



HAL
open science

Compositional statistical mechanics, entropy and variational inference

Grégoire Sergeant-Perthuis

► **To cite this version:**

Grégoire Sergeant-Perthuis. Compositional statistical mechanics, entropy and variational inference. Twelfth Symposium on Compositional Structures (SYCO 12), Apr 2024, Birmingham (UK), United Kingdom. hal-04518736

HAL Id: hal-04518736

<https://hal.sorbonne-universite.fr/hal-04518736v1>

Submitted on 24 Mar 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

Compositional statistical mechanics, entropy and variational inference

Grégoire Sergeant-Perthuis^{*1}

¹Laboratoire de Biologie Computationnelle et Quantitative, Sorbonne Université, Paris , France , ORCID: 0000-0002-2079-3410

March 24, 2024

Abstract

In this document, accepted for the Twelfth Symposium on Compositional Structures (SYCO 12) [1], we aim to gather various results related to a compositional/categorical approach to *rigorous* Statistical Mechanics [2–9]. Rigorous Statistical Mechanics is centered on the mathematical study of statistical systems. Central concepts in this field have a natural expression in terms of diagrams in a category that couples measurable maps and Markov kernels [8]. We showed that statistical systems are particular representations of partially ordered sets (posets), that we call \mathcal{A} -specifications, and expressed their phases, i.e., Gibbs measures, as invariants of these representations. It opens the way to the use of homological algebra to compute phases of statistical systems. Two central results of rigorous Statistical Mechanics are, firstly, the characterization of extreme Gibbs measure as it relates to the zero–one law for extreme Gibbs measures, and, secondly, their variational principle which states that for translation invariant Hamiltonians, Gibbs measures are the minima of the Gibbs free energy. We showed in [9] how the characterization of extreme Gibbs measures extends to \mathcal{A} -specifications; we proposed in [10] an Entropy functional for \mathcal{A} -specifications and gave a message-passing algorithm, that generalized the belief propagation algorithm of graphical models (see O. Peltre’s SYCO 8 talk, [11] and [12]), to minimize variational free energy.

1 Introduction

Context: Statistical physics is a framework that focuses on the probabilistic description of complex systems: a collection of interacting ‘particles’

^{*}gregoireserper@gmail.com

or components of a whole in most generality [13]. Its main feature is to introduce an energy function H for the system, that associates to any of its configurations a real value; the probability p of a configuration is given in terms of the Boltzmann distribution ($p = e^{-\beta H} / \int dx e^{-\beta H(x)}$), which associates high-energy configurations to unlikely events. We will call ‘statistical system’ the complex system provided with the Boltzmann distribution. Statistical physics serves as a rich framework for probabilistic modeling. It has several names depending on the community [14]; for example, it is called ‘energy-based modeling’ in machine learning, the probabilistic model is called a Gibbs Random Field in graphical modeling and is a particular case of a statistical system. It is widely used in engineering, for example: computational structural biology (in computational statistical physics: Hamiltonian Monte Carlo) [15, 16], robotics (reinforcement learning: Markov Chains, Markov Decision Processes) [17], and more generally for modeling interaction and dependencies of random variables using graphical models [18].

Applied Category Theory: New foundations, based on topology and geometry, were proposed for probability theory, information theory, and deep learning [19–23]. More generally, they fall in a field of research that has recently emerged, Applied Category Theory, which focuses on applying these principles to engineering [20–22, 24–26]. Recent results give a characterization of the zero–one law for independent random variables and for Markov chains in a categorical formulation [27, 28]. The zero–one law for extreme Gibbs measures is known to extend the ones of independent random variables and Markov chains [29], so it would be expected that the categorical formulation of extreme Gibbs measures we propose may also relate to the categorical formulation developed in the case of independent random variables and Markov chains.

Motivation: A landmark rigorous formulation of Statistical Mechanics can be found in Georgii’s *Gibbs Measures and Phase Transition* [29]. Such a framework, which revolves around the concepts of ‘specification’ and Hamiltonian, is needed to define rigorously pure phases or, more generally, Gibbs measures of statistical systems. For a given Hamiltonian, there can be multiple phases only for infinitely many interacting particles. The main constructions of such a framework rely on the necessity to have a *universe*, denoted Ω , that encompasses all possible configurations of the system, i.e., all the possible joint configurations of the particles $\Omega := \prod_{i \in \mathbb{N}} E_i$, with E_i the state space of the particle i . There are important limitations to making references to a ‘global’ universe [19]. Let us state two of them. The first one is that it can be difficult to compute with Hamiltonians (e.g., computing free energies [15], expectations of observables [30], infinite volume Gibbs measures [29]) as it requires summing over a set which ‘size’ increases exponentially with the number of components (variables) involved in the system. The second one is that they don’t account for heterogeneity, incompleteness, or incompatibility in the description of the system to model: by definition,

each configuration of the system corresponds to observing simultaneously and without loss of information the state of all its components. This second point is related to the first one as to model possible configurations of a system it is required to weight all possible simultaneous configurations of its particles; this forbids granularity of the description of the system: invariably, the dimension of the space of possible models increases exponentially with the number of (elementary) particles. Mathematically, Hamiltonians are sums of local potentials, however, the relation between potentials and their phases is highly nontrivial and noncompositional. In other words, the Hamiltonian point of view does not allow for building complex statistical systems from simple ones in a way that allows controlling and computing the phases of the associated statistical systems from the phases of the simpler one. The relation between the complexity of the interactions that appear in the statistical system and its number of (pure) phases (more than 2) is known only for very specific systems (e.g. Ising model); results in this direction are considered difficult problems [31, 32]. In our approach, we want to forget about the ‘ Ω ’ and we want to study statistical systems ‘locally’ focusing on their local interactions and how they compose. We propose a constructive approach: ‘composing’ statistical systems together by building in a controlled and computable manner nontrivial statistical systems from simpler ones. This is what reformulating specifications as presheaves allows us to do, in particular, by making use of operations such as coproducts and products but also geometric morphisms. The approach we propose also has the advantage to model statistical systems for which it is not possible to have complete knowledge of the states of the particles at the same time and that could also have conflicting local descriptions inconsistent with a global description.

Contribution: In this paper, we gather various results related to a compositional/categorical approach to ‘rigorous’ Statistical Mechanics [2–9].

Our motivation in [7] was to build a bridge between geometry and rigorous statistical physics by showing that the standard formal definition of a statistical system, a ‘specification’ (Definition 1.23 [29]), can be identified with a particular representation of a partially ordered set (poset), i.e., a functor of a poset in a certain category. For a given poset \mathcal{A} , we call such representation an \mathcal{A} -specification. Representations of posets have a precise geometric interpretation [33–37], and there is a rich literature coming from algebra, geometry, and topology to study them. We showed that phases of statistical systems are geometric invariants of these representations and computed them for ‘projective’ poset representations. We remarked that, in our setting, systems of ‘finite size’ can have multiple phases, which is not possible in the current formalism for phase transition, but we do not claim that those finite size systems capture correctly the possible phase transitions of infinite size systems. However, we believe that the algebraic treatment of statistical systems we propose to be fairly similar for finite size systems and

infinite size systems.

In [4], we gave a characterization of independent variables in terms of projective objects in the category of presheaves over a poset and an easy-to-verify condition that characterizes such objects. Injective representations were characterized in [3], and their relationship with the marginal extension problem is studied in [5]; a unifying perspective for Hilbert spaces can be found in [6]. In [10], we proposed a general framework for optimization of presheaves which in particular allows defining an Entropy functional for \mathcal{A} -specifications. In [9], we showed how the characterization of extreme Gibbs measures, one of the steps for proving a zero-one law for the extreme Gibbs measures, transfers to \mathcal{A} -specifications.

2 Structure of the Paper and Contribution

We start by recalling some important notions of rigorous statistical mechanics [29] (specifications, Gibbs measures). We then show that the standard concept that encodes the statistical system, which is the notion of 'specifications', can be extended into a poset representation; we also show that the Gibbs measures of a statistical system are geometric invariants of this representation. From there, we propose a novel categorical formulation for statistical mechanics; we define 'generalized specifications' and 'generalized Gibbs measures' for those specifications: the \mathcal{A} -specifications. We characterize the generalized Gibbs measures for projective \mathcal{A} -specifications and their Gibbs measures. We then recall what extreme Gibbs measures are and their characterization on the tail σ -algebra. We extend such characterization to extreme Gibbs measures of \mathcal{A} -specifications. Finally, we recall what the Bethe free energy is in the context of graphical models and what the Belief Propagation algorithm is; when the graphical model is a Hidden Markov model (HMM), the Belief Propagation relates to Kalman filtering. We explain how we can define such variational free energy for \mathcal{A} -specifications. We then give a message-passing algorithm for finding critical points of the variational free energy of \mathcal{A} -specifications.

3 Background: Rigorous Statistical Mechanics

We will follow the presentation of Georgii's reference book *Gibbs Measures and Phase Transitions* [29].

Definition 1 (Markov Kernel). A Markov kernel k from the measurable space (E, \mathcal{E}) to the measurable space (E_1, \mathcal{E}_1) is a function $k : \mathcal{E}_1 \times E \rightarrow [0, 1]$ such that

1. $\forall \omega \in E, k(\cdot|\omega)$ is a measure on E_1

2. $\forall A \in \mathcal{E}_1$, $k(A|\cdot)$ is a measurable map from E to \mathbb{R}
3. $\forall \omega \in E$, $k(E_1|\omega) = 1$, i.e., $k(\cdot|\omega)$ is a probability measure.

We will denote a Markov kernel k from (E, \mathcal{E}) to (E_1, \mathcal{E}_1) as $k : (E, \mathcal{E}) \rightarrow (E_1, \mathcal{E}_1)$. We denote $[(E, \mathcal{E}), (E_1, \mathcal{E}_1)]_K$, the set of kernels from (E, \mathcal{E}) to (E_1, \mathcal{E}_1) ; if there is no ambiguity on the σ -algebras the spaces are provided with, we will simply denote it as $[E, E_1]_K$. We denote $\mathbb{P}(E)$, the space of probability distributions over E ; it is a measurable space for the smallest σ -algebra that makes the evaluation maps, on measurable sets of E , measurable.

Markov kernels can be composed as follows. Let $k : E \rightarrow E_1$, $k_1 : E_1 \rightarrow E_2$ be two Markov kernels, then the composition $k_1 \circ k : E \rightarrow E_2$ is the following Markov kernel: for any $A \in \mathcal{E}_2$ and $\omega \in E$,

$$k_1 \circ k(A|\omega) = \int k_1(A|\omega_1)k(d\omega_1|\omega) \quad (3.1)$$

A measurable map $f : (E, \mathcal{E}) \rightarrow (E_1, \mathcal{E}_1)$ between two measurable spaces can be extended into the following Markov kernel: for any $A \in \mathcal{E}_1$ and $\omega \in E$,

$$k_f(A|\omega) = 1[f(\omega) \in A] \quad (3.2)$$

To avoid having too many notations, we will denote k_f as f and the context will specify if f refers to the measurable map or its extension; for example, a composition $k \circ f$ between a Markov kernel k and f necessarily means that, here, f refers to k_f . We will also denote $k_1 \circ k$ as $k_1 k$.

A probability measure $p \in \mathbb{P}(E)$ can also be identified with the following Markov kernel k_p from $*$, the measurable space with one element, to (E, \mathcal{E}) ; for any $A \in \mathcal{E}$, $k_p(A|*) = p(A)$. Similarly, we identify p and k_p .

Remark 1. Measurable spaces and measurable maps form a category. Giry [38] and Lawvere [39] are the first to have remarked that measurable spaces and Markov kernels also form a category; the latter category is the Kleisli category of the first.

Definition 2 (Proper Kernel, Section 1.1. [29]). Let $\mathcal{E}_1 \subseteq \mathcal{E}$ be two σ -algebras of a set E , a kernel $k \in [(E, \mathcal{E}_1), (E, \mathcal{E})]_K$ is proper if and only if, for any $A \in \mathcal{E}$, $B \in \mathcal{E}_1$,

$$k(A \cap B|\cdot) = k(A|\cdot)1_B \quad (3.3)$$

Let us set the notations. I is the set of components of a complex system. $(E_i, \mathcal{E}_i, i \in I)$ is a collection of measurable spaces, with each E_i being the space of configuration (state space) of the component $i \in I$. (E, \mathcal{E}) denotes the state space of the system; it is the product $E := \prod_{i \in I} E_i$ with the

product σ -algebra. For a sub-collection of components $a \subseteq I$, we denote $E_a := \prod_{i \in a} E_i$ the associated state space, and \mathcal{E}_a the associated product σ -algebra. $i^a : E \rightarrow E_a$ is the projection that sends a configuration $\omega := (\omega_i, i \in I)$ to the configuration of the sub-collection a , $(\omega_i, i \in a)$. Finally, let us denote $\mathcal{P}_f(I)$ the set of finite subsets of I .

Definition 3 (Specification, Adaptation of Def. 1.23 [29]). A specification γ with state space (E, \mathcal{E}) is a collection $(\gamma_a, a \in \mathcal{P}_f(I))$ of proper Markov kernels such that for any $a \in \mathcal{P}_f(I)$, $\gamma_a \in [(E_{\bar{a}}, \mathcal{E}_{\bar{a}}), (E, \mathcal{E})]_K$ and which satisfies that for any $b \subseteq a$, i.e. $\bar{a} \subseteq \bar{b}$ and any $A \in \mathcal{E}$,

$$\gamma_b \circ i^{\bar{b}} \circ \gamma_a(A|\cdot) = \gamma_a(A|\cdot) \quad (3.4)$$

Remark 2. In Definition (1.23) [29], E is the product over the same measurable space X over I , and I is a countably infinite set.

In the standard definition of a specification (Definition 3), the Markov kernels encode border conditions for experiments that only involve finite numbers of components ($a \in \mathcal{P}_f(I)$). It is what formally encodes the statistical system. The next definition defines the ‘phases’ of the system.

Definition 4 (Gibbs measures, Def. 1.23 [29]). Let γ be a specification with state space E ; the set of probability measures,

$$\mathcal{G}(\gamma) := \{p \in \mathbb{P}(E) : \mathbb{E}_p[A|\mathcal{E}_{\bar{a}}](\omega) = \gamma_a(A|\omega_a) \text{ p a.s.}\} \quad (3.5)$$

is the set of Gibbs measures of γ .

One of the central problems of (rigorous) statistical mechanics is to understand the relationship between a specification γ (statistical system) and its set of Gibbs measures $\mathcal{G}(\gamma)$ (its phases).

4 Statistical Systems as poset representations

Theorem 1. *Let γ be a specification with state space E . For any $a, b \in \mathcal{P}_f(I)$ such that $b \subseteq a$, there is a unique Markov kernel $F_b^a : E_{\bar{a}} \rightarrow E_{\bar{b}}$ such that the following diagram commutes,*

$$\begin{array}{ccc} E_{\bar{a}} & \xrightarrow{\gamma_a} & E \\ F_b^a \downarrow & \nearrow \gamma_b & \\ E_{\bar{b}} & & \end{array} \quad (4.1)$$

i.e. such that $\gamma_b \circ F_b^a = \gamma_a$. Furthermore for any collection $a, b, c \in \mathcal{P}_f(I)$ with $a \subseteq b \subseteq c$,

$$F_c^b \circ F_b^a = F_c^a \quad (4.2)$$

Proof. Let $b \subseteq a$, let F_b^a satisfy the commutative diagram 4.1, then,

$$i^{\bar{b}} \gamma_b F_b^a = i^{\bar{b}} \gamma_a \quad (4.3)$$

therefore,

$$F_b^a = i^{\bar{b}} \gamma_a \quad (4.4)$$

For any $a, b \in \mathcal{A}$ such that $b \subseteq a$ let $F_b^a = i^{\bar{b}} \gamma_a$ then for $c \subseteq b \subseteq a$,

$$F_c^b F_b^a = i^{\bar{c}} \gamma_b i^{\bar{b}} \gamma_a \quad (4.5)$$

Equation 3.4 can be rewritten as, for $b \subseteq a$,

$$\gamma_b i^{\bar{b}} \gamma_a = \gamma_a \quad (4.6)$$

therefore,

$$F_c^b F_b^a = F_c^a \quad (4.7)$$

□

Definition 5 (Partially ordered set). A partially ordered set (poset), (\mathcal{A}, \leq) , is a set \mathcal{A} provided with a binary relation $\leq: \mathcal{A} \times \mathcal{A} \rightarrow \{0, 1\}$, such that

1. reflexive: $\forall a \in \mathcal{A}, a \leq a$
2. transitive: if $c \leq b$ and $b \leq a$ then $c \leq a$
3. antisymmetric: if $b \leq a$ and $a \leq b$ then $a = b$

$(\mathcal{P}_f(I), \subseteq)$ is a poset for the inclusion relation; $(\mathcal{P}_f(I), \supseteq)$ is also a poset for the reversed inclusion relation: $b \supseteq a \iff a \subseteq b$. The convention is to denote $(\mathcal{P}_f(I), \supseteq)$ as $\mathcal{P}_f(I)^{op}$ because what relates $(\mathcal{P}_f(I), \supseteq)$ to $\mathcal{P}_f(I)$ is the fact that the order is ‘opposed’; the same convention holds for any poset; \mathcal{A}^{op} is the set \mathcal{A} with opposed order.

We call a representation of the poset \mathcal{A} : a collection of ‘spaces’ $(G(a), a \in \mathcal{A})$ and a collection of ‘maps’ $(G_a^b: G(b) \rightarrow G(a); b, a \in \mathcal{A}, b \leq a)$ which satisfies for any $c \leq b \leq a$, $G_c^b \circ G_b^a = G_c^a$. We keep the notion of ‘poset representation’ a bit vague for now (we do not say what ‘spaces’ or ‘maps’ are); we keep this notion vague at this stage but in the next section will make this notion formal by introducing the concept of category and of functor.

We call a representation of the poset \mathcal{A} : a collection of ‘spaces’ $(G(a), a \in \mathcal{A})$ and a collection of ‘maps’ $(G_a^b: G(b) \rightarrow G(a); b, a \in \mathcal{A}, b \leq a)$ which

satisfies for any $c \leq b \leq a$, $G_c^b \circ G_b^a = G_c^a$. We keep the notion of ‘poset representation’ a bit vague for now (we do not say what ‘spaces’ or ‘maps’ are); we keep this notion vague at this stage but in the next section will make this notion formal by introducing the concept of category and of functor.

Theorem 1 implies that a specification γ can be promoted to a representation F of $(\mathcal{P}_f(I), \supseteq)$.

Theorem 2. *Let γ be a specification with state space E , let $p \in \mathcal{G}(\gamma)$. For any $a \in \mathcal{A}$ let $p_a := i^{\bar{a}} \circ p$; it is the marginal distribution on $E_{\bar{a}}$ of p . Then, for any $a, b \in \mathcal{P}_f(I)$ such that $b \subseteq a$,*

$$F_b^a \circ p_a = p_b \quad (4.8)$$

Proof. For any $a, b \in \mathcal{P}_f(I)$ such that $b \subseteq a$,

$$\gamma_b \circ (F_b^a p_a) = p \quad (4.9)$$

therefore, $F_b^a p_a = i^{\bar{b}} p$ and

$$F_b^a p_a = p_b \quad (4.10)$$

□

5 Categorical formulation of specification and Gibbs measures

We denote categories in bold, e.g. **C**. We will denote **Mes** as the category that has as objects measurable spaces and as morphisms measurable maps (Section 1 [38]); we will denote **Kern** as the category that has as objects measurable spaces and as morphisms Markov kernels (the Kleisli category of the monad **Mes**).

A poset, (\mathcal{A}, \leq) , can be seen as a category, **A**, with at most one morphism between two objects: the objects of **A** are the elements of \mathcal{A} and for any two elements $b, a \in \mathcal{A}$ there is one morphism $b \rightarrow a$ when $b \leq a$. From now on, we will drop the bold notation for the category **A** and denote it simply as \mathcal{A} . A functor G from a poset \mathcal{A} to a category **C** is precisely a collection of maps G_a^b for $b, a \in \mathcal{A}$ such that $b \leq a$, which satisfy $G_a^b \circ G_b^c = G_a^c$ for any three elements $c \leq b \leq a$. A functor from \mathcal{A} to some target category is what we will call a representation of the poset \mathcal{A} ; in general, the target category is the category of vector spaces or modules [33]. For this article, the target category will be **Mes** and **Kern**.

Consider a functor $G : \mathcal{A} \rightarrow \mathbf{Set}$ from a poset \mathcal{A} to the category of sets. A collection $(\omega_a, a \in \mathcal{A})$ is called a section of G if for any $b \leq a$, $G_a^b(\omega_b) = \omega_a$; the set of sections of a poset representation is called the limit of G (III.4 [40]) and denoted $\lim G$. It is an ‘invariant’ of G that can

be computed using homological algebra when the target category of G is enriched with some algebraic structure (Chapter 13 [41]).

In Theorem 1 we showed that we can associate to a specification γ , a functor from $(\mathcal{P}_f(I), \supseteq)$ to **Kern**. The convention is to call a functor with source \mathcal{A}^{op} , a *presheaf*. In Theorem 2 we showed that Gibbs measures of γ are ‘sections’ of F . We will denote this set of ‘sections’ as $[\ast, F]_{K, \mathcal{A}}$; more precisely for a functor $F : \mathcal{A} \rightarrow \mathbf{Kern}$,

$$[\ast, F]_{K, \mathcal{A}} := \{(p_a \in \mathbb{P}(F(a)), a \in \mathcal{A}) \mid \forall b \leq a, F_a^b p_b = p_a\} \quad (5.1)$$

Introducing the presheaf F and $[\ast, F]_{K, \mathcal{A}}$ is our way to emphasize that the compatible measures $p \in [\ast, F]_{K, \mathcal{A}}$ don’t have to be measures over the whole space E . There is no need to define statistical systems ‘globally’, one can also define them ‘locally’.

Let us now introduce the more general, categorical setting we propose for statistical systems.

Definition 6 (Generalized Specification, \mathcal{A} -Specifications). Let \mathcal{A} be a poset, a generalized specification over \mathcal{A} , or simply \mathcal{A} -specification, is a couple (G, F) of a presheaf and a functor where $G : \mathcal{A}^{op} \rightarrow \mathbf{Mes}$ and $F : \mathcal{A} \rightarrow \mathbf{Kern}$ are such that for any $a, b \in \mathcal{A}$ with $b \leq a$,

$$G_b^a F_a^b = \text{id} \quad (5.2)$$

In the previous definition, G , in the particular case of $E = \prod_{i \in I} E_i$, encodes the collection of projections $i_b^a : E_a \rightarrow E_b$ for $b \subseteq a$; it is in some way the ‘skeleton’ of the spaces of observables of the statistical system. It is a key ingredient for the generalization of (rigorous) statistical mechanics to a categorical framework.

Definition 7 (Gibbs measures for \mathcal{A} -specifications). Let $\gamma = (G, F)$ be an \mathcal{A} -specification, we call the Gibbs measures of γ the sections of F ,

$$\mathcal{G}_g(\gamma) := [\ast, F]_{K, \mathcal{A}} \quad (5.3)$$

6 Gibbs measure of projective \mathcal{A} -specifications

Let E be a measurable space; we denote $L^\infty(E)$ the set of bounded, real-valued, measurable functions over E . One associates to a Markov kernel $F : E_b \rightarrow E_a$ a linear map $\pi : L^\infty(E_a) \rightarrow L^\infty(E_b)$ defined as follows: for any $f \in L^\infty(E_a)$,

$$\forall \omega_b \in E_b, \quad \pi(f)(\omega_b) = \int f(\omega_a) F(d\omega_a | \omega_b) \quad (6.1)$$

This association is ‘functorial’, we may denote the underlying functor $L^\infty : \mathbf{Kern}^{op} \rightarrow \mathbf{Vect}$ which is presheaf from the category of Markov kernels to the category of vector spaces. It is the presheaf that associates spaces to their space of observables. Let us denote $L^\infty \circ G : \mathcal{A} \rightarrow \mathbf{Vect}$ as i and $L^\infty \circ F : \mathcal{A}^{op} \rightarrow \mathbf{Vect}$ as π . In these notation one has that for any $a, b \in \mathcal{A}$ such that $b \leq a$ then $\pi_b^a \circ i_a^b = \text{id}$.

For the definition and characterization of projective presheaf over a poset see [34–36].

Definition 8 (Projective \mathcal{A} -specifications). An \mathcal{A} -specification (G, F) is called projective when $L^\infty \circ F$ is a projective presheaf (in \mathbf{Vect}). In other words, there is a collection of presheaves $(S_a, a \in \mathcal{A})$ such that, $L^\infty \circ F \cong \bigoplus_{a \in \mathcal{A}} S_a$ where for any $b \geq a$, $S_a(b)$ is a constant vector space denoted S_a and $S_{a_c}^b = \text{id}$ for any $a \leq c \leq b$ and $S_a(b) = 0$ if $b \not\geq a$. The collection of presheaves $(S_a, a \in \mathcal{A})$ is called the decomposition of (G, F) .

For any poset \mathcal{A} , symmetrizing the order defines the following equivalence relation,

$$\forall a, b \in \mathcal{A}, a \sim b \iff a \leq b \text{ or } b \leq a \quad (6.2)$$

The equivalence classes of this equivalence relation are the connected components of \mathcal{A} that we will denote as $\mathcal{C}(\mathcal{A})$. To each element of $a \in \mathcal{A}$ one can associate its connected component $\mathcal{C}(a)$. If each connected component has a minimum element, in other words, if for any $C \in \mathcal{C}(\mathcal{A})$, and any $b \in C$, there is $c \in C$ such that, $c \leq b$, then we shall denote, $\mathcal{C}_*(\mathcal{A})$ as the collection of these minimum elements; if not $\mathcal{C}_*(\mathcal{A}) = \emptyset$.

To conclude this article let us characterize Gibbs measures of projective \mathcal{A} -specifications.

Theorem 3. *Let $\gamma = (G, F)$ be a projective \mathcal{A} -specification. If at least one of the connected components of \mathcal{A} does not have a minimum element, i.e. when,*

$$\mathcal{C}_*(\mathcal{A}) = \emptyset \quad (6.3)$$

then,

$$\mathcal{G}_g(\gamma) = \emptyset \quad (6.4)$$

if not,

$$\mathcal{G}_g(\gamma) = \prod_{a \in \mathcal{C}_*(\mathcal{A})} \mathbb{P}(G(a)) \quad (6.5)$$

Proof. Let us denote $L^\infty G$ as i and, $L^\infty F$ as π ; (i, π) is decomposable, let $(S_a, a \in \mathcal{A})$ be its decomposition. For any $a, b \in \mathcal{A}$ such that $b \leq a$ and $\mu \in \mathcal{G}(\gamma)$, let us denote $L^\infty \mu$ as ν . For any $v \in L^\infty F(a)$,

$$\nu_b \pi_b^a \left(\sum_{c \leq a} S_c(a)(v) \right) = \nu_a \left(\sum_{c \leq a} S_c(a)(v) \right) \quad (6.6)$$

where $S_c(a)(v)$ denotes the projection of v on $S_c(a)$; therefore,

$$\nu_b \left(\sum_{c \leq b} S_{cb}^a(v) \right) = \nu_a i_a^b \left(\sum_{c \leq b} S_{cb}^a(v) \right) = \nu_a \left(\sum_{c \leq a} S_c(a)(v) \right) \quad (6.7)$$

and so,

$$\nu_a \left(\sum_{\substack{c \leq a \\ c \not\leq b}} S_c(a) \right) = 0 \quad (6.8)$$

Therefore for any $a \notin \mathcal{C}_*(\mathcal{A})$, $\nu_a|_{S_a(a)} = 0$. Furthermore,

$$\operatorname{colim} i \cong \bigoplus_{a \in \mathcal{A}} S_a(a). \quad (6.9)$$

ν is uniquely determined by $(\nu_a|_{S_a(a)}, a \in \mathcal{A})$; if there is a connected component $C \in \mathcal{C}(\mathcal{A})$ that does not have a minimal element, for any $a \in C$,

$$\nu_a|_{S_a} = 0 \quad (6.10)$$

Therefore for any $a \in C$, $\nu_a = 0$; this is contradictory with the fact that $\mu_a \in \mathbb{P}(\gamma(a))$ and so,

$$\mathcal{G}_g(\gamma) = \emptyset \quad (6.11)$$

When $\mathcal{C}_*(\mathcal{A})$ is non empty for any functor, H , from \mathcal{A} to **Set**,

$$\lim H \cong \prod_{a \in \mathcal{C}_*(\mathcal{A})} H(a) \quad (6.12)$$

therefore,

$$\mathcal{G}_g(\gamma) = \prod_{a \in \mathcal{C}_*(\mathcal{A})} \mathbb{P}(G(a)) \quad (6.13)$$

□

7 Characterization of extreme Gibbs measures of \mathcal{A} -specifications

We now turn to the characterization of extreme Gibbs measures (Theorem 7.7 [29]), which is one of the steps for proving a zero-one law for the extreme Gibbs measures. We will show how it transfers to \mathcal{A} -specifications. In the classical theory of rigorous statistical mechanics, the tail σ -algebra generates the observables for which a ‘generalized’ law of large numbers (zero–one law) holds. Following [9], we give a candidate for such a σ -algebra in the categorical setting and show the associated extreme Gibbs measures decomposition.

Let us recall the definition of the tail σ -algebra in the classical formulation. Let us first consider the case of time series, i.e. $\Omega = \prod_{i \in \mathbb{N}} E_i$ with E_i are measurable spaces; let us denote $\mathcal{E}_{\geq k}$ as the σ -algebra generated by the cylinders $\prod_{n \geq k} E_n$. The tail σ -algebra is defined as $\bigcap_{k \in \mathbb{N}} \mathcal{E}_{\geq k}$. For specifications, I is any set and $\mathcal{E}_{\bar{a}}$ is indexed by a subset $a \subseteq I$ that is finite. We defined $\mathcal{E}_{\bar{a}}$ to be a σ -algebra of $E_{\bar{a}}$, however it can also be identified with the smallest σ -algebra on E that make $i^{\bar{a}} : E, \mathcal{E} \rightarrow E_{\bar{a}}, \mathcal{E}_{\bar{a}}$ measurable. Through this identification one defined $\mathcal{E}_{\infty} := \bigcap_{a \in \mathcal{P}_f(I)} \mathcal{E}_{\bar{a}}$.

For a functor from \mathcal{A} to **Mes**, let us denote $\sigma(G(a))$ as the σ -algebra of the measurable space $G(a)$, where $a \in \mathcal{A}$, and $\sigma(G)$ as the underlying functor defined as $\sigma(G)_a^b A_b := G_b^{a-1} A_b$, with $b \leq a$. We propose that one candidate that plays the role of the tail σ -algebra for a given specification $\gamma = (G, F)$ is $\lim \sigma(G)$ defined as,

$$\lim \sigma(G) := \{(A_a \in \sigma(G(a)), a \in \mathcal{A}) \mid \forall a, b \in \mathcal{A}, \quad A_a = G_b^{a-1} A_b\} \quad (7.1)$$

Let us denote $1_A : E \rightarrow \{0, 1\}$ the indicator function over the set A that sends $\omega \in A$ to 1 and $\omega \notin A$ to 0. Let us remark that $1_{A_b} \circ G_b^a = 1[G_b^a(\omega_a) \in A_b] = 1_{G_b^{a-1} A_b}$. Remark that $A \in \lim \sigma(G)$ is equivalent to $1_A \in \lim i$; in other words, $\lim \sigma(G)$ is the restriction of $\lim i$ to indicator functions of the form $1_{A_a}, a \in \mathcal{A}$.

Finally, we also need to recall that for any $f \in L^\infty(E)$ and $\mu \in \mathbb{P}(E)$, one can define a measure $f \cdot \mu$ as $f \cdot \mu(d\omega) = f(\omega)d\omega$.

The key proposition of this section is Proposition 1; the proof of that proposition is given for $G(a), a \in \mathcal{A}$ finite measurable sets.

Assumption: Therefore, we assume in what follows that the measurable sets $G(a)$ are finite.

However, there is no finiteness constraint on \mathcal{A} . A weaker version holds when $G(a)$ is not finite. We will say that $F > 0$ when for any $a, b \in \mathcal{A}$, such

that $b \leq a$, $F(\omega_a|\omega_b) > 0$ for any ω_b such that $G_b^a(\omega_a) = \omega_b$; $G \circ F = \text{id}$ requires that $F(\omega_a|\omega_b) = 0$ when $G_b^a(\omega_a) \neq \omega_b$.

The following lemma is an extension of the classical result that states that conditioning over a σ -subalgebra $\mathcal{F}_1 \subseteq \mathcal{F}$ defines a morphism of modules when finer (\mathcal{F} measurable) observables are seen as modules over the coarser (\mathcal{F}_1 measurable observables).

Lemma 1. *Let E_1, E_2 be two measurable spaces, let $g : E_2 \rightarrow E_1$ be a measurable map and $f : E_1 \rightarrow E_2$ be a Markov kernel so that, $f \circ g = \text{id}$. Let us denote respectively i and π the induced linear maps on $L^\infty(E_1), L^\infty(E_2)$. Let $h \in L^\infty(E_2)$ and $k \in L^\infty(E_1)$, then,*

$$\pi(h).k = \pi(h.i(k)) \quad (7.2)$$

Proof. Let us first prove the result in the particular case when $h = 1_B$ with $B \in \sigma(E_2)$ and $k = 1_A$ with $A \in \sigma(E_1)$. Let us denote \overline{A} the complement of A , then $1_A + 1_{\overline{A}} = 1$ and $1_A.1_{\overline{A}} = 0$. Furthermore $i(1_A) = 1_{g^{-1}A}$

$$\pi(1_B) = \pi(1_B.1_{g^{-1}A}) + \pi(1_B.1_{\overline{g^{-1}A}}) \quad (7.3)$$

and

$$\pi(1_B.1_{g^{-1}A}) \leq \pi(1_{g^{-1}A}) = \pi \circ i(1_A) = 1_A \quad (7.4)$$

$$\pi(1_B.1_{\overline{g^{-1}A}}) \leq \pi(1_{\overline{g^{-1}A}}) = 1_{\overline{A}} \quad (7.5)$$

Therefore,

$$\pi(1_B)1_A = \pi(1_B.1_{g^{-1}A})1_A + \pi(1_B.1_{\overline{g^{-1}A}})1_A \quad (7.6)$$

But, $\pi(1_B.1_{g^{-1}A})1_A \leq 1_{\overline{A}}.1_A = 0$ so, $\pi(1_B)1_A = \pi(1_B.1_{g^{-1}A})1_A$. Furthermore, $\pi(1_B.1_{g^{-1}A}) = \pi(1_B.1_{g^{-1}A})1_A + \pi(1_B.1_{g^{-1}A})1_{\overline{A}}$ therefore,

$$\pi(1_B.1_{g^{-1}A})1_{\overline{A}} \leq \pi(i(1_A)).1_{\overline{A}} = 0 \quad (7.7)$$

We just showed that,

$$\pi(1_B.1_{g^{-1}A}) = \pi(1_B.1_{g^{-1}A})1_{\overline{A}} \quad (7.8)$$

So $\pi(1_B).1_A = \pi(1_B.i(1_A))$. The result then extends by linearity directly to $h = \sum_{k \leq n} 1_{B_k}$ and $k = \sum_{k \leq n_1} 1_{A_k}$, which ends the proof. \square

Let us remark that if $A \in \lim \sigma(G)$ then $\overline{A} := (\overline{A}_a, a \in \mathcal{A})$ is also in $\lim \sigma(G)$.

Proposition 1. Let $\gamma = (G, F)$ be a specification, let $G(a)$ be finite sets for any $a \in \mathcal{A}$, let $F > 0$. Let $\mu \in \mathcal{G}(\gamma)$, for any $f \in \prod_{a \in \mathcal{A}} L^\infty(G(a))$, such that $\forall a \in \mathcal{A}, \mu_a(f_a) = 1$,

$$f \cdot \mu \in \mathcal{G}_g(\gamma) \iff \exists \tilde{f} \in \lim i, \text{ s.t. } f \cdot \mu = \tilde{f} \cdot \mu \quad (7.9)$$

Proof. Let us assume that $f \cdot \mu \in \mathcal{G}_g(\gamma)$, then for any $a, b \in \mathcal{A}$ such that $b \leq a$, and for any $g_a \in L^\infty(G(a))$, by hypothesis, $(f \cdot \mu)_b \pi_b^a(g_a) = (f \cdot \mu)_a(g_a)$; it can be rewritten as,

$$\mu_b(\pi_b^a(g_a) \cdot f_b) = \mu_a(f_a \cdot g_a) \quad (7.10)$$

By Lemma 1, $\pi_b^a(g_a) \cdot f_b = \pi_b^a(g_a \cdot i_a^b f_b)$; therefore,

$$\mu_b \pi_b^a(g_a \cdot i_a^b f_b) = \mu_a(f_a \cdot g_a) \quad (7.11)$$

Therefore $f_a = i_a^b f_b$ μ_a -almost surely.

We will now show that there is $\tilde{f} \in \lim i$ such that $\tilde{f} \cdot \mu = f \cdot \mu$. It is in this part of the proof that we assume that $G(a)$ are finite sets and that $F > 0$. Let's call $S_a := \text{supp} \mu_a$, the support of μ_a , i.e., the set $\text{supp} \mu_a := \{\omega_a \in G(a) \mid \mu_a(\omega_a) > 0\}$. Let us denote $N_a := \overline{S_a}$ its complement and $M_a = 1_{N_a}$. We will now show that $(N_a, a \in \mathcal{A}) \in \lim \sigma(G)$.

For any $b, a \in \mathcal{A}$ such that $b \leq a$, $\mu_a i_a^b = \mu_b$; therefore, as $\mu_b(M_b) = 0$, one has that $\mu_a i_a^b(M_b) = 0$. Recall that $i_a^b(M_b)$ is the indicator function of the set $G_b^{a-1} N_b$; the previous remark implies that $i_a^b(M_b) \leq M_a$. Furthermore, $\mu_b \pi_b^a(M_a) = \mu_a(M_a) = 0$; therefore, $\pi_b^a(M_a) \leq M_b$.

Hence, as $\pi_b^a(i_a^b M_b) = M_b$ and $i_a^b M_b \leq M_a$, then by applying π_b^a on both sides, $M_b \leq \pi_b^a(M_a)$. And so $\pi_b^a(M_a) = M_b$.

Recall that we showed that $\pi_b^a(M_a) = M_b$ and $i_a^b(M_b) \leq M_a$. In particular, $M_a - i_a^b(M_b) \geq 0$; furthermore, $\pi_b^a(M_a - i_a^b(M_b)) = 0$ so $M_a = i_a^b(M_b)$. To be more explicit: $\forall \omega_b \in G(b)$,

$$\pi_b^a(M_a - i_a^b(M_b))(\omega_b) = \sum_{\omega_a \in G(a)} F(\omega_a | \omega_b) [M_a - i_a^b(M_b)](\omega_a) \quad (7.12)$$

As, by hypothesis, for any ω_b such that $G_b^a(\omega_a) = \omega_b$ one has that $F(\omega_a | \omega_b) > 0$, then $M_a = i_a^b(M_b)$. This implies that $M \in \lim i$ and $N \in \lim \sigma(G)$. This also implies that $S \in \lim \sigma(G)$.

Let $\tilde{f} = f 1_S$. Then $f 1_S \in \lim F$ and for any $a \in \mathcal{A}$, $f_a = \tilde{f}_a$ μ_a -a.s., which ends the proof. \square

One remarks that $\lim i$ is a subset of $\lim \pi$: for any $f \in \lim i$, by definition for any $a, b \in \mathcal{A}$ such that $b \leq a$, $i_a^b f_b = f_a$ and so $\pi_b^a i_a^b f_b = \pi_b^a f_a$ so $f_b = \pi_b^a f_a$.

Let us also remark that for any $b \leq a \in \mathcal{A}$, $\mu_a(A_a) = \mu_b(A_b)$ for $\mu \in \mathcal{G}(\gamma)$, $A \in \lim G$.

Theorem 4 (Extreme measure characterisation (Generalisation of Theorem 7.7 [29])). *Let $\gamma = (G, F)$ be a specification, let $G(a)$ be finite sets for any $a \in \mathcal{A}$, let $F > 0$. $\mathcal{G}_g(\gamma)$ is a convex set. Each $\mu \in \mathcal{G}_g(\gamma)$ is uniquely determined by its restriction to $\lim \sigma(G)$. Furthermore μ is extreme in $\mathcal{G}(\gamma)$ if and only if for any $A \in \lim \sigma(G)$, $\forall a \in \mathcal{A}$, $\mu_a(A_a) = 0$ or 1 .*

Proof. Let us denote π^* the functor from \mathcal{A} to **Vect** for which for any $b \leq a$, $\pi_a^b : (L^\infty F(b))^* \rightarrow (L^\infty F(a))^*$ is the dual of π_b^a that send linear forms to linear forms. Then $\mathcal{G}_g(\gamma)$ is a subspace of the vector space $\lim F^*$ and furthermore for any $a \in \mathcal{A}$ and $p \in [0, 1]$, $p\mu_a + (1-p)\nu_a \in \mathbb{P}(G(a))$ whenever $\mu_a, \nu_a \in \mathbb{P}(G(a))$. Therefore $\mathcal{G}_g(\gamma)$ is a convex set.

Proposition 1, allows us to apply a similar proof, when done with caution, to the one found of Theorem 7.7 in [29]. Let us recall the proof. Let $\mu, \nu \in \mathbb{G}(\gamma)$ such that $\mu|_{\lim i} = \nu|_{\lim i}$. Let $\bar{\mu} = \frac{\mu + \nu}{2}$, then $\bar{\mu} \in \mathcal{G}(\gamma)$. But μ and ν are absolutely continuous with respect to $\bar{\mu}$ therefore for any $a \in \mathcal{A}$ there is $f_a, g_a \in L^\infty(G(a))$ such that $\mu_a = f_a \bar{\mu}_a$ and $\nu_a = g_a \bar{\mu}_a$. By Proposition 1, $f, g \in \lim i$. By hypothesis, for any $h \in \lim i$ $\bar{\mu}_a(h_a) = \mu_a(h_a) = \nu_a(h_a)$. Importantly i_a^b is a ring morphism of $L^\infty(G(b))$, i.e. $i_a^b(k_b \cdot h_b) = i_a^b(k_b) \cdot i_a^b(h_b)$. Therefore for any $h, k \in \lim i$, $k \cdot h \in \lim i$; as $f - g \in \lim i$, then it is also the case that $(f - g)^2 \in \lim i$; but for any $a \in \mathcal{A}$,

$$\bar{\mu}_a[(f_a - g_a)^2] = 0 \quad (7.13)$$

so $f_a = g_a$ $\bar{\mu}_a$ - a.s. Therefore $f\bar{\mu} = g\bar{\mu}$ and $\mu = \nu$.

Showing that $\mu \in \mathcal{G}_g(\gamma)$ is extreme is equivalent to μ being trivial on $\lim i$ is a direct generalization of Corollary 7.4 [29] thanks to Proposition 1. Let $\mu \in \mathcal{G}(\gamma)$ be not trivial on $\lim \sigma(G)$ then there is $A = (A_a, a \in \mathcal{A}) \in \lim \sigma(G)$ such that,

$$\forall a \in \mathcal{A}, \quad 0 < \mu_a(A_a) < 1 \quad (7.14)$$

Therefore for any $a \in \mathcal{A}$,

$$\mu_a = \mu_a(A_a) \frac{1_{A_a} \cdot \mu_a}{\mu_a(A_a)} + \mu_a(\bar{A}_a) \frac{1_{\bar{A}_a} \cdot \mu_a}{\mu_a(\bar{A}_a)} \quad (7.15)$$

Furthermore, for any $b \leq a$,

$$i_a^b \frac{1_{A_b}}{\mu_a(A_b)} = \frac{1_{A_a}}{\mu_a(A_b)} = \frac{1_{A_a}}{\mu_a(A_a)} \quad (7.16)$$

Therefore $\frac{1_{A_a}}{\mu_a(A_a)}, a \in \mathcal{A}$ is in $\lim i$ and so by Lemma 1, $\left(\frac{1_{A_a} \cdot \mu_a}{\mu_a(A_a)}, a \in \mathcal{A}\right) \in \mathcal{G}_g(\gamma)$. Similarly $\left(\frac{1_{\bar{A}_a} \cdot \mu_a}{\mu_a(\bar{A}_a)}, a \in \mathcal{A}\right) \in \mathcal{G}_g(\gamma)$. In particular there is $0 < p < 1$

so that $\mu = p\nu + (1-p)\nu_1$ with $\nu, \nu_1 \in \mathcal{G}_g(\gamma)$. Therefore μ is not an extreme measure.

Assume now that $\mu \in \mathcal{G}_g(\gamma)$ is such that for any $A \in \lim \sigma(A)$, and any $a \in \mathcal{A}$, $\mu_a(A_a) = 0$ or 1 . Suppose that there is $0 < p < 1$ such that $\mu = p\nu + (1-p)\nu_1$ with $\nu, \nu_1 \in \mathcal{G}_g(\gamma)$. Then for any $a \in \mathcal{A}$, ν_a, ν_{1a} is absolutely continuous with respect to μ_a . Therefore, there are $(f_a \geq 0, a \in \mathcal{A}), (g_a \geq 0, a \in \mathcal{A})$ both in $\prod_{a \in \mathcal{A}} L^\infty(G(a))$ such that $\nu = f\mu$ and $\nu_1 = g\mu$. As $\nu, \nu_1 \in \mathcal{G}_g(\gamma)$, then by Lemma 1, $f, g \in \lim i$. Therefore, for all $a \in \mathcal{A}$, $\mu_a(f_a) = 0$ or for all $a \in \mathcal{A}$, $\mu_a(g_a) = 0$. So, $\mu = \nu$ or $\mu = \nu_1$ and μ is extreme in $\mathcal{G}_g(\gamma)$. □

Let us remark that if \mathcal{A} has only one connected component, then for $A \in \lim \sigma(G)$, satisfying $\forall a \in \mathcal{A}, \mu_a(A_a) = 0$ or 1 is equivalent to $\exists a \in \mathcal{A}, \mu_a(A_a) = 0$ or 1 . Indeed, if a, b are in the same connected component, i.e., $a \leq b$ or $b \leq a$, then $\mu_a(A_a) = \mu_b(A_b)$.

8 Background: Variational inference for graphical model and more

Consider a joint distribution $P_{X,Y} \in \mathbb{P}(E \times E_1)$ over two random variables $X \in E, Y \in E_1$. A classical problem is, given an observation ω_1 on Y , to compute the posterior $P_{X|Y}(\omega, \omega_1) = \frac{P_{X,Y}(\omega, \omega_1)}{P_Y(\omega_1)}$ with $P_Y(\omega_1) = \sum_{\omega \in E} P_{X|Y}(\omega, \omega_1)$, the marginal distribution of Y . However, doing so requires summing over all possible configurations of X , which can be computationally too costly. This is the case, for example, when $X = X_0, \dots, X_T$ with $X_i \in S$ and $E = \prod_{i \leq T} B$. Instead, one resorts to variational inference to compute $P_{X|Y}$ approximately [30]. We will now explain what variational inference is, but first let us introduce entropy and Gibbs free energy. When E is a finite set, the entropy of a probability distribution Q on E is defined as:

$$S(Q) = - \sum_{\omega \in E} Q(\omega) \ln Q(\omega) \quad (8.1)$$

Let H be a measurable function $H : E \rightarrow \mathbb{R}$. For $Q \in \mathbb{P}(E)$, one calls $\mathbb{E}_Q[H] - \frac{1}{\beta} S(Q)$ the Gibbs free energy; in general $\beta = 1$. An important property is that,

$$- \ln \sum_{\omega \in E} e^{-\beta H(\omega)} = \inf_{Q \in \mathbb{P}(E)} \mathbb{E}_Q[\beta H] - S(Q) \quad (8.2)$$

The optimal solution to Equation 8.2 is given by the Boltzmann distri-

bution

$$Q^*(\omega) = \frac{e^{-H(\omega)}}{\sum_{\omega \in E} e^{-\beta H(\omega)}} \quad (8.3)$$

Let $H(\omega) = -\ln P_{X,Y}(\omega, \omega_1)$ and $\beta = 1$, then $Q^*(\omega) = P_{X|Y}(\omega|\omega_1)$. Therefore, solving the optimization problem of Equation 8.2 is equivalent to computing the posterior $P_{X|Y}(\omega|\omega_1)$. Solving Equation 8.2 over a subset of distributions $Q \in \Theta \subseteq \mathbb{P}(E)$ is called variational inference. If furthermore the Gibbs free energy is replaced by an approximation, we call it approximate variational inference.

One remarks that $\inf_{Q \in \mathbb{P}(E)} \mathbb{E}_Q[\beta H] - S(Q)$ is equivalent to $\sup_{Q \in \mathbb{P}(E)} S(Q) - \mathbb{E}_Q[\beta H]$. And this last optimization problem relates, through Lagrange multipliers, to maximizing entropy under energy constraints $U \in \mathbb{R}$,

$$\sup_{\substack{Q \in \mathbb{P}(E) \\ \mathbb{E}_Q[H]=U}} - \sum_{x \in E} Q(x) \ln Q(x) \quad (8.4)$$

In the physics literature, one refers to Equation 8.4 as MaxEnt [42], which stands for the principle of maximum entropy and such principle has many application see [42, 43]. Variational inference is called the variational principle.

Graphical models translate relations on graphs into conditional independence relations between variables [2, 44, 45]. The Hammersley-Clifford Theorem states that for strictly positive distributions, these conditional independence relations are equivalent to factorizing the joint distributions on the cliques of the graph. When a joint distribution factors according to an acyclic graph, computing the posterior $P_{X|Y}$ over some nodes of the graph can be done efficiently through dynamic programming with an algorithm called Belief Propagation. In the particular case where the graphical model is a Gaussian Hidden Markov Model, this dynamic programming algorithm is the smoother Kalman filter [46].

The Belief Propagation algorithm solves an approximate variational inference problem for a variational free energy called the Bethe free energy [10, 11, 18, 47].

We now specify the previous statement and explain what we mean by a probability distribution that factors according to a collection of subsets of variables. Then, we recall the remarkable property that for probability distributions that factor according to an acyclic graph, entropy can be decomposed into sums of 'local' entropies. This will allow us to introduce what the Bethe free energy is and how it relates to entropy.

It is important to consider the case of graphical models to understand the more general setting of *factor graphs* or *factorization models* Yedidia, Freeman, Weiss consider in their seminal article, *Constructing Free Energy Approximations and Generalized Belief Propagation Algorithms*, extending the correspondence between a Generalized Belief Propagation algorithm and

the associated variational free energy that we will call the *Generalized Bethe free energy*. We present their work and O. Peltre's complemented version and extension of their results (see SYCO 8 talk, [11,48]) in the next section. Doing so will then motivate the entropy we introduce for \mathcal{A} -specifications and give intuition on the message passing algorithm we introduce for such specifications.

Let $E = \prod_{i \in I} E_i$ where I is finite. \mathbb{P} is said to factor according to a subset \mathcal{A} of $\mathcal{P}(I)$ when for any $a \in \mathcal{A}$ there is $f_a : E_a \rightarrow \mathbb{R}$ such that,

$$P(\omega) = \prod_{a \in \mathcal{A}} f_a(\omega_a) \quad (8.5)$$

We will denote the space of probability distributions that factor according to \mathcal{A} as $\text{Fac}_{\mathcal{A}}$. We will call such space the space of (\mathcal{A} -)factorization models.

A graph $G = (I, A)$ is a collection of vertices I and edges A ; one can see it as a poset $\mathcal{A}(G)$ with the relation of inclusion. In other words, for any $v \in I$ and $e \in A$, $v \leq e$ if and only if $v \in e$.

Let us now recall a proposition that states that a strictly positive probability distribution on a finite set E , i.e., $Q(\omega) > 0$ for $\omega \in E$, factors according to its marginal distribution on the edges and vertices [12,49]. We denote $\mathbb{P}_{>0}$ as the space of strictly positive probability distributions.

Proposition 2 (Factorization on acyclic graphs). *Let I be a finite set and let $E = \prod_{i \in I} E_i$ be a product of finite sets, and let $G = (I, A)$ be a finite acyclic graph. $Q \in \mathbb{P}_{>0}(E)$ factors accordingly to $\mathcal{A}(G)$, i.e., $Q \in G_{\mathcal{A}(G)}$ if and only if for any $\omega \in E$,*

$$Q(\omega) = \frac{\prod_{e \in A} i_e^* Q(\omega_e)}{\prod_{v \in I} i_v^* Q^{d(v)-1}(\omega_v)}, \quad (8.6)$$

where $i^e : E \rightarrow E_e$ is the projection onto the state space $E_{v_1} \times E_{v_2}$ associated with the edge $e = (v_1, v_2)$, and similarly for $i^v : E \rightarrow E_v$.

For a given sub-poset $\mathcal{A} \subseteq \mathcal{P}(I)$, let us denote

$$[* , i]_{\mathcal{A}} := \{(Q_a \in \mathcal{P}(E_a), a \in \mathcal{A}) \mid \forall b \leq a, i_a^b Q_b = Q_a\} \quad (8.7)$$

as the collection of 'local' probability measures that have compatible marginals with respect to the projections $i_b^a : E_a \rightarrow E_b$ when $b \subseteq a$ in $\mathcal{P}(I)$. The next proposition expresses the free energy of probability distributions that factor according to an acyclic graph as a weighted sum of free energies of its marginals on edges and vertices. The weights are given by the inclusion-exclusion formula of the graph when seen as a poset. First, let us recall the inclusion-exclusion formula for a poset.

Definition 9 (Zeta operator of a poset). Let \mathcal{A} be a finite poset. We call the ‘zeta-operator’ of a poset \mathcal{A} , denoted ζ , the operator from $\bigoplus_{a \in \mathcal{A}} \mathbb{R} \rightarrow \bigoplus_{a \in \mathcal{A}} \mathbb{R}$ defined as, for any $\lambda \in \bigoplus_{a \in \mathcal{A}} \mathbb{R}$ and any $a \in \mathcal{A}$,

$$\zeta(\lambda)(a) = \sum_{b \leq a} \lambda_b \quad (8.8)$$

Proposition 3 (Reformulation of Proposition 2 [50], Rota’ 64). *Let \mathcal{A} be a finite poset. The zeta-operator of \mathcal{A} is invertible. We will call its inverse the Möbius inversion of \mathcal{A} , denoted μ . Furthermore, there is a collection $(\mu(a, b); b, a \in \mathcal{A} \text{ s.t. } b \leq a)$ such that, for any $\lambda \in \bigoplus_{a \in \mathcal{A}} \mathbb{R}$ and $a \in \mathcal{A}$,*

$$\mu(\lambda)(a) = \sum_{b \leq a} \mu(a, b) \lambda_b \quad (8.9)$$

We call the coefficient $(\mu(a, b), b, a \text{ s.t. } b \leq a)$ the Möbius coefficients of \mathcal{A} . In particular Proposition 3 implies that, for any $b, a \in \mathcal{A}$ such that $b \leq a$,

$$\sum_{c: b \leq c \leq a} \mu(a, c) = 1[b = a] \quad (8.10)$$

$$\sum_{c: b \leq c \leq a} \mu(c, b) = 1[b = a] \quad (8.11)$$

For a collection of values $\lambda_a \in \mathbb{R}$, $a \in \mathcal{A}$, we call the following expression $\sum_{a \in \mathcal{A}} \sum_{b \leq a} \mu(a, b) \lambda_b$ the inclusion-exclusion formula over a poset \mathcal{A} ; this formula corresponds to the value one would attribute to a maximal element, denoted 1, added to \mathcal{A} :

$$\lambda_1 := \sum_{a \in \mathcal{A}} \sum_{b \leq a} \mu(a, b) \lambda_b \quad (8.12)$$

which can be rewritten as $\lambda_1 = \sum_{a \in \mathcal{A}} [\sum_{b \geq a} \mu(b, a)] \lambda_a$. We will denote $c(a) = \sum_{b \leq a} \mu(a, b)$ the weighted coefficients

To find the classical inclusion-exclusion formula, consider I a finite set; the poset \mathcal{A} is $(\mathcal{P}(I), \supseteq)$ with the reversed order. The quantities λ_a , for $a \in \mathcal{A}$, are $|A_i|$ when $a = i \in I$, and represent the cardinality of the sets A_i , as well as the cardinality of all possible intersections $|\cap_{i \in a} A_i|$ when $a \subseteq I$. In this setting, the maximal element 1 has a value $\lambda_1 = |\cup_{i \in I} A_i|$.

Proposition 4. *Let I be a finite set, and let $E = \prod_{i \in I} E_i$ be a product of finite sets. Consider a finite acyclic graph $G = (I, A)$. Let $(H_a : E_a \rightarrow \mathbb{R}, a \in \mathcal{A})$ be a collection of Hamiltonians (measurable maps) that respectively factor through the projection $i^a : E \rightarrow E_a$. The following map,*

$$\begin{aligned} \phi : \quad & [*, i]_{\mathcal{A}} \quad \rightarrow \quad \text{Fac}_{\mathcal{A}(G)} \\ & (Q_a, a \in \mathcal{A}(G)) \quad \mapsto \quad \frac{\text{Fac}_{\mathcal{A}(G)}}{\prod_{e \in A} Q_e} \\ & \quad \quad \quad \quad \quad \mapsto \quad \frac{\text{Fac}_{\mathcal{A}(G)}}{\prod_{v \in I} Q_v^{d(v)-1}} \end{aligned} \quad (8.13)$$

is a bijection; here $d(v)$ is the degree of node v . Furthermore,

$$\mathbb{E}_{\phi(Q)} \left[\sum_{e \in A} H_e + \sum_{v \in I} H_v \right] - S(\phi(Q)) = F(Q_a, a \in \mathcal{A}(G)) \quad (8.14)$$

where

$$F(Q_a, a \in \mathcal{A}(G)) = \sum_{a \in \mathcal{A}(G)} c(a) [S(i_*^a Q) - \mathbb{E}_{i_*^a Q}[H_a]] \quad (8.15)$$

Proof. Let for any $p \in \mathbb{P}_{>0}(E) \cap \text{Fac}_{\mathcal{A}(G)}$,

$$\begin{aligned} \psi : \quad & \mathbb{P}(E) \cap \text{Fac}_{\mathcal{A}(G)} \quad \rightarrow \quad \lim F \\ & P \quad \mapsto \quad (\pi_{a,*} P, a \in \mathcal{A}(G)) \end{aligned} \quad (8.16)$$

Then $\phi\psi = \text{id}$ and $\psi\phi = \text{id}$, furthermore for any $P \in \mathbb{P}(E) \cap \text{Fac}_{\mathcal{A}(G)}$,

$$\mathbb{E}_Q \left[\sum_{a \in \mathcal{A}} H_a \right] = \mathbb{E}_Q \left[\sum_{a \in \mathcal{A}} \sum_{b \leq a} \mu(a, b) H_a \right] = \sum_{a \in \mathcal{A}} \sum_{b \leq a} \mu(a, b) \mathbb{E}_{i_*^a Q} [H_a] \quad (8.17)$$

and by Proposition 2,

$$S(Q) = \sum_{e \in A} S(i_*^e Q) - \sum_{v \in I} (d(v) - 1) S(i_*^v Q) = \sum_{a \in \mathcal{A}(G)} c(a) S(i_*^a Q) \quad (8.18)$$

which ends the proof. \square

$F(Q_a, a \in \mathcal{A}(G)) = \sum_{a \in \mathcal{A}(G)} c(a) [S(i_*^a Q) - \mathbb{E}_{i_*^a Q}[H_a]]$ is called the Bethe free energy and depends only on the marginal distributions $i_*^a Q$. Proposition 4 implies that variational inference on acyclic graphical models can be done over marginal distributions and still give the exact posterior. This remark is very important as, for a graph (G, E) , the state space of $(E_e, e \in A)$ is of size $|A| \times N^2$ with N the maximal size of the $E_v, v \in I$, whereas the cardinality of E is $O(N^{|I|})$. The collection of marginal distributions serves therefore as a compressed representation of the joint distribution. When G is acyclic, the Belief Propagation algorithm finds the minima of the Bethe free energy exploiting the compressed reformulation of $\text{Fac}_{\mathcal{A}(G)}$; the complexity of the algorithm is $O(|A|)$.

The Bethe free energy can be defined on any graph, even those that may contain cycles; in such cases, the equality between the Bethe free energy and the Gibbs free energy does not hold in general. The key idea is to replace the Gibbs free energy with the Bethe free energy in variational inference and solve the associated optimization problem, hoping that the true posterior will be correctly approximated.

We introduced the Bethe free energy in a similar fashion to [11]. This also allows us to see how one can define the Bethe free energy for any factorization model. In this context, the Bethe free energy is called the Generalized Bethe free Energy or region-based approximation of free energy in [12]. In the next section, we define the (Generalized) Belief Propagation algorithm for the Generalized Bethe free Energy.

9 Background: Generalized Bethe free energy

The Generalized Bethe free energy is an approximation of the Gibbs free energy that generalizes the Bethe free energy [12]. Consider the configuration space $E = \prod_{i \in I} E_i$ of finite sets over a finite set I . Let $\mathcal{A} \subseteq \mathcal{P}(I)$, and let $(H_a \in E_a \rightarrow \mathbb{R}, a \in \mathcal{A})$ be a collection of (measurable) Hamiltonians. For $Q = (Q_a \in \mathbb{P}(E_a), a \in \mathcal{A})$, the General Bethe free energy is defined as:

$$F_{\text{Bethe}}(Q) = \sum_{a \in \mathcal{A}} \sum_{b \geq a} \mu(b, a) (\mathbb{E}_{Q_a}[H_a] - S_a(Q_a)) \quad (9.1)$$

Here we considered $Q = (Q_a \in \mathbb{P}(E_a), a \in \mathcal{A})$ to be any collection of distributions over the local variables $X_a \in E_a$; however, let's keep in mind that the collection we are interested in are $Q \in [*, i]_{\mathcal{A}}$ that are compatible with respect to marginalization.

In the Bethe Free energy, the entropy functional is replaced by a reconstruction of the 'local' entropies that only makes use of the entropy of the marginals. This is the term $S_{GB}(p) = \sum_{a \in \mathcal{A}} c(a) S(p_a)$ in F_{Bethe} . It is economical in the sense that it is a non-redundant way of computing entropy, similar to counting the cardinality of the union of sets using the Inclusion-Exclusion principle (see [51] for a detailed presentation of this idea).

The General Belief Propagation is an algorithm that enables us to find the critical points of the Generalized Bethe Free Energy. A classical result states that fixed points of this algorithm correspond to critical points of that free energy. Let us now recall the expression of this algorithm.

For $\mathcal{A} \subseteq \mathcal{P}(I)$, E a finite product of finite sets, and $(H_a, a \in \mathcal{A})$ a collection of Hamiltonians. For two elements of \mathcal{A} , a and b such that $b \subseteq a$, two types of messages are considered: top-down messages $m_{a \rightarrow b} \in \mathbb{R}^{E_b}$ and bottom-up messages $n_{b \rightarrow a} \in \mathbb{R}^{E_a}$. The update rule is as follows: consider

messages at times t , $(n_{b \rightarrow a}^t, m_{a \rightarrow b}^t | b, a \in \mathcal{A} \text{ s. t. } b \leq a)$, they are related by the following relation,

$$\forall a, b \in \mathcal{A}, \text{ s.t. } b \leq a, \quad n_{b \rightarrow a}^t = \prod_{\substack{c: b \leq c \\ c \not\leq a}} m_{c \rightarrow b}^t \quad (9.2)$$

One can define beliefs as ,

$$\forall a \in \mathcal{A}, \quad b_a^t \propto e^{-H_a} \prod_{\substack{b \in \mathcal{A}: \\ b \leq a}} n_{b \rightarrow a}^t \quad (9.3)$$

The beliefs are sometimes interpreted as probability distributions up to a multiplicative constant; here to make the presentation clearer, we require that b_a is a probability distribution. The update rule is given by,

$$\forall a, b \in \mathcal{A}, \text{ s.t. } b \leq a \quad m_{a \rightarrow b}^{t+1}(x_b) = m_{a \rightarrow b}^t(x_b) \frac{\sum_{y_a: \pi_b^a(y_a) = x_b} b_a^t(y_a)}{b_b^t(x_b)} \quad (9.4)$$

The multiplication of function $n_{b \rightarrow a}$ that have different domains is made possible because there is an the embedding of \mathbb{R}^{E_b} into \mathbb{R}^{E_a} implicitly implied in the last equation.

The algorithm can be rewritten in a more condensed manner, updating only the top-down messages, for all $a, b \in \mathcal{A}$, such that $b \leq a$,

$$m_{a \rightarrow b}^{t+1}(x_b) = m_{a \rightarrow b}^t(x_b) \frac{\sum_{y_a: i_b^a(y_a) = x_b} e^{-H_a(y_a)} \prod_{\substack{c \in \mathcal{A}: \\ c \subseteq a}} \prod_{\substack{d: c \subseteq d \\ d \not\subseteq a}} m_{d \rightarrow c}^t(x_c)}{e^{-H_b(x_b)} \prod_{\substack{c \in \mathcal{A}: \\ c \subseteq b}} \prod_{\substack{d: c \subseteq d \\ d \not\subseteq b}} m_{d \rightarrow c}^t(x_c)} \quad (9.5)$$

We will denote this update rule as GBP , $m(t+1) = GBP(m(t))$

Theorem 5 (Yedidia, Freeman, Weiss, Peltre). *Let $(m_{a \rightarrow b}, a, b \in \mathcal{A} \text{ s.t. } b \subseteq a)$ be a fix point of the Generalized Belief Propagation, i.e.*

$$m = GBP(m) \quad (9.6)$$

Let $(b_a, a \in \mathcal{A})$ be the associated beliefs and let, for $a \in \mathcal{A}$, $p_a = b_a / \sum_{x \in E_a} b_a(x)$ be the associated normalized beliefs. Then $(p_a, a \in \mathcal{A})$ is a critical point of F_{Bethe} under the constraint that $p \in \lim F$.

Proof. Theorem 5.15 [11], Theorem 5 [12].

□

10 Entropy of \mathcal{A} -specifications and variational principle

Let \mathcal{A} be a finite poset and $\gamma = (G, F)$ be a specifications with $G(a)$ being a finite set for any $a \in \mathcal{A}$. We propose that the entropy of $Q \in \mathcal{G}_g(\gamma)$ to be $S_{GB}(Q) = \sum_{a \in \mathcal{A}} c(a)S(Q_a)$ and that the free energy is $F_{\text{Bethe}}(Q) = \sum_{a \in \mathcal{A}} c(a) (\mathbb{E}_{Q_a}[H_a] - S_a(Q_a))$.

It might seem at first glance that it is sufficient to apply GBP to find the critical points of the free energy for specifications, but it is actually much trickier than that. In GBP the free energy is constrained over $Q \in [* , i]_{\mathcal{A}}$ but for $Q \in \mathcal{G}_g(\gamma)$, $Q \in [* , F]_{K, \mathcal{A}}$. The first significant difference is that for GBP the presheaf is prescribed to the one associated with marginalization whereas here F can be anything. In Section 3.2 [10] (version 1) the case where F is any presheaf is treated in detail and corresponds to a fair amount of generalization with respect to GBP as presented in [11, 12]. The second difference is that F is a functor and not a presheaf; it is an essential difference tackled in the most general setting by Theorem 2.3 [10] (version 1). The intuition behind such a difference is as follows. The dual of the F acts on Lagrange multipliers; the Lagrange multipliers are analogous to the messages ‘ $m_{a \rightarrow b}$ ’ of GBP. When F is a presheaf its dual is a functor and one can send the Lagrange multipliers in a cell $F(a)$ but applying F_b^{a*} to $m_{a \rightarrow b}$ for any $b \leq a$ and so one can do the product of such messages (or the sum if one takes the logarithm). One cannot send the Lagrange multipliers for the b ’s smaller than a into the cell $F(a)^*$ when the F is a functor. However, when there is G such that F is a section of G , G can send the multipliers into $F(a)$ and one can build a message-passing algorithm that finds the critical points of the Bethe free energy for specifications. We now explicitly state Theorem 2.3 [10] (version 1) in the particular case of specifications and detail the message-passing algorithms for finding optimal Gibbs measures.

Problem to solve: The optimization problem we want to solve is the following,

$$\inf_{Q \in \mathcal{G}_g(\gamma)} F_{\text{Bethe}}(Q) \quad (10.1)$$

We will need to generalize the ζ and Möbius inversion of a poset to the ones for functors and presheaves.

Definition 10 (Möbius inversion associated to a functor). Let $G : \mathcal{A} \rightarrow \mathbf{Mod}$ be a functor from a finite poset to the category of R -modules with R a ring; let $\zeta_G : \bigoplus_{a \in \mathcal{A}} G(a) \rightarrow \bigoplus_{a \in \mathcal{A}} G(a)$ be such that for any $a \in \mathcal{A}$ and $v \in \bigoplus_{a \in \mathcal{A}} G(a)$,

$$\mu_G(v)(a) = \sum_{b \leq a} \mu(a, b) G_a^b(v_b) \quad (10.2)$$

Proposition 5. Let $G : \mathcal{A} \rightarrow \mathbf{Mod}$ be a functor from a finite poset to the category of modules, μ_G is invertible and its inverse, denoted ζ_G , is defined as follows, for any $a \in \mathcal{A}$ and $v \in \bigoplus_{a \in \mathcal{A}} G(a)$,

$$\zeta_G(v)(a) = \sum_{b \leq a} G_a^b(v_b) \quad (10.3)$$

Proof. Let $v \in \bigoplus_{a \in \mathcal{A}} G(a)$ and $a \in \mathcal{A}$,

$$\zeta_G \mu_G(v)(a) = \sum_{b \leq a} \sum_{c \leq b} \mu(b, c) G_a^b G_b^c(v_c) \quad (10.4)$$

therefore,

$$\zeta_G \mu_G(v)(a) = \sum_{c \leq a} \left(\sum_{b: c \leq b \leq a} \mu(b, c) \right) G_a^c(v_c) = G_a^a(v_a) \quad (10.5)$$

Furthermore,

$$\mu_G \zeta_G(v)(a) = \sum_{b \leq a} \mu(a, b) \sum_{c \leq b} G_a^c(v_c) = v_a \quad (10.6)$$

□

Remark 3. Let us remark for any poset (\mathcal{A}, \leq) one can reverse the relations, in other words, for any $a, b \in \mathcal{A}$,

$$a \leq_{op} b \iff b \leq a \quad (10.7)$$

We shall also denote \leq_{op} as \geq and the corresponding poset as \mathcal{A}^{op} or (\mathcal{A}, \geq) . One has that, for any $a, b \in \mathcal{A}$ such that $a \geq b$,

$$\zeta_{\mathcal{A}^{op}}(b, a) = \zeta_{\mathcal{A}}(a, b) \quad (10.8)$$

$$\mu_{\mathcal{A}^{op}}(b, a) = \mu_{\mathcal{A}}(a, b) \quad (10.9)$$

In particular for any $G : \mathcal{A} \rightarrow \mathbf{Vect}$ functor from a finite poset to the category of modules,

$$\mu_{G^*} = (\mu_G)^* \quad (10.10)$$

as for any $(l_a \in G(a)^*, a \in \mathcal{A})$,

$$\sum_{a \in \mathcal{A}} \sum_{b \leq a} \mu(a, b) l_a G_a^b = \sum_{b \in \mathcal{A}} \sum_{a \geq b} \mu(a, b) G_b^a(l_a) \quad (10.11)$$

In what follows, we go back to the convention for specification, which is that G is a presheaf (to the category of measurable spaces). Let us define the function $FE(Q) : \prod_{a \in \mathcal{A}} \mathbb{P}(E_a) \rightarrow \prod_{a \in \mathcal{A}} \mathbb{R}$ as $FE(Q) = (\mathbb{E}_{Q_a}[H_a] - S_a(Q_a), a \in \mathcal{A})$, which sends a collection of probability measures over \mathcal{A} to their Gibbs free energies. For any $Q \in \prod_{a \in \mathcal{A}} \mathbb{P}(E_a)$, let us denote $d_Q FE$ as the differential of FE at the point Q .

Theorem 6. *Let \mathcal{A} be a finite poset, let $\gamma = (G, F)$ be a \mathcal{A} -specification such that $G(a)$ is a finite set for any $a \in \mathcal{A}$. Let $H_a : G(a) \rightarrow \mathbb{R}$ be a collection of (measurable) Hamiltonians. The critical points of the Generalized Bethe free energy are the $Q \in [*, F]_{K, \mathcal{A}}$ such that,*

$$\mu_{G^*} d_Q FE|_{[*, F]_{K, \mathcal{A}}} = 0 \quad (10.12)$$

Let us now present the message-passing algorithm we consider. For two elements of \mathcal{A} , a and b such that $b \leq a$, two types of messages are considered: top-down messages $m_{a \rightarrow b} \in \mathbb{R}^{G(a)}$ and bottom-up messages $n_{b \rightarrow a} \in \mathbb{R}^{G(b)}$. Consider messages at times t , $(n_{b \rightarrow a}^t, m_{a \rightarrow b}^t; b, a \in \mathcal{A} \text{ s.t. } b \leq a)$, they are related by the following relation,

$$\forall b \leq a, \forall \omega_1 \in G(a) \quad n_{b \rightarrow a}^t(\omega_1) = \prod_{\substack{c: b \leq c \\ c \not\leq a}} \sum_{\omega \in G(b)} m_{c \rightarrow b}^t(\omega) G_b^a(\omega | \omega_1) \quad (10.13)$$

One then defines beliefs as ,

$$\forall a \in \mathcal{A}, \quad b_a^t \propto e^{-H_a} \prod_{\substack{b \in \mathcal{A}: \\ b \leq a}} n_{b \rightarrow a}^t \quad (10.14)$$

where $b_a \in \mathbb{P}(G(a))$. The update rule is,

$$\forall a, b \in \mathcal{A}, \text{ s.t. } b \leq a \quad m_{a \rightarrow b}^{t+1} = m_{a \rightarrow b}^t \frac{F_a^b b_b^t}{b_a^t} \quad (10.15)$$

Theorem 2.2 [10] (version 2) applied to specifications implies that fixed points of the previous message-passing algorithm are in correspondence with critical points of the Bethe free energy over the space of Gibbs measures.

11 Acknowledgement

We would like to thank the reviewers for their comments.

References

- [1] G. Sergeant-Perthuis, “Compositional statistical mechanics, entropy and variational inference,” in *Twelfth Symposium on Compositional Structures*, 2024.
- [2] G. Sergeant-Perthuis, “Bayesian/graphoid intersection property for factorisation spaces,” 2021. <https://arxiv.org/abs/1903.06026>.
- [3] G. Sergeant-Perthuis, “Intersection property and interaction decomposition,” Apr. 2019. <https://arxiv.org/abs/1904.09017>.
- [4] G. Sergeant-Perthuis, “Interaction decomposition for presheaves,” Aug. 2020. <https://arxiv.org/abs/2008.09029>.
- [5] D. Bennequin, O. Peltre, G. Sergeant-Perthuis, and J. P. Vigneaux, “Extra-fine sheaves and interaction decompositions,” Sept. 2020. <https://arxiv.org/abs/2009.12646>.
- [6] G. Sergeant-Perthuis, “Interaction decomposition for Hilbert spaces,” 2021. <https://arxiv.org/abs/2107.06444>.
- [7] G. Sergeant-Perthuis, *Intersection property, interaction decomposition, regionalized optimization and applications*. PhD thesis, Université de Paris, 2021. Link.
- [8] G. Sergeant-Perthuis, “A categorical approach to statistical mechanics,” in *Geometric Science of Information* (F. Nielsen and F. Barbaresco, eds.), (Cham), pp. 258–267, Springer Nature Switzerland, 2023.
- [9] G. Sergeant-Perthuis, “Characterization of extreme Gibbs measures for a Categorical Approach to Statistical Mechanics,” Feb. 2024. <https://hal.sorbonne-universite.fr/hal-04456412>.
- [10] G. Sergeant-Perthuis, “Regionalized optimization,” 2022. <https://arxiv.org/abs/2201.11876>.
- [11] O. Peltre, “Message passing algorithms and homology,” 2020. Ph.D. thesis, Link.
- [12] J. Yedidia, W. Freeman, and Y. Weiss, “Constructing Free-Energy Approximations and Generalized Belief Propagation Algorithms,” *IEEE Transactions on Information Theory*, vol. 51, pp. 2282–2312, July 2005.
- [13] D. Ruelle, *Statistical Mechanics*. Imperial College Press, 1999.
- [14] Y. Lecun, S. Chopra, R. Hadsell, M. Ranzato, and F. Huang, *A tutorial on energy-based learning*. MIT Press, 2006.

- [15] T. Lelièvre, M. Rousset, and G. Stoltz, *Free Energy Computations*. Imperial College Press, 2010.
- [16] C. Chipot and A. Pohorille, eds., *Free energy calculations: theory and applications in chemistry and biology*. No. 86 in Springer series in chemical physics, Berlin ; New York: Springer, 2007. OCLC: ocm79447449.
- [17] M. Wiering and M. van Otterlo, eds., *Reinforcement Learning: State-of-the-Art*, vol. 12 of *Adaptation, Learning, and Optimization*. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012.
- [18] M. Mezard and A. Montanari, *Information, Physics, and Computation*. USA: Oxford University Press, Inc., 2009.
- [19] M. Gromov, “In a search for a structure, part 1: On entropy,,” 2013.
- [20] T. Fritz, “A synthetic approach to markov kernels, conditional independence and theorems on sufficient statistics,” *Advances in Mathematics*, 2020.
- [21] J. P. Vigneaux, “Information structures and their cohomology,” *Theory and Applications of Categories*, vol. 35, no. 38, pp. 1476–1529, 2020.
- [22] T. Fritz and D. I. Spivak, “Internal probability valuations.” Workshop Categorical Probability and Statistics, 2020.
- [23] J.-C. Belfiore and D. Bennequin, “Topos and stacks of deep neural networks,” *ArXiv*, 2021.
- [24] T. Fritz and P. Perrone, “A probability monad as the colimit of spaces of finite samples,” *Theory and Applications of Categories*, 2019.
- [25] M. Marcolli, “Gamma spaces and information,” *Journal of Geometry and Physics*, 2019.
- [26] P. Baudot and D. Bennequin, “The homological nature of entropy,” *Entropy*, vol. 17, no. 5, pp. 3253–3318, 2015.
- [27] T. Fritz, T. Gonda, and P. Perrone, “De finetti’s theorem in categorical probability,” *arXiv preprint arXiv:2105.02639*, 2021.
- [28] S. Moss and P. Perrone, “A category-theoretic proof of the ergodic decomposition theorem,” *Ergodic Theory and Dynamical Systems*, vol. 43, no. 12, p. 4166–4192, 2023.
- [29] H.-O. Georgii, *Gibbs Measures and Phase Transitions*. Berlin, New York: De Gruyter, 2011.
- [30] P. Alquier, “Approximate Bayesian Inference,” *Entropy*, vol. 22, p. 1272, Nov. 2020.

- [31] L. Onsager, “Crystal statistics. i. a two-dimensional model with an order-disorder transition,” *Phys. Rev.*, vol. 65, pp. 117–149, Feb 1944.
- [32] S. El-Showk, M. F. Paulos, D. Poland, S. Rychkov, D. Simmons-Duffin, and A. Vichi, “Solving the 3d ising model with the conformal bootstrap,” *Phys. Rev. D*, vol. 86, p. 025022, Jul 2012.
- [33] D. Simson, “Linear representations of partially ordered sets and vector space categories,” 1993.
- [34] C.-S. Hu, “A brief note for sheaf structures on posets,” 2020.
- [35] K. Yanagawa, “Sheaves on finite posets and modules over normal semi-group rings,” *Journal of Pure and Applied Algebra*, vol. 161, no. 3, pp. 341–366, 2001.
- [36] A. Brown and O. Draganov, “Computing minimal injective resolutions of sheaves on finite posets,” 2021.
- [37] E. Spiegel and C. J. O’Donnell, *Incidence Algebras*. CRC Press, 1997.
- [38] M. Giry, “A categorical approach to probability theory,” in *Categorical aspects of topology and analysis*, vol. 915 of *Lecture notes in Mathematics*, pp. 68–85, Springer, 1982.
- [39] E. W. Lawvere, “The category of probabilistic mappings,” 1962. [Link](#).
- [40] S. Mac Lane, *Categories for the working mathematician*, vol. 5. Springer Science & Business Media, 2013.
- [41] J. Gallier and J. Quaintance, *Homology, Cohomology, and Sheaf Cohomology for Algebraic Topology, Algebraic Geometry, and Differential Geometry*. World Scientific, 2022.
- [42] H. K. Kesavan, *Jaynes’ maximum entropy principle*, pp. 1779–1782. Boston, MA: Springer US, 2009.
- [43] A. De Martino and D. De Martino, “An introduction to the maximum entropy approach and its application to inference problems in biology,” *Heliyon*, vol. 4, no. 4, p. e00596, 2018.
- [44] J. Pearl, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Inference*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1988.
- [45] S. L. Lauritzen, *Graphical Models*. Oxford Science Publications, 1996.

- [46] E. Sudderth and J. Pacheco, “Probabilistic graphical models.” Brown University CSCI 2950-P, Spring 2013, Lecture 12: Gaussian Belief Propagation, State Space Models and Kalman Filters, Guest Kalman Filter Lecture by Jason Pacheco, 2013.
- [47] T. Morita, “Cluster variation method of cooperative phenomena and its generalization i,” *Journal of the Physical Society of Japan*, vol. 12, no. 7, pp. 753–755, 1957.
- [48] O. Peltre, “A homological approach to belief propagation and bethe approximations,” in *Geometric Science of Information* (F. Nielsen and F. Barbaresco, eds.), (Cham), pp. 218–227, Springer International Publishing, 2019.
- [49] T. P. Speed, “A note on nearest-neighbour gibbs and markov probabilities,” *Sankhyā: The Indian Journal of Statistics, Series A*, 1979.
- [50] G.-C. Rota, “On the foundations of combinatorial theory I. Theory of Möbius functions,” *Probability theory and related fields*, vol. 2, no. 4, pp. 340–368, 1964.
- [51] H. K. T., “On the amount of information,” *Theory of Probability and its Applications*, 1962.