

# Information-Theoretic Analysis of MAXCUT Algorithms

An Bian, Alexey Gronskiy and Joachim M. Buhmann  
Department of Computer Science, ETH Zurich  
{ybian, alexeygr, jbuhmann}@inf.ethz.ch

July 12, 2016

## Abstract

NP-hard combinatorial optimization algorithms are often characterized by their approximation ratios. In real world applications, the resilience of algorithms to input fluctuations and to modeling errors pose important robustness requirements. Motivated by this problem, we propose an information-theoretic algorithmic regularization and validation strategy based on *posterior agreement*, and further theoretically justify it by presenting the “coding by posterior” framework. The strategy regularizes algorithms and ranks them according to the informativeness of their output given noisy input. To illustrate this strategy, we develop methods to evaluate the posterior distribution of the Goemans-Williamson’s MAXCUT algorithm using semidefinite programming relaxation (MAXCUT-SDP by [Goemans and Williamson \(1995\)](#)). Experimental comparison with representative greedy MAXCUT algorithms shows that MAXCUT-SDP with the best known approximation ratio generalizes worse than greedy MAXCUT algorithms under high noise level.

## 1. Introduction

Algorithms are usually characterized by time and space complexity in the worst-case setting, and for specific instance distributions, in the average case. The robustness of an algorithm to input fluctuations is rarely investigated although such a property might often be indispensable in applications. Taking the MAXCUT problem for example, in practice, instead of having the graph  $G$  as input to recover the maximal cut, one usually only have access to multiple noisy observations of the graph  $G$ . Assuming for simplicity, there are two noisy observations of the underlying “master” graph  $G$ :  $G'$  and  $G''$ , and we want to recover the maximal cut with respect to  $G$ . We call this setting the “two-instance” scenario, and the ability of an algorithm to recover the true solutions given only noisy observations is closely related to the robustness/informativeness of the algorithm. According to [Buhmann \(2010\)](#), the informativeness of algorithms sensitively depends on the input distributions, which controls the precision of the output given the input uncertainty.

For an algorithm  $\mathcal{A}$  to survive in this two-instance scenario, e.g., the MAXCUT-SDP ([Goemans and Williamson, 1995](#)), we propose a general information-theoretic regularization and validation strategy, which is based on a provable analogue of information content for algorithms. Classical algorithms usually search for a unique or a randomized solution in the hypothesis class. Input noise often renders such algorithmic solutions highly unstable. Therefore, we require an algorithm to return a posterior distribution of solutions given the noisy input. Such a posterior should concentrate on few solutions but the posterior must be stable for equally likely inputs. We interpret this tradeoff between precise localization in the hypothesis class and stability of posteriors as the *generalization* property of an algorithm. Under this strategy, an algorithm should stop early to recover the stable solutions (posterior distribution of solutions).

It is well-known that when training machine learning models, e.g., training neural network using the stochastic gradient descent algorithm, one should stop the algorithm early to recover generalizable models, which is called the “early-stopping” strategy ([Caruana et al., 2001](#)). With empirical success, few theory has been proposed for this well-utilized strategy. By simple analogue between the generalizable solutions and machine learning models, this work also gives an information-theoretic verification of this “early-stopping” strategy.

Though we use MAXCUT algorithms as an illustrating example in this work, it is noteworthy that the strategy applies generally to any algorithms in the two-instance scenario. Previously, generalization of algorithms as a selection principle has been demonstrated for minimum spanning tree algorithms ([Gronskiy and Buhmann, 2014](#)) and data clustering ([Buhmann, 2010](#)).

## 1.1 MaxCut as an exemplary problem

Given an undirected graph  $G = (V, E)$  with a vertex set  $V = \{1, \dots, n\}$  and an edge set  $E$  with non-negative weights  $w_{ij}, \forall (i, j) \in E$ , MAXCUT determines a cut  $c := (S, \bar{S})$  of the vertex set  $V = S \cup \bar{S}, S \cap \bar{S} = \emptyset$  such that the cut value  $\text{cut}(c) := \sum_{i \in S} \sum_{j \in \bar{S}} w_{ij}$  is maximal. MAXCUT is a prototypical case of an unconstrained submodular maximization problem, and it is utilized in various applications, such as semisupervised learning (Wang et al., 2013) and social networks (Agrawal et al., 2003). Bian et al. (2015) has investigated the robustness of *greedy* MAXCUT algorithms based on the ‘‘approximation set coding’’ framework (Buhmann, 2010; Gronskiy and Buhmann, 2014), the methodology used there can not handle the continuous MAXCUT algorithms, e.g., the MAXCUT-SDP (Goemans and Williamson, 1995). This work aims to analyze the generalization performance of non-greedy, continuous MAXCUT algorithm based on strategy derived from a general framework called ‘‘coding by posterior’’.

## 1.2 Algorithm analysis by algorithmic information content

We investigate the generalization ability of an algorithm  $\mathcal{A}$  under the *two-instance scenario*. Assume there are two noisy instances  $G', G''$  drawn from the same underlying distribution  $\mathcal{D}$ . The algorithm then calculates a sequence of posteriors  $\{\mathbb{P}_t^{\mathcal{A}}(c|G')\}, \{\mathbb{P}_t^{\mathcal{A}}(c|G'')\}$  as a function of time  $t$ ,  $c$  is the solution in the hypothesis/solution space  $\mathcal{C}$ . The *posterior agreement* is defined to measure the overlap between the two posteriors at time  $t$ ,

$$k_t^{\mathcal{A}}(G', G'') := \sum_{c \in \mathcal{C}} \mathbb{P}_t^{\mathcal{A}}(c|G') \mathbb{P}_t^{\mathcal{A}}(c|G''). \quad (\text{posterior agreement}) \quad (1)$$

We define the *information content* of an algorithm  $\mathcal{A}$  as the maximal *temporal information content*  $I_t^{\mathcal{A}}$  at time  $t$ :

$$I^{\mathcal{A}}(G'; G'') := \max_t I_t^{\mathcal{A}}(G'; G'') = \max_t \mathbb{E}_{G', G''} [\log(|\mathcal{C}| k_t^{\mathcal{A}}(G', G''))]. \quad (2)$$

It generalizes the algorithmic information content in Gronskiy and Buhmann (2014).  $I_t^{\mathcal{A}}$  measures how much information is extracted by  $\mathcal{A}$  at time  $t$  from the input distribution that is relevant to the output distribution, thus reflecting the generalization ability. It naturally suggests the following algorithmic regularization and validation strategy based on posterior agreement:

- *Regularize* an algorithm  $\mathcal{A}$  by stopping it at the optimal time, which is defined as  $t^* = \arg \max_t \mathbb{E}_{G', G''} [k_t^{\mathcal{A}}(G', G'')]$ . It corresponds to the early-stopping strategy;
- *Validation*: Rank two algorithms  $\mathcal{A}_1$  and  $\mathcal{A}_2$  by comparing  $I^{\mathcal{A}_1}$  and  $I^{\mathcal{A}_2}$  under a specific input distribution  $\mathcal{D}$ .

The rest of this paper is organized as: Section 2 verifies definition of the algorithmic information content (Equation 2) by detailing and proving the “coding by posterior” framework; Section 3 interprets the MAXCUT-SDP algorithm; Section 4 describes the provable methods to evaluate MAXCUT-SDP posteriors; Section 5 contains experimental results; Section 6 discusses and concludes the paper.

## 2. The coding by posterior framework

We introduce the “coding by posterior” framework in this section, which is based on an analogue to the noisy communication channel in information theory.

We denote as  $G, G', G'' \in \mathcal{G}$  different data instances generated from the underlying distribution  $\mathcal{D}$ , e.g.,  $G$  may be a graph instance for the MAXCUT problem. Often, a computational problem is associated with some cost function  $R(c, G)$ , which measures how well a hypothesis  $c$  in the hypothesis space  $\mathcal{C}$  will solve the problem on input  $G$ . An algorithm  $\mathcal{A}$  maps the input space to the hypothesis space  $\mathcal{A} : \mathcal{G} \rightarrow \mathcal{C}$ . We introduce parameters  $\theta \in \Theta$  to enumerate a set of algorithms. In optimization for example,  $\theta$  might denote the approximation precision. In general, we assume that algorithm  $\mathcal{A}$  assigns non-negative weights  $w_\theta(c, G)$  to all hypotheses dependent on the input and the parameters, i.e.,

$$w : \mathcal{C} \times \mathcal{G} \times \Theta \rightarrow [0, +\infty), \quad (c, G, \theta) \mapsto w_\theta(c, G). \quad (3)$$

Gibbs weights  $w_\beta(c, G) = \exp(-\beta R(c, G))$ , for example, rank different hypotheses according to how well they solve the problem in terms of costs  $R(c, G)$ .

Such a weighting of hypotheses given data can be interpreted as a posterior distribution  $\mathbb{P}_\theta(c|G)$  induced by algorithm  $\mathcal{A}$ , and is defined as

$$\mathbb{P}_\theta(c|G) := w_\theta(c, G) / \sum_{c' \in \mathcal{C}} w_\theta(c', G), \quad \forall c \in \mathcal{C}. \quad (4)$$

For example, if we choose an indicator function as weights

$$w_\theta(c, G) = \mathbb{1}\{R(c, G) \leq R(c^\perp, G) + \gamma(\theta)\}, \quad (5)$$

where  $\gamma(\theta)$  denotes a precision value determined by a specific  $\mathcal{A}$  and the empirical risk minimizer  $c^\perp(G) = \arg \min_{c \in \mathcal{C}} R(c, G)$  centers an “approximation set” of size  $\gamma(\theta)$  in  $\mathcal{C}$ . In this manner we can recover the “approximation set coding” in [Buhmann \(2010\)](#).

The posterior  $\mathbb{P}_\theta(c|G)$  effectively partitions the hypothesis class into statistically equivalent solutions with high weight values and discards hypotheses with vanishing weights.  $\mathbb{P}_\theta(c|G)$  plays the role of a codebook vector with the associated Voronoi cell. To generate alternative posteriors for a coding protocol we have to use the given data  $G$  and have to transform the mapping from  $\mathcal{G}$  to  $\mathcal{C}$ . Such transformations should not change the measurements represented by  $G$  but the algorithmic mapping. Given data  $G$  and algorithm  $\mathcal{A}$  with posterior  $\mathbb{P}_\theta(c|G)$ , we define the *transformation set*  $\mathbb{T}$  as a set of mappings  $\tau : \mathcal{G} \rightarrow \mathcal{G}$  s.t. the following two conditions are satisfied,

1.  $\mathcal{A}(\tau \circ G), \tau \in \mathbb{T}$  generates an “approximately uniform cover” of the hypothesis space  $\mathcal{C}$ , i.e.,  $\sum_{\tau \in \mathbb{T}} \mathbb{P}_{\theta}(c|\tau \circ G) \in \left[ \frac{|\mathbb{T}|}{|\mathcal{C}|}(1 - \rho), \frac{|\mathbb{T}|}{|\mathcal{C}|}(1 + \rho) \right]$ , for  $0 < \rho < 1$ ;
2. For every transformation  $\tau \in \mathbb{T}$  there exists an associated transformation  $\tau^{\mathcal{C}}$  such that  $w_{\theta}(c, \tau \circ G) = w_{\theta}(\tau^{\mathcal{C}} \circ c, G)$ .

Given a posterior and transformations, we can define a *virtual communication scenario*. It requires a sender  $\mathfrak{S}$ , a receiver  $\mathfrak{R}$ , and a problem generator  $\mathfrak{PG}$  as a noisy channel between  $\mathfrak{S}$  and  $\mathfrak{R}$ . Sender and receiver agree on algorithm  $\mathcal{A}$  and its induced posterior.

## 2.1 Code book generation

The communication code is generated by the procedure:

- 1) Sender  $\mathfrak{S}$  and receiver  $\mathfrak{R}$  obtain the sample set  $G'$  from the problem generator  $\mathfrak{PG}$ .
- 2) Sender  $\mathfrak{S}$  and receiver  $\mathfrak{R}$  calculate the posterior  $\mathbb{P}_{\theta}(c|G')$ .
- 3) A set of transformations  $T = \{\tau_1, \dots, \tau_M\} \subseteq \mathbb{T}$  is generated uniformly with associated posteriors  $\mathbb{P}_{\theta}(c|\tau_j \circ G'), 1 \leq j \leq M$ .
- 4)  $\mathfrak{S}$  and  $\mathfrak{R}$  agree on a transformation set  $T$  and posteriors  $\mathbb{P}_{\theta}(c|\tau_j \circ G'), 1 \leq j \leq M$ .

The posteriors  $\mathbb{P}_{\theta}(c|\tau_j \circ G'), \tau_j \in T$  play the role of codebook vectors in Shannon’s theory of communication.

## 2.2 Communication protocol

- 1) The sender  $\mathfrak{S}$  selects a transformation  $\tau_s \in T$  as message and sends it to the problem generator  $\mathfrak{PG}$ .
- 2)  $\mathfrak{PG}$  generates the test sample set  $G''$  and applies the transformation  $\tau_s$  to  $G''$ , yielding  $\tilde{G} := \tau_s \circ G''$ .
- 3)  $\mathfrak{PG}$  sends  $\tilde{G}$  to  $\mathfrak{R}$  without revealing  $\tau_s$ .
- 4)  $\mathfrak{R}$  calculates the posterior  $\mathbb{P}_{\theta}(c|\tilde{G})$ .
- 5)  $\mathfrak{R}$  estimates the message  $\tau_s$  by using the decoding rule:

$$\hat{\tau} = \arg \max_{\tau \in T} \mathbb{E}_{c \sim \mathbb{P}_{\theta}(c|\tau \circ G')} \mathbb{P}_{\theta}(c|\tau \circ G') = \arg \max_{\tau \in T} \sum_{c \in \mathcal{C}} \mathbb{P}_{\theta}(c|\tau \circ G') \mathbb{P}_{\theta}(c|\tilde{G}) \quad (6)$$

### 2.3 Error analysis of the virtual communication protocol

The probability of a communication error amounts to

$$\begin{aligned}
\mathbb{P}(\hat{\tau} \neq \tau_s | \tau_s) &= \mathbb{P}\left(\max_{\tau_j \in T \setminus \tau_s} \mathbb{E}_{\mathbb{P}(c|\tau_j \circ G')}[\mathbb{P}(c|\tilde{G})] \geq \mathbb{E}_{\mathbb{P}(c|\tau_s \circ G')}[\mathbb{P}(c|\tilde{G})]\right) \\
&\stackrel{(a)}{\leq} \sum_{\tau_j \in T \setminus \tau_s} \mathbb{P}\left(\mathbb{E}_{\mathbb{P}(c|\tau_j \circ G')}[\mathbb{P}(c|\tilde{G})] \geq \mathbb{E}_{\mathbb{P}(c|G')}[\mathbb{P}(c|G'')]\right) \\
&\stackrel{(b)}{\leq} \sum_{\tau_j \in T \setminus \tau_s} \mathbb{E}_{G', G''} \frac{\mathbb{E}_{\tau_j} \mathbb{E}_{\mathbb{P}(c|\tau_j \circ G')}[\mathbb{P}(c|\tilde{G})]}{\mathbb{E}_{\mathbb{P}(c|G')}[\mathbb{P}(c|G'')]},
\end{aligned} \tag{7}$$

by applying the union bound (a) and Markov's inequality (b).

Abbreviating  $Z_{\mathbb{T}} := \mathbb{E}_{\tau_j} \mathbb{E}_{\mathbb{P}(c|\tau_j \circ G')}[\mathbb{P}(c|\tilde{G})]$ , we derive

$$\begin{aligned}
Z_{\mathbb{T}} &= \mathbb{E}_{\tau_j} \mathbb{E}_{\mathbb{P}(c|\tau_j \circ G')} \mathbb{P}(c|\tilde{G}) = \mathbb{E}_{\tau_j} \sum_{c \in \mathcal{C}} \mathbb{P}(c|\tau_j \circ G') \mathbb{P}(c|\tilde{G}) \\
&= \sum_{c \in \mathcal{C}} \mathbb{P}(c|\tilde{G}) \mathbb{E}_{\tau_j} \mathbb{P}(c|\tau_j \circ G') = \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(\tau_j) \mathbb{P}(c|\tau_j \circ G) \leq (1 + \rho) |\mathcal{C}|^{-1}
\end{aligned} \tag{8}$$

where Equation 8 arises from the approximately uniform coverage of  $\mathcal{C}$  by the posteriors, i.e.,  $\sum_{\tau \in \mathbb{T}} \mathbb{P}(c|\tau \circ G) \in [\frac{|\mathbb{T}|}{|\mathcal{C}|}(1 - \rho), \frac{|\mathbb{T}|}{|\mathcal{C}|}(1 + \rho)]$  and  $\mathbb{P}(\tau) = 1/|\mathbb{T}|$ . Substituting Equation 8 into Equation 7 we derive the error bound

$$\begin{aligned}
\mathbb{P}(\hat{\tau} \neq \tau_s | \tau_s) &\leq \sum_{\tau_j \in T \setminus \tau_s} \mathbb{E}_{G', G''} \left[ \left( \frac{|\mathcal{C}|}{1 + \rho} \mathbb{E}_{\mathbb{P}(c|G')}[\mathbb{P}(c|G'')] \right)^{-1} \right] \\
&= (M - 1) \mathbb{E}_{G', G''} \left[ \left( \frac{|\mathcal{C}|}{1 + \rho} k(G', G'') \right)^{-1} \right] \\
&\leq M \mathbb{E}_{G', G''} \left[ \exp(-\log(\frac{|\mathcal{C}|}{1 + \rho} k(G', G''))) \right].
\end{aligned} \tag{9}$$

$$\leq M \mathbb{E}_{G', G''} \left[ \exp(-\log(\frac{|\mathcal{C}|}{1 + \rho} k(G', G''))) \right]. \tag{10}$$

Where Equation 9 comes from the definition of posterior agreement in Equation 1. We analyze  $\hat{I} := \log(|\mathcal{C}|k(G', G''))$  in Equation 10 at its expected value

$$I := \mathbb{E}_{G', G''} [\log(|\mathcal{C}|k(G', G''))]. \tag{11}$$

To control the fluctuations  $\Delta_{G', G''} := \hat{I} - I$ , we assume that for all  $\epsilon > 0, \delta > 0$ , there exists  $n_0 \in \mathbb{N}$  s.t. for all  $n > n_0$

$$\mathbb{P}(|\Delta_{G', G''}| \geq \epsilon I) < \delta. \tag{12}$$

This assumption of asymptotically vanishing fluctuations yields the following upper bound

$$\mathbb{E}_{G', G''} [\exp(-\hat{I})] \leq \exp(-I(1 - \epsilon)). \tag{13}$$

---

**Algorithm 1:** MAXCUT-SDP (Goemans and Williamson, 1995)

---

**Input:** undirected graph  $G = (V, E)$  with non-negative weights  $w$

**Output:** cut  $c = (S, \bar{S})$

- 1 solve  $(R)$ , obtaining an optimal set of vectors  $\mathbf{v}_i \in S_{n-1}$ ;
  - 2 let  $\mathbf{r}$  be a vector *uniformly* distributed on  $S_{n-1}$ ;
  - 3 **return**  $S := \{i \mid \mathbf{v}_i \cdot \mathbf{r} \geq 0, \forall i \in V\}$  and  $\bar{S}$
- 

Since  $\epsilon$  can be chosen arbitrarily small in the asymptotic limit, the error probability is bounded with high probability by

$$\mathbb{P}(\hat{\tau} \neq \tau_s | \tau_s) \leq \exp(-I + \log(M(1 + \rho))). \quad (14)$$

For  $I$  exceeding the effective total rate  $\log(M(1 + \rho))$ , the error vanishes asymptotically since  $I = O(\log |\mathcal{C}|)$ . This bound suggests that we should maximize  $I$  (Equation 11) when searching for informative algorithms, thus verifying our definition of algorithmic information content in Equation 2.

## 2.4 Connection to classical mutual information

We analyse the classical mutual information  $\mathcal{I}(G'; G'')$  by expanding the joint distribution  $\mathbb{P}(G', G'')$  with cut variables  $c \in \mathcal{C}$  and transformations  $\tau \in T$ :

$$\mathcal{I}(G'; G'') = \mathbb{E}_{G', G''} \log \frac{\mathbb{P}(G', G'')}{\mathbb{P}(G')\mathbb{P}(G'')} = \mathbb{E}_{G', G''} \log \frac{\sum_c \sum_\tau \mathbb{P}(G', G'' | c, \tau) \mathbb{P}(c, \tau)}{\mathbb{P}(G')\mathbb{P}(G'')} \quad (15)$$

The conditional distribution  $\mathbb{P}(G', G'' | c, \tau)$  of  $G', G''$  factorizes due to conditioning on  $c$  and  $\tau$ ,

$$\mathbb{P}(G', G'' | c, \tau) = \mathbb{P}(G' | c, \tau) \mathbb{P}(G'' | c, \tau) \stackrel{(a)}{=} \frac{\mathbb{P}(c | G', \tau)}{\mathbb{P}(c | \tau)} \mathbb{P}(G' | \tau) \frac{\mathbb{P}(c | G'', \tau)}{\mathbb{P}(c | \tau)} \mathbb{P}(G'' | \tau), \quad (16)$$

since  $G', G''$  are drawn i.i.d. from the same distribution  $\mathbb{P}(\tau_s \circ G)$ . The transformation  $\tau$  plays the role of a latent variable. Step (a) applies the Bayes rule twice. Substitute Equation 16 into 15 we get (detailed derivation in Appendix A)

$$\mathcal{I}(G'; G'') = \mathbb{E}_{G', G''} \log \sum_c \left[ \frac{\mathbb{P}(c | G') \mathbb{P}(c | G'')}{\mathbb{P}(c)} \right] \leq \mathbb{E}_{G', G''} \log |\mathcal{C}| \sum_c \mathbb{P}(c | G') \mathbb{P}(c | G'') = I(G'; G'').$$

With the uniform distribution  $\mathbb{P}(c) = |\mathcal{C}|^{-1}$ ,  $\mathcal{I}(G'; G'')$  is maximized and we reach the algorithmic information content in Equation 2.

## 3. MaxCut algorithm using SDP relaxation

In this section we give a geometric interpretation of Goemans-Williamson's MAX-CUT algorithm using semidefinite programming relaxation (Goemans and Williamson

(1995), abbreviated as MAXCUT-SDP), which will facilitate deriving methods to calculate the posterior of cuts. Algorithm 1 summarizes the MAXCUT-SDP algorithm: It rounds the solution to a non-linear programming relaxation, which can be interpreted as SDP, then it solves the SDP using standard algorithms, such as interior-point methods (Helmborg et al., 1996), bundle method or block coordinate descent (Waldspurger et al., 2015). Concretely, MAXCUT is formulated as the NP-complete integer program:

$$(Q) \quad \begin{aligned} & \max \frac{1}{2} \sum_{i < j} w_{ij} (1 - v_i v_j) \\ & \text{s.t. } v_i \in \{-1, 1\} \quad \forall i \in V \end{aligned} \quad (17)$$

then (Q) is relaxed to define the following non-linear problem,

$$(R) \quad \begin{aligned} & \max \frac{1}{2} \sum_{i < j} w_{ij} (1 - \mathbf{v}_i \cdot \mathbf{v}_j) \\ & \text{s.t. } \mathbf{v}_i \in S_{n-1} \quad \forall i \in V \end{aligned} \quad (18)$$

where  $S_{n-1}$  is the  $(n - 1)$ -dimensional unit sphere, i.e.,  $S_{n-1} = \{\mathbf{v} \in \mathbb{R}^n \mid \|\mathbf{v}\|_2 = 1\}$ . Arrange the  $n$  vectors  $\mathbf{v}_1, \dots, \mathbf{v}_n$  to be the  $n$  columns of a  $n \times n$  matrix  $D$ , that is,  $D = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n)$ . Let  $X := D^\top D$ , then the  $ij$ -th entry of  $X$  is  $x_{ij} = \mathbf{v}_i \cdot \mathbf{v}_j$ . One can observe that (R) equals to the following SDP problem with only equality constraints:

$$(SDP) \quad \begin{aligned} & \max \frac{1}{2} \sum_{i < j} w_{ij} (1 - x_{ij}) \\ & \text{such that } x_{ii} = 1, \forall i \in V, X \text{ is symmetric positive semidefinite.} \end{aligned} \quad (19)$$

We use one classical interior-point method (Helmborg et al., 1996) to solve the SDP problem in (19).

A geometric view in Figure 1 explains the essence of Algorithm 1: It maps vertices to vectors on the unit sphere. A feasible SDP solution corresponds to a point configuration on the unit sphere, while a feasible solution to MAXCUT assigns a sign variable  $\{\pm 1\}$  with every graph vertex. An optimal solution of SDP tends to send adjacent vertices with heavy edges to antipodal points, thereby maximizing  $(1 - \mathbf{v}_i \cdot \mathbf{v}_j)/2$ . A rounding technique is required that separates most far away pairs, and hence keeps close pairs together.

*Random hyperplane rounding* works gracefully: A random hyperplane through the origin partitions the sphere into two halves, which correspond to cut parts (see Figure 1). The ratio between the expected cut value over the maximum cut value is never worse than  $\alpha \approx .87856$ , which is the expected approximation guarantee.

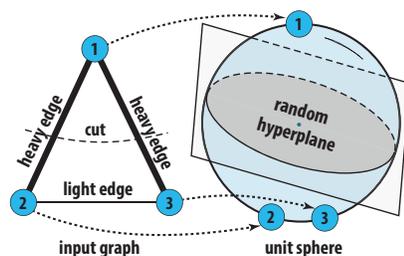


Figure 1: A geometric view of Algorithm 1

## 4. Calculate posterior probability of cuts

Given the geometric interpretation of the MAXCUT-SDP algorithm in Section 3, we can derive the scheme to calculate posterior probability of cuts here, which will be used to evaluate the posterior agreement in Equation 1.

In step  $t$ , the relaxed SDP problem (18) outputs the  $n$  intermediate vectors  $O_t = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$ . Each cut  $c := (S, \bar{S})$ , where  $S = \{1, \dots, \ell\}$ ,  $\bar{S} = \{\ell + 1, \dots, n\}$ , induces a set,

$$B(c) = \{\mathbf{b}_1, \dots, \mathbf{b}_n\} := \{\mathbf{v}_1, \dots, \mathbf{v}_\ell, -\mathbf{v}_{\ell+1}, \dots, -\mathbf{v}_n\}.$$

$B(c)$  is used to define a polygonal intersection cone:

**Definition 1** (Polygonal intersection cone). *The cone  $C$  determined by cut  $c = (S, \bar{S})$  is the intersection of  $n$  half-spaces:*

$$C = C(c) := \{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2 \leq 1 \mid \mathbf{x} \cdot \mathbf{b}_i \geq 0, \forall i \in V\}. \quad (20)$$

It determines the posterior of the corresponding cut by Lemma 1 (all proof of lemmas are in Appendix A) in the following,

**Lemma 1.** *The posterior  $\mathbb{P}(c|G)$  of a cut  $c$  is*

$$\mathbb{P}(c|G) = \frac{2 * \text{unit spherical area of } C(c)}{\text{area of unit sphere}} = \frac{2 * \text{volume of } C(c)}{\text{volume of unit ball}} = \frac{2 * \text{solid angle of } C(c)}{\text{solid angle of unit sphere}}. \quad (21)$$

Ensured by Lemma 1, the cut probabilities are measured either by spherical area, by volume or by solid angle. Without loss of generality, we calculate the solid angle to derive  $\mathbb{P}(c|G)$ . Since it is convenient to express the method of calculating solid angle in terms of the spanning cone definition, let us transform the intersection cone  $C$  into the spanning cone,

**Definition 2** (Polygonal spanning cone). *According to Tiel (1984), a polygonal spanning cone is spanned by a set of  $n$  linearly independent unit vectors  $A = \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n\}$  in  $\mathbb{R}^n$ :*

$$C' := \{\mathbf{x} \in \mathbb{R}^n, \|\mathbf{x}\|_2 \leq 1 \mid \mathbf{x} = \sum_{i=1}^n \lambda_i \mathbf{a}_i, \lambda_i \geq 0, 1 \leq i \leq n\}$$

Given an intersection cone  $C$ , one can get an equivalent spanning cone  $C'$  by taking  $A^\top = B^{-1}$ , which is ensured by,

**Lemma 2.** *Given one intersection cone  $C$  (Definition 1) and one spanning cone  $C'$  (Definition 2), if  $\exists k_1, \dots, k_n > 0$ , s.t.,  $A^\top = \text{diag}(k_1, \dots, k_n)B^{-1}$ , then  $C' = C$ .*

---

**Algorithm 2:** Calculate posterior of each cut

---

**Input:** independent vectors  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  on  $S_{n-1}$

**Output:** posterior of each cut  $\mathbb{P}(c|G), \forall c \in \mathcal{C}$

```
1 for each cut  $c \in \mathcal{C}$  do
2   | get the cut induced set  $B$ 
3   |  $A^\top \leftarrow B^{-1}$ ; //ensured by Lemma 2
4   | compute  $\mathbb{P}(c|G)$  by Equations 23 and 22;
5 return  $\mathbb{P}(c|G), \forall c \in \mathcal{C}$ 
```

---

Now we have the spanning cone  $C'$  associated with the cut  $c$ , we borrow the results of  $n$ -dimensional solid angle calculating (Hajja and Walker, 2002; Ribando, 2006): The solid angle of a spanning cone  $C'$  from Definition 2 is given by:

$$E = |\det(A)| \int_S \|As\|_2^{-n} dS, \quad (22)$$

where the integral is calculated over a unit sphere  $\|s\|_2 = 1$  in the positive orthant given by  $s_i \geq 0$ . Combined with the fact that the solid angle subtended by  $S_{n-1}$  is  $\Omega_n = \frac{2\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2})}$  ( $\Gamma(\cdot)$  is the Gamma function), according to Equation 21,

$$\mathbb{P}(c|G) = 2E/\Omega_n = E \cdot \Gamma(n/2)/\pi^{\frac{n}{2}}. \quad (23)$$

The complete procedure<sup>1</sup> to calculate the posterior probability of each cut is summarized in Algorithm 2.

The way to exactly evaluate the surface integral (Equation 22) is in Appendix B. It involves a  $(n - 1)$ -variate integral, which is computationally intractable, we only use it in the low dimensional case as ground truth.

**Sampling to approximate posterior of cut.** For the high dimensional case, ensured by Lemma 1, we propose one simple and efficient sampling method in Algorithm 3 to approximate the posterior: In each iteration it uniformly samples one hyperplane with normal vector  $\mathbf{r}$  and records the cut  $c$  separated by that hyperplane, then it estimates  $\mathbb{P}(c|G)$  by the statistics of each cut's frequency of occurrence. Theoretical analysis of approximation guarantee of Algorithm 3 and space-efficient implementation of it is in Appendix C and D, respectively.

## 5. Experiments

We compare the MAXCUT-SDP algorithm (abbreviated as “SDP” in the following) with two representative greedy MAXCUT algorithms: The double greedy D2Greedy

---

1. If the  $n$  intermediate vectors  $O_i = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  have mutual dependencies, one can add small perturbations to them in order to make them independent, and the perturbation would still be insignificant w.r.t. vector positions.

---

**Algorithm 3:** Approximate cut’s posterior by sampling

---

**Input:**  $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n\}$  on  $S_{n-1}$ , #samplings

**Output:** approximate posterior of each cut

- 1 initialize  $\text{count}(c) \leftarrow 0, \forall c \in \mathcal{C}$ ;
  - 2 **for** each  $\mathbf{r}$  *uniformly* sampled from  $S_{n-1}$  **do**
  - 3      $\tilde{c} \leftarrow (S, \bar{S})$ , where  $S = \{i \mid \mathbf{v}_i \cdot \mathbf{r} \geq 0, \forall i \in V\}$ ;
  - 4      $\text{count}(\tilde{c}) \leftarrow \text{count}(\tilde{c}) + 1$ ;
  - 5 **return**  $\mathbb{P}(c|G) \approx \text{count}(c)/\#\text{samplings}, \forall c \in \mathcal{C}$
- 

(Deterministic Double Greedy algorithm in Buchbinder et al. (2012)), and the backward greedy EC (Edge Contraction algorithm in Kahruman et al. (2007)). The way to evaluate their posteriors (approximation sets) can be found in Bian et al. (2015). Let  $W_{\mathcal{A}}(G)$  be the cut value generated by an algorithm  $\mathcal{A}$  on graph  $G$ ,  $W_*(G)$  be the optimal cut value of  $G$ . The approximation ratio of an algorithm  $\mathcal{A}$  is the worst-case bound  $\min_G \frac{W_{\mathcal{A}}(G)}{W_*(G)}$ , which ranks the three algorithms as  $\text{SDP} \succ \text{D2Greedy} \succ \text{EC}$ . Since finding  $W_*(G)$  for NP-complete problem is non-trivial, we use  $\frac{W_{\mathcal{A}}(G)}{W(G)}$  ( $W(G)$ : total weight of  $G$ ) as a natural lower bound of  $\frac{W_{\mathcal{A}}(G)}{W_*(G)}$ .

**Experimental setting.** We experiment with the Gaussian edge weights model (Gronskiy and Buhmann, 2014): The graph instances are generated in a two-step fashion: Firstly, a random “master” graph  $G$  is generated with Gaussian distributed edge weights  $w_{ij} \sim N(\mu, \sigma_m^2)$ ,  $\mu = 300, \sigma_m = 50$ , negative edges are set to be  $\mu$ . Secondly, noisy graphs  $G', G''$  are obtained by adding Gaussian distributed noise  $n_{ij} \sim N(0, \sigma^2)$ , negative edges are set to be zero. We perform 1000 repeated noisy samplings to estimate the expectation over  $(G', G'')$  in Equation 2.

**Results and analysis.** Figure 2 shows the temporal information content ( $I_t^{\mathcal{A}}$  in Equation 2) for two  $\sigma$  values: 10 and 58. For all the algorithms,  $I_t^{\mathcal{A}}$  increases at the beginning. After reaching some optimal step  $t^*$ , where the highest  $I_t^{\mathcal{A}}$  ( $I^{\mathcal{A}}$  in Equation 2) is achieved, it decreases and finally vanishes. This observation confirms the principle of regularization by early stopping at time  $t^*$  when maximum  $I_t^{\mathcal{A}}$  is reached.

Figure 3 shows (a) the information content, and (b) the fully overlapping curves  $W_{\mathcal{A}}(G')/W(G'), W_{\mathcal{A}}(G'')/W(G'')$ , respectively.  $\sigma$  controls the noise level, larger  $\sigma$  means larger noise. In the noiseless case,  $G' = G''$ , so  $\mathbb{P}(c|G') = \mathbb{P}(c|G'')$ , and  $I_t^{\mathcal{A}} = \mathbb{E}_{G', G''}[\log(|\mathcal{C}| \sum_{c \in \mathcal{C}} \mathbb{P}^2(c|G'))]$ . All algorithms start with uniform distribution of solutions when  $t = 0$ ; as the algorithm proceeds, the distribution of solutions concentrates more and more on a small support,  $\sum_{c \in \mathcal{C}} \mathbb{P}^2(c|G')$  increases and reaches a maximum in the final step, so all algorithms reach the maximum  $I_t^{\mathcal{A}}$  in the final step. For greedy algorithms (D2Greedy and EC), there is only one final solution

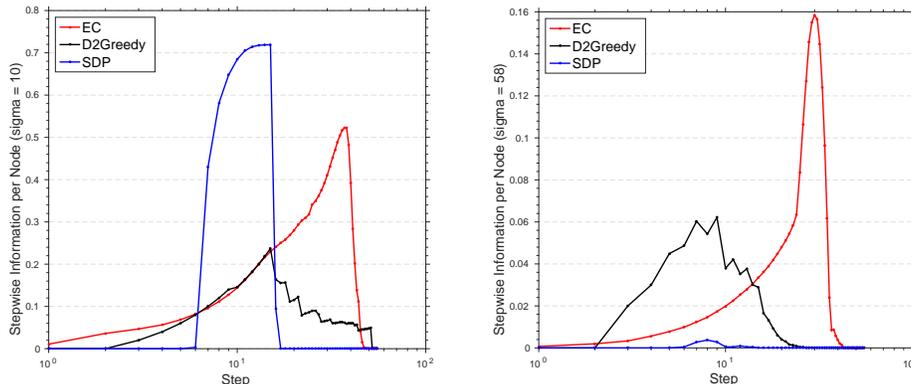


Figure 2:  $I_t^{\mathcal{A}}$  per vertex w.r.t.  $t$  (a)  $\sigma = 10$  (b)  $\sigma = 58$  ( $n = 50$ )

with probability 1 in the last step, so  $\sum_{c \in \mathcal{C}} \mathbb{P}^2(c|G') = 1$  and the maximum  $I_t^{\mathcal{A}}$  is  $\log(|\mathcal{C}|) = \log(2^{n-1} - 1)$ , as shown by Figure 3(a). For SDP, however, when  $\sigma = 0$ , SDP can only approximately solve the input graphs, in the last step there are several solutions with non-zero probability, which renders its information content less than  $\log(|\mathcal{C}|)$ .

It is worth noting that for greedy algorithms (D2Greedy and EC), the higher the approximation ratio is for noisy graphs, the lower is the information content achieved by the algorithm. This behavior is quite intuitive since high approximation ratio means better adaptation to empirical fluctuations and, therefore, overfitting to noisy graphs. Consequently, there will be less agreement between the solutions of the two noisy graphs and the information content of the algorithm drops. A similar conclusion has also been drawn in Bousquet and Bottou (2008).

However, for non-greedy algorithm SDP, there are two factors affecting its information content: The approximation ratio and its probabilistic weighting strategy to down-weight solutions without discarding them. SDP keeps all the possible solutions, instead of removing the bad solutions as greedy algorithms do, it assigns less probabilistic weights to them, so it can capture some uncertainty in the input.

The information content of SDP shows the influence of both factors: For low noise, the probabilistic weighting strategy dominates, SDP outperforms greedy algorithms in information content; while in high noise level, the influence of approximation ratio dominates, and SDP is inferior to greedy algorithms.

## 6. Discussion and conclusion

$I^{\mathcal{A}}$  in Equation (2) measures the information content of an algorithm given a noisy source of instances. We theoretically justify this criterion, and apply it to study the robustness of MAXCUT algorithms with different approximation ratios. Of particular interest is the SDP based algorithm by Goemans and Williamson (1995), since it pursues a non-greedy strategy for MAXCUT.

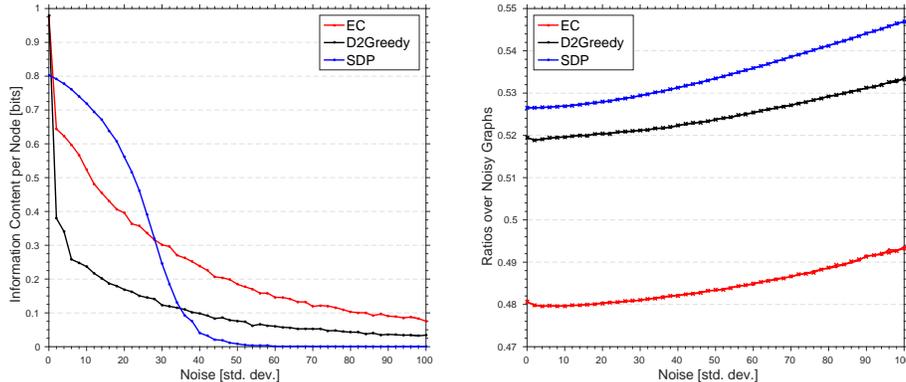


Figure 3: (a)  $I^{\mathcal{A}}$  per vertex w.r.t.  $\sigma$  (b)  $\frac{W_{\mathcal{A}}(G')}{W(G')}$ ,  $\frac{W_{\mathcal{A}}(G'')}{W(G'')}$

Comparison of SDP with two representative greedy MAXCUT algorithms (D2Greedy and EC) demonstrates that the ability of this approximation algorithm to achieve a high approximation ratio might decrease its generalization ability. The property of an algorithm to efficiently find a good empirical minimum might increase its fragility due to noise adaptation. This observation could be generalized or even proved for general approximation algorithms provided that the algorithms operate in a similar settings or use similar optimization strategies.

The posterior agreement based criterion also enables a meta-algorithm to search for more informative algorithms. Algorithms are usually tuned by parameter adaptation or by modifying the algorithmic strategy in the spirit of genetic programming. Thereby, the meta-algorithm will search through the space of algorithms guided by maximal gradient ascent on posterior agreement. With a validation criterion as posterior agreement, we enable algorithm engineering to explore multi-objective optimization of algorithms with respect to time, space and robustness.

#### ACKNOWLEDGMENT

This research was partially supported by the Max Planck ETH Center for Learning Systems.

## References

- Rakesh Agrawal, Sridhar Rajagopalan, Ramakrishnan Srikant, and Yirong Xu. Mining newsgroups using networks arising from social behavior. In *WWW*, pages 529–535, 2003.
- Yatao Bian, Alexey Gronskiy, and Joachim M. Buhmann. Greedy maxcut algorithms and their information content. In *IEEE Information Theory Workshop (ITW)*, pages 1–5, 2015.
- Olivier Bousquet and Léon Bottou. The tradeoffs of large scale learning. In *NIPS*, pages 161–168, 2008.
- Niv Buchbinder, Moran Feldman, Joseph Naor, and Roy Schwartz. A tight linear time  $(1/2)$ -approximation for unconstrained submodular maximization. In *FOCS, 2012*, pages 649–658. IEEE, 2012.
- Joachim M. Buhmann. Information theoretic model validation for clustering. In *ISIT*, pages 1398–1402, 2010.
- Rich Caruana, Steve Lawrence, and Lee Giles. Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping. In *NIPS*, volume 13, page 402, 2001.
- Michel X Goemans and David P Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *Journal of the ACM*, 42(6):1115–1145, 1995.
- Alexey Gronskiy and Joachim M. Buhmann. How informative are minimum spanning tree algorithms. In *ISIT*, 2014.
- Torben Hagerup and Christine Rüb. A guided tour of chernoff bounds. *Information processing letters*, 33(6):305–308, 1990.
- Mowaffaq Hajja and Peter Walker. The measure of solid angles in n-dimensional euclidean space. *International Journal of Mathematical Education in Science and Technology*, 33(5):725–729, 2002.
- Christoph Helmberg, Franz Rendl, Robert J Vanderbei, and Henry Wolkowicz. An interior-point method for semidefinite programming. *SIAM Journal on Optimization*, 6(2):342–361, 1996.
- Sera Kahruman, Elif Kolotoglu, Sergiy Butenko, and Illya V Hicks. On greedy construction heuristics for the maxcut problem. *International Journal of Computational Science Engineering*, 3(3):211–218, 2007.

Jason M Ribando. Measuring solid angles beyond dimension three. *Discrete & Computational Geometry*, 36(3):479–487, 2006.

Jan van Tiel. *Convex analysis*. John Wiley, 1984.

Irène Waldspurger, Alexandre d’Aspremont, and Stéphane Mallat. Phase recovery, maxcut and complex semidefinite programming. *Mathematical Programming*, 149(1-2):47–81, 2015.

Jun Wang, Tony Jebara, and Shih-Fu Chang. Semi-supervised learning using greedy max-cut. *JMLR*, 14(1):771–800, March 2013. ISSN 1532-4435.

# Appendix

## A. Proofs

### A.1 Detailed proof in Section 2.4

The classical mutual information determines a lower bound on information content defined in Equation 2. The classical mutual information  $\mathcal{I}(G'; G'')$  is defined as

$$\mathcal{I}(G'; G'') = \mathbb{E}_{G', G''} \log \frac{\mathbb{P}(G', G'')}{\mathbb{P}(G')\mathbb{P}(G'')} \quad (24)$$

From the definition of virtual communication scenario, the data instances  $G', G''$  can be treated to be drawn from a mixture distribution, as illustrated in Figure 4.

We consider the special transformations that map from input space to output spaces, then the joint probability can be factorized in the following way,

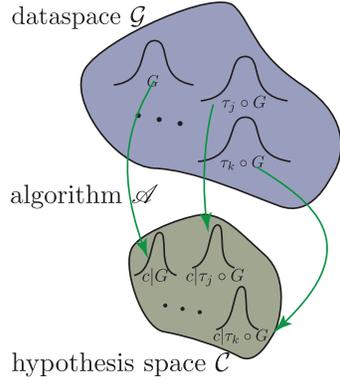


Figure 4: Illustration of the mixture distribution

$$\begin{aligned} \mathbb{P}(G', G'') &= \sum_{c \in \mathcal{C}} \sum_{T \in \mathbb{T}^M} \mathbb{P}(G', G'' | c, T) \mathbb{P}(c, T) \\ &= \sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \sum_{c \in \mathcal{C}} \mathbb{P}(G', G'' | \tau_j, c) \mathbb{P}(\tau_j, c) \\ &\stackrel{(a)}{=} \sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \sum_{c \in \mathcal{C}} \mathbb{P}(G' | \tau_j, c) \mathbb{P}(G'' | \tau_j, c) \mathbb{P}(c | \tau_j) \mathbb{P}(\tau_j) \\ &\stackrel{(b)}{=} \sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(\tau_j) \sum_{c \in \mathcal{C}} \frac{\mathbb{P}(c | G', \tau_j)}{\mathbb{P}(c | \tau_j)} \mathbb{P}(G' | \tau_j) \frac{\mathbb{P}(c | G'', \tau_j)}{\mathbb{P}(c | \tau_j)} \mathbb{P}(G'' | \tau_j) \mathbb{P}(c | \tau_j) \\ &\stackrel{(c)}{=} \sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G' | \tau_j) \mathbb{P}(G'' | \tau_j) \mathbb{P}(\tau_j) \underbrace{\sum_{c \in \mathcal{C}} \frac{\mathbb{P}(c | G', \tau_j) \mathbb{P}(c | G'', \tau_j)}{\mathbb{P}(c | \tau_j)}}_{=\tilde{k}(\tau_j \circ G', \tau_j \circ G'')} \end{aligned}$$

Step (a) exploits the fact that conditioning on the transformation  $\tau_j$  and hypothesis  $c$  renders  $G', G''$  statistically independent since the two instances are drawn i.i.d. from the same component of the mixture distribution; (b) applies Bayes rule twice. In step (c) we define the *generalized posterior agreement* as

$$\tilde{k}(G', G'') := \sum_{c \in \mathcal{C}} \frac{\mathbb{P}(c | G') \mathbb{P}(c | G'')}{\mathbb{P}(c)}.$$

From condition 2) of the definition of the transformation set  $\mathbb{T}$ , one can get that

$$\tilde{k}(\tau_j \circ G', \tau_j \circ G'') = \tilde{k}(G', G'') \quad (25)$$

Combining Equation 25 with (c) one can get,

$$\frac{\mathbb{P}(G', G'')}{\mathbb{P}(G')\mathbb{P}(G'')} = \tilde{k}(G', G'') \frac{\sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_j)\mathbb{P}(\tau_j)}{\underbrace{\sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(\tau_j) \sum_{l=1}^M p_l \sum_{\tau_l \in \mathbb{T}} \mathbb{P}(G''|\tau_l)\mathbb{P}(\tau_l)}_{\text{to be proved} = |\mathbb{T}|}} \quad (26)$$

$$= \tilde{k}(G', G'')|\mathbb{T}| \quad (27)$$

Let us prove Equation 27 first of all. We simplify the term,

$$\begin{aligned} & \frac{\sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_j)\mathbb{P}(\tau_j)}{\sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(\tau_j) \sum_{l=1}^M p_l \sum_{\tau_l \in \mathbb{T}} \mathbb{P}(G''|\tau_l)\mathbb{P}(\tau_l)} = \\ & \frac{\sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_j)\mathbb{P}(\tau_j)}{\sum_{j=1}^M \sum_{l=1}^M \sum_{\tau_j \in \mathbb{T}} \sum_{\tau_l \in \mathbb{T}} p_j p_l \mathbb{P}(\tau_j)\mathbb{P}(\tau_l)\mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l)} = \\ & \frac{\sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_j)\mathbb{P}(\tau_j)}{\sum_{j=1}^M \sum_{\tau_j \in \mathbb{T}} p_j^2 \mathbb{P}(\tau_j)^2 \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l) + \sum_{j \neq l} \sum_{\tau_j \in \mathbb{T}} \sum_{\tau_l \in \mathbb{T}} p_j p_l \mathbb{P}(\tau_j)\mathbb{P}(\tau_l)\mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l)} \end{aligned} \quad (28)$$

We further make the simplifying assumption that  $\mathbb{P}(\tau_j) = \mathbb{P}(\tau_l) = 1/|\mathbb{T}|$ , then

$$\begin{aligned} & \Rightarrow \frac{\frac{1}{|\mathbb{T}|} \sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_j)}{\frac{1}{|\mathbb{T}|^2} \left( \sum_{j=1}^M \sum_{\tau_j \in \mathbb{T}} p_j^2 \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l) + \sum_{j \neq l} \sum_{\tau_j \in \mathbb{T}} \sum_{\tau_l \in \mathbb{T}} p_j p_l \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l) \right)} \\ & = |\mathbb{T}| \frac{\sum_{j=1}^M p_j \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_j)}{\left( \sum_{j=1}^M p_j^2 \sum_{\tau_j \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l) + \sum_{j \neq l} p_j p_l \sum_{\tau_j \in \mathbb{T}} \sum_{\tau_l \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l) \right)} \\ & \stackrel{p_j = \Delta_{js}}{=} |\mathbb{T}| \frac{\sum_{\tau_s \in \mathbb{T}} \mathbb{P}(G'|\tau_s)\mathbb{P}(G''|\tau_s)}{\sum_{\tau_s \in \mathbb{T}} \mathbb{P}(G'|\tau_s)\mathbb{P}(G''|\tau_s) + \underbrace{\sum_{j \neq l} \Delta_{js} \Delta_{ls}}_{= 0 \text{ for } l \neq j} \sum_{\tau_j \in \mathbb{T}} \sum_{\tau_l \in \mathbb{T}} \mathbb{P}(G'|\tau_j)\mathbb{P}(G''|\tau_l)} \\ & = |\mathbb{T}| \end{aligned}$$

Inserting Equation (27) into Equation (24) proves the claim that the mutual information is a lower bound on the information content:

$$\begin{aligned}
\mathcal{I}(G'; G'') &= \mathbb{E}_{G', G''} \log \frac{\mathbb{P}(G', G'')}{\mathbb{P}(G')\mathbb{P}(G'')} \\
&= \mathbb{E}_{G', G''} \log \tilde{k}(G', G'') \\
&= \mathbb{E}_{G', G''} \log \sum_{c \in \mathcal{C}} \frac{\mathbb{P}(c|G')\mathbb{P}(c|G'')}{\mathbb{P}(c)} \\
&\leq \mathbb{E}_{G', G''} \log |\mathcal{C}| \sum_{c \in \mathcal{C}} \mathbb{P}(c|G')\mathbb{P}(c|G'') \\
&= I(G'; G'').
\end{aligned}$$

With the uniform distribution  $\mathbb{P}(c) = |\mathcal{C}|^{-1}$ ,  $\mathcal{I}(G'; G'')$  is maximized and we derive the information content in Equation 2.

## A.2 Proof of Lemma 1

*Proof.* For a specific cut  $c := (S, \bar{S})$ , assume the collections of normal vectors of all hyperplanes that give the cut  $c$  is  $R(c)$ , according to the random hyperplane rounding technique,

$$\begin{aligned}
R(c) &= \{\mathbf{r} \in S_{n-1} \mid \mathbf{r} \cdot \mathbf{b}_i \geq 0, \forall i \in V\} \\
&\cup \{\mathbf{r} \in S_{n-1} \mid \mathbf{r} \cdot \mathbf{b}_i \leq 0, \forall i \in V\}
\end{aligned} \tag{29}$$

So  $R(c)$  is two times the unit spherical surface of  $C(c)$ . Considering the fact that normal vectors of all hyperplanes constitute the surface of unit sphere, we get the first equality in (21). Using simple geometrical knowledge, we can get the second and third equalities. ■

## A.3 Proof of Lemma 2

*Proof.*  $C' = C \Leftrightarrow$  any point in  $C'$  must be in  $C \Leftrightarrow (\sum_1^n \lambda_i \mathbf{a}_i) \cdot \mathbf{b}_j \geq 0, 1 \leq i, j \leq n$  holds  $\forall \lambda_i \geq 0 \Leftrightarrow (\lambda_1, \dots, \lambda_n) \cdot A^T B \geq 0$  holds  $\forall (\lambda_1, \dots, \lambda_n) \geq 0$

So if  $\exists k_1, \dots, k_n > 0$ , s.t.  $A^T = \text{diag}(k_1, \dots, k_n) B^{-1}$ , one can get that  $(\lambda_1, \dots, \lambda_n) \cdot A^T B \geq 0$  holds  $\forall (\lambda_1, \dots, \lambda_n) \geq 0$ , so  $C' = C$ . ■

## B. Exactly evaluate the surface integral (Equation 22)

Exactly calculating probability of cuts involves evaluating the high dimensional surface integral in Equation 22, To do this, we first of all parametrize it using spherical polar coordinates, then transform it to be a multivariate integral. Writing  $\mathbf{s} = \sum_1^n s_i \mathbf{e}_i$ , we get:

$$\|A\mathbf{s}\|_2^2 = \sum_{i=1}^n \mathbf{a}_i \cdot \mathbf{a}_i s_i^2 + 2 \sum_{i<j} \mathbf{a}_i \cdot \mathbf{a}_j s_i s_j = 1 + 2 \sum_{i<j} \mathbf{a}_i \cdot \mathbf{a}_j s_i s_j \tag{30}$$

Plugging Equation 30 into 22 one can express the surface integral in a more manageable form

$$E = |\det(A)| \int_S (1 + 2 \sum_{i < j} \mathbf{a}_i \cdot \mathbf{a}_j s_i s_j)^{-n/2} dS = |\det(A)| \int_S f^{-n/2}(\mathbf{s}) dS \quad (31)$$

where  $f(\mathbf{s}) = 1 + 2 \sum_{i < j} \mathbf{a}_i \cdot \mathbf{a}_j s_i s_j$ . Then parametrizing by spherical polar coordinates  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_{n-1})$ :

$$s_i = \cos(\theta_i) \prod_{j=1}^{i-1} \sin(\theta_j), i = 1, \dots, n-1; s_n = \prod_{i=1}^{n-1} \sin(\theta_i) \quad (32)$$

for  $0 \leq \theta_i \leq \pi/2, 0 \leq i \leq n-1$ , considering that the Jacobian is  $\prod_{i=1}^{n-2} \sin^{n-1-i}(\theta_i)$ , substitute Equation 32 to Equation 31 it reaches the multivariate integral:

$$E = |\det(A)| \int_{\theta_1} \dots \int_{\theta_{n-1}} \frac{\prod_{i=1}^{n-2} \sin^{n-1-i}(\theta_i)}{f^{n/2}(\boldsymbol{\theta})} d\theta_1 \dots d\theta_{n-1} \quad (33)$$

### C. Theoretical analysis of Algorithm 3

We will show that for a cut  $c$  with high ground truth probability  $p_c := \mathbb{P}(c|V)$ , the estimated cut probability  $\hat{p}_c$  by uniform sampling in Algorithm 3 will be close to  $p_c$  with high probability.

Let  $k = \#\text{samplings}$ , random variable  $X_i = 1$  means recovering cut  $c$  in the  $i^{\text{th}}$  sampling,  $X_i = 0$  means not recovering  $c$  in the  $i^{\text{th}}$  sampling. So  $\hat{p}_c = \sum_{i=1}^k X_i/k$ , from the Chernoff-Hoeffding theorem (Hagerup and Rüb, 1990), for  $\epsilon > 0$ ,

$$\begin{aligned} \mathbb{P}(\hat{p}_c \geq p_c + \epsilon) &\leq \left[ \left( \frac{p_c}{p_c + \epsilon} \right)^{p_c + \epsilon} \left( \frac{1 - p_c}{1 - p_c - \epsilon} \right)^{1 - p_c - \epsilon} \right]^k = e^{-D(p_c + \epsilon \| p_c)k} \\ &\leq \left( \frac{p_c}{p_c + \epsilon} \right)^{k(p_c + \epsilon)} \cdot e^{k\epsilon} \\ \mathbb{P}(\hat{p}_c \leq p_c - \epsilon) &\leq \left[ \left( \frac{p_c}{p_c - \epsilon} \right)^{p_c - \epsilon} \left( \frac{1 - p_c}{1 - p_c + \epsilon} \right)^{1 - p_c + \epsilon} \right]^k \\ &= e^{-D(p_c - \epsilon \| p_c)k} \leq \left( \frac{p_c}{p_c - \epsilon} \right)^{k(p_c - \epsilon)} \cdot e^{-k\epsilon} \end{aligned}$$

where  $D(\cdot)$  is the Kullback-Leibler divergence between two Bernoulli random variables.

So to ensure that with probability at most  $\delta < 1$ , the estimated probability  $\hat{p}_c$  is at most  $\epsilon$ -distant from the true probability  $p_c$ , one need to ensure that:

$$\max(e^{-D(p_c + \epsilon \| p_c)k}, e^{-D(p_c - \epsilon \| p_c)k}) \leq \delta$$

---

**Algorithm 4:** Pseudo-code to calculate estimate of  $\sum_c \mathbb{P}(c|G')\mathbb{P}(c|G'')$  when  $k < |\mathcal{C}|$

---

**Input:**  $\text{cutIndices}(G'), \text{cutIndices}(G'') \in \mathbb{R}^k$ , wherein the indices are in ascending order

**Output:** estimate of  $\sum_c \mathbb{P}(c|G')\mathbb{P}(c|G'')$

```

1 initialize  $\text{idx1} = \text{idx2} = 1, \text{sum} = 0$ ;
2 while  $\text{idx1} \leq k \ \&\& \ \text{idx2} \leq k$  do
3   if  $\text{cutIndices}(G')_{\text{idx1}} == \text{cutIndices}(G'')_{\text{idx2}}$  then
4      $\text{commonIdx} = \text{cutIndices}(G')_{\text{idx1}}$ ;
5      $\text{cutNum1} = \text{cutNum2} = 1$ ;
6      $\text{idx1} ++, \text{idx2} ++$ ;
7     while  $\text{idx1} \leq k \ \&\& \ \text{commonIdx} == \text{cutIndices}(G')_{\text{idx1}}$  do
8        $\text{cutNum1} ++, \text{idx1} ++$ ;
9     while  $\text{idx2} \leq k \ \&\& \ \text{commonIdx} == \text{cutIndices}(G'')_{\text{idx2}}$  do
10       $\text{cutNum2} ++, \text{idx2} ++$ ;
11     $\text{sum} += \text{cutNum1} * \text{cutNum2}$ ;
12  else if  $\text{cutIndices}(G')_{\text{idx1}} < \text{cutIndices}(G'')_{\text{idx2}}$  then
13     $\text{idx1} ++$ ;
14  else
15     $\text{idx2} ++$ ;
16 return  $\sum_c \mathbb{P}(c|G')\mathbb{P}(c|G'') \approx \text{sum}/k^2$ 

```

---

which is equivalent to:

$$k \geq \max \left( \frac{-\ln \delta}{D(p_c + \epsilon || p_c)}, \frac{-\ln \delta}{D(p_c - \epsilon || p_c)} \right) \quad (34)$$

which gives the lower bound of the sampling number  $k$  required to recover the ground truth  $p_c$  with probability  $\delta$  at a specific error level  $\epsilon$ .

## D. Space-efficient implementation of Algorithm 3

When sampling number  $k \geq |\mathcal{C}|$ , use array  $\text{cutFrequency} \in \mathbb{R}^{|\mathcal{C}|}$  to record cuts' frequency of occurrence, and the posterior agreement  $\sum_c \mathbb{P}(c|G')\mathbb{P}(c|G'')$  is estimated as the inner product  $\langle \text{cutFrequency}(G'), \text{cutFrequency}(G'') \rangle$ .

When  $k < |\mathcal{C}|$ , use array  $\text{cutIndices} \in \mathbb{R}^k$  to record indices of sampled cuts in each sampling, note that there would be duplicated cuts in  $\text{cutIndices}$ . Then sort the array  $\text{cutIndices}$  to make the indices in it be in ascending order. Finally, use the way described by the pseudo-code in Algorithm 4 to calculate estimate of posterior agreement  $\sum_c \mathbb{P}(c|G')\mathbb{P}(c|G'')$ .