 2\}$ and $\beta := \{2, 3\}$ for the independent sets so that $\mathcal{J} = \{\alpha, \beta\}$. It means that the r-problem has six non-negative variables with the following constraints:

$$r_{\alpha|1} + r_{\beta|1} = 1; \quad r_{\alpha|2} + r_{\beta|2} = 1; \quad r_{\alpha|3} + r_{\beta|3} = 1.$$

Then the mapping A is described by the following coordinate functions:

$$\begin{aligned} A_1(\mathbf{r}) &= \sqrt{r_{\alpha|2} r_{\alpha|3}} \\ A_2(\mathbf{r}) &= \sqrt{r_{\alpha|1} r_{\alpha|3}} + \sqrt{r_{\beta|1} r_{\beta|3}} \\ A_3(\mathbf{r}) &= \sqrt{r_{\beta|1} r_{\beta|2}} \end{aligned}$$

Next we describe the convex corner K_a associated to this example. Note that $A_1^2(\mathbf{r}) + A_3^2(\mathbf{r}) \leq r_{\alpha|2} + r_{\beta|2} = 1$, and $A_2(\mathbf{r}) \leq \sqrt{r_{\alpha|1} + r_{\beta|1}} \cdot \sqrt{r_{\alpha|3} + r_{\beta|3}} = 1$ by Cauchy–Schwarz. It follows that for any $\mathbf{a} \in K_a$ we have $a_1^2 + a_3^2 \leq 1$ and $a_2 \leq 1$. Now fix a_1, a_3 such that $a_1^2 + a_3^2 \leq 1$, and let us try to find the largest possible corresponding a_2 value:

$$\text{let } w := r_{\alpha|2}; \text{ then } r_{\beta|2} = 1 - w; r_{\alpha|3} = \frac{a_1^2}{w}; r_{\beta|1} = \frac{a_3^2}{1 - w}.$$

Therefore

$$a_2 = A_2(\mathbf{r}) = \sqrt{\frac{a_1^2}{w} \left(1 - \frac{a_3^2}{1 - w}\right)} + \sqrt{\frac{a_3^2}{1 - w} \left(1 - \frac{a_1^2}{w}\right)}.$$

We need to maximize this formula in the one free variable w . It is easy to see that when $a_1 + a_3 \leq 1$, the maximum is always 1, meaning that the boundary of K_a includes a triangle whose vertices are $(0, 1, 0)$; $(1, 1, 0)$ and $(0, 1, 1)$. When $a_1 + a_3 > 1$ and $a_1^2 + a_3^2 \leq 1$, we did not find a closed formula, but one can easily plot the maximum as a function of the parameters a_1 and a_3 ; see Figure 2. Note that when $a_1 = a_3$, the maximum is always taken at $w = 1/2$, so we get $a_2 = 2\sqrt{2}a_1\sqrt{1 - 2a_1^2}$ for the boundary of K_a in this cross section.

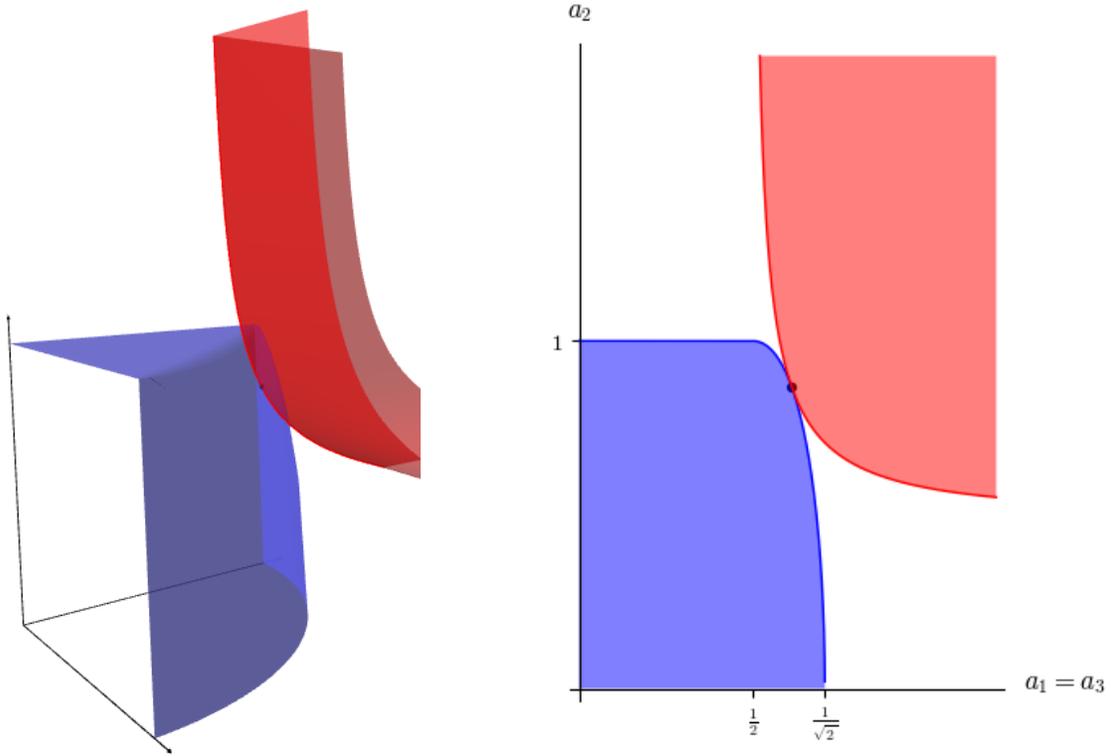


Fig. 2. On the left: plots of the boundaries of K_a (blue) and L^{-1} (red). On the right: the two-dimensional cross section corresponding to the plane $a_1 = a_3$. The black dot marks the unique intersection point, where φ_a takes its minimum over K_a and its maximum over L^{-1} .

As for the convex corner L corresponding to the dual problem, first we need to work out the formulas for τ_α and τ_β :

$$\begin{aligned} \tau_\alpha(\mathbf{b}, \mathbf{t}) &= \frac{1}{3} (b_1\sqrt{t_2t_3} + b_2\sqrt{t_1t_3}); \\ \tau_\beta(\mathbf{b}, \mathbf{t}) &= \frac{1}{3} (b_2\sqrt{t_1t_3} + b_3\sqrt{t_1t_2}). \end{aligned}$$

By Cauchy–Schwarz we have

$$\tau_\alpha(\mathbf{b}, \mathbf{t}) \leq \frac{1}{3} \sqrt{b_1^2 + b_2^2} \sqrt{(t_1 + t_2)t_3} \leq \frac{1}{3} \sqrt{b_1^2 + b_2^2} \frac{t_1 + t_2 + t_3}{2},$$

which shows by (15) that $\mathbf{b} \in L_\alpha$ provided that $b_1^2 + b_2^2 \leq 4$. It is also easy to see that $\mathbf{b} \notin L_\alpha$ if $b_1^2 + b_2^2 > 4$. Similar calculations show that $\mathbf{b} \in L_\beta$ if and only if $b_2^2 + b_3^2 \leq 4$. We conclude that

$$L = \{\mathbf{b} = (b_1, b_2, b_3) : b_1^2 + b_2^2 \leq 4 \text{ and } b_2^2 + b_3^2 \leq 4\}.$$

In Figure 2 we plotted (the boundary of) L^{-1} instead of L to illustrate the fact that K_a and L^{-1} intersect in a single point (marked by a black dot in the figure). This intersection point is where φ_a takes its minimum over K_a . The minimum points of the various problems are as follows.

The minimum of φ_q is attained at the following point \mathbf{q} :

$$\begin{aligned} q_{\alpha|1} &= 1; q_{\beta|1} = 0; \\ q_{\alpha|2} &= q_{\beta|2} = 1/2; \\ q_{\alpha|3} &= 0; q_{\beta|3} = 1. \end{aligned}$$

Then φ_r takes its minimum at the corresponding point $\mathbf{r} = R(\mathbf{q})$:

$$\begin{aligned} r_{\alpha|1} &= 1/4; r_{\beta|1} = 3/4; \\ r_{\alpha|2} &= 1/2; r_{\beta|2} = 1/2; \\ r_{\alpha|3} &= 3/4; r_{\beta|3} = 1/4. \end{aligned}$$

Lastly, φ_a takes its minimum at

$$\mathbf{a} = A(\mathbf{r}) = \left(\sqrt{3/8}, \sqrt{3/4}, \sqrt{3/8} \right).$$

Using our results, one can easily verify that this is the optimal point in K_a by checking that $\mathbf{a}^{-1} \in L$, which indeed holds as

$$\text{for } \mathbf{b} = \mathbf{a}^{-1} = \left(\sqrt{8/3}, \sqrt{4/3}, \sqrt{8/3} \right) \text{ we have } b_1^2 + b_2^2 = b_2^2 + b_3^2 = 4/3 + 8/3 = 4.$$

This confirms the value of $H_G(X|Y)$; see (19).

Finally, the table below shows the values $\varphi_r(\mathbf{r}^{(n)})$ of the iterative process started from a random point $\mathbf{q}^{(0)}$. We also included the error (i.e., the distance from the minimum) and our error bound based on the dual problem (see Theorem 27 in Section V).

n	value $\varphi_r(\mathbf{r}^{(n)})$	error $\varphi_r(\mathbf{r}^{(n)}) - H_G(X Y)$	error bound (see Thm 27) $\delta(A(\mathbf{r}^{(n)}))$
5	0.3749085763210158	$1.8 \cdot 10^{-5}$	$2.5 \cdot 10^{-3}$
10	0.3748904169016328	$3.2 \cdot 10^{-7}$	$3.2 \cdot 10^{-4}$
15	0.3748901019703158	$5.5 \cdot 10^{-9}$	$4.3 \cdot 10^{-5}$
20	0.3748900965089192	$9.6 \cdot 10^{-11}$	$5.6 \cdot 10^{-6}$
25	0.3748900964142102	$1.6 \cdot 10^{-12}$	$7.4 \cdot 10^{-7}$
30	0.3748900964125679	$2.9 \cdot 10^{-14}$	$9.8 \cdot 10^{-8}$
35	0.3748900964125393	$4.6 \cdot 10^{-16}$	$1.2 \cdot 10^{-8}$
$H_G(X Y)$	0.3748900964125389...		

E. Fractional chromatic number

Given a convex corner K , it is natural to ask what the maximum of its entropy function is. That is, by varying the distribution of X , what is the maximal possible $H_K(X)$ we can get for a fixed K ? In general one can say the following about this maximum entropy.

Lemma 25 (see Corollary 1.2.21 in [4]). *Let $K \subset \mathbb{R}^{\mathcal{X}}$ be an arbitrary convex corner. Then*

$$\max_X H_K(X) = \log \tau(K),$$

where $\tau(K)$ denotes the smallest $t \geq 1$ such that the constant $1/t$ vector lies in K . (Note that here X can be any random variable on \mathcal{X} : its support may be a proper subset of \mathcal{X} .)

The question arises: is there a special meaning of $\tau(K)$ in our setting? In the unconditioned case, that is, for the vertex packing polytope $K = \text{VP}(G)$, $\tau(K)$ is known to be equal to the fractional chromatic number of the graph [17, Lemma 4]. Is there a generalization of this result: does $\tau(K_a)$ have a nice graph theoretic meaning in the conditional setting?

Problem 26. Fix a graph G equipped with a distribution on \mathcal{Y} at each vertex x (described by $p_{y|x}$). Note that this determines the convex corner K_a . Is it possible to give a (graph theoretic) description of $\tau(K_a)$? This could lead to a notion generalizing the fractional chromatic number to *measure-labelled graphs*.

V. DISCUSSION OF THE ALGORITHM

As we have seen in the introduction, one may start at any point $\mathbf{q}^{(0)} \in \text{int}(K_q)$ and alternate in applying the mappings R and Q to get a sequence (4) with decreasing φ -values. In fact, Theorem 3 tells us that the values always converge to $\min \varphi(\mathbf{q}, \mathbf{r}) = H_G(X|Y)$.

In this section, we provide an error bound for the algorithm, then propose a tweak for improving the running time, and finally analyze the rate of convergence in the unconditioned case of graph entropy.

A. Error bound

How long should we run the iterations? A natural stopping rule is to terminate the algorithm at a step where the drop in the φ -value gets below some threshold. Is there a way to know how far we are from the actual minimum? Using the dual problem defined in Section IV-C, we can easily get an error bound for any given $\mathbf{r}^{(n)}$ we stop at.

Theorem 27. Let $\mathbf{r} \in K_r$ arbitrary and set $\mathbf{a} = A(\mathbf{r})$. For each j consider the following maximization problem:

$$1 + \delta_j(\mathbf{a}) := \max_{\mathbf{t}} \tau_j(\mathbf{a}^{-1}, \mathbf{t}) = \max_{\mathbf{t}} \sum_{x \in j} \frac{p_x}{a_x} \prod_y (t_y)^{p^{y|x}} \text{ under the constraints } t_y \geq 0; \sum_y p^y t_y = 1. \quad (20)$$

Then $\varphi_a(\mathbf{a}) = \varphi_r(\mathbf{r})$ is at most $\delta(\mathbf{a}) := \max_j \delta_j(\mathbf{a})$ away from the minimum. More precisely,

$$\varphi_r(\mathbf{r}) - H_G(X|Y) \leq \log(1 + \delta(\mathbf{a})) \leq \delta(\mathbf{a}).$$

In particular, \mathbf{a} (and hence \mathbf{r}) is optimal if and only if $\delta(\mathbf{a}) = 0$.

Remark 28. Note that each maximization is a convex optimization problem, whose dimension ($|\mathcal{Y}|$) is small compared to that of the \mathbf{r} -problem ($|\mathcal{Y}| \cdot |\mathcal{J}|$) so we can solve them with high precision relatively fast.

Proof. By definition, $\mathbf{b} := (1 + \delta(\mathbf{a}))^{-1} \mathbf{a}^{-1}$ lies in L . Therefore

$$H_G(X|Y) = -\min_L \varphi_a \leq -\varphi_a(\mathbf{b}) = \log(1 + \delta(\mathbf{a})) + \varphi_a(\mathbf{a}).$$

□

The table at the end of Section IV-D compares this error bound to the true error for the Orlitsky–Roche example.

B. A tweak: deleting redundant sets

The running time of the algorithm depends on two things: the time required to perform a single step and the number of steps required to get within the desired distance of the minimum. With one small tweak we can achieve significant gains for both at the same time.

First of all, note that at each step the algorithm performs $\mathcal{O}(|\mathcal{X}| \cdot |\mathcal{Y}| \cdot |\mathcal{J}|)$ operations when computing $\mathbf{r}^{(n)} = R(\mathbf{q}^{(n)})$ and $\mathbf{q}^{(n+1)} = Q(\mathbf{r}^{(n)})$.

In examples there are often a large number of (independent sets) j that are actually not “used” at the optimal \mathbf{q} and \mathbf{r} in the sense that $q_{j|x} = 0$ and $r_{j|y} = 0$ for all x, y . For any such j , these variables will converge to 0 through the iterations. To speed things up, we may want to detect such *redundant* sets j early and set the corresponding variables to 0. Note that these variables remain to be 0 from this point on, so we may remove such a j from \mathcal{J} and proceed with the iterations using a smaller set \mathcal{J} . This immediately reduces the computational complexity for each subsequent step. Moreover, it typically results in a better rate of convergence as well: without redundant sets, the error usually decays at a faster rate. Consequently, this version of the algorithm often requires considerably fewer steps to reach the desired precision. (This phenomenon will be illustrated for graph entropy both by an example and an analysis.)

However, when the algorithm terminates and outputs an (approximate) minimum point for some subsystem \mathcal{J} of the original \mathcal{J}_{or} , we should justify that all deletions we made along the way were indeed necessary. So we take the corresponding point \mathbf{a} and perform our optimality check/error bound calculations: we compute $\delta_j(\mathbf{a})$ as in (20). For each $j \in \mathcal{J}$ we should get a negative number or a very small positive number, confirming that we are indeed close to the minimum point of the problem corresponding to the subsystem \mathcal{J} . If $\delta_j(\mathbf{a}) \leq 0$ for all deleted sets $j \in \mathcal{J}_{\text{or}} \setminus \mathcal{J}$, it means that we cannot do better even if we used the deleted sets. If, on the other hand, $\delta_j(\mathbf{a}) > 0$ for some of the deleted sets j , then we should “re-activate” them (i.e., add them back to \mathcal{J}).

So we propose the following **tweaked version of the iterative process**.

- Set $r_{j|y}^{(0)} = 1/|\mathcal{J}|$ for each j and y . Note that $\mathbf{r}^{(0)} \in K_r$.
- Set $\varepsilon_{\text{act}} = 10^{-3} |\mathcal{Y}|/|\mathcal{J}|$.
- At step n :
 - compute $\mathbf{r}^{(n-1)} \xrightarrow{Q} \mathbf{q}^{(n)} \xrightarrow{R} \mathbf{r}^{(n)}$;
 - for any j with $\sum_y r_{j|y}^{(n)} < \varepsilon_{\text{act}}$, remove j from \mathcal{J} and delete the corresponding variables $r_{j|y}^{(n)}$ for all y ;
 - for each y , re-normalize the remaining variables $r_{j|y}^{(n)}$, $j \in \mathcal{J}$ such that

$$\text{the constraint } \sum_j r_{j|y}^{(n)} = 1 \text{ is satisfied again.}$$

- Compute the value $\varphi_r(\mathbf{r}^{(n)})$ after every 10 steps, and terminate the iterations when this value, compared to the

previous one, decreases by less than some small $\varepsilon_{\text{prec}}$ (say, 10^{-15}).

- Set $\mathbf{a} = A(\mathbf{r}^{(n)})$.
- Compute $\delta_j(\mathbf{a})$ as in (20) for all $j \in \mathcal{J}$ as well as for all previously deleted sets j .
- If $\delta_j(\mathbf{a}) \leq 0$ for each deleted j , then return $\varphi_{\mathbf{a}}(\mathbf{a})$ with the error bound $\max_{j \in \mathcal{J}} \delta_j(\mathbf{a})$.
- Otherwise, for each deleted $j \in \mathcal{J}_{\text{or}} \setminus \mathcal{J}$ with $\delta_j(\mathbf{a}) > 0$, add j back to \mathcal{J} and create the corresponding variables $r_{j|y}^{(n)}$ for each y , setting them to some small positive values.^a Then re-normalize as before so that $\mathbf{r}^{(n)} \in K_r$ holds again. Finally, restart the iterations, this time with no set-deletions.

^aAny values work but the following choice should guarantee that we get a smaller φ_r -value right after restart: set $r_{j|y}^{(n)} := \varepsilon t_y$ for a sufficiently small ε , where \mathbf{t} denotes the vector at which (20) takes its maximum.

Normally, we set ε_{act} to be fairly small so that it is extremely unlikely that we unjustifiably delete a set j , and the check at the end should (essentially always) confirm this.

Our implementation in Python is available on GitHub [11].

C. An example

The next example shows how detecting redundant sets can speed the convergence up.

Example 29. Let G be the dodecahedral graph: a 3-regular graph with 20 vertices and 30 edges; see Figure 3. It has 295 maximal independent sets. For a uniform X we have

$$H_G(X) = \log \frac{5}{2}.$$

This can be seen easily using that $|j| \leq 8$ for each $j \in \mathcal{J}$ and that one can find five independent sets j_1, \dots, j_5 such that each vertex is contained in exactly two of them (and hence $(2/5)\mathbf{e} \in K_a = \text{VP}(G)$, where \mathbf{e} is the all-ones vector).

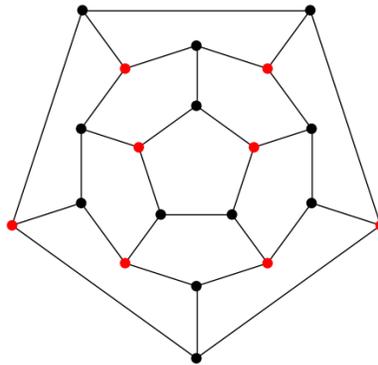
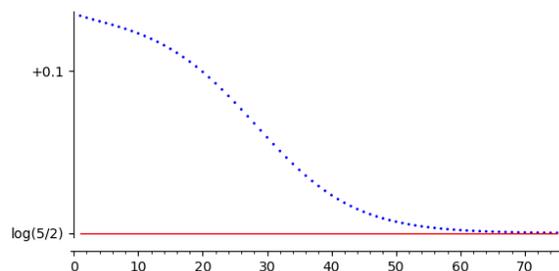
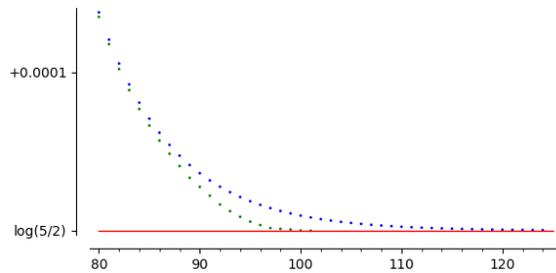


Fig. 3. The graph of the dodecahedron. The red vertices form an independent set of size 8. By “rotation” one can get five independent sets in a way that each vertex is contained in exactly two of them.

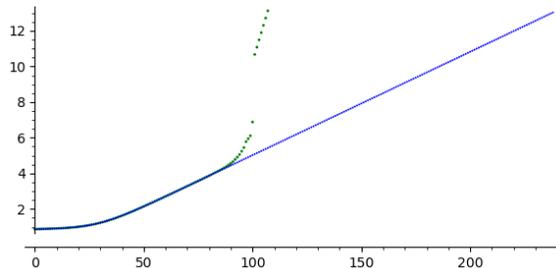
Starting from a random point $\mathbf{q}^{(0)}$, the blue dots below show the value $\varphi_r(\mathbf{r}^{(n)})$ for each iteration $n = 1, \dots, 75$.



For comparison, we run the process from the same starting point, but this time deleting a set j if $r_j^{(n)}$ gets below $\varepsilon_{\text{act}} := 2^{-20} \approx 10^{-6}$. Up to $n = 75$ only 23 sets were deleted and there was little difference in the value compared to the plot above. Afterwards the deletion rate accelerated and by step $n = 101$ all but the five independent sets of size 8 were deleted. The figure below compares the values in the two cases after step 80. (We used blue dots for the original process with no set-deletions and green dots for the one with set-deletions.)



The original algorithm (blue dots) needed 239 steps to get within distance 10^{-13} of the true minimum $\log(5/2)$, while the refined process (green dots) reached this threshold after only 108 iterations. We plotted the distance to the minimum in a logarithmic scale below: the horizontal axis shows the number of steps, while the vertical axis shows $-\log_{10}$ of the distance (i.e., the number of precise decimal digits essentially).



We see that there is a considerable leap in precision at the point when all 290 “redundant” independent sets have been deleted. Both versions eventually settle into a phase where the precision (“number of precise digits”) grows at a linear rate. The tweaked version clearly exhibits a faster rate. In fact, the analysis in the next section will reveal that this faster rate is $2 \cdot \lg(8/5) \approx 0.408$ compared to the rate $\lg(8/7) \approx 0.058$ of the original version.

D. Rate of convergence for graph entropy

As we have seen in the example of the previous section, the precision of the iterative algorithm appears to grow at some linear rate (for steps $n \geq n_0$). The following analysis confirms this observation and explains how one can determine this (“eventual”) rate in the unconditioned setting (i.e., graph entropy). Rigorous proofs would make the analysis undesirably long and technical so in this section we settle for only sketching the arguments.

Formulas for graph entropy: In the special case of graph entropy (i.e., when $|\mathcal{Y}| = 1$ so there is only one y) the formulas simplify considerably. First of all, we have $p^y = p^{y|x} = 1$ and $p_{x,y} = p_{x|y} = p_x$, and we may omit y in the indices. So \mathbf{r} now denotes a point $(r_j)_{j \in \mathcal{J}}$ in the set

$$K_r = \left\{ \mathbf{r} = (r_j) : r_j \geq 0; \sum_j r_j = 1 \right\} \subset \mathbb{R}^{\mathcal{J}}.$$

Furthermore, we have the following simple formulas:

$$\begin{aligned} R_j(\mathbf{q}) &= \sum_{x \in j} p_x q_{j|x}; \\ A_x(\mathbf{r}) &= \sum_{j \ni x} r_j; \\ Q_{j|x}(\mathbf{r}) &= \begin{cases} 0 & \text{if } x \notin j; \\ r_j / A_x(\mathbf{r}) & \text{if } x \in j; \end{cases} \\ \varphi_r(\mathbf{r}) &= - \sum_x p_x \log A_x(\mathbf{r}) = - \sum_x p_x \log \sum_{j \ni x} r_j. \end{aligned}$$

So A is simply a linear $\mathbb{R}^{\mathcal{J}} \rightarrow \mathbb{R}^{\mathcal{X}}$ map corresponding to the following matrix $M \in \mathbb{R}^{\mathcal{J} \times \mathcal{X}}$:

$$M_{x,j} := \begin{cases} 1 & \text{if } x \in j; \\ 0 & \text{if } x \notin j. \end{cases}$$

That is, the columns of M are the indicators functions of the sets j , and we have $A(\mathbf{r}) = M\mathbf{r}$.

As for the stepping map $F_r: \mathbf{r} \mapsto \mathbf{r}'$ for the r -problem, we have

$$r'_j = \left(\underbrace{\sum_{x \in j} \frac{p_x}{A_x(\mathbf{r})}}_{\Delta_j(\mathbf{r})} \right) r_j. \quad (21)$$

It follows that if \mathbf{r} is a fixed point of F_r (i.e., $r'_j = r_j$ for each j), if and only if $\Delta_j(\mathbf{r}) = 1$ for any j with $r_j > 0$.

It is worth mentioning that if ∂_j denotes the partial derivative w.r.t. the variable r_j , then we have

$$\partial_j \varphi_r(\mathbf{r}) = - \sum_x p_x \frac{\partial_j A_x(\mathbf{r})}{A_x(\mathbf{r})} = - \sum_{x \in j} \frac{p_x}{A_x(\mathbf{r})} = -\Delta_j(\mathbf{r}).$$

So what the stepping map $F_r(\mathbf{r})$ does in this unconditioned setting is simply multiply \mathbf{r} (coordinate-wise) by the negative of the gradient $\nabla \varphi_r(\mathbf{r})$.

Case of no redundant sets: We start our analysis with the case when each j is “used” ($r_j > 0$) at the minimum point \mathbf{r} of φ_r . This is always the case in the tweaked version of the algorithm which ensures that all redundant sets are eventually deleted. Note that $\Delta_j(\mathbf{r}) = 1$ for all j in this case, and hence the gradient $\nabla \varphi_r(\mathbf{r}) = -\mathbf{e}$ for the all-ones vector \mathbf{e} .

For a vector \mathbf{v} we will use the notation $D(\mathbf{v})$ for the corresponding diagonal matrix. In particular, $D(\mathbf{r}) \in \mathbb{R}^{\mathcal{J} \times \mathcal{J}}$ is the diagonal matrix with entries r_j , while $D(\mathbf{p}\mathbf{a}^{-2}) \in \mathbb{R}^{\mathcal{X} \times \mathcal{X}}$ is the diagonal matrix with entries p_x/a_x^2 .

Lemma 30. *Assume that \mathbf{r} is a minimum point of φ_r and that each $r_j > 0$. Set $\mathbf{a} := A(\mathbf{r}) = M\mathbf{r}$ so that \mathbf{a} is the minimum point of φ_a . Let*

$$N := D(\mathbf{r})M^\top D(\mathbf{p}\mathbf{a}^{-2})M.$$

Then $N \in \mathbb{R}^{\mathcal{J} \times \mathcal{J}}$ is a square matrix with nonnegative entries and with the following properties:

- in each column the sum of the entries is 1 (and hence 1 is an eigenvalue);
- N is diagonalizable with eigenvalues in $[0, 1]$;
- $\ker N = \ker M$.

The proof of the lemma can be found at the end of the section.

Claim. *The rate of convergence is governed by the smallest nonzero eigenvalue λ_{\min} of N :*

$$\varphi_r(\mathbf{r}^{(n)}) = \min_{K_r} \varphi_r + O((1 - \lambda_{\min})^{2n}). \quad (22)$$

So the rate of growth for the precision is $-2 \lg(1 - \lambda_{\min})$.

Example. The dodecahedral graph of Example 29 has five independent sets of size 8. Note that their pairwise intersections are of size 2. Independent sets of smaller size are all redundant so let \mathcal{J} be the set of these five sets. Then we have $r_j = 1/5$ for all j and $a_x = 2/5$ for all x . It follows that each diagonal entry of N is equal to $1/2$, while all other entries are equal to $1/8$. Therefore, the eigenvalues (with multiplicity) are 1; $3/8$; $3/8$; $3/8$; $3/8$. So $\lambda_{\min} = 3/8$ and we get that the rate of growth for the precision is $-2 \lg(5/8) + o(1)$, which is consistent with our numerical findings presented earlier.

Now we will sketch the proof of the claim. For the sake of simplicity we assume that $\ker M = \{0\}$. In this case \mathbf{r} is the unique minimum point of φ_r and hence $\mathbf{r}^{(n)} \rightarrow \mathbf{r}$ as $n \rightarrow \infty$. So difference vector $\boldsymbol{\varrho}^{(n)} := \mathbf{r}^{(n)} - \mathbf{r}$ converges to 0. Set $\boldsymbol{\alpha}^{(n)} := M\boldsymbol{\varrho}^{(n)}$. Then $\|\boldsymbol{\alpha}^{(n)}\| = O(\|\boldsymbol{\varrho}^{(n)}\|)$ converges to 0 as well.

With these notations, we compute the coordinates of the next point $\mathbf{r}^{(n+1)} = F_r(\mathbf{r}^{(n)})$ of our sequence:

$$r_j^{(n+1)} = \left(\sum_{x \in j} \frac{p_x}{a_x + \alpha_x^{(n)}} \right) (r_j + \varrho_j^{(n)}) = \left(\underbrace{\sum_{x \in j} \frac{p_x}{a_x}}_{=1} - \sum_{x \in j} \frac{p_x \alpha_x^{(n)}}{a_x^2} + \sum_{x \in j} \frac{p_x (\alpha_x^{(n)})^2}{a_x^2 (a_x + \alpha_x^{(n)})} \right) (r_j + \varrho_j^{(n)}).$$

It follows that

$$\varrho_j^{(n+1)} = \varrho_j^{(n)} - r_j \sum_{x \in j} \frac{p_x}{a_x^2} \alpha_x^{(n)} + O(\|\boldsymbol{\varrho}^{(n)}\|^2).$$

Since $\boldsymbol{\alpha}^{(n)} = M\boldsymbol{\varrho}^{(n)}$, we conclude that

$$\boldsymbol{\varrho}^{(n+1)} = (I - N)\boldsymbol{\varrho}^{(n)} + O(\|\boldsymbol{\varrho}^{(n)}\|^2). \quad (23)$$

Recall that λ_{\min} is the smallest nonzero eigenvalue of N . Under our assumption $\ker N = \ker M = \{0\}$, so 0 is not an eigenvalue now, meaning that the largest eigenvalue of the diagonalizable matrix $I - N$ is $1 - \lambda_{\min}$. Then it is not hard to deduce from (23) that

$$\|\boldsymbol{\varrho}^{(n)}\| = O((1 - \lambda_{\min})^n).$$

As for the φ_r -value,

$$\varphi_r(\mathbf{r}^{(n)}) = \varphi_r(\mathbf{r} + \boldsymbol{\varrho}^{(n)}) = \varphi_r(\mathbf{r}) + \underbrace{\nabla \varphi_r(\mathbf{r}) \cdot \boldsymbol{\varrho}^{(n)}}_{=-\mathbf{e}} + O(\|\boldsymbol{\varrho}^{(n)}\|^2),$$

where the dot product $\mathbf{e} \cdot \boldsymbol{\varrho}^{(n)}$, which is simply the sum of the coordinates of $\boldsymbol{\varrho}^{(n)}$, is equal to 0 because this sum is 1 both for $\mathbf{r} \in K_r$ and for $\mathbf{r}^{(n)} \in K_r$. Then (22) clearly follows.

In fact, heuristically, $\boldsymbol{\varrho}^{(n+1)} \approx (I - N)\boldsymbol{\varrho}^{(n)}$ means that if we write $\boldsymbol{\varrho}^{(n)}$ in an eigenbasis, then the parts corresponding to smaller eigenvalues will become negligible and $\boldsymbol{\varrho}^{(n)}$ will be close to an eigenvector with the maximal eigenvalue $1 - \lambda_{\min}$, and hence $\boldsymbol{\varrho}^{(n+1)} \approx (1 - \lambda_{\min})\boldsymbol{\varrho}^{(n)}$ for large n (at least for typical starting points).

When $\ker M = \ker N$ has positive dimension, 1 is an eigenvalue of $I - N$ and it seems that we do not necessarily have exponential decay. Note, however, that vectors from $\ker M$ do not make a difference from the point of view of φ -value because for any $\mathbf{v} \in \ker M$ we have $A(\mathbf{r}' + \mathbf{v}) = A(\mathbf{r}')$, and hence $\varphi_r(\mathbf{r}' + \mathbf{v}) = \varphi_r(\mathbf{r}')$.

We close this section by proving the required properties of N .

Proof of Lemma 30. Since \mathbf{r} is a minimum point of φ_r , by Proposition 24 \mathbf{r} is a fixed point of F_r , and hence $\Delta_j(\mathbf{r}) = 1$ for each j . This means that $M^\top(\mathbf{p}\mathbf{a}^{-1})$ is the all-ones vector \mathbf{e} .

We also have $M\mathbf{r} = A(\mathbf{r}) = \mathbf{a}$. Then

$$N^\top \mathbf{e} = M^\top D(\mathbf{p}\mathbf{a}^{-2})M \underbrace{D(\mathbf{r})\mathbf{e}}_{=\mathbf{r}} = M^\top D(\mathbf{p}\mathbf{a}^{-2}) \underbrace{M\mathbf{r}}_{=\mathbf{a}} = M^\top \underbrace{D(\mathbf{p}\mathbf{a}^{-2})\mathbf{a}}_{=\mathbf{p}\mathbf{a}^{-1}} = M^\top(\mathbf{p}\mathbf{a}^{-1}) = \mathbf{e},$$

confirming that 1 is an eigenvalue and that each column sum of N is 1. Since all entries are nonnegative, it follows that N is a (left) stochastic matrix, and hence $|\lambda| \leq 1$ for each eigenvalue λ .

Furthermore, N is similar to a positive semidefinite matrix:

$$D(\mathbf{r}^{-1/2})ND(\mathbf{r}^{1/2}) = D(\mathbf{r}^{1/2})M^\top D(\mathbf{p}\mathbf{a}^{-2})MD(\mathbf{r}^{1/2}),$$

so N is diagonalizable with nonnegative eigenvalues.

Finally, let $B = D(\mathbf{p}^{1/2}\mathbf{a}^{-1})M$. Then

$$\ker M = \ker B = \ker(B^\top B) = \ker(M^\top D(\mathbf{p}\mathbf{a}^{-2})M) = \ker N. \quad \square$$

Convergence for a redundant set: If $r_j = 0$ for a given j at the limiting point $\mathbf{r} = \lim_{n \rightarrow \infty} \mathbf{r}^{(n)}$, then we must have $\Delta_j(\mathbf{r}) < 1$. We then eventually see an exponential decay in the j -coordinate:

$$r_j^{(n+1)} = (\Delta_j(\mathbf{r}) + o(1))r_j^{(n)}.$$

Then the growth rate for the precision of $\varphi_r(\mathbf{r}^{(n)})$ is at most $-\log \Delta$, where Δ denotes the largest value of $\Delta_j(\mathbf{r})$ among all redundant sets j . For the dodecahedral graph we have $\Delta = 7/8$, which is consistent with our previous numerical findings.

VI. CONCLUSION

The optimal rate of lossless functional compression with side information at the receiver can be characterized by conditional graph entropy. However, little can be found in the literature about this entropy notion. So we set out to study conditional graph entropy in more detail. Our starting point was the original formula (2) which can also be formulated as the q -problem. Our first step was the discovery of the r -problem and the stepping maps Q, R between the two problems. This interaction was reminiscent of the alternating optimization in the EM algorithm, which made us realize that there might be an underlying alternating minimization problem. This, in turn, helped us to analyze the iterative algorithm because we could turn to the general theory of Csiszár and Tusnády: we verified that the 3-point, 4-point, and 5-point properties hold in our setting, and showed that the iterations always converge to the minimum. Our theoretical results lead to a practical algorithm for computing conditional graph entropy that also comes with an error bound based on a dual problem.

Alternating optimization has a vast and growing literature. The fact that (conditional) graph entropy is part of this family of problems will hopefully inspire future research in the area.

ACKNOWLEDGMENTS

We are grateful to two anonymous reviewers for their numerous valuable remarks and suggestions that helped us tremendously in improving the paper.

REFERENCES

- [1] N. Alon and A. Orlitsky. Source coding and graph entropies. *IEEE Transactions on Information Theory*, 42(5):1329–1339, 1996.
- [2] S. Arimoto. An algorithm for computing the capacity of arbitrary discrete memoryless channels. *IEEE Transactions on Information Theory*, 18(1):14–20, 1972.
- [3] R. Blahut. Computation of channel capacity and rate-distortion functions. *IEEE Transactions on Information Theory*, 18(4):460–473, 1972.
- [4] Gareth Boreland. *Information theoretic parameters for graphs and operator systems*. PhD thesis, Queen’s University Belfast, 2020.
- [5] Imre Csiszár, János Körner, László Lovász, Katalin Marton, and Gábor Simonyi. Entropy splitting for antiblocking corners and perfect graphs. *Combinatorica*, 10(1):27–40, 1990.
- [6] Imre Csiszár and Gábor Tusnády. Information geometry and alternating minimization procedures. *Statistics and Decisions*, Supp. 1:205–237, 1984.
- [7] Vishal Doshi, Devavrat Shah, Muriel Médard, and Michelle Effros. Functional compression through graph coloring. *IEEE Transactions on Information Theory*, 56(8):3901–3917, 2010.
- [8] F. Dupuis, W. Yu, and F.M.J. Willems. Blahut-arimoto algorithms for computing channel capacity and rate-distortion with side information. In *International Symposium on Information Theory, 2004. ISIT 2004. Proceedings.*, pages 179–, 2004.
- [9] S.I. Gel’fand and M.S. Pinsker. Coding for channel with random parameters. *Probl. Contr. and Inf. Theory*, 1980.
- [10] M. Grötschel, L. Lovász, and A. Schrijver. Relaxations of vertex packing. *J. Combin. Theory Ser. B*, 40(3):330–343, 1986.
- [11] Viktor Harangi. Graph entropy program code. <https://github.com/harangi/graphentropy>, 2023.
- [12] János Körner. Coding of an information source having ambiguous alphabet and the entropy of graphs. In *6th Prague conference on information theory*, pages 411–425, 1973.
- [13] László Lovász. On the Shannon capacity of a graph. *IEEE Trans. Inform. Theory*, 25(1):1–7, 1979.
- [14] Alon Orlitsky and James R. Roche. Coding for computing. *IEEE Transactions on Information Theory*, 47:903–917, 1998.
- [15] Claude E. Shannon and Warren Weaver. *The Mathematical Theory of Communication*. University of Illinois Press, Urbana and Chicago, 1949.
- [16] Gábor Simonyi. Graph entropy: A survey. In William Cook, László Lovász, and Paul Seymour, editors, *Combinatorial Optimization*, volume 20 of *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, pages 399–441, 1993.
- [17] Gábor Simonyi. Perfect graphs and graph entropy. An updated survey. In Jorge Ramirez-Alfonsín and Bruce Reed, editors, *Perfect Graphs*, pages 293–328. John Wiley and Sons, 2001.
- [18] Pascal O. Vontobel, Aleksandar Kavcic, Dieter M. Arnold, and Hans-Andrea Loeliger. A generalization of the Blahut–Arimoto algorithm to finite-state channels. *IEEE Transactions on Information Theory*, 54(5):1887–1918, 2008.
- [19] Péter Vrana. Probabilistic refinement of the asymptotic spectrum of graphs. *Combinatorica*, 41(6):873–904, 2021.
- [20] H. Witsenhausen. The zero-error side information problem and chromatic numbers (corresp.). *IEEE Transactions on Information Theory*, 22(5):592–593, 1976.