Ergod. Th. & Dynam. Sys., page 1 of 38 © The Author(s), 2022. Published by Cambridge University Press. This is an Open Access article, distributed under the terms of the Creative Commons Attribution licence (http://creativecommons.org/licenses/by/4.0), which permits unrestricted re-use, distribution and reproduction, provided the original article is properly cited. doi:10.1017/etds.2022.27

The relative f-invariant and non-uniform random sofic approximations

CHRISTOPHER SHRIVER

Department of Mathematics, University of California, Los Angeles, Los Angeles 90095, CA, USA (e-mail: christopher.shriver@math.utexas.edu)

(Received 29 May 2021 and accepted in revised form 1 March 2022)

Abstract. The f-invariant is an isomorphism invariant of free-group measure-preserving actions introduced by Lewis Bowen, who first used it to show that two finite-entropy Bernoulli shifts over a finitely generated free group can be isomorphic only if their base measures have the same Shannon entropy. Bowen also showed that the f-invariant is a variant of sofic entropy; in particular, it is the exponential growth rate of the expected number of good models over a uniform random homomorphism. In this paper we present an analogous formula for the relative f-invariant and use it to prove a formula for the exponential growth rate of the expected number of good models over a random sofic approximation which is a type of stochastic block model.

Key words: sofic entropy, random regular graphs, stochastic block model 2020 Mathematics Subject Classification: 37A35 (Primary); 28D20 (Secondary)

1. Introduction and main results

Let $G = \langle S \rangle$ denote the rank-r free group with generating set $S = \{s_1, \ldots, s_r\}$ and identity e, and let (X, μ, T) be a measure-preserving G-system, that is, T is a homomorphism from G to the automorphism group of the standard probability space (X, μ) . We will not need to make explicit use of the σ -algebra on X, so we leave it unnamed.

An *observable* on *X* is a measurable map with domain *X*. In this paper the codomain will be a finite set endowed with the discrete sigma algebra; in this case we call the map a *finite observable* and the codomain an *alphabet*.

Any observable $\alpha: X \to A$ induces a map $\alpha^G: X \to A^G$ by setting

$$(\alpha^G(x))_g = \alpha(T_g x)$$
 for all $g \in G$.

We call the A-coloring $\alpha^G(x)$ of G the *itinerary* of x, since it records the observations that will be made over the entire orbit of x under the action of G. We also similarly define

the map $\alpha^H: X \to \mathbb{A}^H$ for any subset H of G. We abbreviate $\alpha^n := \alpha^{\mathbb{B}(e,n)}$, where $\mathbb{B}(e,n)$ is the closed ball of radius n centered at the identity in G, which is endowed with the word-length metric. If $\beta: X \to \mathbb{B}$ is a second finite observable, we denote by $\alpha\beta: X \to \mathbb{A} \times \mathbb{B}$ the map $\alpha\beta(x) = (\alpha(x), \beta(x))$.

The (Shannon) entropy of a finite observable $\alpha: X \to A$ is defined by

$$H_{\mu}(\alpha) = -\sum_{a \in \mathbb{A}} \alpha_* \mu(a) \log \alpha_* \mu(a),$$

where $\alpha_*\mu \in \operatorname{Prob}(A)$ is the pushforward measure, with the convention $0 \log 0 = 0$. The entropy of α can be interpreted as the expected amount of information revealed by observing α , assuming its distribution $\alpha_*\mu$ is known.

An early application of Shannon's entropy to ergodic theory was its use by Kolmogorov and Sinai to show that there exist non-isomorphic Bernoulli shifts over \mathbb{Z} . A Bernoulli shift over \mathbb{Z} is a system of the form $(\mathbb{A}^{\mathbb{Z}}, \mu^{\mathbb{Z}}, S)$ for some alphabet \mathbb{A} and $\mu \in \operatorname{Prob}(\mathbb{A})$; S is the shift action of \mathbb{Z} . They did this by defining an *entropy rate* for \mathbb{Z} -systems, which can be interpreted as the average information per unit time revealed by observing the system. For a Bernoulli shift $(\mathbb{A}^{\mathbb{Z}}, \mu^{\mathbb{Z}}, S)$, the entropy rate is simply the 'base entropy' $H_{\mu}(\alpha)$, where $\alpha : \mathbb{A}^n \to \mathbb{A}$ is the 'time zero' observable.

Isomorphism invariance of the Kolmogorov–Sinai entropy rate is typically proven using the fact that entropy rate is non-increasing under factor maps (which are surjective homomorphisms of measure-preserving systems). This fact can be interpreted as stating that a system cannot simulate another system that is 'more random'.

The entropy rate was soon generalized to systems acted on by an arbitrary amenable group (such as \mathbb{Z}^d). Extending beyond amenable groups proved more difficult, and in fact it was found to be impossible for such an extension to preserve all desirable properties of the Kolmogorov–Sinai entropy rate. In particular, an entropy rate for non-amenable group actions which assigns Bernoulli shifts their base entropy cannot be non-increasing under factor maps [13, Appendix C].

The first invariant to distinguish between Bernoulli shifts over free groups is Lewis Bowen's f-invariant. Following [2], this can be defined by

$$F_{\mu}(T,\alpha) = (1-2r)H_{\mu}(\alpha) + \sum_{i=1}^{r} H_{\mu}(\alpha^{\{e,s_i\}}),$$

$$f_{\mu}(T,\alpha) = \inf_{n} F_{\mu}(T,\alpha^{n}) = \lim_{n \to \infty} F_{\mu}(T,\alpha^{n}).$$

The main theorem of [3] is that $f_{\mu}(T,\alpha)$ depends on the observable α only through the σ -algebra it generates. In particular, the common value of $f_{\mu}(T,\alpha)$ among all α which generate the σ -algebra of the measurable space X (assuming such α exist) is a measure-conjugacy invariant of the system (X,μ,T) . In the same paper, Bowen showed that the f-invariant of a Bernoulli shift is the Shannon entropy of the base measure; in particular, Bernoulli shifts with different base entropies are non-isomorphic.

In [2], Bowen gave an alternate formula for the f-invariant, which we now introduce.

For any homomorphism $\sigma: G \to \operatorname{Sym}(n)$ we have a G-system ([n], $\operatorname{Unif}(n)$, σ), and we can consider a labeling $\mathbf{x} \in \mathbb{A}^n$ as an A-valued observable on this system. We denote

the law of its itinerary by $P_{\mathbf{x}}^{\sigma} = \mathbf{x}_{*}^{G}$ Unif(n) and call this the *empirical distribution* of \mathbf{x} . We say that \mathbf{x} is a good model for α over σ if it is difficult to distinguish the *G*-systems (X, μ, T) and $([n], \mathrm{Unif}(n), \sigma)$ via their respective observables α and \mathbf{x} . To make this precise, we denote

$$\Omega(\sigma, O) := \{ \mathbf{x} \in \mathbb{A}^n : P_{\mathbf{x}}^{\sigma} \in O \},$$

which is a set of good models for α over σ if O is a weak*-open neighborhood of $\alpha_*^G \mu \in \operatorname{Prob}(\mathbb{A}^G)$; the particular set O quantifies how good the models are. The alphabet \mathbb{A} is given the discrete topology and \mathbb{A}^G the product topology, so 'weak*-close' means marginals on some finite sets are close in total variation norm.

For each $n \in \mathbb{N}$, let $s_n = \text{Unif}(\text{Hom}(G, \text{Sym}(n)))$. Bowen showed in [2] that the f-invariant is given by

$$f_{\mu}(T,\alpha) = \inf_{O \ni \alpha_*^G \mu} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_n}{\mathbb{E}} |\Omega(\sigma, O)|.$$

To make an analogy with statistical physics, we can think of $\alpha_*^G \mu$ as a macroscopic statistical distribution of the state of a system; then the f-invariant is the exponential growth rate of the expected number of 'microstates' that are consistent with these statistics. What we here call good models are often called microstates for this reason.

More generally, given any random or deterministic sofic approximation $\Sigma = \{s_n\}_{n=1}^{\infty}$, we can define the sofic entropy relative to Σ by

$$\mathbf{h}_{\Sigma,\mu}(T,\alpha) = \inf_{O\ni\alpha_*^G\mu} \limsup_{n\to\infty} \frac{1}{n}\log \underset{\sigma\sim \mathbf{s}_n}{\mathbb{E}} |\Omega(\sigma,O)|.$$

Here each s_n is a probability measure on the set of functions $G \to \operatorname{Sym}(n)$ which is supported on functions which are approximately free homomorphisms.

This paper is motivated by a desire to better understand the dependence of sofic entropy on the sofic approximation Σ . For any choice of Σ , the sofic entropy agrees with Kolmogorov–Sinai entropy if the acting group is amenable [6] and with the Shannon entropy of the base if the system is a Bernoulli shift [4]. For some systems, the sofic entropy can be finite relative to some sofic approximations and $-\infty$ relative to others. It is unknown whether two deterministic sofic approximations can yield different finite entropy values for the same system.

In this paper, we express the entropy relative to a type of stochastic block model in terms of the relative *f*-invariant, which we now introduce.

If α , β are two finite observables with codomains A, B, the conditional entropy is

$$H_{\mu}(\alpha|\beta) = H_{\mu}(\alpha\beta) - H_{\mu}(\beta).$$

This can be interpreted as the expected amount of information revealed by observing α if both the value of β and the joint distribution of α and β are known. The relative f-invariant

is defined by

$$\begin{split} F_{\mu}(T,\alpha|\beta) &= F_{\mu}(T,\alpha\beta) - F_{\mu}(T,\beta) \\ &= (1-2r)\mathrm{H}_{\mu}(\alpha|\beta) + \sum_{i=1}^{r} \mathrm{H}_{\mu}(\alpha^{\{e,s_i\}} \mid \beta^{\{e,s_i\}}), \\ f_{\mu}(T,\alpha|\beta) &= \inf_{k_1 \in \mathbb{N}} \sup_{k_2 \in \mathbb{N}} F_{\mu}(T,\alpha^{k_1} \mid \beta^{k_2}). \end{split}$$

Both the infimum and supremum can be replaced by limits; this follows from Lemma 3.2 below. It follows from Corollary 3.5 that we could also directly define

$$f_{\mu}(T, \alpha|\beta) = f_{\mu}(T, \alpha\beta) - f_{\mu}(T, \beta),$$

as long as $f_{\mu}(T, \beta) > -\infty$.

We now define the relevant type of stochastic block model. If H is a finite subset of G, we denote by $d^H(\mu, \nu)$ the total variation distance between the marginals of μ and ν on \mathbb{A}^H . Our convention for the total variation distance between measures $\mu, \nu \in \operatorname{Prob}(\mathbb{A})$ is

$$\|\mu - \nu\|_{\text{TV}} = \frac{1}{2} \sum_{a \in \lambda} |\mu\{a\} - \nu\{a\}|.$$

For each $k \in \mathbb{N}$ we define a pseudometric on Prob(\mathbb{A}^G) by

$$d_k^*(\mu, \nu) = \sum_{i \in [r]} d^{B(e,k) \cup B(s_i,k)}(\mu, \nu).$$

Note that $\{d_k^*\}_{k\in\mathbb{N}}$ together generate the weak* topology on $\operatorname{Prob}(\mathbb{A}^G)$. These generalize the function d_σ^* from [2], which corresponds to the case k=0. For $O=\{\nu\in\operatorname{Prob}(\mathbb{A}^G):d_k^*(\alpha_*^G\mu,\nu)<\varepsilon\}$ we write

$$\Omega(\sigma, O) =: \Omega_{\iota}^*(\sigma, \alpha, \varepsilon) \subseteq \mathbb{A}^n.$$

Our stochastic block model is now defined as follows: given $\mathbf{y}_0 \in \mathbb{B}^n$, $\sigma_0 \in \text{Hom}(G, \text{Sym}(n))$, and $k \in \mathbb{N}$, let

$$\mathrm{SBM}(\sigma_0, \mathbf{y}_0, k) := \mathrm{Unif}(\{\sigma \in \mathrm{Hom}(G, \mathrm{Sym}(n)) : d_k^*(P_{\mathbf{y}_0}^{\sigma}, P_{\mathbf{y}_0}^{\sigma_0}) = 0\}).$$

The labeling y_0 partitions the elements of [n] into |B| communities, and we can think of the random homomorphism σ as a random choice of directed edges between and within the communities. Certain statistics of these random edge choices are determined by the reference homomorphism σ_0 ; note that for k > 0 these statistics are more precise than those specified by a standard stochastic block model. In §2 we define weights, which are the objects used to record the relevant statistics.

1.1. Main results. Our main theorems show that the relative f-invariant can be interpreted as the growth rate of the expected number of ways to extend a planted good model for β to a good model for $\alpha\beta$, over a stochastic block model which has statistics determined by β and its planted model.

We first prove that if $\beta_*^G \mu$ is Markov then we can use a stochastic block model which only takes into account 'one-step statistics'.

THEOREM A. Let $\alpha: X \to \mathbb{A}$ and $\beta: X \to \mathbb{B}$ be finite observables, and for each n let $\mathbf{y}_n \in \mathbb{B}^n$ and $\sigma_n \in \text{Hom}(G, \text{Sym}(n))$ be such that

$$\lim_{n\to\infty} d_0^*(P_{\mathbf{y}_n}^{\sigma_n}, \beta_*^G \mu) = 0.$$

Suppose that $\beta_*^G \mu$ is a Markov measure. With $s_n = SBM(\sigma_n, \mathbf{y}_n, 0)$, we have

$$f_{\mu}(T, \alpha \mid \beta) = \inf_{O \ni (\alpha\beta)_{*,\mu}^{G}} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_{n}}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^{n} : (\mathbf{x}, \mathbf{y}_{n}) \in \Omega(\sigma, O)\}|.$$

PROPOSITION A. The assumptions of Theorem A are non-vacuous; that is, for any finite observable $\beta: X \to \mathbb{B}$ there exist sequences $\{\mathbf{y}_n \in \mathbb{B}^n\}_{n=1}^{\infty}$ and $\{\sigma_n \in \operatorname{Hom}(G,\operatorname{Sym}(n))\}_{n=1}^{\infty}$ such that $\lim_{n\to\infty} d_0^*(P_{\mathbf{y}_n}^{\sigma_n},\beta_{\mu}^G\mu) = 0$.

This follows from the fact that free group actions are 'sofic', which is proven for example in [10, 14, 15]. A more elementary proof is given in §4 below.

If $\beta_*^G \mu$ is not Markov, then the same formula holds with a more precise type of stochastic block model.

THEOREM B. Let $\alpha: X \to \mathbb{A}$ and $\beta: X \to \mathbb{B}$ be finite observables. Let m_n approach infinity as n goes to infinity while satisfying $m_n = o(\log \log n)$. For each n, let $\mathbf{y}_n \in \mathbb{B}^n$ and $\sigma_n \in \text{Hom}(G, \text{Sym}(n))$ be such that

$$d_{m_n}^*(P_{\mathbf{y}_n}^{\sigma_n}, \beta_*^G \mu) = O\left(\frac{1}{\log n}\right).$$

Suppose that $f_{\mu}(T, \beta) > -\infty$. With $s_n = SBM(\sigma_n, \mathbf{y}_n, m_n)$,

$$f_{\mu}(T, \alpha \mid \beta) = \inf_{\substack{O \ni (\alpha\beta)^{G}_{u} \\ n \to \infty}} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_{n}}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^{n} : (\mathbf{x}, \mathbf{y}_{n}) \in \Omega(\sigma, O)\}|.$$

PROPOSITION B. The assumptions of Theorem B are non-vacuous; that is, for any finite observable $\beta: X \to \mathbb{B}$ and any sequence $\{m_n \in \mathbb{N}\}_{n=1}^{\infty}$ approaching infinity while satisfying $m_n = o(\log \log n)$, there exist sequences $\{\mathbf{y}_n \in \mathbb{B}^n\}_{n=1}^{\infty}$ and $\{\sigma_n \in \operatorname{Hom}(G, \operatorname{Sym}(n))\}_{n=1}^{\infty}$ such that $\lim_{n \to \infty} d_{m_n}^*(P_{\mathbf{y}_n}^{\sigma_n}, \beta_{\mathbf{y}_n}^G \mu) = O(1/\log n)$.

Using Theorem B, we prove the following formula for the growth rate of the expected number of good models over a stochastic block model. This can be compared to the variational principle in [12], and has a similar proof.

THEOREM C. Let s_n , α , β be as in the statement of Theorem B. Then

$$\inf_{O\ni\alpha_*^G\mu}\limsup_{n\to\infty}\frac{1}{n}\log\underset{\sigma\sim s_n}{\mathbb{E}}|\Omega(\sigma,O)|=\sup_{\lambda\in\mathsf{J}(\alpha_*^G\mu,\,\beta_*^G\mu)}f_\lambda(S,\,\mathsf{a}\mid\mathsf{b}).$$

Here $J(\alpha_*^G\mu, \beta_*^G\mu)$ is the set of joinings of the G-systems $(A^G, \alpha_*^G\mu, S)$ and $(B^G, \beta_*^G\mu, S)$, that is, shift-invariant probability measures on $(A \times B)^G$ whose A^G, B^G marginals are $\alpha_*^G\mu$, $\beta_*^G\mu$, respectively. S denotes the shift action of G. We use a, b to denote the maps

$$a: (A \times B)^G \to A$$

 $((a_g, b_g))_{g \in G} \mapsto a_e$

and

b:
$$(A \times B)^G \to B$$

 $((a_g, b_g))_{g \in G} \mapsto b_e$,

which observe the A (respectively, B) label at the identity.

Note that the supremum is always greater than or equal to $f_{\mu}(T,\alpha)$, with equality attained by the product joining; this means that the expected number of good models for α over a block model with built-in good models for any β is at least the expected number of good models over a uniformly random homomorphism. It is possible for the supremum to be strictly larger, however. For example, suppose $f_{\mu}(T,\alpha) < 0$ and $\alpha = \beta$, and let λ be the diagonal joining. Then

$$f_{\lambda}(S, \mathbf{a} \mid \mathbf{b}) = 0 > f_{\mu}(T, \alpha).$$

1.2. Related work. The expressions appearing on the right-hand sides of Theorems A and B are very closely related to Ben Hayes' definition of 'relative sofic entropy in the presence' [11, Definition 2.5]. Some differences are that we consider *expected* numbers of good models over *random* sofic approximations, and that Hayes takes a supremum inside the logarithm over which good model is to be extended, while we fix a sequence $\{y_n\}$ of planted good models. Hayes also does not restrict to shift systems as we do here.

In [8], the free energy (that is, the limit of normalized log partition functions) over a type of stochastic block model is shown to satisfy a variational principle; see Propositions 3.6 and 3.7 of that paper.

1.3. Random sofic approximations. As noted above, the f-invariant is closely related to another invariant of measure-preserving systems called sofic entropy, which was introduced by Bowen in [4].

A homomorphism $\sigma \in \text{Hom}(G, \text{Sym}(n))$ is called (D, δ) -sofic for some finite $D \subset G$ and $\delta > 0$ if

$$|\{j \in [n] : \sigma(\gamma) | j \neq j \text{ for all } \gamma \in D \setminus \{e\}\}| > (1 - \delta)n.$$

A sequence of homomorphisms $\Sigma = (\sigma_n \in \text{Hom}(G, \text{Sym}(n)))_{n \in \mathbb{N}}$ is called a sofic approximation if, for every (D, δ) , the homomorphism σ_n is (D, δ) -sofic for all large enough n.

The sofic entropy relative to Σ is the exponential growth rate of the number of good models over σ_n . Specifically, for any finite observable α on X we have

$$h_{\Sigma,\mu}(T,\alpha) = \inf_{O \ni \alpha_*^G \mu} \limsup_{n \to \infty} \frac{1}{n} \log |\Omega(\sigma_n, O)|.$$

This is an isomorphism invariant of the system (X, μ, T) if α is any *generating* observable, that is if the σ -algebra of the measurable space X is the coarsest one which is shift-invariant and α -measurable.

By analogy with this expression, we might call the sequences of random homomorphisms appearing in expressions above 'random sofic approximations'. The following proposition provides further justification for this terminology.

PROPOSITION 1.1. If (s_n) is any of the sequences appearing in Theorems A, B, and C, then for any (D, δ) there exists $\varepsilon > 0$ such that

$$\mathbb{P}_{\sigma \sim \mathsf{s}_n}(\sigma \text{ is } (D, \delta) \text{-sofic}) \geq 1 - n^{-\varepsilon n}$$

for all large enough n.

In particular, if $\sigma_1 \sim s_1$, $\sigma_2 \sim s_2$ etc. are independent then (σ_n) is a sofic approximation with probability 1.

1.4. Organization. In §2 we define weights and discuss some of their useful properties. In §3 we prove a few basic results about the functions f and F. Some of the results of these two sections are used in §4 to show that the assumptions of the main theorems are not vacuous. In §5 we show how the function F is related to the number of homomorphism-labeling pairs (σ, \mathbf{y}) that realize a given weight, which is the main ingredient of the proofs of Theorems A and B given in the next two sections. In §8 we show how to deduce Theorem C from Theorem B. Section 9 contains a proof of Proposition 1.1. The final section contains a proof of Lemma 2.3, which asserts that a weight can be approximated by a denominator-n weight with a specified marginal.

2. Weights

If $\alpha: X \to A$ is a finite observable, for $a, a' \in A$ and $i \in [r]$ let

$$W_{\alpha}(a, a'; i) = \alpha_*^{\{e, s_i\}} \mu(a, a') = \mu\{x \in X : \alpha(x) = a, \alpha(T_{s_i} x) = a'\}$$

and also denote

$$W_{\alpha}(a) = \alpha_* \mu(a).$$

For $\mathbf{x} \in \mathbb{A}^n$ and $\sigma \in \text{Hom}(G, \text{Sym}(n))$ let

$$W_{\sigma,\mathbf{x}}(a,a';i) = P_{\mathbf{x}}^{\sigma,\{e,s_i\}}(a,a')$$

and $W_{\sigma,\mathbf{x}}(a) = P_{\mathbf{x}}^{\sigma,\{e\}}(a)$. This could equivalently be defined as a special case of the previous construction, with σ specifying an action on X = [n] with an observable $\mathbf{x}:[n] \to \mathbb{A}$.

More abstractly, any $W \in (\text{Prob}(\mathbb{A}^2))^r$ is called an \mathbb{A} -weight if

$$\sum_{a' \in \mathbb{A}} W(a, a'; i) = \sum_{a' \in \mathbb{A}} W(a', a; j)$$

for all $i, j \in [r]$ and $a \in A$. For each $a \in A$ we denote this common value by W(a). Note that the objects W_{α} and $W_{\sigma, \mathbf{x}}$ defined above satisfy this condition.

We say that W has denominator n if $n \cdot W(a, a'; i) \in \mathbb{N}$ for all a, a', i.

The measures $W(\cdot, \cdot; i)$ for $i \in [r]$ are called the *edge measures* of W, and $W(\cdot)$ is called the *vertex measure*.

For any alphabet A, we use the metric on A-weights defined by

$$\begin{split} d(W_1, W_2) &:= \sum_{i \in [r]} \|W_1(\cdot, \cdot; i) - W_2(\cdot, \cdot; i)\|_{\text{TV}} \\ &= \frac{1}{2} \sum_{i \in [r]} \sum_{a, a' \in \mathbb{A}} |W_1(a, a'; i) - W_2(a, a'; i)|. \end{split}$$

We can use weights to count good models up to equivalence under the pseudometrics d_k^* using the following proposition.

PROPOSITION 2.1. If $\sigma \in \text{Hom}(G, \text{Sym}(n))$ and $\mathbf{x} \in \mathbb{A}^n$, then for any observable $\alpha: X \to \mathbb{A}$,

$$d(W_{\sigma,\mathbf{X}^k},W_{\alpha^k}) = d_k^*(P_{\mathbf{X}}^{\sigma},\alpha_*^G\mu).$$

Note this implies also that

$$d_k^*(P_{\mathbf{x}}^{\sigma}, \alpha_*^G \mu) = d_0^*(P_{\mathbf{x}^k}^{\sigma}, (\alpha^k)_*^G \mu).$$

Proof. By definition of the distance between weights,

$$d(W_{\sigma,\mathbf{x}^k}, W_{\alpha^k}) = \frac{1}{2} \sum_{i \in [r]} \sum_{\mathbf{a}, \mathbf{a}' \in \mathbb{A}^{\mathrm{B}(e,k)}} |W_{\sigma,\mathbf{x}^k}(\mathbf{a}, \mathbf{a}'; i) - W_{\alpha^k}(\mathbf{a}, \mathbf{a}'; i)|$$

$$= \frac{1}{2} \sum_{i \in [r]} \sum_{\mathbf{a}, \mathbf{a}' \in \mathbb{A}^{\mathrm{B}(e,k)}} \left| \frac{1}{n} \right| \left\{ j \in [n] : \frac{(\mathbf{x}^k)_j = \mathbf{a}}{(\mathbf{x}^k)_{\sigma(s_i)j} = \mathbf{a}'} \right\} \left|$$

$$- \mu \left\{ x \in X : \frac{\alpha^k(x) = \mathbf{a}}{\alpha^k(T_{s_i}x) = \mathbf{a}'} \right\} \right|.$$

For many 'incompatible' pairs \mathbf{a} , \mathbf{a}' , both terms will be zero: suppose $g \in B(e, k) \cap B(s_i, k)$, so that $gs_i^{-1} \in B(e, k)$. If the second term in the absolute value is non-zero, then for some $x \in X$ we have $\alpha^k(x) = \mathbf{a}$ and $\alpha^k(T_{s_i}x) = \mathbf{a}'$, and therefore

$$\mathbf{a}'_{gs_i^{-1}} = (\alpha^k(T_{s_i}x))_{gs_i^{-1}} = \alpha(T_{gs_i^{-1}}T_{s_i}x) = \alpha(T_gx) = (\alpha^k(x))_g = \mathbf{a}_g.$$

The same argument shows that $\mathbf{a}'_{gs_i^{-1}} = \mathbf{a}_g$ for all $g \in B(e, k) \cap B(s_i, k)$ whenever the first term is non-zero. Therefore we can restrict the sum to pairs \mathbf{a} , \mathbf{a}' with $\mathbf{a}'_{gs_i^{-1}} = \mathbf{a}_g$ for all $g \in B(e, k) \cap B(s_i, k)$. Equivalently, we can sum over all $\mathbf{A} \in \mathbb{A}^{B(e,k) \cup B(s_i,k)}$ to get

$$d(W_{\sigma,\mathbf{x}^{k}}, W_{\alpha^{k}}) = \frac{1}{2} \sum_{i \in [r]} \sum_{\mathbf{A} \in \mathbb{A}^{B(e,k) \cup B(s_{i},k)}} \left| \frac{1}{n} | \{j \in [n] : (\mathbf{x}^{B(e,k) \cup B(s_{i},k)})_{j} = \mathbf{A} \} \right|$$

$$- \mu \{x \in X : \alpha^{B(e,k) \cup B(s_{i},k)}(x) = \mathbf{A} \}$$

$$= \sum_{i \in [r]} d^{B(e,k) \cup B(s_{i},k)}(P_{\mathbf{x}}^{\sigma}, \alpha_{*}^{G}\mu).$$

It will be useful to consider the pushforward map induced by a map between alphabets: if $\pi : A \to B$ is a measurable map and W is an A-weight, then πW is the B-weight given by

$$\pi W(b,b';i) = \sum_{a \in \pi^{-1}\{b\}} \sum_{a' \in \pi^{-1}\{b'\}} W(a,a';i).$$

Note that this implies that the vertex measure of W is

$$\pi W(b) = \sum_{a \in \pi^{-1}\{b\}} W(a).$$

For example, let $\pi_B: A \times B \to B$ be the projection map. If W is an $A \times B$ -weight then $\pi_{\rm B}W$ is given by

$$\pi_{\mathrm{B}}W(b_1) = \sum_{a \in \mathbb{A}} W((a,b_1)) \quad \pi_{\mathrm{B}}W(b_1,b_2;i) = \sum_{a_1,a_2 \in \mathbb{A}} W((a_1,b_1),(a_2,b_2);i).$$

We call this the B-marginal of W.

All weights in the present paper will be over alphabets of the form $A^{B(e,k)} \times B^{B(e,k')}$. We use this fact to introduce some simplified notation for projections.

- π_A denotes projection onto the entire A factor $A^{B(e,k)}$; π_B is used similarly.
- For m < k and m' < k', π_{m,m'} denotes projection onto A^{B(e,m)} × B^{B(e,m')}.
 π_m denotes the projection A^{B(e,k)} → A^{B(e,m)}, except that if m = 0 we write π_e. We define F(W) for an abstract weight W by

$$F(W) = (1 - 2r)H(W(\cdot)) + \sum_{i \in [r]} H(W(\cdot, \cdot; i))$$

where H is the Shannon entropy. Note that this is consistent with the above definitions in that, for example,

$$F(W_{\alpha}) = F_{\mu}(T, \alpha).$$

We can revisit the definition of our version of the stochastic block model using weights. Let $H \subset G$ and let W be a denominator-n $B^{B(e,k)}$ -weight. Suppose there exist $\mathbf{y} \in B^n$ and $\sigma \in \text{Hom}(G, \text{Sym}(n))$ such that $W = W_{\sigma, \mathbf{v}^k}$. Then

$$\mathrm{SBM}(\sigma,\mathbf{y},k) = \mathrm{Unif}(\{\sigma' \in \mathrm{Hom}(G,\mathrm{Sym}(n)): W_{\sigma',\mathbf{y}^k} = W\}),$$

so we can also denote this distribution by SBM(y, W). Specifying the distribution by a weight rather than a specific homomorphism will occasionally be more convenient.

2.1. Constructing weights and good models. We borrow the first result of this type from [2]; it allows us to find a denominator-*n* approximation to a given weight.

LEMMA 2.2. (Lemma 2.3 of [2]) There is a constant C such that for any A-weight W there is a denominator-n A-weight within distance $C|A|^2r/n$ of W.

The following lemma allows us not only to construct a denominator-n approximation to a given weight, but also to specify a marginal of this approximation:

LEMMA 2.3. Let W be an $A \times B$ -weight. If W_B is a B-weight of denominator n with $d(W_{\rm B},\pi_{\rm B}W)<\delta$ then there is an ${\rm A}\times{\rm B}$ -weight $W_{\rm AB}$ with denominator n such that $\pi_{\rm B}W_{\rm AB} = W_{\rm B}$ and $d(W_{\rm AB}, W) < 265r(|A \times B|^2/n + \delta)$.

The construction is fairly involved, so it is postponed to §10. The constant 265 is not intended to be optimal.

The definition of a weight W_{σ,\mathbf{x}^k} in terms of a homomorphism σ and a labeling \mathbf{x} is straightforward. However, we will also need to know whether a given weight can be realized in this way. The next two results address this inverse problem.

PROPOSITION 2.4. If W is a denominator-n A-weight, then there exist $\mathbf{x} \in \mathbb{A}^n$ and $\sigma \in \text{Hom}(G, \text{Sym}(n))$ such that $W = W_{\sigma, \mathbf{x}}$.

П

Proof. This is implied by Proposition 2.1 of [2].

Unfortunately, this does not imply that for every denominator- $n \ \mathbb{A}^{B(e,k)}$ -weight W there is some $\sigma \in \text{Hom}(G, \text{Sym}(n))$ and $\mathbf{x} \in \mathbb{A}^n$ such that $W = W_{\sigma,\mathbf{x}^k}$; instead it provides $\mathbf{X} \in (\mathbb{A}^{B(e,k)})^n$ such that $W = W_{\sigma,\mathbf{X}}$.

However, if we already know that W is close to a weight of the form W_{α^k} for some observable α , then the following proposition shows that W is also close to a weight of the form W_{σ,\mathbf{x}^k} .

PROPOSITION 2.5. Let $\alpha: X \to \mathbb{A}$, $\sigma \in \text{Hom}(G, \text{Sym}(n))$, and $\mathbf{X} \in (\mathbb{A}^{B(e,k)})^n$ be such that $d(W_{\sigma,\mathbf{X}}, W_{\sigma^k}) \leq \varepsilon$ for some $\varepsilon \geq 0$. Writing $\mathbf{x} = \pi_e \mathbf{X} \in \mathbb{A}^n$, we have

$$d(W_{\sigma,\mathbf{X}}, W_{\sigma,\mathbf{X}^k}) \le 2r|\mathbf{B}(e,k)|\varepsilon.$$

An immediate consequence is that $\mathbf{X} \in \Omega_0^*(\sigma, \alpha^k, \varepsilon)$ implies $\pi_e \mathbf{X} \in \Omega_k^*(\sigma, \alpha, c\varepsilon)$ where $c = 1 + 2r|\mathbf{B}(e, k)|$; cf. Claim 2 in the proof of Proposition 3.2 of [2].

Proof. Claim 4 in the proof of Proposition 3.2 of [2] implies that

$$|\{j \in [n] : \mathbf{X}(j) \neq \mathbf{x}^k(j)\}| \le n|\mathbf{B}(e,k)|\varepsilon.$$

It follows that for any $i \in [r]$,

$$\begin{aligned} |\{j \in [n] : \mathbf{X}^{\{e, s_i\}}(j) \neq (\mathbf{x}^k)^{\{e, s_i\}}(j)\}| \\ &\leq |\{j \in [n] : \mathbf{X}(j) \neq \mathbf{x}^k(j)\}| + |\{j \in [n] : \mathbf{X}(\sigma(s_i)j) \neq \mathbf{x}^k(\sigma(s_i)j)\}| \\ &< 2n|\mathbf{B}(e, k)|\varepsilon, \end{aligned}$$

so

$$d(W_{\sigma,\mathbf{X}}, W_{\sigma,\mathbf{x}^k}) = \sum_{i \in [r]} \|(\mathbf{X}^{\{e,s_i\}})_* \operatorname{Unif}(n) - ((\mathbf{x}^k)^{\{e,s_i\}})_* \operatorname{Unif}(n)\|_{\mathrm{TV}}$$

$$\leq \sum_{i \in [r]} 2|B(e,k)|\varepsilon = 2r|B(e,k)|\varepsilon.$$

The following corollary of the first part of the proof will be useful later. It says that if the weight $W_{\sigma,\mathbf{X}}$ generated by some $\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n$ and $\sigma \in \mathrm{Hom}(G,\mathrm{Sym}(n))$ is exactly attainable in some sense, then \mathbf{X} can be exactly recovered from σ and the projection $\pi_e\mathbf{X} \in \mathbb{A}^n$.

COROLLARY 2.6. Suppose that $\sigma \in \text{Hom}(G, \text{Sym}(n))$ and $\mathbf{X} \in (\mathbb{A}^{B(e,k)})^n$ are such that either

- (1) $W_{\sigma,\mathbf{X}} = W_{\alpha^k}$ for some $\alpha: X \to A$, or
- (2) $W_{\sigma,\mathbf{X}} = W_{\sigma_0,\mathbf{X}_0^k}$ for some $\sigma_0 \in \text{Hom}(G, \text{Sym}(m))$ and $\mathbf{x}_0 \in \mathbb{A}^m$.

Then $(\pi_e \mathbf{X})^k = \mathbf{X}$.

Note that $(\pi_e \mathbf{X})^k$ is the k-neighborhood labeling generated from $\pi_e \mathbf{X}$ using σ , rather than σ_0 or some other homomorphism.

Proof. In the first case, we are in the setting of the previous proposition with $\varepsilon = 0$, so the first inequality of its proof gives the claimed result.

The second case is actually the same; this is only obscured somewhat by the notation. We are in the setting of the previous proposition with the space X = [m] having a G-action specified by σ_0 and a finite observable $\mathbf{x}_0:[m] \to \mathbb{A}$.

3. *Properties of F and f*

LEMMA 3.1. (Continuity as weight function) If W_1 , W_2 are A-weights with $d(W_1, W_2) \le \varepsilon < 1$ then

$$|F(W_1) - F(W_2)| \le 4r(H(\varepsilon) + \varepsilon \log_2 |A|),$$

where H(p) denotes the entropy of the probability measure $(p, 1 - p) \in Prob(\{0, 1\})$.

Proof. We use Fano's inequality in the following form (equation (2.139) of [9]). Suppose X, Y are A-valued random variables defined on the same probability space and let $p_e = \mathbb{P}(X \neq Y)$ be their probability of disagreement. Then

$$H(X \mid Y) \le H(p_e) + p_e \log |A|$$
.

Using the chain rule and non-negativity of Shannon entropy, we can deduce that

$$|H(X) - H(Y)| \le H(p_e) + p_e \log |A|$$
.

Let $\mu_1, \mu_2 \in \text{Prob}(A)$ be the respective distributions of X_1, X_2 . Because $\|\mu_1 - \mu_2\|_{\text{TV}}$ is the minimum value of $\mathbb{P}(X_1 \neq X_2)$ over all possible couplings, if $\|\mu_1 - \mu_2\|_{\text{TV}} < \varepsilon$ then

$$|H(\mu_1) - H(\mu_2)| \le H(\varepsilon) + \varepsilon \log |A|$$
.

The assumed bound $d(W_1, W_2) \le \varepsilon$ implies that each vertex and edge measure of W_1 is within total variation distance ε of its counterpart in W_2 , so

$$\begin{split} |F(W_1) - F(W_2)| &\leq |1 - 2r| \cdot |\mathrm{H}(W_1(\cdot)) - \mathrm{H}(W_2(\cdot))| \\ &+ \sum_{i \in [r]} |\mathrm{H}(W_1(\cdot, \cdot; i)) - \mathrm{H}(W_2(\cdot, \cdot; i))| \\ &\leq (2r - 1)(\mathrm{H}(\varepsilon) + \varepsilon \log |\mathbb{A}|) \\ &+ r \cdot (\mathrm{H}(\varepsilon) + \varepsilon \log |\mathbb{A}|^2) \\ &\leq 4r(\mathrm{H}(\varepsilon) + \varepsilon \log |\mathbb{A}|). \end{split}$$

Let $\alpha: X \to A$ and $\beta: X \to B$ be observables. We say that β is a *coarsening* of α if each part of the partition of X induced by β is a union of parts of the partition induced by α (up to null sets). Equivalently, there is some function $g: A \to B$ such that $\beta = g \circ \alpha$ almost surely. In this situation we can also call α a refinement of β .

A useful property of the Shannon entropy $H_{\mu}(\alpha)$ is monotonicity under refinement. The function F does not share this property, but it is monotone under the following particular kind of refinement introduced in [3].

We say that β is a *simple splitting* of α if there is some $s \in \{s_1^{\pm 1}, \ldots, s_r^{\pm 1}\}$ and a coarsening $\tilde{\alpha}$ of α such that, up to null sets, the partition induced by β is the coarsest common refinement of the partitions induced by α and $\tilde{\alpha} \circ T_s$.

We say that β is a *splitting* of α if there are observables $\alpha = \beta_0, \beta_1, \ldots, \beta_n = \beta$ such that β_i is a simple splitting of β_{i-1} for $i = 1, 2, \ldots, n$. We will use the following monotonicity properties of the relative version of F.

LEMMA 3.2. (Monotonicity under splitting)

- (1) If α_1 is a splitting of α_2 then $F(\alpha_1|\beta) < F(\alpha_2|\beta)$.
- (2) If β_1 is a splitting of β_2 then $F(\alpha|\beta_1) \geq F(\alpha|\beta_2)$.

Proof. (1) This is essentially Proposition 5.1 of [3]; conditioning on β makes no difference to the proof.

(2) The proof is based on the proof of part (1), but in place of the chain rule for conditional entropy we use the following bound:

$$\begin{aligned} H(\alpha \mid \beta_2) &\leq H(\alpha, \beta_1 \mid \beta_2) \quad \text{(monotonicity)} \\ &= H(\beta_1 \mid \beta_2) + H(\alpha \mid \beta_1, \beta_2) \quad \text{(chain rule)} \\ &\leq H(\beta_1 \mid \beta_2) + H(\alpha \mid \beta_1) \quad \text{(monotonicity)}. \end{aligned}$$

We will also use the following consequence of the previous bound:

$$\begin{split} & \text{H}(\alpha^{\{e,s_i\}} \mid \beta_1^{\{e,s_i\}}) - \text{H}(\alpha^{\{e,s_i\}} \mid \beta_2^{\{e,s_i\}}) \\ & \geq -\text{H}(\beta_1^{\{e,s_i\}} \mid \beta_2^{\{e,s_i\}}) \quad \text{(previous bound)} \\ & \geq -(\text{H}(\beta_1^{\{s_i\}} \mid \beta_2^{\{e,s_i\}}) + \text{H}(\beta_1 \mid \beta_2^{\{e,s_i\}})) \quad \text{(subadditivity)} \\ & = -(\text{H}(\beta_1 \mid \beta_2^{\{e,s_i^{-1}\}}) + \text{H}(\beta_1 \mid \beta_2^{\{e,s_i\}})) \quad (\textit{T-invariance of μ)}. \end{split}$$

It suffices to check the case where β_1 is a simple splitting of β_2 . Let $t \in \{s_1^{\pm 1}, \ldots, s_r^{\pm 1}\}$ and let $\tilde{\beta}$ be a coarsening of β_2 such that the partition induced by β_1 is the same as the coarsest common refinement of the partitions induced by β_2 and $\tilde{\beta} \circ T_t$ up to null sets. Then, using the two bounds just derived,

$$\begin{split} F(\alpha|\beta_1) - F(\alpha|\beta_2) &= (1 - 2r)(\mathrm{H}(\alpha|\beta_1) - \mathrm{H}(\alpha|\beta_2)) \\ &+ \sum_{i \in [r]} (\mathrm{H}(\alpha^{\{e,s_i\}}|\beta_1^{\{e,s_i\}}) - \mathrm{H}(\alpha^{\{e,s_i\}}|\beta_1^{\{e,s_i\}})) \end{split}$$

$$\geq (1 - 2r)(-H(\beta_1|\beta_2)) - \sum_{i \in [r]} (H(\beta_1 \mid \beta_2^{\{e, s_i^{-1}\}})$$

$$+ H(\beta_1 \mid \beta_2^{\{e, s_i\}}))$$

$$= (2r - 1)H(\beta_1|\beta_2) - \sum_{s \in \{s_1^{\pm 1} \dots s_r^{\pm 1}\}} H(\beta_1 \mid \beta_2^{\{e, s\}})$$

But

$$H(\beta_1 \mid \beta_2^{\{e,t\}}) \le H(\beta_1 \mid \beta_2 \tilde{\beta}^{\{t\}}) = 0,$$

so we can remove the t term from the sum to get

$$\begin{split} F(\alpha|\beta_1) - F(\alpha|\beta_2) &\geq (2r-1) \mathrm{H}(\beta_1|\beta_2) - \sum_{s \in \{s_1^{\pm 1} \dots s_r^{\pm 1}\} \setminus \{t\}} \mathrm{H}(\beta_1 \mid \beta_2^{\{e,s\}}) \\ &= \sum_{s \in \{s_1^{\pm 1} \dots s_r^{\pm 1}\} \setminus \{t\}} (\mathrm{H}(\beta_1|\beta_2) - \mathrm{H}(\beta_1 \mid \beta_2^{\{e,s\}})) \\ &\geq 0. \end{split}$$

One corollary is the following convenient formula.

COROLLARY 3.3. Let α , β be finite observables such that $\beta_*^G \mu$ is a Markov measure. Then $F_{\mu}(T, \alpha^{k_1} | \beta^{k_2})$ is independent of k_2 . In particular,

$$f_{\mu}(T, \alpha \mid \beta) = \inf_{k} F_{\mu}(T, \alpha^{k} \mid \beta).$$

Proof. By the previous proposition, for any $k \le k_2$ we have

$$F_{\mu}(T, \alpha^{k_1} | \beta^k) \leq F_{\mu}(T, \alpha^{k_1} | \beta^{k_2}).$$

On the other hand, by Theorem 6.1 of [5] $F_{\mu}(T, \beta^k) = F_{\mu}(T, \beta^{k_2})$ so

$$F_{\mu}(T, \alpha^{k_1} \mid \beta^k) = F_{\mu}(T, \alpha^{k_1} \beta^k) - F_{\mu}(T, \beta^{k_2}).$$

Applying monotonicity under splitting to the first term on the right gives

$$F_{\mu}(T, \alpha^{k_1} \mid \beta^k) \ge F_{\mu}(T, \alpha^{k_1} \beta^{k_2}) - F_{\mu}(T, \beta^{k_2}) = F_{\mu}(T, \alpha^{k_1} \mid \beta^{k_2}).$$

This establishes independence of k_2 ; the formula for f follows.

PROPOSITION 3.4. Let α , β be finite observables. Then for any $k \in \mathbb{N}$,

$$F_{\mu}(T, \alpha^k \mid \beta) \leq \mathrm{H}_{\mu}(\alpha \mid \beta).$$

It follows that

$$f_{\mu}(T, \alpha \mid \beta) \leq \mathrm{H}_{\mu}(\alpha \mid \beta).$$

Proof. By Lemma 3.2, $F_{\mu}(T, \alpha^k \mid \beta) \leq F_{\mu}(T, \alpha \mid \beta)$. Using elementary properties of Shannon entropy, we have

$$\begin{split} F_{\mu}(T,\alpha\mid\beta) &= (1-2r)\mathrm{H}_{\mu}(\alpha\mid\beta) + \sum_{i\in[r]} \mathrm{H}_{\mu}(\alpha^{\{e,s_i\}}\mid\beta^{\{e,s_i\}}) \\ &\leq (1-2r)\mathrm{H}_{\mu}(\alpha\mid\beta) + \sum_{i\in[r]} [\mathrm{H}_{\mu}(\alpha\mid\beta^{\{e,s_i\}}) + \mathrm{H}_{\mu}(\alpha^{\{s_i\}}\mid\beta^{\{e,s_i\}})] \\ &\leq (1-2r)\mathrm{H}_{\mu}(\alpha\mid\beta) + \sum_{i\in[r]} [\mathrm{H}_{\mu}(\alpha\mid\beta) + \mathrm{H}_{\mu}(\alpha^{\{s_i\}}\mid\beta^{\{s_i\}})]. \end{split}$$

By T-invariance of μ we have

$$H_{\mu}(\alpha^{\{s_i\}} \mid \beta^{\{s_i\}}) = H_{\mu}(\alpha \mid \beta),$$

so the first inequality follows.

For any $k_1, k_2 \in \mathbb{N}$ this gives

$$F_{\mu}(T, \alpha^{k_1} \mid \beta^{k_2}) \leq H_{\mu}(\alpha \mid \beta^{k_2}) \leq H_{\mu}(\alpha \mid \beta),$$

so the second inequality follows upon taking the supremum over k_2 then the infimum over k_1 .

We can use this bound to give a proof of the chain rule for the relative f-invariant, a version of which first appeared in [5] (there it is called the Abramov–Rokhlin formula; see also [7]).

COROLLARY 3.5. (Chain rule)

$$f_{\mu}(T, \alpha\beta) = f_{\mu}(T, \alpha \mid \beta) + f_{\mu}(T, \beta).$$

Proof. By definition of the relative version of F and the chain rule for conditional entropy, for each k_1 , k_2 we have

$$F_{\mu}(T, \alpha^{k_1} \beta^{k_2}) = F_{\mu}(T, \alpha^{k_1} \mid \beta^{k_2}) + F_{\mu}(T, \beta^{k_2}).$$

By Lemma 3.2 each term is monotone in k_2 , so the limits as $k_2 \to \infty$ exist. By Proposition 3.4 all terms are bounded above (recall we only consider finite observables, so in particular all observables have finite entropy), so we can split the limit across the sum on the right to get

$$\lim_{k_2 \to \infty} F_{\mu}(T, \alpha^{k_1} \beta^{k_2}) = \lim_{k_2 \to \infty} F_{\mu}(T, \alpha^{k_1} \mid \beta^{k_2}) + f_{\mu}(T, \beta).$$

Taking k_1 to infinity gives the result.

4. Non-vacuity of main theorems

4.1. Theorem A. Here we prove Proposition A, which asserts the non-vacuity of Theorem A. Given $\beta: X \to \mathbb{B}$, we need to show that there exist $\mathbf{y}_n \in \mathbb{B}^n$ and $\sigma_n \in \operatorname{Hom}(G,\operatorname{Sym}(n))$ such that $\lim_{n\to\infty} d_0^*(P_{\mathbf{y}_n}^{\sigma_n},\beta_*^G\mu)=0$.

By Lemma 2.2, there is a sequence $\{W_n\}_{n=1}^{\infty}$ of B-weights such that W_n has denominator n for each n and $d(W_n, W_\beta) = o(1)$. By Proposition 2.4, for each n we can pick \mathbf{y}_n , σ_n such that $W_{\sigma_n,\mathbf{y}_n} = W_n$. Since $d_0^*(P_{\mathbf{y}_n}^{\sigma_n}, \beta_*^G \mu) = d(W_{\sigma_n,\mathbf{y}_n}, W_\beta)$, these suffice.

4.2. *Theorems B and C*. Here we prove Proposition B, which asserts the non-vacuity of Theorem B (and by extension Theorem C, since the assumptions are the same).

Let m_n approach infinity as n approaches infinity while satisfying $m_n = o(\log \log n)$ and let $\beta: X \to \mathbb{B}$ be a finite observable. We need to show that there exist $\mathbf{y}_n \in \mathbb{B}^n$ and $\sigma_n \in \text{Hom}(G, \text{Sym}(n))$ such that $d_{m_n}^*(P_{\mathbf{y}_n}^{\sigma_n}, \beta_*^G \mu) = O(1/\log n)$.

By Lemma 2.2, there is a sequence $\{W_n\}_{n=1}^{\infty}$ of weights such that W_n is a denominator-n $\mathbb{B}^{\mathrm{B}(e,m_n)}$ -weight for each n and $d(W_n,W_{\beta^{m_n}})=O(|\mathbb{B}^{\mathrm{B}(e,m_n)}|^2/n)$. By Proposition 2.4, for each n we can pick \mathbf{Y}_n , σ_n such that $W_{\sigma_n,\mathbf{Y}_n}=W_n$. Let $\mathbf{y}_n=\pi_e\mathbf{Y}_n$. By Proposition 2.5,

$$d_{m_n}^*(P_{\mathbf{y}_n}^{\sigma_n}, \beta_*^G \mu) = d(W_{\sigma_n, \mathbf{y}_n^{m_n}}, W_{\beta^{m_n}}) = O\left(|\mathsf{B}(e, m_n)| \cdot \frac{|\mathsf{B}^{\mathsf{B}(e, m_n)}|^2}{n}\right) = O\left(\frac{1}{\log n}\right).$$

5. Counting lemmas

For a B-weight W, let $Z_n(W)$ denote the number of pairs $(\sigma, \mathbf{y}) \in \text{Hom}(G, \text{Sym}(n)) \times \mathbb{B}^n$ such that $W_{\sigma, \mathbf{y}} = W$.

PROPOSITION 5.1. *If W is a* B-weight with denominator n then

$$(3\sqrt{n})^{-r|\mathbf{B}|^2} \le \frac{Z_n(W)}{e^{F(W)n}(n!)^r n^{(1-r)/2}} \le (3\sqrt{n})^{r|\mathbf{B}|^2}.$$

Proof. We write

$$Z_n(W) = \sum_{\sigma} |\{\mathbf{y} \in \mathbb{B}^n : W_{\sigma,\mathbf{y}} = W\}| = (n!)^r \mathbb{E}_{\sigma} |\{\mathbf{y} \in \mathbb{B}^n : W_{\sigma,\mathbf{y}} = W\}|.$$

where \mathbb{E}_{σ} denotes the expectation over a uniform choice of $\sigma \in \text{Hom}(G, \text{Sym}(n))$. Proposition 2.1 of [2] states that

$$\mathbb{E}_{\sigma}|\{\mathbf{y} \in \mathbb{B}^{n} : W_{\sigma,\mathbf{y}} = W\}| = \frac{n!^{1-r} \prod_{b \in \mathbb{B}} (nW(b))!^{2r-1}}{\prod_{i=1}^{r} \prod_{b,b' \in \mathbb{B}} (nW(b,b';i))!}.$$

Lemma 2.2 of the same paper gives an estimate of this quantity, but for our purposes we need to be more careful about how the estimate depends on the size of the alphabet.

We use the version of Stirling's approximation given by

$$k^{k+1/2}e^{-k} \le k! \le 3 \cdot k^{k+1/2}e^{-k}$$

valid for $k \ge 1$. To estimate the products that appear in the expectation, we will need to omit all factors which equal 0! = 1 since Stirling's approximation is not valid for these. To do this carefully, let

$$\mathbf{B}' = \{ b \in \mathbf{B} : W(b) \neq 0 \}$$

and for each $i \in [r]$ let

$$B'_i = \{(b, b') \in B^2 : W(b, b'; i) \neq 0\}.$$

For the numerator of the above expectation we get

$$\begin{split} n!^{1-r} \prod_{b \in \mathbb{B}'} (nW(b))!^{2r-1} &\leq (3n^{n+1/2} \ e^{-n})^{1-r} \prod_{b \in \mathbb{B}'} (3(nW(b))^{nW(b)+1/2} e^{-nW(b)})^{2r-1} \\ &= 3^{1-r+|\mathbb{B}'|(2r-1)} \ n^{rn+1/2-r/2+(2r-1)|\mathbb{B}'|/2} \\ &\times e^{-rn+(2r-1)[n \sum_{b \in \mathbb{B}'} W(b) \log W(b) + 1/2 \sum_{b \in \mathbb{B}'} \log W(b)]} \end{split}$$

and a lower bound which is identical except missing the first factor. For the denominator, let $S = \sum_{i \in [r]} |B'_i|$. We get

$$\begin{split} \prod_{i=1}^{r} \prod_{(b,b') \in \mathbb{B}_{i}'} (nW(b,b';i))! &\leq \prod_{i=1}^{r} \prod_{(b,b') \in \mathbb{B}_{i}'} 3(nW(b,b';i))^{nW(b,b';i)+1/2} e^{-nW(b,b';i)} \\ &= 3^{S} n^{nr+S/2} \\ &\times e^{n \sum_{i} \sum_{b,b'} W(b,b';i) \log W(b,b';i) + 1/2 \sum_{i,b,b'} \log W(b,b';i) - nr} \end{split}$$

and again we have a lower bound which is identical except missing the first factor 3^S. Therefore the quotient is bounded above by

$$3^{1-r+|\mathsf{B}'|(2r-1)} \; n^{(1-r)/2+(2r-1)|\mathsf{B}'|/2-S/2} \; e^{nF(W)+(2r-1)\frac{1}{2} \sum_b \log W(b) - \frac{1}{2} \sum_{i,b,b'} \log W(b,b';i)}$$

and below by

$$3^{-S} \; n^{(1-r)/2 + (2r-1)|\mathbb{B}'|/2 - S/2} \; e^{nF(W) + (2r-1)\frac{1}{2} \sum_b \log W(b) - \frac{1}{2} \sum_{i,b,b'} \log W(b,b';i)}.$$

Since W has denominator n, we have

$$0 \ge (2r - 1)\frac{1}{2} \sum_{b \in \mathbb{B}'} \log W(b) \ge (2r - 1)\frac{1}{2} \sum_{b \in \mathbb{B}'} \log \frac{1}{n} = -\frac{2r - 1}{2} |\mathbb{B}'| \log n$$

and

$$0 \le -\frac{1}{2} \sum_{i} \sum_{(b,b') \in B'_i} \log W(b,b';i) \le -\frac{1}{2} \sum_{i} \sum_{(b,b') \in B'_i} \log \frac{1}{n} = \frac{S}{2} \log n.$$

Therefore $Z_n(W)$ satisfies

$$3^{-S}n^{((1-r)-S)/2}e^{F(W)n}(n!)^r < Z_n(W) < 3^{1-r+|B'|(2r-1)}n^{((1-r)+(2r-1)|B'|)/2}e^{F(W)n}(n!)^r.$$

Since $S \le r|B|^2$ and $|B'| \le |B|$, we conclude that

$$3^{-r|\mathbf{B}|^2} n^{((1-r)-r|\mathbf{B}|^2)/2} e^{F(W)n} (n!)^r$$

$$\leq Z_n(W) \leq 3^{1-r+|\mathbf{B}|(2r-1)} n^{((1-r)+(2r-1)|\mathbf{B}|)/2} e^{F(W)n} (n!)^r,$$

П

and the stated inequality follows.

The following proposition establishes the connection between the relative version of *F* and expected numbers of good models over stochastic block models.

PROPOSITION 5.2. Given any denominator-n ($\mathbb{A} \times \mathbb{B}^{B(e,k)}$)-weight W_{AB} , let W_B denote the $\mathbb{B}^{B(e,k)}$ -weight $\pi_B W_{AB}$. Let $\mathbf{y} \in \mathbb{B}^n$ be a fixed labeling with $p_{\mathbf{y}} = \pi_e W_B(\cdot)$, and let

$$\mu = \operatorname{SBM}(\mathbf{y}, W_{\mathrm{B}}) = \operatorname{Unif}(\{\sigma \in \operatorname{Hom}(G, \operatorname{Sym}(n)) : W_{\sigma, \mathbf{y}^k} = W_{\mathrm{B}}\}),$$

assuming W_B is such that the desired support is non-empty. Then

$$\mathcal{E} := \underset{\sigma \sim \mu}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^n : W_{\sigma, (\mathbf{x}, \mathbf{y}^k)} = W_{\text{AB}}\}| = \frac{Z_n(W_{\text{AB}})}{Z_n(W_{\text{B}})}.$$

In particular,

$$\frac{\mathcal{E}}{e^{n(F(W_{\mathbb{A}\mathbb{B}})-F(W_{\mathbb{B}}))}} \in ((9n)^{-r|\mathbb{B}|^2(|\mathbb{A}|^2+1)}, (9n)^{r|\mathbb{B}|^2(|\mathbb{A}|^2+1)}).$$

LEMMA 5.3. Let W_{AB} be an $A \times B^{B(e,k)}$ -weight of denominator n. Then

$$|\{(\sigma, \mathbf{x}, \mathbf{y}) : W_{\sigma, (\mathbf{x}, \mathbf{y}^k)} = W_{AB}\}| \in \{0, |\{(\sigma, \mathbf{x}, \mathbf{Y}) : W_{\sigma, (\mathbf{x}, \mathbf{Y})} = W_{AB}\}|\}.$$

Proof. Suppose $|\{(\sigma, \mathbf{x}, \mathbf{y}) : W_{\sigma, (\mathbf{x}, \mathbf{y}^k)} = W_{AB}\}| \neq 0$; we then need to show

$$|\{(\sigma,\mathbf{x},\mathbf{y}):W_{\sigma,(\mathbf{x},\mathbf{y}^k)}=W_{\mathrm{AB}}\}|=|\{(\sigma,\mathbf{x},\mathbf{Y}):W_{\sigma,(\mathbf{x},\mathbf{Y})}=W_{\mathrm{AB}}\}|.$$

The inequality \leq is clear, since we have an injection $(\sigma, \mathbf{x}, \mathbf{y}) \mapsto (\sigma, \mathbf{x}, \mathbf{y}^k)$.

The converse inequality holds because $(\sigma, \mathbf{x}, \mathbf{Y}) \mapsto (\sigma, \mathbf{x}, \pi_e \mathbf{Y})$ in an injection from the set on the right to the set on the left. This follows from Corollary 2.6.

Proof of Proposition 5.2. Let

$$\tilde{\mu} = \mathrm{Unif}(\{(\sigma, \tilde{\mathbf{y}}) : W_{\sigma, \tilde{\mathbf{y}}^k} = W_{\mathrm{B}}\}),$$

and let $\tilde{\mu}_2$ be its marginal on the ' $\tilde{\mathbf{y}}$ '-coordinate. This marginal is supported on { $\tilde{\mathbf{y}}$: $p_{\tilde{\mathbf{y}}} = \pi_e W_{\mathrm{B}}(\cdot)$ }. Note that $\tilde{\mu}$ conditioned on any particular $\tilde{\mathbf{y}}$ in the support of $\tilde{\mu}_2$ is SBM($\tilde{\mathbf{y}}$, W_{B}), and that

$$\underset{\sigma \sim \text{SBM}(\tilde{\mathbf{y}}, W_{\text{B}})}{\mathbb{E}} |\{\mathbf{x} \in \text{A}^n : W_{\sigma, (\mathbf{x}, \tilde{\mathbf{y}}^k)} = W_{\text{AB}}\}|$$

is the same for each $\tilde{\mathbf{y}}$ in the support of $\tilde{\mu}_2$, with one choice being \mathbf{y} from the proposition statement. This is because any two choices have the same label frequencies and hence are related by a permutation of [n]. With the choice $\tilde{\mathbf{y}} = \mathbf{y}$ the expectation is \mathcal{E} , so

$$\begin{split} \mathcal{E} &= \underset{\tilde{\mathbf{y}} \sim \tilde{\mu}_{2}}{\mathbb{E}} \mathcal{E} \\ &= \underset{\tilde{\mathbf{y}} \sim \tilde{\mu}_{2}}{\mathbb{E}} [\mathbb{E}_{\sigma \sim \text{SBM}(\tilde{\mathbf{y}}, W_{\mathbb{B}})} | \{\mathbf{x} \in \mathbb{A}^{n} : W_{\sigma, (\mathbf{x}, \tilde{\mathbf{y}}^{k})} = W_{\mathbb{A}\mathbb{B}} \} |] \\ &= \underset{(\sigma, \tilde{\mathbf{y}}) \sim \tilde{\mu}}{\mathbb{E}} | \{\mathbf{x} \in \mathbb{A}^{n} : W_{\sigma, (\mathbf{x}, \tilde{\mathbf{y}}^{k})} = W_{\mathbb{A}\mathbb{B}} \} | \\ &= \frac{\sum_{\sigma, \tilde{\mathbf{y}}} | \{\mathbf{x} \in \mathbb{A}^{n} : W_{\sigma, (\mathbf{x}, \tilde{\mathbf{y}}^{k})} = W_{\mathbb{A}\mathbb{B}} \} | }{| \{(\sigma, \tilde{\mathbf{y}}) : W_{\sigma, \tilde{\mathbf{y}}^{k}} = W_{\mathbb{B}} \} | } \\ &= \frac{| \{(\sigma, \mathbf{x}, \tilde{\mathbf{y}}) : W_{\sigma, (\mathbf{x}, \tilde{\mathbf{y}}^{k})} = W_{\mathbb{A}\mathbb{B}} \} | }{| \{(\sigma, \tilde{\mathbf{y}}) : W_{\sigma, (\tilde{\mathbf{y}}, \tilde{\mathbf{y}}^{k})} = W_{\mathbb{B}} \} | } \end{split}$$

$$= \frac{|\{(\sigma, \mathbf{x}, \mathbf{Y}) : W_{\sigma, (\mathbf{x}, \mathbf{Y})} = W_{\mathtt{AB}}\}|}{|\{(\sigma, \mathbf{Y}) : W_{\sigma, \mathbf{Y}} = W_{\mathtt{B}}\}|} \quad \text{(previous lemma)}$$

$$= \frac{Z_n(W_{\mathtt{AB}})}{Z_n(W_{\mathtt{B}})}.$$

Note that our assumption that the intended support of μ is non-empty allows us to rule out the '0' case in the application of Lemma 5.3.

The rest of the result then follows from our estimates on Z_n in Proposition 5.1.

6. Proof of Theorem A

6.1. *Upper bound*. Note that we will not rely on the Markov assumption for the upper bound.

For each $k \in \mathbb{N}$,

$$\begin{split} &\inf_{O\ni(\alpha\beta)_*^G\mu} \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim s_n}{\mathbb{E}} |\{\mathbf{x}\in \mathbb{A}^n: (\mathbf{x},\mathbf{y}_n)\in \Omega(\sigma,O)\}| \\ &\leq \inf_{\varepsilon} \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim s_n}{\mathbb{E}} |\{\mathbf{x}\in \mathbb{A}^n: (\mathbf{x},\mathbf{y}_n)\in \Omega_k^*(\sigma,\alpha\beta,\varepsilon)\}| \\ &= \inf_{\varepsilon} \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim s_n}{\mathbb{E}} |\{\mathbf{x}\in \mathbb{A}^n: (\mathbf{x}^k,\mathbf{y}_n^k)\in \Omega_0^*(\sigma,(\alpha\beta)^k,\varepsilon)\}| \\ &\leq \inf_{\varepsilon} \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim s_n}{\mathbb{E}} |\{\mathbf{X}\in (\mathbb{A}^{\mathrm{B}(e,k)})^n: (\mathbf{X},\mathbf{y}_n^k)\in \Omega_0^*(\sigma,(\alpha\beta)^k,\varepsilon)\}|. \end{split}$$

Write

$$\begin{split} \mathcal{E}_k(n,\varepsilon) &:= \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : (\mathbf{X},\mathbf{y}_n^k) \in \Omega_0^*(\sigma,(\alpha\beta)^k,\varepsilon)\}| \\ &= \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : d(W_{\sigma,(\mathbf{X},\mathbf{y}_n^k)},W_{(\alpha\beta)^k}) < \varepsilon)\}| \end{split}$$

and assume that *n* is large enough that $m_n \ge k$.

Writing $W_n(\alpha\beta, k, \varepsilon)$ for the set of all denominator-n weights W with $d(W, W_{(\alpha\beta)^k}) < \varepsilon$,

$$\begin{split} \mathcal{E}_k(n,\varepsilon) &= \underset{\sigma \sim_{\mathbb{S}_n}}{\mathbb{E}} \sum_{W \in \mathcal{W}_n(\alpha\beta,k,\varepsilon)} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : W_{\sigma,(\mathbf{X},\mathbf{y}_n^k)} = W\}| \\ &= \sum_{W \in \mathcal{W}_n(\alpha\beta,k,\varepsilon)} \underset{\sigma \sim_{\mathbb{S}_n}}{\mathbb{E}} [|\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : W_{\sigma,(\mathbf{X},\mathbf{y}_n^k)} = W\}||W_{\sigma,\mathbf{y}_n^k} = \pi_{\mathbb{B}}W] \\ &\cdot \underset{\sigma \sim_{\mathbb{S}_n}}{\mathbb{P}} (W_{\sigma,\mathbf{y}_n^k} = \pi_{\mathbb{B}}W) \end{split}$$

since if $W_{\sigma,\mathbf{y}_n^k} \neq \pi_B W$ then $W_{\sigma,(\mathbf{X},\mathbf{y}_n^k)} \neq W$. But s_n conditioned on $\{W_{\sigma,\mathbf{y}_n^k} = \pi_B W\}$ is SBM $(\mathbf{y}_n, \pi_B W)$, so we can bound the expectation above using Proposition 5.2, getting

$$\mathcal{E}_k(n,\varepsilon) \leq (9n)^{r|\mathsf{B}^{\mathsf{B}(e,k)}|^2(|\mathsf{A}^{\mathsf{B}(e,k)}|+1)} \sum_{W \in \mathcal{W}_n(\alpha\beta,k,\varepsilon)} e^{n(F(W)-F(\pi_\mathsf{B}W))} \mathop{\mathbb{P}}_{\sigma \sim \mathsf{s}_n}(W_{\sigma,\mathbf{y}_n^k} = \pi_\mathsf{B}W).$$

Note $(9n)^{r|\mathbb{B}^{\mathrm{B}(e,k)}|^2(|\mathbb{A}^{\mathrm{B}(e,k)}|+1)} \le e^{o_{n\to\infty}(n)}$. Fix $\delta > 0$. By continuity of F (Lemma 3.1), for all small enough ε (possibly depending on k), we have

$$\mathcal{E}_k(n,\varepsilon) \leq e^{n(F_\mu(T,\alpha^k|\beta^k) + \delta + o_{n \to \infty}(1))} \sum_{W \in \mathcal{W}_n(\alpha\beta,k,\varepsilon)} \mathbb{P}_{\sigma \sim s_n}(W_{\sigma,\mathbf{y}_n^k} = \pi_{\mathsf{B}}W).$$

Bounding each probability by 1, we get

$$\mathcal{E}_k(n,\varepsilon) \leq e^{n(F_{\mu}(T,\alpha^k|\beta^k) + \delta + o_{n\to\infty}(1))} |\mathcal{W}_n(\alpha\beta,k,\varepsilon)|.$$

But

$$|\mathcal{W}_n(\alpha\beta, k, \varepsilon)| \le n^{r|(\mathbb{A}\times\mathbb{B})^{\mathrm{B}(e,k)}|^2} \le e^{o_{n\to\infty}(n)},$$

so this implies

$$\limsup_{n \to \infty} \frac{1}{n} \log \mathcal{E}_k(n, \varepsilon) \le F_{\mu}(T, \alpha^k \mid \beta^k) + \delta$$
$$\le F_{\mu}(T, \alpha^k \mid \beta^{k_2}) + \delta$$

for any $k_2 \ge k$, by monotonicity under splitting. Taking the limit as $k_2 \to \infty$ followed by the infimum over ε (which takes δ to 0) and k gives

$$\inf_{\varepsilon,k} \limsup_{n \to \infty} \frac{1}{n} \log \mathcal{E}_k(n, \varepsilon) \le f_{\mu}(T, \alpha \mid \beta).$$

Since

$$\inf_{O\ni(\alpha\beta)_*^G\mu} \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim s_n}{\mathbb{E}} |\{\mathbf{x}\in \mathbb{A}^n : (\mathbf{x},\mathbf{y}_n)\in \Omega(\sigma,O)\}|$$

$$\leq \inf_{\varepsilon} \limsup_{n\to\infty} \frac{1}{n} \log \mathcal{E}_k(n,\varepsilon)$$

for every k, this completes the upper bound.

6.2. Lower bound. Fix $k \in \mathbb{N}$. To estimate

$$\mathcal{E} := \underset{\sigma \sim S_n}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^n : (\mathbf{x}, \mathbf{y}_n) \in \Omega_k^*(\sigma, \alpha\beta, \varepsilon)\}|$$

we bound below using the expected size of

$$X_k(\sigma,\alpha\beta,\varepsilon\mid\mathbf{y}_n):=\{\mathbf{X}\in(\mathbb{A}^{\mathrm{B}(e,k)})^n:(\mathbf{X},\mathbf{y}_n^k)\in\Omega_0^*(\sigma,(\alpha\beta)^k,\varepsilon)\}.$$

This is not a true lower bound but, by equation (7.1) below, there are constants C, d, c independent of n such that

$$|\mathcal{X}_k(\sigma, \alpha\beta, \varepsilon \mid \mathbf{y}_n)| \le C \exp(nd\varepsilon + nH(2|B(e, k)|\varepsilon))$$
$$\cdot |\{\mathbf{x} \in \mathbb{A}^n : (\mathbf{x}, \mathbf{y}_n) \in \Omega_k^*(\sigma, \alpha\beta, \varepsilon)\}|.$$

The 'error' factor has an exponential growth rate which vanishes as $\varepsilon \to 0$, so will not be a problem.

We now find a lower bound for the expectation of $|X_k|$. Applying Proposition 5.2 as above, we have

$$\begin{split} & \underset{\sigma \sim s_n}{\mathbb{E}} |\mathcal{X}_k(\sigma, \alpha\beta, \varepsilon \mid \mathbf{y}_n)| \\ & = \sum_{W \in \mathcal{W}_n(\alpha\beta, k, \varepsilon)} \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : W_{\sigma, (\mathbf{X}, \mathbf{y}_n^k)} = W\}| \\ & \geq \sum_{W \in \mathcal{W}_n(\alpha\beta, k, \varepsilon)} \exp[n(F(W) - F(\pi_{\mathrm{B}}W) - o_n(1))] \underset{\sigma \sim s_n}{\mathbb{P}} (\pi_{\mathrm{B}}W = W_{\sigma, \mathbf{y}_n^k}). \end{split}$$

For any $\delta > 0$, for small enough $\varepsilon > 0$ (independent of n), by continuity of F this is at least

$$\exp[n(F_{\mu}(\alpha^{k} \mid \beta^{k}) - \delta - o_{n}(1))] \sum_{W \in W_{n}(\alpha\beta,k,\varepsilon)} \mathbb{P}_{\alpha \sim S_{n}}(\pi_{\mathsf{B}}W = W_{\sigma,\mathbf{y}_{n}^{k}}).$$

We give a lower bound for the sum by first rewriting it as

$$\sum_{W_{\mathbb{B}} \text{ denom.-} n} |\{W \in \mathcal{W}_n(\alpha\beta, k, \varepsilon) : \pi_{\mathbb{B}}W = W_{\mathbb{B}}\}| \cdot \underset{\sigma \sim_{\mathbb{S}_n}}{\mathbb{P}} (W_{\sigma,\mathbf{y}_n^k} = W_{\mathbb{B}}).$$

Fix $\eta > 0$. By Lemma 2.3, for all large enough n the B-weight $W_{\sigma_n,\mathbf{y}_n}$ can be extended to a $\mathsf{B}^{\mathsf{B}(e,k)}$ -weight W_B with $d(W_\mathsf{B},W_{\beta^k}) \leq \eta$; to apply the lemma we can think of the extended weight W_B as having alphabet $\mathsf{B}^{\mathsf{B}(e,k)\setminus\{e\}}\times\mathsf{B}$, and recall that we assume $\lim_{n\to\infty}d(W_{\sigma_n,\mathbf{y}_n},W_\beta)=0$. Choose σ , \mathbf{Y} such that $W_{\sigma,\mathbf{Y}}=W_\mathsf{B}$. Since $\pi_eW_\mathsf{B}=W_{\sigma_n,\mathbf{y}_n}$, it must be that $\pi_e\mathbf{Y}$ is a permutation of \mathbf{y}_n : they must assign labels with the same frequencies since

$$p_{\pi_e \mathbf{Y}}(\cdot) = (\pi_e W_{\mathrm{B}})(\cdot) = W_{\sigma_n, \mathbf{V}_n}(\cdot) = p_{\mathbf{V}_n}(\cdot).$$

Therefore we can choose σ , **Y** such that π_e **Y** = \mathbf{y}_n .

Let
$$\widetilde{W}_{B} = W_{\sigma, \mathbf{y}_{n}^{k}} = W_{\sigma, (\pi_{e}\mathbf{Y})^{k}}$$
. By Proposition 2.5,

$$d(\widetilde{W}_{\mathrm{B}}, W_{\beta^k}) \le d(\widetilde{W}_{\mathrm{B}}, W_{\mathrm{B}}) + d(W_{\mathrm{B}}, W_{\beta^k}) \le 2r|\mathrm{B}(e, k)|\eta + \eta.$$

So, as long as η is small enough and n is large enough (depending on ε , k), by Lemma 2.3,

$$|\{W \in \mathcal{W}_n(\alpha\beta, k, \varepsilon) : \pi_B W = W_B\}| > 1.$$

Now consider the probability appearing in the $\widetilde{W}_{\rm B}$ term:

$$\mathbb{P}_{\sigma \sim \mathbf{s}_n}(W_{\sigma, \mathbf{y}_n^k} = \widetilde{W}_{\mathbf{B}}) = \frac{|\{\sigma : W_{\sigma, \mathbf{y}_n^k} = \widetilde{W}_{\mathbf{B}}\}|}{|\{\sigma : W_{\sigma, \mathbf{y}_n} = W_{\sigma, \mathbf{y}_n}\}|}.$$

By symmetry in the choice of y with the correct letter frequencies (any two y with the same p_y are related by a permutation of [n], so have the same number of σ which give a particular weight), we can write this as

$$\mathbb{P}_{\sigma \sim s_n}(W_{\sigma, \mathbf{y}_n^k} = \widetilde{W}_{\mathsf{B}}) = \frac{|\{(\sigma, \mathbf{y}) : W_{\sigma, \mathbf{y}^k} = \widetilde{W}_{\mathsf{B}}\}|}{|\{(\sigma, \mathbf{y}) : W_{\sigma, \mathbf{y}} = W_{\sigma_n, \mathbf{y}_n}\}|}$$

$$= \frac{|\{(\sigma, \mathbf{Y}) : W_{\sigma, \mathbf{Y}} = \widetilde{W}_{\mathsf{B}}\}|}{|\{(\sigma, \mathbf{y}) : W_{\sigma, \mathbf{Y}} = W_{\sigma_n, \mathbf{y}_n}\}\}|} \quad \text{(Lemma 5.3)}$$

$$= \frac{Z_n(\widetilde{W}_{\mathsf{B}})}{Z_n(W_{\sigma_n,\mathbf{y}_n})} \quad \text{(definition of } Z_n)$$

$$\geq \exp(n[F(\widetilde{W}_{\mathsf{B}}) - F(W_{\sigma_n,\mathbf{y}_n})]) \cdot (3\sqrt{n})^{-r(|\mathsf{B}^{\mathsf{B}(e,k)}|^2 - |\mathsf{B}|)} \quad \text{(Prop. 5.1)}$$

$$= \exp(n[F(\widetilde{W}_{\mathsf{B}}) - F(W_{\sigma_n,\mathbf{y}_n}) - o(1)]).$$

By continuity of F, we then get

$$\underset{\sigma \sim s_n}{\mathbb{P}} (W_{\sigma, \mathbf{y}_n^k} = \widetilde{W}_{\mathrm{B}}) \ge \exp n(F_{\mu}(\beta^k) - F_{\mu}(\beta) - 2\delta + o(1))$$

for all large enough n and small enough η (again depending on k, ε), with $\delta > 0$ the same as chosen above. Since $\beta_*^G \mu$ is a Markov chain, $F_{\mu}(\beta^k) = F_{\mu}(\beta)$ [5, Theorem 6.1].

Putting this all together, for any $k \in \mathbb{N}$, for all $\delta > 0$, we have

$$\underset{\sigma \sim_{\mathbf{S}_n}}{\mathbb{E}} |X_k(\sigma, \alpha\beta, \varepsilon \mid \mathbf{y}_n)| \ge \exp[n(F_{\mu}(\alpha^k \mid \beta^k) - 3\delta - o(1))]$$

for all large enough n and small enough $\varepsilon > 0$.

It follows that for any $k \in \mathbb{N}$,

$$\inf_{\varepsilon} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{x} \in \mathbf{A}^n : (\mathbf{x}, \mathbf{y}_n) \in \Omega_k^*(\sigma, \alpha\beta, \varepsilon)\}| \ge F_{\mu}(T, \alpha^k \mid \beta^k).$$

Taking the limit as $k \to \infty$ gives the desired bound, using Corollary 3.3 and that the family of pseudometrics $\{d_k^* : k \in \mathbb{N}\}$ generates the weak* topology.

7. Proof of Theorem B

Let $W_n = W_{\sigma_n, \mathbf{y}_n^{m_n}}$, so that

$$s_n = SBM(\mathbf{y}_n, W_n).$$

Note that, by definition of s_n ,

$$\mathbb{P}_{\sigma \sim S_n}(W_{\sigma,\mathbf{y}_n^{m_n}} = W_n) = 1.$$

LEMMA 7.1. With W_n as just defined in terms of m_n , σ_n , and \mathbf{y}_n , we have

$$\lim_{n\to\infty} F(W_n) = f_{\mu}(T,\beta).$$

Proof. The assumption in the theorem statement that $d_{m_n}^*(P_{y_n}^{\sigma_n}, \beta_*^G \mu) = O(1/\log n)$ implies the existence of a constant C such that

$$d(W_n, W_{\beta^{m_n}}) \leq \frac{C}{\log n}.$$

By Lemma 3.1 we have

$$|F(W_{\sigma,\mathbf{y}^{m_n}}) - F(W_{\beta^{m_n}})| \le 4r \left(\mathbf{H}\left(\frac{C}{\log n}\right) + \frac{C}{\log n} |\mathbf{B}(e,m_n)| \log |\mathbf{B}| \right) = o(1)$$

using that $m_n = o(\log \log n)$. Since m_n approaches infinity as n goes to infinity we have $f_{\mu}(T, \beta) = \lim_{n \to \infty} F(W_{\beta^{m_n}})$, so the result follows.

LEMMA 7.2. If $m_n = o(\log \log n)$, then for any k > 0 and $\varepsilon > 0$ we have $|B^{B(e,m_n)}|^k = o(n^{\varepsilon})$.

Proof. This is certainly true if |B| = 1; assume therefore that $|B| \ge 2$.

Our assumption $m_n = o(\log \log n)$ guarantees that

$$(2r-1)^{m_n} < \frac{r-1}{r} \frac{\varepsilon}{k \log|\mathbf{B}|} \log n$$

for all large enough n. Therefore

$$|B(e, m_n)| = \frac{r(2r-1)^{m_n} - 1}{r-1} < \frac{\varepsilon}{k \log |B|} \log n.$$

This inequality can be rearranged to give

$$|\mathsf{B}^{\mathrm{B}(e,m_n)}|^k < n^{\varepsilon}.$$

Since $\varepsilon > 0$ is arbitrary, the result follows.

In the remainder of this section we prove Theorem B by first proving the right-hand side is an upper bound for the left, then proving it is also lower bound.

7.1. *Upper bound*. Just as in the proof of the upper bound in Theorem A, for each $k \in \mathbb{N}$ and $\varepsilon > 0$ we have

$$\inf_{\boldsymbol{O}\ni(\alpha\boldsymbol{\beta})_{s}^{G}\mu}\limsup_{n\to\infty}\frac{1}{n}\log\underset{\sigma\sim s_{n}}{\mathbb{E}}|\{\mathbf{x}\in\boldsymbol{\mathbb{A}}^{n}:(\mathbf{x},\mathbf{y}_{n})\in\Omega(\sigma,\boldsymbol{O})\}|\leq \limsup_{n\to\infty}\frac{1}{n}\log\mathcal{E}_{k}(n,\varepsilon),$$

where

$$\begin{split} \mathcal{E}_k(n,\varepsilon) &:= \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : (\mathbf{X},\mathbf{y}_n^k) \in \Omega_0^*(\sigma,(\alpha\beta)^k,\varepsilon)\}| \\ &= \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : d(W_{\sigma,(\mathbf{X},\mathbf{y}_n^k)},W_{(\alpha\beta)^k}) < \varepsilon)\}|. \end{split}$$

We assume that *n* is large enough that $m_n \ge k$.

Since s_n is SBM(σ_n , \mathbf{y}_n , m_n) rather than SBM(σ_n , \mathbf{y}_n , k), we cannot apply Proposition 5.2 directly to this expression. We get around this as follows. Let

$$\mathcal{W}_n(m, m') := \{W_{\sigma, (\mathbf{X}, \mathbf{y}^{m'})} : \sigma \in \text{Hom}(G, \text{Sym}(n)), \mathbf{X} \in (\mathbb{A}^{B(e, m)})^n, \mathbf{y} \in \mathbb{B}^n\}.$$

All elements of this set are denominator-n $\mathbb{A}^{B(e,m)} \times \mathbb{B}^{B(e,m')}$ -weights; we avoid the question of exactly which weights are in this set, but call such weights *attainable*. For k < m and k' < m' let

$$\mathcal{W}_n(m,m';\alpha\beta,k,k';\varepsilon) = \{W \in \mathcal{W}_n(m,m') : d(\pi_{k,k'}W,W_{\alpha^k\beta^{k'}}) < \varepsilon\}$$

denote the set of such weights whose appropriate marginal is within ε of the $(\mathbb{A}^{\mathrm{B}(e,k)} \times \mathbb{B}^{\mathrm{B}(e,k')})$ -weight $W_{\alpha^k\beta^{k'}}$. For now we take m=k=k', but we will need more generality below. Then

$$\mathcal{E}_k(n,\varepsilon) = \underset{\sigma \sim s_n}{\mathbb{E}} \sum_{W \in \mathcal{W}_n(k,m,\gamma,\sigma,\beta,k,k,\varepsilon)} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k)})^n : W_{\sigma,(\mathbf{X},\mathbf{y}_n^{m_n})} = W\}|,$$

so we can apply Proposition 5.2 to get

$$\mathcal{E}_k(n,\varepsilon) \leq (9n)^{r|\mathbb{B}^{\mathrm{B}(e,m_n)}|^2(|\mathbb{A}^{\mathrm{B}(e,k)}|+1)} \sum_{W \in \mathcal{W}_n(k,m_n;\alpha\beta,k,k;\varepsilon)} e^{n(F(W)-F(\pi_{\mathbb{B}}W))} \mathbf{1}_{\{\pi_{\mathbb{B}}W = W_n\}}.$$

By Lemma 7.2 we have $(9n)^{r|\mathbb{B}^{\mathrm{B}(e,m_n)}|^2(|\mathbb{A}^{\mathrm{B}(e,k)}|+1)} \leq e^{o_{n\to\infty}(n)}$. Using this and Lemma 7.1, we have

$$\mathcal{E}_k(n,\varepsilon) \leq \sum_{W \in \mathcal{W}_n(k,m_n;\alpha\beta,k,k;\varepsilon)} e^{n(F(W) - f(T,\beta) + o_{n \to \infty}(1))} \mathbf{1}_{\{\pi_{\mathbb{B}}W = W_n\}},$$

where the little o is uniform over all terms in the sum. Here we use the assumption that $f_{\mu}(T, \beta)$ is finite.

By definition of $W_n(k, m_n)$, for any $W \in W_n(k, m_n; \alpha\beta, k, k; \varepsilon)$ we can pick $\sigma \in \text{Hom}(G, \text{Sym}(n))$, $\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(\varepsilon,k)})^n$, and $\mathbf{y} \in \mathbb{B}^n$ so that $W = W_{\sigma,(\mathbf{X},\mathbf{y}^{m_n})}$. Then since $\mathbf{X}\mathbf{y}^{m_n}$ is a splitting of $\mathbf{X}\mathbf{y}^k$, by Lemma 3.2 we have

$$F(W) = F_{\text{Unif}([n])}(\sigma, \mathbf{X}\mathbf{y}^{m_n}) \le F_{\text{Unif}([n])}(\sigma, \mathbf{X}\mathbf{y}^k) = F(\pi_{k,k}W),$$

where here $F_{\text{Unif}([n])}(\sigma, \mathbf{X}\mathbf{y}^{m_n})$ is F of the observable $\mathbf{X}\mathbf{y}^{m_n}$ on the measure-preserving system $([n], \text{Unif}([n]), \sigma)$ (we shift to this notation from weights in order to apply the splitting lemma). By continuity of F, for all small enough ε (depending on k) we have

$$F(\pi_{k,k}W) \le F(W_{(\alpha\beta)^k}) + \delta = F_{\mu}(T, (\alpha\beta)^k) + \delta.$$

Along with the above, this implies that

$$\mathcal{E}_k(n,\varepsilon) \leq e^{n(F(T,(\alpha\beta)^k) - f(T,\beta) + o_n(1) + \delta)} \sum_{W \in \mathcal{W}_n(k,m_n;\alpha\beta,k,k;\varepsilon)} \mathbf{1}_{\{\pi_{\mathbb{B}}W = W_n\}}.$$

Bounding all terms in the sum by 1, we get

$$\mathcal{E}_k(n,\varepsilon) \leq e^{n(F(T,(\alpha\beta)^k) - f_\mu(T,\beta) + o_n(1) + \delta)} \, |\mathcal{W}_n(k,m_n;\alpha\beta,k,k;\varepsilon)|.$$

Using Lemma 7.2, we have

$$|\mathcal{W}_n(k, m_n; \alpha\beta, k, k; \varepsilon)| \leq |\mathcal{W}_n(k, m_n)| \leq n^{r|\mathbb{A}^{\mathrm{B}(e,k)} \times \mathbb{B}^{\mathrm{B}(e,m_n)}|^2} \leq e^{o_{n \to \infty}(n)},$$

so this implies

$$\limsup_{n\to\infty} \frac{1}{n} \log \mathcal{E}_k(n,\varepsilon) \leq F_{\mu}(T,(\alpha\beta)^k) - f_{\mu}(T,\beta) + \delta.$$

Taking the infimum over ε and k, and using the chain rule for f (Corollary 3.5, again using the assumption that $f_{\mu}(T, \beta)$ is finite), gives

$$\inf_{\varepsilon,k} \limsup_{n \to \infty} \frac{1}{n} \log \mathcal{E}_k(n,\varepsilon) \leq f_{\mu}(T,\alpha\beta) - f_{\mu}(T,\beta) = f_{\mu}(T,\alpha \mid \beta).$$

Since

$$\inf_{O\ni(\alpha\beta)_*^G\mu} \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim_{\mathbb{S}_n}}{\mathbb{E}} |\{\mathbf{x}\in\mathbb{A}^n: (\mathbf{x},\mathbf{y}_n)\in\Omega(\sigma,O)\}|$$

$$\leq \inf_{\varepsilon} \limsup_{n\to\infty} \frac{1}{n} \log \mathcal{E}_k(n,\varepsilon),$$

for every k, this completes the upper bound.

7.2. Lower bound. In this section we denote

$$\mathcal{X}_{k_1,k_2}(\sigma,\alpha\beta,\varepsilon\mid\mathbf{y}):=\{\mathbf{X}\in(\mathbb{A}^{\mathrm{B}(e,k_1)})^n:(\mathbf{X},\mathbf{y}^{k_2})\in\Omega_0^*(\sigma,\alpha^{k_1}\beta^{k_2},\varepsilon)\},$$

$$\Omega_k^*(\sigma,\alpha\beta,\varepsilon\mid\mathbf{y}):=\{\mathbf{x}\in\mathbb{A}^n:(\mathbf{x},\mathbf{y})\in\Omega_k^*(\sigma,\alpha\beta,\varepsilon)\}$$

(note the dependence on n is implicitly specified by $\sigma \in \text{Hom}(G, \text{Sym}(n))$ and $\mathbf{y} \in \mathbb{B}^n$), and with $\Sigma = \{s_n\}_{n=1}^{\infty}$,

$$\begin{split} \mathbf{h}_{\Sigma,\mu}(T,\alpha\mid\beta:k,\varepsilon) &:= \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim\mathbf{s}_n}{\mathbb{E}} |\{\mathbf{x}\in\mathbf{A}^n: (\mathbf{x},\mathbf{y})\in\Omega_k^*(\sigma,\alpha\beta,\varepsilon)\}| \\ &= \limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim\mathbf{s}_n}{\mathbb{E}} |\Omega_k^*(\sigma,\alpha\beta,\varepsilon\mid\mathbf{y})|. \end{split}$$

The following two claims are used to relate the sizes of the sets defined above.

CLAIM 1. Let $k \leq \min(k_1, k_2)$. For any σ , \mathbf{y} , we have

$$\pi_e[X_{k_1,k_2}(\sigma,\alpha\beta,\varepsilon\mid\mathbf{y})]\subseteq\Omega_k^*(\sigma,\alpha\beta,c\varepsilon\mid\mathbf{y})$$

where c = 1 + |B(e, k)|.

Proof. If $(\mathbf{X}, \mathbf{y}^{k_2}) \in \Omega_0^*(\sigma, \alpha^{k_1} \beta^{k_2}, \varepsilon)$, then

$$\pi_{k,k}(\mathbf{X}, \mathbf{y}^{k_2}) \in \Omega_0^*(\sigma, (\alpha\beta)^k, \varepsilon);$$

this follows from the fact that total variation distance is non-increasing under pushforwards. Applying Proposition 2.5, we get

$$(\pi_e \mathbf{X}, \mathbf{y}) = \pi_e(\pi_{k,k}(\mathbf{X}, \mathbf{y}^{k_2})) \in \Omega_k^*(\sigma, \alpha\beta, c\varepsilon).$$

CLAIM 2. Fix σ , \mathbf{y} , and $k \leq \min(k_1, k_2)$. As established in the previous claim, we can consider π_e as a map from $X_{k_1,k_2}(\sigma, \alpha\beta, \varepsilon \mid \mathbf{y})$ to $\Omega_k^*(\sigma, \alpha\beta, c\varepsilon \mid \mathbf{y})$. There are constants C, d independent of n such that π_e is at most C exp $(nd\varepsilon + nH(2|B(e, k)|\varepsilon))$ -to-one.

Proof. If $\Omega_k^*(\sigma, \alpha\beta, c\varepsilon \mid \mathbf{y})$ is empty, then the claim is vacuously true. Otherwise, fix $\mathbf{x} \in \Omega_k^*(\sigma, \alpha\beta, c\varepsilon \mid \mathbf{y})$. If $\mathbf{X} \in \pi_e^{-1}\{\mathbf{x}\}$, then $\pi_e(\mathbf{X}, \mathbf{y}^k) = (\mathbf{x}, \mathbf{y})$. Claim 3 in the proof of Proposition 3.2 of [2] gives an upper bound of the desired form for the number of such pairs $(\mathbf{X}, \mathbf{y}^k)$, and therefore the number of such \mathbf{X} .

Claim 2 implies that

$$|X_{k_1,k_2}(\sigma,\alpha\beta,\varepsilon\mid\mathbf{y})| \leq C \exp(nd\varepsilon + n\mathrm{H}(2|\mathrm{B}(e,k)|\varepsilon)) \cdot |\Omega_k^*(\sigma,\alpha\beta,c\varepsilon\mid\mathbf{y})|, \tag{7.1}$$
 where C,d are independent of n .

We now find a lower bound for the expectation of $|\mathcal{X}|$. Fix $k_1, k_2 \in \mathbb{N}$, and suppose n is large enough that $m_n \ge \max(k_1, k_2)$. Using Proposition 5.2 and Lemma 7.2, we have

$$\begin{split} & \underset{\sigma \sim \mathbb{S}_{n}}{\mathbb{E}} |\mathcal{X}_{k_{1},k_{2}}(\sigma,\alpha\beta,\varepsilon \mid \mathbf{y}_{n})| \\ & = \sum_{W \in \mathcal{W}_{n}(k_{1},m_{n};\alpha\beta,k_{1},k_{2};\varepsilon)} \underset{\sigma \sim \mathbb{S}_{n}}{\mathbb{E}} |\{\mathbf{X} \in (\mathbb{A}^{\mathrm{B}(e,k_{1})})^{n} : W_{\sigma,(\mathbf{X},\mathbf{y}_{n}^{m_{n}})} = W\}| \\ & \geq \sum_{W \in \mathcal{W}_{n}(k_{1},m_{n};\alpha\beta,k_{1},k_{2};\varepsilon)} \exp[n(F(W) - F(\pi_{\mathrm{B}}W) - o_{n}(1))] \mathbf{1}_{\{\pi_{\mathrm{B}}W = W_{\sigma,\mathbf{y}_{n}^{m_{n}}}\}} \\ & \geq \inf_{W \in \mathcal{W}_{n}(k_{1},m_{n};\alpha\beta,k_{1},k_{2};\varepsilon)} \exp[n(F(W) - F(\pi_{\mathrm{B}}W) - o_{n}(1))] \\ & \times \sum_{W \in \mathcal{W}_{n}(k_{1},m_{n};\alpha\beta,k_{1},k_{2};\varepsilon)} \mathbf{1}_{\{\pi_{\mathrm{B}}W = W_{\sigma,\mathbf{y}_{n}^{m_{n}}}\}}. \end{split}$$

We bound the infimum below as follows. Given any $W \in W_n(k_1, m_n; \alpha\beta, k_1, k_2; \varepsilon)$, we can let X, y, σ be such that $W = W_{\sigma,(X,y^{m_n})}$. Then by Lemma 3.2 and continuity of F,

$$\begin{split} F(W) - F(\pi_{\mathbb{B}}W) &= F(\sigma, \mathbf{X}|\mathbf{y}^{m_n}) \\ &\geq F(\sigma, \mathbf{X}|\mathbf{y}^{k_2}) \\ &= F(\pi_{k_1, k_2}W) - F(\pi_{\mathbb{B}}\pi_{k_1, k_2}W) \\ &\geq F_{\mu}(T, \alpha^{k_1}|\beta^{k_2}) - \delta \end{split}$$

for any $\delta > 0$ for all small enough ε (with 'small enough' dependent only on k_1, k_2). This implies that the infimum is bounded below by

$$\exp[n(F_u(T, \alpha^{k_1} | \beta^{k_2}) - o_n(1) - \delta)].$$

We bound the sum below by first rewriting it as

$$|\{W\in\mathcal{W}_n(k_1,m_n;\alpha\beta,k_1,k_2;\varepsilon):\pi_{\mathbb{B}}W=W_{\sigma,\mathbf{y}_n^{m_n}}\}|.$$

The following claim, then, implies that the sum is bounded below by 1.

CLAIM 3. For all large enough n,

$$\{W \in \mathcal{W}_n(k_1, m_n; \alpha\beta, k_1, k_2; \varepsilon) : \pi_{\mathbb{B}}W = W_{\sigma, \mathbf{y}_n^{m_n}}\} \neq \varnothing.$$

Proof. By Lemma 2.3, if

$$n > 680 |\mathbf{A}^{\mathbf{B}(e,k_1)} \times \mathbf{B}^{\mathbf{B}(e,m_n)}|^2 r/\varepsilon$$

and $d(W_{\sigma,\mathbf{y}_n^{m_n}},W_{\beta^{m_n}}) < \varepsilon/530r$ then there is a $(\mathbb{A}^{\mathrm{B}(e,k_1)} \times \mathbb{B}^{\mathrm{B}(e,m_n)})$ -weight W with $\pi_{\mathrm{B}}W = W_{\sigma,\mathbf{y}_n^{m_n}}$ and $d(W,W_{\alpha^{k_1}\beta^{m_n}}) < \varepsilon$. By definition of s_n and Lemma 7.2, both conditions are met for all large enough n.

The claim will follow if we show that W is attainable.

With W as chosen above, by Proposition 2.4 we can choose $\tilde{\sigma} \in \text{Hom}(G, \text{Sym}(n))$, $\tilde{\mathbf{X}} \in (\mathbb{A}^{B(e,k_1)})^n$, and $\tilde{\mathbf{Y}} \in (\mathbb{B}^{B(e,m_n)})^n$ such that $W = W_{\tilde{\sigma},(\tilde{\mathbf{X}}\tilde{\mathbf{Y}})}$.

Let $\tilde{\mathbf{y}} = \pi_e \tilde{\mathbf{Y}} \in \mathbb{B}^n$. To complete the proof we show that $\tilde{\mathbf{y}}^{m_n} = \tilde{\mathbf{Y}}$, that is,

$$\tilde{\mathbf{y}}(\tilde{\sigma}(g)i) = (\tilde{\mathbf{Y}}(i))_g$$

for all $i \in [n]$ and $g \in B(e, m_n)$. We prove this by induction on the word length |g|.

The base case |g| = 0 (that is, g = e) follows immediately from the definition of \tilde{y} .

For the inductive step, write g = ht with |h| = |g| - 1 and $t \in \{s_1^{\pm 1}, \dots, s_r^{\pm 1}\}$. Then, assuming the result holds for h,

$$\tilde{\mathbf{y}}(\tilde{\sigma}(g)i) = \tilde{\mathbf{y}}(\tilde{\sigma}(h)\tilde{\sigma}(t)i) = (\tilde{\mathbf{Y}}(\tilde{\sigma}(t)i))_h.$$

Now since $W_{\tilde{\sigma},\tilde{\mathbf{Y}}} = W_{\sigma_n,\mathbf{y}_n^{m_n}}$, we can pick $j \in [n]$ such that

$$\tilde{\mathbf{Y}}(i) = \mathbf{y}_n^{m_n}(j)$$
 and $\tilde{\mathbf{Y}}(\tilde{\sigma}(t)i) = \mathbf{y}_n^{m_n}(\sigma(t)j)$.

This implies

$$(\tilde{\mathbf{Y}}(\tilde{\sigma}(t)i))_h = (\mathbf{y}_n^{m_n}(\sigma(t)j))_h = \mathbf{y}_n(\sigma(g)j) = (\mathbf{y}_n^{m_n}(j))_g = (\tilde{\mathbf{Y}}(i))_g.$$

Hence, for all large enough n, we have

$$\underset{\sigma \sim s_n}{\mathbb{E}} |X_{k_1,k_2}(\sigma, \alpha\beta, \varepsilon \mid \mathbf{y}_n)| \ge \exp[n(F_{\mu}(T, \alpha^{k_1} \mid \beta^{k_2}) - o_n(1) - \delta)],$$

and therefore

$$\limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma\sim s_n}{\mathbb{E}} |X_{k_1,k_2}(\sigma,\alpha\beta,\varepsilon \mid \mathbf{y}_n)| \geq F_{\mu}(T,\alpha^{k_1} \mid \beta^{k_2}) - \delta.$$

Combining this lower bound with equation (7.1) and the definition of $h_{\Sigma,\mu}(T,\alpha \mid \beta : k, c\varepsilon)$, we get

$$d\varepsilon + H(2|B(e,k)|\varepsilon) + h_{\Sigma,\mu}(T,\alpha \mid \beta : k, c\varepsilon) \ge F_{\mu}(T,\alpha^{k_1} \mid \beta^{k_2}) - \delta.$$

Taking the infimum in ε then letting δ go to zero gives

$$\inf_{\varepsilon} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^n : (\mathbf{x}, \mathbf{y}_n) \in \Omega_k^*(\sigma, \alpha\beta, \varepsilon)\}| \ge F_{\mu}(T, \alpha^{k_1} \mid \beta^{k_2})$$

for $k \leq \min(k_1, k_2)$. First take $k_2 \to \infty$, then $k_1 \to \infty$, then take the infimum over k. We get

$$f_{\mu}(T, \alpha \mid \beta) \leq \inf_{\varepsilon, k} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_{n}}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^{n} : (\mathbf{x}, \mathbf{y}_{n}) \in \Omega_{k}^{*}(\sigma, \alpha\beta, \varepsilon)\}|$$

$$= \inf_{O \ni (\alpha\beta)_{*}^{G} \mu} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_{n}}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^{n} : (\mathbf{x}, \mathbf{y}_{n}) \in \Omega(\sigma, O)\}|$$

where the last line follows because the collection of pseudometrics $\{d_k^*: k \in \mathbb{N}\}$ generates the weak* topology on $\operatorname{Prob}((\mathbb{A} \times \mathbb{B})^G)$.

8. Proof of Theorem C

By analogy with sofic entropy, we denote $\Sigma := \{s_n\}_{n=1}^{\infty}$ and denote the left-hand side of the formula in the theorem statement by $h_{\Sigma,\mu}(T,\alpha)$.

Endow $Prob(A^G)$ with the metric

$$d(\lambda, \nu) := \sum_{r=1}^{\infty} 2^{-r} d^{\mathbf{B}(e,r)}(\lambda, \nu).$$

Note that this induces the weak* topology (where A is given the discrete topology and A^G the product topology).

Writing $\mu_{\mathbb{A}} = \alpha_*^G \mu \in \text{Prob}(\mathbb{A}^G)$, we then have

$$h_{\Sigma,\mu}(T,\alpha) = \inf_{\varepsilon > 0} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^n : d(P_{\mathbf{x}}^{\sigma}, \mu_{\mathbb{A}}) < \varepsilon\}|.$$

We will similarly denote $\mu_{\mathbb{B}} = \beta_*^G \mu \in \text{Prob}(\mathbb{B}^G)$.

8.1. Lower bound. Let $\lambda \in \text{Prob}((\mathbb{A} \times \mathbb{B})^G)$ be any joining of (the shift systems with respective measures) $\mu_{\mathbb{A}}$ and $\mu_{\mathbb{B}}$. Then, for any $\mathbf{x} \in \mathbb{A}^n$ and $\mathbf{y} \in \mathbb{B}^n$, we have

$$d(P_{\mathbf{X}}^{\sigma}, \mu_{\mathbb{A}}) \leq d(P_{(\mathbf{X},\mathbf{Y})}^{\sigma}, \lambda),$$

where d is defined on $Prob((A \times B)^G)$ analogously to the definition given on $Prob(A^G)$ above. This inequality holds because total variation distance is non-increasing under pushforwards. Consequently,

$$\mathrm{h}_{\Sigma,\mu}(T,\alpha) \geq \inf_{\varepsilon>0} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim \mathrm{s}_n}{\mathbb{E}} |\{\mathbf{x} \in \mathrm{A}^n : d(P^{\sigma}_{(\mathbf{x},\mathbf{y}_n)},\lambda) < \varepsilon\}| = f_{\lambda}(S,\mathrm{a} \mid \mathrm{b}).$$

Taking the supremum over joinings λ gives the lower bound.

8.2. *Upper bound.* For $\varepsilon > 0$, let

$$\mathsf{J}_{\varepsilon} := \{ \lambda \in \mathsf{Prob}^{S}((\mathsf{A} \times \mathsf{B})^{G}) : d(\mathsf{a}_{*}^{G}\lambda, \mu_{\mathsf{A}}) < \varepsilon \text{ and } d(\mathsf{b}_{*}^{G}\lambda, \mu_{\mathsf{B}}) < \varepsilon \}$$

be the set of shift-invariant 'approximate joinings' of $\mu_{\mathbb{A}}$ and $\mu_{\mathbb{B}}$. Since $\operatorname{Prob}((\mathbb{A} \times \mathbb{B})^G)$ is compact, for each $\varepsilon > 0$ there exist $\lambda_1, \ldots, \lambda_m \in J_{\varepsilon}$ such that

$$\mathsf{J}_{\varepsilon} \subseteq \bigcup_{i=1}^m \mathsf{B}(\lambda_i, \varepsilon).$$

By definition of s_n we have $\mathbb{P}_{\sigma \sim s_n}(d(P_{\mathbf{y}_n}^{\sigma}, \mu_{\mathbb{B}}) < \varepsilon) = 1$ for all large enough n. Therefore,

$$\begin{split} \mathbf{h}_{\Sigma,\mu}(T,\alpha) &= \inf_{\varepsilon} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim \mathbf{s}_n}{\mathbb{E}} |\{\mathbf{x} \in \mathbf{A}^n : P_{(\mathbf{x},\mathbf{y}_n)}^{\sigma} \in \mathsf{J}_{\varepsilon}\}| \\ &\leq \inf_{\varepsilon} \limsup_{n \to \infty} \frac{1}{n} \log \underset{i=1}{\overset{m}{\sum}} \underset{\sigma \sim \mathbf{s}_n}{\mathbb{E}} |\{\mathbf{x} \in \mathbf{A}^n : P_{(\mathbf{x},\mathbf{y}_n)}^{\sigma} \in \mathsf{B}(\lambda_i,\varepsilon)\}| \\ &= \inf_{\varepsilon} \max_{1 \le i \le m} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim \mathbf{s}_n}{\mathbb{E}} |\{\mathbf{x} \in \mathsf{A}^n : P_{(\mathbf{x},\mathbf{y}_n)}^{\sigma} \in \mathsf{B}(\lambda_i,\varepsilon)\}| \\ &\leq \inf_{\varepsilon} \sup_{\lambda \in \mathsf{J}_{\varepsilon}} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim \mathbf{s}_n}{\mathbb{E}} |\{\mathbf{x} \in \mathsf{A}^n : P_{(\mathbf{x},\mathbf{y}_n)}^{\sigma} \in \mathsf{B}(\lambda,\varepsilon)\}|. \end{split}$$

Note that the entire expression in the infimum is decreasing as $\varepsilon \to 0$, so we may replace the infimum with a limit. Rather than taking a continuous limit we write

$$h_{\Sigma,\mu}(T,\alpha) \leq \lim_{m \to \infty} \sup_{\lambda \in \mathsf{J}_{1/m}} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim \mathsf{s}_n}{\mathbb{E}} |\{\mathbf{x} \in \mathsf{A}^n : P^{\sigma}_{(\mathbf{x},\mathbf{y}_n)} \in \mathsf{B}(\lambda,1/m)\}|.$$

For each m, pick $\lambda_m \in \mathsf{J}_{1/m}$ to get within 1/m of the supremum. Then the right-hand side is equal to

$$\lim_{m \to \infty} \limsup_{n \to \infty} \frac{1}{n} \log \underset{\sigma \sim s_n}{\mathbb{E}} |\{ \mathbf{x} \in A^n : P_{(\mathbf{x}, \mathbf{y}_n)}^{\sigma} \in B(\lambda_m, 1/m) \}|. \tag{*}$$

Let λ_{m_j} be a subsequence with weak* limit λ_0 . By weak* continuity of pushforwards under projection we have $\lambda_0 \in J(\mu_A, \mu_B)$. Now, for any $\delta > 0$, for all large enough j we have both $1/m_j < \delta/2$ and $d(\lambda_{m_j}, \lambda_0) < \delta/2$, so by the triangle inequality

$$B(\lambda_{m_i}, 1/m_i) \subseteq B(\lambda_0, \delta).$$

It follows that the expression in (*), and hence $h_{\Sigma}(\alpha)$, is bounded above by

$$\limsup_{n\to\infty} \frac{1}{n} \log \underset{\sigma \sim s_n}{\mathbb{E}} |\{\mathbf{x} \in \mathbb{A}^n : P^{\sigma}_{(\mathbf{x},\mathbf{y}_n)} \in \mathrm{B}(\lambda_0,\delta)\}|.$$

Taking the infimum over δ shows that

$$h_{\Sigma}(\mu, \alpha) \leq f_{\lambda_0}(S, \mathbf{a} \mid \mathbf{b}) \leq \sup_{\lambda \in \mathsf{J}(\mu_h, \mu_B)} f_{\lambda}(S, \mathbf{a} \mid \mathbf{b}).$$

9. Proof of Proposition 1.1

All sequences of interest are of the form

$$\mathbf{s}_n = \mathrm{SBM}(\sigma_n, \mathbf{y}_n, m_n) = \mathrm{Unif}(\{\sigma \in \mathrm{Hom}(G, \mathrm{Sym}(n)) : W_{\sigma, \mathbf{y}_n^{m_n}} = W_n\})$$

with $\mathbf{y}_n \in \mathbb{B}^n$, $\sigma_n \in \operatorname{Sym}(n)$, $m_n = o(\log \log n)$, and where W_n is the $\mathbb{B}^{\mathrm{B}(e,m_n)}$ -weight $W_{\sigma_n,\mathbf{y}_n^{m_n}}$. In the case of Theorem A we simply have $m_n = 0$ for all n.

The theorem will follow from the following lemma.

LEMMA 9.1. Let ζ_n denote the uniform measure on $\operatorname{Hom}(G,\operatorname{Sym}(n))$. Then, for any finite $D\subset G$ and $\delta>0$, there exists $\varepsilon>0$ such that

$$\underset{\sigma \sim \xi_n}{\mathbb{P}} (\sigma \text{ is } (D, \delta) \text{-softc}) \geq 1 - n^{-\varepsilon n}$$

for all large enough n.

This can be proven by making superficial changes to the proof of the similar result [1, Lemma 3.1].

To prove Proposition 1.1, it now suffices to show that, for any $\varepsilon > 0$,

$$\mathbb{P}_{\sigma \sim \zeta_n}(W_{\sigma,\mathbf{y}_n^{m_n}} = W_n) \ge n^{-\varepsilon n}$$

for all large enough n. To do this, first note that the left-hand side here depends only on the vector $p_{\mathbf{y}_n} \in \operatorname{Prob}(\mathbb{B})$ of letter frequencies. Therefore,

$$\mathbb{P}_{\sigma \sim \zeta_n} \text{ (there exists } \mathbf{y} \in \mathbb{B}^n \text{ s.t. } W_{\sigma, \mathbf{y}^{m_n}} = W_n) \leq \sum_{\mathbf{y}: p_{\mathbf{y}} = p_{\mathbf{y}_n}} \mathbb{P}_{\sigma \sim \zeta_n} (W_{\sigma, \mathbf{y}^{m_n}} = W_n) \\
= \exp\{n \mathbf{H}(p_{\mathbf{y}_n}) + o(n)\} \mathbb{P}_{\sigma \sim \zeta_n} (W_{\sigma, \mathbf{y}_n^{m_n}} = W_n).$$

But by Proposition 2.5, if $\sigma \in \text{Hom}(G, \text{Sym}(n))$ and $\mathbf{Y} \in (\mathbb{B}^{B(e,m_n)})^n$ are such that $W_{\sigma,\mathbf{Y}} = W_n = W_{\sigma_n,\mathbf{V}_n^{m_n}}$, then the projection $\mathbf{Y}_e \in \mathbb{B}^n$ satisfies $(\mathbf{Y}_e)^{m_n} = \mathbf{Y}$. Therefore, for each σ ,

$$|\{\mathbf{Y} \in (\mathbf{B}^{\mathbf{B}(e,m_n)})^n : W_{\sigma,\mathbf{Y}} = W_n\}| = |\{\mathbf{y} \in \mathbf{B}^n : W_{\sigma,\mathbf{y}^{m_n}} = W_n\}|.$$

Hence.

$$\mathbb{E}_{\sigma \sim \zeta_n} |\{ \mathbf{Y} \in (\mathbb{B}^{\mathbf{B}(e, m_n)})^n : W_{\sigma, \mathbf{Y}} = W_n \}| = \mathbb{E}_{\sigma \sim \zeta_n} |\{ \mathbf{y} \in \mathbb{B}^n : W_{\sigma, \mathbf{y}^{m_n}} = W_n \}| \\
\leq |\mathbb{B}|^n \mathbb{P}_{\sigma \sim \zeta_n} \text{ (there exists } \mathbf{y} \in \mathbb{B}^n \text{ s.t. } W_{\sigma, \mathbf{y}^{m_n}} = W_n).$$

Combining these last few statements, we see that

$$\mathbb{P}_{\sigma \sim \zeta_n}(W_{\sigma,\mathbf{y}_n^{m_n}} = W_n) \ge \exp\{-2n \log|\mathbf{B}| + o(n)\} \mathbb{E}_{\sigma \sim \zeta_n}[\{\mathbf{Y} \in (\mathbf{B}^{\mathbf{B}(e,m_n)})^n : W_{\sigma,\mathbf{Y}} = W_n\}].$$

We can ignore the first factor here since it only decays exponentially fast. By Proposition 5.1,

$$\mathbb{E}_{\sigma \sim \zeta_n} |\{ \mathbf{Y} \in (\mathsf{B}^{\mathsf{B}(e,m_n)})^n : W_{\sigma,\mathbf{Y}} = W_n \}| = \frac{Z_n(W_n)}{(n!)^r} \ge (3\sqrt{n})^{-r|\mathsf{B}^{\mathsf{B}(e,m_n)}|^2} e^{F(W_n)n} n^{(1-r)/2}.$$

The third factor is clearly not a problem and can also be ignored. For the first factor,

$$\frac{1}{n \log n} \log(3\sqrt{n})^{-r|\mathbf{B}^{\mathrm{B}(e,m_n)}|^2} = -r \frac{|\mathbf{B}^{\mathrm{B}(e,m_n)}|^2}{n} \frac{\log 3\sqrt{n}}{\log n} \to 0 \quad \text{as } n \to \infty$$

using Lemma 7.2. For the second factor, first note that by definition of $F(W_n)$ we have

$$F(W_n) = (1 - 2r)H(W_n(\cdot)) + \sum_{i \in [r]} H(W_n(\cdot, \cdot; i))$$

$$\geq -2rH(W_n(\cdot))$$

$$> -2r \log|\mathsf{B}^{\mathsf{B}(e, m_n)}|.$$

So

$$\frac{1}{n \log n} \log e^{F(W_n)n} = \frac{F(W_n)}{\log n} \ge -2r \frac{\log |\mathbb{B}^{\mathrm{B}(e, m_n)}|}{\log n} \to 0 \quad \text{as } n \to \infty,$$

again using Lemma 7.2. This implies that for every $\varepsilon > 0$ we have

$$(3\sqrt{n})^{-r|\mathbf{B}^{\mathbf{B}(e,m_n)}|^2}e^{F(W_n)n} \ge n^{-\varepsilon n}$$

for all large enough n, which implies the result.

10. Proof of Lemma 2.3

We show how to construct a denominator-n weight W_{AB} that has a given B-marginal W_B and is close to a given (A \times B)-weight W whose B-marginal $\pi_B W$ is close to W_B . As in the theorem statement, we assume

$$d(\pi_{\rm B}W, W_{\rm B}) < \delta$$
.

To minimize the appearance of factors of $\frac{1}{2}$, in this section we work with the ℓ^1 distance on weights, which is twice the distance defined above. Therefore the previous assumption becomes

$$d_1(\pi_{\mathsf{B}} W, W_{\mathsf{B}}) = \sum_{i \in [r]} \sum_{b, b' \in \mathsf{B}} |\pi_{\mathsf{B}} W(b, b'; i) - W_{\mathsf{B}}(b, b'; i)| < 2\delta.$$

We fix distinguished elements $a_0 \in \mathbb{A}$ and $b_0 \in \mathbb{B}$ which will be referred to throughout this section.

10.1. The vertex measure. We first define the weight's vertex measure by

$$\begin{split} W_{\mathrm{AB}}((a,b)) &= \frac{1}{n} \lfloor n \cdot W((a,b)) \rfloor \quad a \in \mathrm{A} \setminus \{a_0\}, \ b \in \mathrm{B}, \\ W_{\mathrm{AB}}((a_0,b)) &= W_{\mathrm{B}}(b) - \sum_{a \neq a_0} W_{\mathrm{AB}}((a,b)) \quad b \in \mathrm{B}. \end{split}$$

See Table 1.

Note that
$$|W_{AB}((a,b)) - W((a,b))| \le 1/n$$
 for $a \ne a_0$ and

$$|W_{AB}((a_0, b)) - W((a_0, b))| \le |W_{B}(b) - \pi_{B}W(b)| + |A|/n.$$

Therefore the ℓ^1 distance between the vertex measures is

$$\begin{split} \sum_{a,b} &|W_{\mathrm{AB}}((a,b)) - W((a,b))| \leq |\mathbb{A}||\mathbb{B}|/n + \sum_{b \in \mathbb{B}} (|W_{\mathrm{B}}(b) - \pi_{\mathbb{B}}W(b)| + |\mathbb{A}|/n) \\ &\leq 2\delta + 2|\mathbb{A}||\mathbb{B}|/n. \end{split}$$

10.1.1. *Nonnegativity*. The terms defined by rounding down *W* using the floor function are guaranteed to be non-negative, but the others are not. In the following we show how to repair any negativity.

TABLE 1. Picking entries of the vertex measure $W_{\rm AB}(\cdot)$. First choose entries of the form $W_{\rm AB}((a,b))$ for $a \neq a_0$ by rounding down W((a,b)), then fill in the first column in a way that guarantees the correct B-marginal.

	a_0	a_1	
b_0 b_1	$\overset{\rightarrow}{\rightarrow}$	[·] [·]	[·]
:	\rightarrow	$\lfloor \cdot floor$	<u>[·]</u>

Let -R/n denote the sum of all negative terms in the vertex measure. Since W contains only non-negative terms, we have

$$\mathbf{1}_{\{W_{AB}((a,b))<0\}} \cdot |W_{AB}((a,b))| \le |W_{AB}((a,b)) - W((a,b))|$$
 for all a, b .

Therefore

$$R/n \le \sum_{b \in \mathbb{R}} |W_{AB}((a_0, b)) - W((a_0, b))| \le 2\delta + |A||B|/n.$$

Suppose there is some $b \in \mathbb{B}$ such that $W_{AB}((a_0, b)) < 0$. Since W_{AB} has denominator n, we must have $W_{AB}((a_0, b)) \le -1/n$. By construction, we have

$$\sum_{a \in A} W_{AB}((a, b)) = W_{B}(b) \ge 0,$$

so there exists some $a^+ \in A$ with $W_{AB}((a^+, b)) \ge 1/n$. Increase $W_{AB}((a_0, b))$ by 1/n and decrease $W_{AB}((a^+, b))$ by 1/n.

The number of times we must repeat this step before all terms are non-negative is exactly R, and each step moves the measure by ℓ^1 distance 2/n; therefore the final edited vertex measure is distance at most 2R/n from the original W_{AB} . If we now let W_{AB} denote the new, non-negative vertex measure, by the above bound on R/n we get

$$\sum_{a,b} |W_{AB}((a,b)) - W((a,b))| \le 6\delta + 4|A||B|/n.$$

10.2. *The* B *half-marginal*. For the purposes of this construction we use the B 'half-marginal', which we denote by

$$W(b, (a', b'); i) := \sum_{a \in \mathbb{A}} W((a, b), (a', b'); i).$$

This is an element of Prob($(B \times (A \times B))^r$).

Before constructing the edge measure of W_{AB} , in this section we first construct what will be its half-marginal.

For each $i \in [r]$, $b, b' \in \mathbb{B}$, and $a' \in \mathbb{A}$, we define

$$W_{AB}(b, (a', b'); i) = \frac{1}{n} [n \cdot W(b, (a', b'); i)] \quad \text{for } a' \neq a_0, b \neq b_0,$$
 (10.1)

$$W_{AB}(b, (a_0, b'); i) = W_B(b, b'; i) - \sum_{a' \neq a_0} W_{AB}(b, (a', b'); i) \quad \text{for } b \neq b_0,$$
 (10.2)

$$W_{AB}(b_0, (a', b'); i) = W_{AB}((a', b')) - \sum_{b \neq b_0} W_{AB}(b, (a', b'); i).$$
(10.3)

See Table 2 for a representation of which terms are defined by each equation.

The definition of the terms in (10.3) ensures that

$$\sum_{b \in \mathbb{B}} W_{AB}(b, (a', b'); i) = W_{AB}((a', b')) \quad \text{for all } a', b', i.$$

TABLE 2. A diagram of how the half-marginal $W_{AB}(\cdot, (\cdot, \cdot); i)$ is chosen if $A = \{a_0, a_1, a_2\}$ and $B = \{b_0, b_1, b_2\}$. First obtain the entries marked $\lfloor \cdot \rfloor$ by rounding down W. Then choose the entries marked \rightarrow according to equation (10.2) which ensures that the B-marginal is W_B . Then choose the entries marked \downarrow according to equation (10.3) which ensures that the vertex weight is the one we chose above.

	(a_0, b_0)	(a_1, b_0)	(a_2, b_0)	(a_0, b_1)	(a_1, b_1)	(a_2, b_1)	(a_0, b_2)	(a_1, b_2)	(a_2, b_2)
b_0	↓	↓	↓	↓	↓	↓	↓	↓	
b_1	\rightarrow	$\lfloor \cdot \rfloor$	$\lfloor \cdot \rfloor$	\rightarrow	$\lfloor \cdot \rfloor$	$\lfloor \cdot floor$	\rightarrow	$\lfloor \cdot \rfloor$	$\lfloor \cdot \rfloor$
b_2	\rightarrow	$\lfloor \cdot \rfloor$	$\lfloor \cdot \rfloor$	\rightarrow	$\lfloor \cdot \rfloor$	$\lfloor \cdot \rfloor$	\rightarrow	$\lfloor \cdot \rfloor$	[·]

This will ensure that W_{AB} has the correct vertex measure. Note also that, by line (10.2),

$$\sum_{a'\in\mathbb{A}}W_{\mathbb{A}\mathbb{B}}(b,(a',b');i)=W_{\mathbb{B}}(b,b';i)\quad\text{ for all }b\in\mathbb{B}\setminus\{b_0\}\text{ and }b'\in\mathbb{B}.$$

Using this and definition (10.3), we also get

$$\sum_{a' \in \mathbb{A}} W_{\mathbb{AB}}(b_0, (a', b'); i) = W_{\mathbb{B}}(b_0, b'; i).$$

This will ensure that the B-marginal of W_{AB} is W_{B} .

We show now that the half-marginal $W_{AB}(\cdot, (\cdot, \cdot); i)$ is ℓ^1 -close to $W(\cdot, (\cdot, \cdot); i)$ by considering separately the contributions to the ℓ^1 distance from terms defined using equations (10.1), (10.2), and (10.3).

(10.1) terms: Each of the terms of W_{AB} defined using the floor in equation (10.1) is distance at most 1/n from the corresponding term of W; therefore the total contribution of these terms to the ℓ^1 distance is

s to the
$$\ell^1$$
 distance is
$$\sum_{\substack{b\in\mathbb{B}\setminus\{b_0\}\\a'\in\mathbb{A}\setminus\{a_0\},b'\in\mathbb{B}\\i\in[r]}}|W_{\mathbb{A}\mathbb{B}}(b,(a',b');i)-W(b,(a',b');i)|\leq |\mathbb{A}||\mathbb{B}|^2r/n.$$

(10.2) terms: By the triangle inequality,

$$\begin{split} |W_{\mathtt{AB}}(b,(a_0,b');i) - W(b,(a_0,b');i)| \\ &= \left| \left(W_{\mathtt{B}}(b,b';i) - \sum_{a' \neq a_0} W_{\mathtt{AB}}(b,(a',b');i) \right) \right. \\ &- \left(\pi_{\mathtt{B}} W(b,b';i) - \sum_{a' \neq a_0} W(b,(a',b');i) \right) \right| \\ &\leq |W_{\mathtt{B}}(b,b';i) - \pi_{\mathtt{B}} W(b,b';i)| \\ &+ \sum_{a' \neq a_0} |W_{\mathtt{AB}}(b,(a',b');i) - W(b,(a',b');i)|. \end{split}$$

The total contribution of such terms is therefore

$$\begin{split} \sum_{b \in \mathbb{B} \setminus \{b_0\}, \ b' \in \mathbb{B}} & |W_{\mathbb{A}\mathbb{B}}(b, (a_0, b'); i) - W(b, (a_0, b'); i)| \\ & \leq \sum_{i \in [r]} & |W_{\mathbb{B}}(b, b'; i) - (\pi_{\mathbb{B}})_* W(b, b'; i)| \\ & \leq \sum_{b \in \mathbb{B} \setminus \{b_0\}, \ b' \in \mathbb{B}} & |W_{\mathbb{B}}(b, b'; i) - (\pi_{\mathbb{B}})_* W(b, b'; i)| \\ & + \sum_{b \in \mathbb{B} \setminus \{b_0\}, \ b' \in \mathbb{B}} & |W_{\mathbb{A}\mathbb{B}}(b, (a', b'); i) - W(b, (a', b'); i)| \\ & \leq 2\delta + |\mathbb{A}| |\mathbb{B}|^2 r/n. \end{split}$$

(10.3) terms: Again applying the triangle inequality,

$$\begin{split} |W_{\mathbb{A}\mathbb{B}}(b_0, (a, b'); i) - W(b_0, (a, b'); i)| \\ &\leq |W_{\mathbb{A}\mathbb{B}}((a, b')) - W((a, b'))| \\ &+ \sum_{b \neq b_0} |W_{\mathbb{A}\mathbb{B}}(b, (a, b'); i) - W(b, (a, b'); i)|. \end{split}$$

Summing over all $a \in A$, $b' \in B$ and $i \in [r]$, we see that the total contribution of such terms is bounded by

$$\begin{split} \sum_{\substack{a \in \mathbb{A}, b' \in \mathbb{B} \\ i \in [r]}} \left[|W_{\mathbb{A}\mathbb{B}}((a,b')) - W((a,b'))| \\ &+ \sum_{\substack{b \neq b_0}} |W_{\mathbb{A}\mathbb{B}}(b,(a,b');i) - W(b,(a,b');i)| \right] \\ &+ \sum_{\substack{vertex \text{ measure} \\ b \in \mathbb{B}}} |W_{\mathbb{A}\mathbb{B}}(b,(a,b)) - W((a,b))| \\ &+ \sum_{\substack{i \in [r] \\ a' \in \mathbb{A} \setminus \{a_0\}, \ b' \in \mathbb{B} \\ i \in [r]}} |W_{\mathbb{A}\mathbb{B}}(b,(a',b');i) - W(b,(a',b');i)| \\ &+ \sum_{\substack{b \in \mathbb{B} \setminus \{b_0\}, \ b' \in \mathbb{B} \\ i \in [r]}} |W_{\mathbb{A}\mathbb{B}}(b,(a_0,b');i) - W(b,(a_0,b');i)| \\ &\leq r \cdot [6\delta + 4|\mathbb{A}||\mathbb{B}|/n] + [|\mathbb{A}||\mathbb{B}|^2r/n] + [2\delta + |\mathbb{A}||\mathbb{B}|^2r/n] \\ &\leq 8r\delta + 6|\mathbb{A}||\mathbb{B}|^2r/n. \end{split}$$

Adding up the contributions of the three types of terms, we see that the ℓ^1 distance between the half-marginals of W and W_{AB} is bounded by

$$10r\delta + 8|\mathbf{A}||\mathbf{B}|^2r/n.$$

10.2.1. *Nonnegativity*. Again, the preceding construction does not guarantee that all terms are non-negative. In the following we describe how to correct negativity.

Let -R/n be the sum of all negative terms of the half-marginal. As above, we get

$$R/n < 10r\delta + 7|\mathbf{A}||\mathbf{B}|^2 r/n$$
.

Suppose there is some $b_{-} \in B$, $(a'_{-}, b'_{-}) \in A \times B$, and $i \in [r]$ such that $W_{AB}(b_{-}, (a'_{-}, b'_{-}); i) < 0$. Then $W_{AB}(b_{-}, (a'_{-}, b'_{-}); i) \le -1/n$. Since

$$\sum_{a' \in \mathbb{A}} W_{\mathbb{A}\mathbb{B}}(b_{-}, (a', b'_{-}); i) = W_{\mathbb{B}}(b_{-}, b'_{-}; i) \ge 0$$

and

$$\sum_{b \in \mathbb{B}} W_{\text{AB}}(b, (a'_-, b'_-); i) = W_{\text{AB}}((a'_-, b'_-)) \ge 0$$

there exist $a'_+ \in A$ and $b_+ \in B$ such that

$$W_{AB}(b_-, (a'_+, b'_-); i) \ge 1/n$$
 and $W_{AB}(b_+, (a'_-, b'_-); i) \ge 1/n$.

Decrease both of these terms by 1/n, and increase both $W_{AB}(b_-, (a'_-, b'_-); i)$ and $W_{AB}(b_+, (a'_+, b'_-); i)$ by 1/n. This moves the half-marginal by ℓ^1 distance 4/n:

$$\sum_{a' \in \mathbb{N}} W_{AB}(b, (a', b'); i) = W_{B}(b, b'; i) \quad \text{and} \quad \sum_{b \in \mathbb{R}} W_{AB}(b, (a', b'); i) = W_{AB}((a', b')).$$

This step must be done at most R times to eliminate all negative entries, so the final half-marginal satisfies

$$\begin{split} & \sum_{i \in [r]} \sum_{b \in \mathbb{B}} \sum_{(a',b') \in \mathbb{A} \times \mathbb{B}} |W_{\mathbb{A}\mathbb{B}}(b,(a',b');i) - W(b,(a',b');i)| \\ & \leq (10r\delta + 8|\mathbb{A}||\mathbb{B}|^2 r/n) + R \cdot 4/n \leq 50r\delta + 36|\mathbb{A}||\mathbb{B}|^2 r/n. \end{split}$$

10.3. The edge measure. Finally, we define the edge measure of W_{AB} by

$$W_{AB}((a, b), (a', b'); i) = \frac{1}{n} \lfloor n \cdot W((a, b), (a', b'); i) \rfloor$$
for $a \neq a_0$ and $(a', b') \neq (a_0, b_0)$, (10.4)

$$W_{\text{AB}}((a_0,b),(a',b');i) = W_{\text{AB}}(b,(a',b');i) - \sum_{a \neq a_0} W_{\text{AB}}((a,b),(a',b');i)$$

for
$$(a', b') \neq (a_0, b_0),$$
 (10.5)

$$W_{AB}((a,b),(a_0,b_0);i) = W_{AB}((a,b)) - \sum_{(a',b')\neq(a_0,b_0)} W_{AB}((a,b),(a',b');i). \quad (10.6)$$

See Table 3.

	(a_0, b_0)	(a_1, b_0)	(a_2, b_0)	(a_0,b_1)	(a_1,b_1)	(a_2,b_1)	(a_0, b_2)	(a_1, b_2)	(a_2, b_2)
(a_0, b_0)	\rightarrow	\downarrow							
(a_1, b_0)	\rightarrow	$\lfloor \cdot \rfloor$							
(a_2, b_0)	\rightarrow	$\lfloor \cdot \rfloor$							
(a_0, b_1)	\rightarrow	\downarrow							
(a_1, b_1)	\rightarrow	$\lfloor \cdot \rfloor$							
(a_2, b_1)	\rightarrow	$\lfloor \cdot floor$	$\lfloor \cdot floor$	$\lfloor \cdot \rfloor$	$\lfloor \cdot floor$	$\lfloor \cdot \rfloor$	$\lfloor \cdot floor$	$\lfloor \cdot floor$	$\lfloor \cdot \rfloor$
(a_0, b_2)	\rightarrow	\downarrow							
(a_1, b_2)	\rightarrow	$\lfloor \cdot \rfloor$							
(a_2, b_2)	\rightarrow	$\lfloor \cdot \rfloor$							

TABLE 3. A diagram of how the edge measure $W_{AB}((\cdot, \cdot), (\cdot, \cdot); i)$ is chosen if $A = \{a_0, a_1, a_2\}$ and $B = \{b_0, b_1, b_2\}$. First obtain the entries marked $\lfloor \cdot \rfloor$ by rounding down entries of W. Then choose entries marked \downarrow according to equation (10.5), which ensures that the B half-marginal is the one chosen above. Then choose entries marked \rightarrow according to equation (10.6), which ensures that the vertex measure is the one chosen above.

It follows from this definition that W_{AB} is a (signed) weight with B-marginal W_{B} .

We now check that W_{AB} is ℓ^1 -close to W. We consider separately the contribution to the ℓ^1 distance of terms defined in equations (10.4), (10.5), and (10.6).

(10.4) terms: Each term of W_{AB} defined using the floor function in equation (10.4) is distance at most 1/n from the corresponding W term. The total contribution of these terms to the ℓ^1 distance is therefore at most $|A|^2|B|^2r/n$.

(10.5) terms: Applying the triangle inequality to terms defined in equation (10.5),

$$\begin{split} |W_{\mathbb{A}\mathbb{B}}((a_0,b),(a',b');i) - W((a_0,b),(a',b');i)| \\ & \leq |W_{\mathbb{A}\mathbb{B}}(b,(a',b');i) - W(b,(a',b');i)| \\ & + \sum_{a \neq a_0} |W_{\mathbb{A}\mathbb{B}}((a,b),(a',b');i) - W((a,b),(a',b');i)| \\ & \leq |W_{\mathbb{A}\mathbb{B}}(b,(a',b');i) - W(b,(a',b');i)| + |\mathbb{A}|/n. \end{split}$$

By the ℓ^1 bound on the distance between the half-marginals, the total contribution of all such terms is therefore at most

$$\sum_{i \in [r]} \sum_{b} \sum_{(a',b') \neq (a_0,b_0)} (|W_{AB}(b,(a',b');i) - W(b,(a',b');i)| + |A|/n)$$

$$\leq [50r\delta + 36|A|^2|B|^2r/n] + |A|^2|B|^2r/n$$

$$= 50r\delta + 37|A|^2|B|^2r/n.$$

(10.6) terms: Applying the triangle inequality to terms defined in equation (10.6),

$$\begin{split} |W_{\mathtt{AB}}((a,b),(a_0,b_0);i) - W_{\mathtt{AB}}((a,b),(a_0,b_0);i)| \\ & \leq |W_{\mathtt{AB}}((a,b)) - W((a,b))| \\ & + \sum_{(a',b') \neq (a_0,b_0)} |W_{\mathtt{AB}}((a,b),(a',b');i) - W((a,b),(a',b');i)|. \end{split}$$

Therefore the total contribution of all such terms is

$$\begin{split} &\sum_{i \in [r]} \sum_{a,b} |W_{\text{AB}}((a,b),(a_0,b_0);i) - W_{\text{AB}}((a,b),(a_0,b_0);i)| \\ &= \sum_{i \in [r]} \sum_{a,b} \left[|W_{\text{AB}}((a,b)) - W((a,b))| \right. \\ &+ \sum_{(a',b') \neq (a_0,b_0)} |W_{\text{AB}}((a,b),(a',b');i) - W((a,b),(a',b');i)| \right] \\ &= \underbrace{\sum_{i \in [r]} \sum_{a,b} |W_{\text{AB}}((a,b)) - W((a,b))|}_{\text{vertex measure}} \\ &+ \underbrace{\sum_{i \in [r]} \sum_{a \neq a_0} \sum_{b} \sum_{(a',b') \neq (a_0,b_0)} |W_{\text{AB}}((a,b),(a',b');i) - W((a,b),(a',b');i)|}_{(10.5) \text{ terms}} \\ &+ \underbrace{\sum_{i \in [r]} \sum_{b} \sum_{(a',b') \neq (a_0,b_0)} |W_{\text{AB}}((a_0,b),(a',b');i) - W((a_0,b),(a',b');i)|}_{(10.5) \text{ terms}} \\ &+ \underbrace{\sum_{i \in [r]} \sum_{b} \sum_{(a',b') \neq (a_0,b_0)} |W_{\text{AB}}((a_0,b),(a',b');i) - W((a_0,b),(a',b');i)|}_{(10.5) \text{ terms}} \\ &+ \underbrace{\sum_{i \in [r]} \sum_{b} \sum_{(a',b') \neq (a_0,b_0)} |W_{\text{AB}}((a_0,b),(a',b');i) - W((a_0,b),(a',b');i)|}_{(10.5) \text{ terms}} \\ &\leq 56r\delta + 41|\mathbf{A}|^2|\mathbf{B}|^2r/n. \end{split}$$

Summing up the contributions from terms of all three types, we get that

$$d_1(W_{AB}, W) \le 106r\delta + 79|A|^2|B|^2r/n.$$

10.3.1. *Nonnegativity.* We can modify a solution with negative entries to get a non-negative one similarly to above. Let -R/n be the sum of all negative entries; then

$$R/n \le 106r\delta + 78|\mathbf{A}|^2|\mathbf{B}|^2r/n.$$

Suppose there is some entry

$$W_{AB}((a_-, b_-), (a'_-, b'_-); i) \le -1/n.$$

We want to increment this term by 1/n without affecting the vertex measure or the B marginal. Since

$$\sum_{(a',b')\in\mathbb{A}\times\mathbb{B}} W_{\mathbb{AB}}((a_{-},b_{-}),(a',b');i) = W_{\mathbb{AB}}((a_{-},b_{-})) \geq 0$$

there exists some $(a'_+,b'_+) \in \mathbb{A} \times \mathbb{B}$ such that $W_{\mathbb{AB}}((a_-,b_-),(a'_+,b'_+);i) \geq 1/n;$ similarly, since

$$\sum_{a \in A} W_{AB}((a, b_{-}), (a', b'_{-}); i) = W_{AB}(b_{-}, (a'_{-}, b'_{-}); i) \ge 0$$

there exists some a_+ such that $W_{AB}((a_+, b_-), (a'_-, b'_-); i) \ge 1/n$. Increase

$$W_{\mathrm{AB}}((a_{-},b_{-}),(a_{-}',b_{-}');i) \quad \text{and} \quad W_{\mathrm{AB}}((a_{+},b_{-}),(a_{+}',b_{+}');i)$$

by 1/n, and decrease

$$W_{AB}((a_-, b_-), (a'_+, b'_+); i)$$
 and $W_{AB}((a_+, b_-), (a'_-, b'_-); i)$

by 1/n. This moves the weight by ℓ^1 distance 4/n.

Since R is the maximum number of times we need to do this before there are no more negative entries, the final weight satisfies

$$d_1(W_{AB}, W) \le 106r\delta + 79|A|^2|B|^2r/n + 4R/n \le 530r\delta + 391|A|^2|B|^2r/n.$$

To simplify, we write

$$d_1(W_{AB}, W) \le 530r(\delta + |A \times B|^2/n),$$

or

$$d(W_{AB}, W) \le 265r(\delta + |A \times B|^2/n).$$

Acknowledgements. Thanks go to Tim Austin for suggesting that results like Theorems B and C should hold, for many helpful discussions, and for providing comments on earlier drafts. Thanks also go to Ben Hayes for sharing helpful references to the operator algebras literature, and to Lewis Bowen and Brandon Seward for helpful conversations. Thanks go to an anonymous referee for many helpful comments on an earlier draft. This material is based upon work supported by the National Science Foundation under grant no. DMS-1855694.

REFERENCES

- [1] D. Airey, L. Bowen and F. Lin. A topological dynamical system with two different positive sofic entropies. *Trans. Amer. Math. Soc.* B **9**(2) (2022), 35–98.
- [2] L. Bowen. The ergodic theory of free group actions: entropy and the f-invariant. Groups Geom. Dyn. 4 (2010), 419–432.
- [3] L. Bowen. A measure-conjugacy invariant for free group actions. Ann. of Math. (2) 171(2) (2010), 1387–1400.
- [4] L. Bowen. Measure conjugacy invariants for actions of countable sofic groups. *J. Amer. Math. Soc.* 23(1) (2010), 217–217.
- [5] L. Bowen. Non-abelian free group actions: Markov processes, the Abramov–Rohlin formula and Yuzvinskii's formula. Ergod. Th. & Dynam. Sys. 30(6) (2010), 1629–1663.
- [6] L. Bowen. Sofic entropy and amenable groups. Ergod. Th. & Dynam. Sys. 32(2) (2012), 427–466.
- [7] L. Bowen and Y. Gutman. Nonabelian free group actions: Markov processes, the Abramov–Rohlin formula and Yuzvinskii's formula – Corrigendum. Ergod. Th. & Dynam. Sys. 33(2) (2013), 643–645.
- [8] A. Coja-Oghlan, M. Hahn-Klimroth, P. Loick, N. Müller, K. Panagiotou and M. Pasch. Inference and mutual information on random factor graphs. 38th International Symposium on Theoretical Aspects of Computer Science (STACS 2021) (Leibniz International Proceedings in Informatics (LIPIcs), 187). Ed. M. Bläser and B. Monmege. Schloss Dagstuhl Leibniz-Zentrum für Informatik, Dagstuhl, Germany, 2021, pp. 24:1–24:15.
- [9] T. M. Cover and J. A. Thomas. *Elements of Information Theory*, 2nd edn. Wiley-Interscience, Hoboken, NJ, 2006.
- [10] K. Dykema, D. Kerr and M. Pichot. Sofic dimension for discrete measured groupoids. *Trans. Amer. Math. Soc.* 366(2) (2013), 707–748.
- [11] B. Hayes. Relative entropy and the Pinsker product formula for sofic groups. Groups Geom. Dyn. 15 (2016), 413–463.
- [12] D. Kerr and H. Li. Entropy and the variational principle for actions of sofic groups. *Invent. Math.* 186(3) (2011), 501–558.

- [13] D. S. Ornstein and B. Weiss. Entropy and isomorphism theorems for actions of amenable groups. *J. Anal. Math.* 48(1) (1987), 1–141.
- [14] L. Păunescu. On sofic actions and equivalence relations. J. Funct. Anal. 261(9) (2011), 2461–2485.
- [15] S. Popa. Independence properties in subalgebras of ultraproduct Π_1 factors. *J. Funct. Anal.* **266**(9) (2014), 5818–5846.