# UNCERTAINTY QUANTIFICATION FOR MARKOV RANDOM FIELDS

PANAGIOTA BIRMPA\* AND MARKOS A. KATSOULAKIS<sup>†</sup>

We present an information-based uncertainty quantification method for general Markov Random Fields, also known as Markov Networks. Markov Random Fields (MRFs) are structured, probabilistic graphical models over undirected graphs, and provide a fundamental unifying modeling tool for statistical mechanics, probabilistic machine learning, and artificial intelligence. Typically MRFs are complex and high-dimensional with nodes and edges (connections) built in a modular fashion from simpler, low-dimensional probabilistic models and their local connections; in turn, this modularity allows to incorporate available data to MRFs and efficiently simulate them by leveraging their graph-theoretic structure. Learning graphical models from data and/or constructing them from physical modeling and constraints necessarily involves uncertainties inherited from data, modeling choices, or numerical approximations. These uncertainties in the MRF can be manifested either in the graph structure or the probability distribution functions, and necessarily will propagate in predictions for quantities of interest. Here we quantify such uncertainties using tight, information-based bounds on the predictions of quantities of interest; these bounds take advantage of the graphical structure of MRFs and are capable of handling the inherent high-dimensionality of such graphical models. We demonstrate our methods in MRFs for medical diagnostics and statistical mechanics models. In the latter, we develop uncertainty quantification bounds for finite-size effects and phase diagrams, which constitute two of the typical predictions goals of statistical mechanics

**Key words.** Markov Random Fields, Uncertainty Quantification, Information Theory, Probabilistic Inequalities, Ising model, Long range interactions

AMS subject classifications. 62H22, 82B20, 94A17

1. Introduction. Probabilistic graphical models (PGM) constitute one of the fundamental tools for Probabilistic Machine Learning (ML) and Artificial Intelligence (AI), allowing for systematic and scalable modeling of uncertainty, causality, domain knowledge, and data assimilation, [38, 49, 36]. The main idea behind PGMs is to represent complex models and associated learning processes using random variables and their interdependence through a graph. We achieve it by constructing structured, high-dimensional probabilistic models, involving many parameters, nodes, and edges, from simpler ones with few parameters, nodes, and edges, thus allowing for distributed probability computations, and by incorporating available data, exploiting graphtheoretic model representations. PGMs are generally classified into Markov Random Fields (MRF) defined over undirected graphs, and Bayesian Networks, defined over Directed Acyclical Graphs [49] that represent conditional independencies between random variables, as well as mixtures of those two classes, [36]. Furthermore, the modeling flexibility of PGMs also allows to combine dynamics, data, and deep learning in Hidden Markov Models [38, 50, 47], as well as in recent work brings together multiscale modeling, physical constraints, and neural networks, [69, 41, 30].

Although the term random field may also refer to continuously indexed processes (e.g. gaussian random fields), in this paper MRFs refer to structured probabilistic models defined on undirected graphs; such PGMs are also referred to as Markov Networks. MRFs arise in statistical mechanics where interactions between particles are usually bi-directional, or when there may be no inherent evidence for causality

<sup>\*</sup>Department of Mathematics and Statistics, University of Massachusetts, Amherst, U.S.A (birmpa@math.umass.edu).

<sup>&</sup>lt;sup>†</sup>Department of Mathematics and Statistics, University of Massachusetts, Amherst, U.S.A (markos@math.umass.edu.

(directionality) and thus undirected graphs are the appropriate structure for such probabilistic models, [38, 49, 71]. Other applications of MRFs include image segmentation, image denoising [49, Sec. 4.2], text processing [64, 56], bioinformatics [61], computer vision [45], Markov logic networks, [21], Gaussian Markov networks [49, Sec. 7.3], artificial intelligence [36], and statistical mechanics [55, Sec. 19.4]. Overall, MRFs provide a fundamental unifying modeling tool for statistical mechanics, probabilistic machine learning, and artificial intelligence, [3, 38].

Learning MRFs can be based on available data, e.g. for learning the graph we refer to [49, 32, 43] for score-based methods, [55, 44] for independence tests on the graph, while maximum likelihood or Bayesian methods can be used for parameter identification, [49]. On the other hand, MRFs in statistical mechanics can be constructed from physical modeling and related constraints, [68, 55]. Therefore, the learning stage of MRFs necessarily involves uncertainties inherited from data, modeling choices, compromises on model complexity, or numerical approximations. These uncertainties in the MRF can be manifested either in the graph structure or the probability distribution functions, and necessarily will propagate through the graph structure and the corresponding structured probabilistic model in the predictions for quantities of interest (QoIs). To understand and quantify the impact of such uncertainties on model predictions, in this paper we present an information-based uncertainty quantification (UQ) method for general MRFs.

Model Uncertainty in Probabilistic Models: In general probabilistic models, uncertainties arising just from the fluctuations of the QoIs, associated with a given probabilistic model p, are referred to as aleatoric and occur when sampling p, [15]. They are handled by well-known tools, e.g. central limit theorems, concentration inequalities, Bayesian posteriors, MCMC, generalized Polynomial Chaos, etc. In contrast to this more standard type of uncertainty quantification, in MRFs, due to the learning process described earlier, we have model uncertainties (also known as epistemic), both in the structure (graph) and the probabilistic model itself—including parametric ones.

Next, we briefly describe the information-theoretic formulation of model uncertainty for general probabilistic models, without assuming any graphical model structures, see [39] for more details. To practically address model uncertainty, we typically compromise by constructing a surrogate or approximation or baseline model p. We construct families Q of (non-parametric) alternative models  $\tilde{p}$  to compare to p, while the "true" model  $p^*$ , which may be intractable or partly unknown, should belong to Q; for this reason we can refer to Q as the ambiguity set, typically defined as a neighborhood of alternative models around the baseline p:

(1.1) 
$$Q = Q^{\eta} = \left\{ \tilde{p} : d(\tilde{p}, p) \le \eta \right\},\,$$

where  $\eta > 0$  corresponds to the size of the ambiguity set and  $d = d(\tilde{p}, p)$  denotes a probability metric or divergence. The next natural mathematical goal is to assess the baseline model "compromise" and understand the resulting biases for QoIs f when we use p for predictions instead of the real model  $p^* \in \mathcal{Q}$ . We define the predictive uncertainty (or bias) for the QoI f when we use the baseline model p instead of any alternative model  $\tilde{p} \in \mathcal{Q}$  (including the real one  $p^*$ ) as the two worst case scenarios:

(1.2) 
$$\sup_{\tilde{p} \in \mathcal{Q}^{\eta}} \{ E_{\tilde{p}} f - E_{p} f \}$$

where  $E_{\tilde{p}}f$  denotes the expected value of the QoI f. Therefore, (1.2) provides a robust performance guarantee for the predictions of the baseline model p for f within

the ambiguity set  $Q^{\eta}$ . This robust perspective for general probabilistic models p is known in Operations Research as Distributionally Robust Optimization (DRO), e.g. [35, 37]. While the definition (1.2) is rather natural and intuitive, it is not obvious that it is practically computable since the neighborhood  $Q^{\eta}$  is infinite-dimensional. However it becomes tractable if we use for metric d in (1.1) the Kullback-Leibler (KL) divergence  $R(\tilde{p}||p)$ . Accordingly,  $\eta$  is a measure of the confidence in KL we put in the baseline model p. In recent work [15, 23, 39], it has been demonstrated that (1.2) (an infinite dimensional optimization problem) is directly computable using the variational formula (follows directly from the Donsker Varadhan variational principle, [23]):

(1.3) 
$$\sup_{\tilde{p} \in \mathcal{Q}^{\eta}} \left\{ E_{\tilde{p}} f - E_{p} f \right\} = \pm \inf_{c > 0} \left[ \frac{1}{c} \log \int e^{\pm c(f - E_{p} f)} p(dx) + \frac{\eta}{c} \right].$$

In this formula we recognize two main ingredients:  $\eta$  is model uncertainty from (1.1) while the Moment Generating Function (MGF)  $\int e^{\pm cf} p(dx)$  encodes the QoI f at the baseline model p. In [23, 39] the authors have developed techniques to compute (exactly or approximately via asymptotics [23]) as well as to provide explicitly upper and lower bounds on (1.2) in terms of concentration inequalities [39]. Tightness, i.e when the sup and inf in (1.2) are attained by an appropriate measure  $\tilde{p}$  have also been studied in [39]. Finally, related UQ bounds have been derived for Markov processes using variational principles and functional inequalities [6], and in rare events [2, 24].

Main results: The main thrust of our results here is to build on the aforementioned perspective for information-based UQ, in order to develop UQ methods for MRFs, and to address their specific UQ challenges. In particular, here we address both structure (graph) and probabilistic uncertainties—including parametric ones—using tight, information-based bounds on the predictions of QoIs; although these new UQ bounds rely on (1.2), they specifically, (a) take advantage of the graphical structure of MRFs, and (b) are capable of handling the inherent high-dimensionality of such graphical models, i.e. there is a necessity for scalable UQ in the size of the system, namely the number of nodes in MRFs such as in the thermodynamic limit of statistical mechanics models.

Regarding the scalability issue, in [46] the authors tested various model uncertainty metrics in defining  $d(\tilde{p}, p)$  in (1.1) such as the Hellinger distance and  $\chi^2$  divergence and inequalities, such as Csiszar-Kullback-Pinsker and the Hammersley-Chapman-Robbins inequalities, [67], in order to bound the model bias with respect of a QoI in the spirit of (1.3). It was shown that among these bounds the only one that scales with the dimension of the model p is (1.3) and  $d(\tilde{p}, p)$  should be the KL divergence.

Once we have settled to the use of the KL divergence for the aforementioned scalability reasons, we turn our attention to the baseline MRF p, the ambiguity set (1.1) and the corresponding alternative MRFs  $\tilde{p}$ . Based on the earlier discussion on model uncertainty for MRFs arising from statistical learning of graph models or physical modeling, we introduce a unifying perspective of three general types of alternative models  $\tilde{p}$ , based on their relative structure to the baseline p: Type I MRFs where the graph structures (nodes and edges) are identical to the baseline p and the parameters of probability distributions are different, Type II where the nodes are the same, but the edges and parameters are different. Finally, Type III where the nodes, structure, and parameters are all different.

In general, MRFs satisfy the specific conditional independence properties dis-

cussed in subsection 2.2. Contrary to Bayesian Networks, their distributions cannot always be factorized by a product of local conditional distributions or local functions over the graph. The celebrated Hammersley-Clifford Theorem, also known as Fundamental Theorem of Random Fields [49, 42, 55], guarantees such a factorization along maximal cliques of the graph under the assumption that p > 0. Here, we make such an assumption for both baseline and Type I-III MRFs. Consequently, the KL divergence is finite without requiring absolute continuity with respect to p.

We take advantage of all the above and we study UQ problems by developing a unified strategy for Type I and II MRFs while Type III is not covered here as explained in Section 3. We focus on the two primary ingredients of (1.2), namely the KL divergence and the MGF, and how they manifest themselves on MRFs. In KL divergence, the factorization discussed earlier is a crucial tool for its simplification and numerical calculation. It allows us to compare local discrepancies in parameters and structure between the baseline p and alternative models. We call these discrepancies excess factors of Type I-II given p. We develop a unifying method for computing the excess factors by interrelating the maximal cliques of alternative MRFs and the baseline MRF p. As for the MGF, the choice of QoIs is determinant. We focus on two different QoIs; those that are involved in the models (e.g. sufficient statistics) as well as characteristic functions defined on events of interest.

Regarding tightness of UQ bounds discussed earlier, we find specific distributions that the derived UQ bounds for MRFs are attainable. In addition, we go beyond that, and pose the question: Given a QoI and a baseline MRF p, what are the possible associated undirected graphs such that the conditional independence properties implied by the graphs are satisfied by the distributions? Such a question introduces the concept of tightness at a graph level. There are cases where we can explicitly determine the associated graphs and others (when the structure is different than the baseline) that depend on the QoI. In the latter case, we give an example that points out a unifying method to construct the right graph or at least, a set of possible graphs.

**Demonstration of UQ for MRFs:** We first demonstrate all the above concepts and UQ methods in a fairly simple and low dimensional MRF example from medical diagnostics. Subsequently, we implement our approach on several high-dimensional statistical mechanics models as they are fundamental in ML [3, 38]. We develop UQ bounds for finite size effects and phase diagrams, which constitute two of the typical predictions goals of statistical mechanics modeling and both require scalable UQ methods.

Specifically, we consider as a baseline model p an Ising-spin system with Kac-type interactions, see [57]. Such a model combines sufficient complexity—since it is not a mean field model—but it is still analytically fairly tractable to serve as a good benchmark problem for high-dimensional MRF. Alternative models  $\tilde{p}$  considered here are 1) Ising models with perturbed interaction potentials with respect to the baseline, 2) models with truncated interactions to facilitate computational implementations, [68], and 3) perturbations by a long-range interaction (even longer than a Kac interaction). As we discuss in Section 6, these systems are typically defined in bounded domains with boundary conditions being a given configuration outside of the domain. To have a graph description of these systems, MRFs need to be modified to account for conditioning a Gibbs distribution on an eliminated set of nodes identified as a configuration defined outside of the domain by using reduced Markov Random Fields (rMRFs) (see [49]). Typical questions we address in these examples include the following: (i) How to capture the phase diagram of a perturbed model through its comparison with the

baseline phase diagram by bounding the model bias. (ii) How to truncate an interaction so as the phase diagram of the baseline model and the truncated one are close within a prescribed tolerance. Note that an extensive analysis on the intersection between other concepts and methods from statistical mechanics—also including non-equilibrium statistical mechanics—and deep learning have been reviewed in [3].

Related methods: We note that existing general-purpose UQ & sensitivity analysis methods, e.g., gradient and ANOVA-based methods, [63, 60, 29] cannot handle UQ with model uncertainties, due to their inherently parametric nature, while it is not clear how they can take advantage of the graphical, causal structure in MRFs. Furthermore, there is earlier work on model uncertainty that represents missing physics with a stochastic noise but without the detailed structure of a graphical model, [51, 65]. In our work, there is a natural structure embedded in the model uncertainty, arising through the graph structure of the MRFs.

Sensitivity analysis has also a long history in statistical mechanics, known as linear and nonlinear response theory, [59, 4], addressing the impact of small and larger parametric perturbations respectively. These types of methods are covered by our approach, as models with perturbed weights are clearly of Type I.

Furthermore, in contrast to these results, a key point in our work here, also immediately clear from (1.3), is that the model perturbations we can consider are not necessarily small. For instance, the parameter  $\eta$  in (1.3) does not need to be small, allowing for global and non-parametric sensitivity analysis; the latter since the KL divergence allows us to consider models outside a specific parametric family, e.g. comparing statistical mechanics models with different potentials. Similarly, we explicitly compute the UQ bounds for large perturbations in a medical diagnostics example.

Sensitivity analysis in MRFs has been also studied in [14]. The authors tackle fundamental questions such as bounding belief change between Markov networks with the same structure but different parameter values. They propose a distance measure and bound the relative change in probability queries by the relative change in parameters (Type I). Global sensitivity in parameters has been studied in [17]. In particular, the authors developed an algorithm that checks the robustness of a MAP configuration i.e. the most likely configuration, in discrete probabilistic graphical models under global perturbations. The present work goes beyond local or global parametric sensitivity analysis that allows us to consider perturbations in both parameters and edges of the graph of the MRF and examines their impact on the prediction of specific QoIs. Special cases of our results for mean field and nearest neighbor Ising models were considered earlier in [46]. Finally, we note that parametric sensitivity analysis for the other class of (directed) probabilistic graphical models, namely Bayesian Networks, was developed in [16] using similar tools to [14]. Parametric sensitivity analysis based on mutual information for multi-scale partial differential equations and neural networks informed by Bayesian Network priors was developed in [70] and [66]. Model uncertainty quantification based on information theory inequalities in the spirit of (1.3) were recently introduced for Bayesian Networks arising in chemical sciences, [30].

This article is organized as follows: We start with some concepts from graph theory to fix notation and then we give a brief background of MRFs/rMRFs (Section 2). Supplementary background behind rMRFs is provided in the Appendix A. We formally introduce the idea of graph interconnections, the impact on distributions and alternative models in Section 3. The main results are presented in Section 4 and

provide UQ bounds for rMRFs, preparing the ground for applications to statistical mechanics models. In Section 5, we present a simple example from medical diagnostics. Section 6 is devoted to UQ for finite size effects, scalability, and finally UQ for phase diagrams for generic interactions and the Ising-Kac model. In the remaining sections of the Appendix, we further discuss the Ising-Kac model, we provide the technical background required for the UQ analysis of Section 6 (e.g Lebowitz-Penrose (LP) limit), we include the proofs of the main results, and explicit calculations of the UQ for medical diagnostics example and statistical mechanics.

#### 2. Preliminaries.

- **2.1. Definitions from Graph Theory.** We start with some notation and terminology from graph theory. A **graph** is a data structure  $\mathcal{G}$  consisting of a set of **nodes**,  $\mathcal{V} = \{1, 2, \dots, N\}$  and a set of **edges**  $\mathcal{E}$ , i.e. all pairs of nodes  $i, j \in \mathcal{V}$  which are connected by an edge, denoted by (i, j). An edge can be directed, denoted by  $i \to j$  or undirected, denoted by i j. A graph is **directed** [resp. **undirected**] if all the edges are directed [resp. **undirected**]. The nodes  $i, j \in \mathcal{V}$  are **adjacent** if and only if  $(i, j) \in \mathcal{E}$ . The **neighborhood** of node i, denoted by  $\mathcal{N}_i$  is the set of nodes which i is adjacent. For sets of nodes A, B and C, C **separates** A **from** B, denoted by  $\{i \in A\} \perp_{\mathcal{G}} \{j \in B\} \mid \{k : k \in C\}$ , if and only if when we remove all the nodes in C there is no path connecting any node in A to any node in B. Lastly, if  $\mathcal{M} \subset \mathcal{V}$ , the **induced subgraph** of  $\mathcal{G}$  is defined as  $\mathcal{G}[\mathcal{M}] = (\mathcal{M}, \mathcal{E}')$  where  $\mathcal{E}'$  includes all the edges  $(i,j) \in \mathcal{E}$  such that  $i,j \in \mathcal{M}$ .
- **2.2.** Conditional Independence Properties and MRFs. In this subsection, we define three conditional independence properties that are necessary for MRFs.
- Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and let  $\mathbf{Y} = \{Y_i\}_{i=1}^{|\mathcal{V}|}$  be a set of random variables that each one is attached to a node and  $|\mathcal{V}|$  denotes the cardinality of  $\mathcal{V}$ .
- Pairwise Markov property (P): Any two non adjacent variables are conditionally independent (CI) given the rest, i.e. a conditional joint can be written as a product of conditional marginals; CI is denoted by  $Y_i \perp Y_j \mid \{Y_k : k \neq i, j\}$ ,
- Local Markov property (L): Any variable  $Y_i$  is conditionally independent of all the others given its neighbors, that is  $Y_i \perp \{Y_k : k \notin \mathcal{N}_i\} \mid \{Y_k : k \in \mathcal{N}_i\}$ ,
- Global Markov property (G): If A, B, C are sets of nodes then any two sets of variables,  $\mathbf{Y}_A = \{Y_i : i \in A\}$  and  $\mathbf{Y}_B = \{Y_i : i \in B\}$  are conditionally independent given a separating set of variables  $\mathbf{Y}_C = \{Y_i : i \in C\}$ , that is  $\mathbf{Y}_A \perp \mathbf{Y}_B \mid \mathbf{Y}_C$ .

It is obvious that (G) implies (L) which implies (P).

DEFINITION 2.1. Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be an undirected graph where  $\mathcal{V} = \{1, 2, ..., N\}$  is the set of nodes and  $\mathcal{E}$  is the set of edges. Let also consider a set of random variables  $\mathbf{Y} = (Y_i)_{i \in \mathcal{V}}$  indexed by  $\mathcal{V}$  where each  $Y_i$  takes values on a finite set  $\mathcal{S}$ . Their joint probability distribution is denoted by p. We say that  $(\mathbf{Y}, p)$  is a Markov Random Field (MRF) iff (G) is satisfied.

As MRFs are defined on an undirected graph, it does not allow to use chain rule of conditional probabilities and further describe the probability distribution  $p(\mathbf{y})$ . A factorization rule for MRFs (i.e. for undirected graphs and the conditional independencies) is important and is provided by Hammersley and Clifford in their unpublished work [42, 40]. To state their result, we need a few more definitions. Let  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  be a graph and let  $c \subset \mathcal{V}$ .

(i) c is called **clique** if any pair of nodes in c is connected by some edge.

(ii) c is called **maximal clique** if any superset c' of c (i.e  $c' \supset c$ ) is not a clique any more. The set of all maximal cliques of graph  $\mathcal{G}$  is denoted by  $\mathcal{C}_{\mathcal{G}}$ .

**Hammersley-Clifford Theorem** A positive distribution  $p(\mathbf{y}) > 0$  satisfies one of (P), (L) and (G) of an undirected graph G iff p parametrized by some parameters  $\mathbf{w} = \{\mathbf{w}_c\}_{c \in \mathcal{C}_{\mathcal{G}}}$  can be represented as a product of clique potentials, i.e

(2.1) 
$$P_{\Psi}^{\mathbf{w}}(\mathbf{y}) \equiv p(\mathbf{y} \mid \mathbf{w}) = \frac{1}{Z(\mathbf{w})} \prod_{c \in \mathcal{C}_{\mathcal{G}}} \Psi_{c}(\mathbf{y}_{c} \mid \mathbf{w}_{c})$$

where  $\Psi_c(\mathbf{y}_c \mid \mathbf{w}_c)$  is a positive function defined on the random variables in clique c and parametrized by some parameters  $\mathbf{w}_c$ , and is called **clique potential**. Also  $Z(\mathbf{w})$  is the partition function given by

(2.2) 
$$Z(\mathbf{w}) = \sum_{\mathbf{y}} \prod_{c \in \mathcal{C}_{\mathcal{G}}} \Psi_c(\mathbf{y}_c \mid \mathbf{w}_c)$$

The theorem states that the set of all joint distributions on an undirected graph  $\mathcal{G}$  that can be factorized as in (2.1) is identical to the set of joint distributions that satisfy the conditional independence properties, under the restriction of strictly positive distributions.

Remark 2.2. Without the assumption of strict positiveness of the joint distribution p, the theorem is not valid. A counterexample has been obtained in [54].

Remark 2.3. The KL divergence or any other f-divergences between a baseline MRF that is assumed nonnegative and alternative MRFs of Type II-III ( different structure, see Introduction) could be infinite due to the loss of absolute continuity. In that case, the Wasserstein metric or the  $\Gamma$ -Divergence, [25], could potentially be good alternatives for the KL divergence in defining (1.1). The implementation of the Wasserstein metric or the  $\Gamma$ -Divergence is still unexplored in the context of such MRFs. For this purpose, the development of new methods constitutes an important step towards comparing MRFs with different structures and nonnegative distributions. In this article, we restrict our attention to the Hammersley-Clifford Theorem and we assume strictly positive probability distributions.

Given a MRF  $(\mathbf{Y}, p)$ , a reduced Markov Random Field (rMRFs) is obtained by conditioning p on some observation  $\mathbf{U} = \mathbf{u}$  with  $\mathbf{U} \subset \mathbf{Y}$ . Hence, the distribution of the resulting rMRF has a reduced number of clique potentials. As we discuss in Section 6, rMRFs are appropriate for formulating statistical mechanics models defined on bounded domains with a given configuration outside of the domain in a graph language. Next, we formally introduce rMRFs.

**2.2.1. Reduced Markov Random Fields (rMRFs).** Let  $\mathbf{Y} = \{Y_i\}_{i \in \mathcal{V}}$  be a collection of random variables indexed by a set of nodes  $\mathcal{V}$  of a graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ , taking values in some space  $\mathcal{Y}^{\mathcal{V}} = \bigotimes_{i=1}^{\mathcal{V}} \mathcal{Y}_i$ . Let  $p \equiv p(\cdot|\mathbf{w})$  be a strictly positive joint probability distribution of  $\mathbf{Y}$  parametrized by  $\mathbf{w}$  such that  $(\mathbf{Y}, p(\cdot|\mathbf{w}))$  is a MRF.

Let **u** be a context and  $\mathcal{M} \subset \mathcal{V}$ . If  $\mathbf{U} = \{Y_i\}_{i \in \mathcal{M}}$  with  $\mathbf{U} = \mathbf{u}$ , we construct the corresponding rMRF as follows: let  $\mathbf{Z} = \{Y_i\}_{i \in \mathcal{V} \setminus \mathcal{M}}$  and  $q(\mathbf{z}|\mathbf{w})$  be the probability distribution factorized according to Proposition A.2 (the analogue of the Hammersley-Clifford Theorem for rMRFs):  $q(\mathbf{z}) \equiv q(\mathbf{z}|\mathbf{w}) = \frac{1}{Z_{\mathbf{u}}(\mathbf{w})} \prod_{c \in \mathcal{C}_{\mathcal{G}}} \Psi_c[\mathbf{u}](\mathbf{z}_c \mid \mathbf{w}_c)$ . More details on rMRFs are given in Appendix A.

The next two sections are presented for rMRFs as we can then recall formulas and the main results directly in the UQ analysis of statistical mechanics models in Section 6.

Their formulation and analysis hold for MRFs and when required, we will be providing more details for their implementation to MRFs.

3. Mathematical Formulation of UQ on MRFs/rMRFs. Let q be a rMRF constructed by learning from available data or from physical modeling and related constraints. Constructing such a model involves uncertainties either in the graph structure or the probability distribution functions, and necessarily will propagate through the graph structure and the corresponding structured probabilistic model in the predictions for QoIs. We quantify the impact of such uncertainties on model predictions by constructing ambiguity sets such as (1.1) consisting of alternative rMRFs given by

(3.1) 
$$Q^{\eta} = \left\{ \text{rMRFs } \tilde{q} : R(\tilde{q}||q) \le \eta \right\},$$

where  $\eta > 0$  corresponds to the size of the ambiguity set. The alternative models  $\tilde{q}$  in (3.1) can be classified into: Type I MRFs where the graph structures (nodes and edges) are identical to the baseline q and the parameters of probability distributions are different, Type II where the nodes are the same, but the edges and parameters are different. Finally, Type III where the nodes, structure, and parameters are all different. Next, we mathematically formulate the alternative models.

**3.1. Alternative models.** Let  $(\mathcal{G}, \mathbf{w}, p)$  and  $(\tilde{\mathcal{G}}, \tilde{\mathbf{w}}, \tilde{p})$  be two MRFs with  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and  $\tilde{\mathcal{G}} = (\tilde{\mathcal{V}}, \tilde{\mathcal{E}})$  being the associated graphs, where  $\mathcal{V}$  and  $\tilde{\mathcal{V}}$  are the sets of nodes and  $\mathcal{E}$  and  $\tilde{\mathcal{E}}$  are the sets of edges.

DEFINITION 3.1.  $(\tilde{\mathcal{G}}, \tilde{\mathbf{w}}, \tilde{p})$  and  $(\mathcal{G}, \mathbf{w}, p)$  can have one of the following interconnections:

Type I:  $\tilde{\mathcal{V}} = \mathcal{V}, \ \tilde{\mathcal{E}} = \mathcal{E} \ and \ \tilde{\mathbf{w}} \neq \mathbf{w}, \ or$ 

Type II:  $\tilde{\mathcal{V}} = \mathcal{V}, \; \mathcal{E} \subset \tilde{\mathcal{E}} \; \text{and} \; \tilde{\mathbf{w}} \neq \mathbf{w} \; \text{or}$ 

Type III:  $\tilde{\mathcal{V}} \neq \mathcal{V}, \; \mathcal{E} \neq \tilde{\mathcal{E}} \; and \; \tilde{\mathbf{w}} \neq \mathbf{w}.$ 

From now on, we refer to the baseline model when we use the notation  $(\mathcal{G}, \mathbf{w}, p)$  and without loss of generality we assume  $\mathcal{E} \subset \tilde{\mathcal{E}}$ . This assumption simplifies the presentation of our approach but intuitively speaking, the fewer edges a rMRF has, the more information it provides, since in a sparser graph, there are more conditional independencies specified.

Based on that, we interrelate the maximal cliques of Type I-II models with those of p. In particular, for Type I there is one to one correspondence between maximal cliques. Changes on the set of edges of a Type II model lead to different sets of maximal cliques and one needs to examine the nature of the new edges and their impact on the maximal cliques of p. Finally, the new set of nodes of a Type III model leads to a drastically new structure that makes such interrelation of maximal cliques hard to achieve. Therefore, this case is not examined here.

Let **u** be a context and  $\mathcal{M} \subset \mathcal{V} \cap \tilde{\mathcal{V}}$ . For  $\mathbf{U} = \{Y_i\}_{i \in \mathcal{M}}$  with  $\mathbf{U} = \mathbf{u}$ , we construct the corresponding rMRFs  $(\mathbf{Z}, q(\cdot|\mathbf{w}))$  and  $(\tilde{\mathbf{Z}}, \tilde{q}(\cdot|\tilde{\mathbf{w}}))$  parametrized by **w** and  $\tilde{\mathbf{w}}$  respectively. Based on the structural classification Type I-III, the probability distributions of  $\tilde{q}$  are treated as follows:

**Type I.** Let  $\mathcal{B} \subset \mathcal{C}_{\mathcal{G}}$  be the set of maximal cliques whose weights differ, i.e for each  $c \in \mathcal{B}$ ,  $\tilde{\mathbf{w}}_c \neq \mathbf{w}_c$ . The clique potentials of  $\tilde{q}(\cdot|\tilde{\mathbf{w}})$  can be rewritten as

$$(3.2) \tilde{\Psi}_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \tilde{\mathbf{w}}_{c}) = \begin{cases} \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c})\Phi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c}, \tilde{\mathbf{w}}_{c}) & , \text{if } c \in \mathcal{B} \\ \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c}) & , \text{otherwise.} \end{cases}$$

We call  $\Phi_c[\mathbf{u}](\cdot \mid \tilde{\mathbf{w}}_c, \mathbf{w}_c) > 0$   $\tilde{q}$ -excess factor of type I relative to q on c and is defined on variables  $\mathbf{z}_c$  in clique  $c \in \mathcal{B}$ . Cliques where no change on weights has occurred, remain the same.

**Type II.** In this type, the class of maximal cliques  $C_{\tilde{\mathcal{G}}}$  is different. The analysis becomes more complicated and clique potentials need to be carefully considered. We look into the nature of one or more new edges by categorizing it as one of the following types: a new edge (i) can create a totally new maximal clique, see Figure 1, third graph, (ii) can connect two or more already existing maximal cliques, see Figure 1, second graph, and (iii) can enlarge an already existing maximal clique, see Figure 1, forth graph. By adding more than one new edges, the new maximal cliques of  $\tilde{\mathcal{G}}$  can









Fig. 1. (First) Baseline MRF model p demonstrated by graph  $\mathcal{G}$ . (Second) Alternative model  $\tilde{p}$  with the associated graph obtained by adding the yellow edge (4-7) and connecting two maximal cliques of p model,  $\{3,4,6\}$  and  $\{3,6,7\}$ , thus  $\tilde{p}$  has a new maximal clique  $\{3,4,6,7\}$ . (Third) Alternative model  $\tilde{p}$  with the associated graph obtained by adding the red edge (6-10) and thus  $\tilde{p}$  has a totally new maximal clique  $\{6,10\}$ . (Forth) Alternative model  $\tilde{p}$  with the associated graph obtained by adding the blue edge (5-10) and enlarging the already existing clique,  $\{5,8\}$ , to  $\{5,8,10\}$ .

be obtained by a combination of (i), (ii), and (iii). We introduce the following sets:

$$(3.3) \mathcal{B}_{\cup} = \{ \tilde{c} \in \mathcal{C}_{\tilde{G}} \setminus \mathcal{C}_{\mathcal{G}} : \tilde{c} = \cup_{i} c_{i}, \text{ for } c_{i} \in \mathcal{C}_{\mathcal{G}} \}$$

(3.4) 
$$\mathcal{B}_{\subseteq} = \{ \tilde{c} \in \mathcal{C}_{\tilde{\mathcal{G}}} \setminus \mathcal{C}_{\mathcal{G}} : \text{there exists } c \in \mathcal{C}_{\mathcal{G}} \text{ s.t. } c \subseteq \tilde{c} \}$$

$$(3.5) \mathcal{B}_{\text{new}} = (\mathcal{C}_{\mathcal{G}} \cup \mathcal{B}_{\sqcup} \cup \mathcal{B}_{\subset})^{c}$$

Then the clique potentials of  $\tilde{q}$  can be rewritten as:

$$(3.6) \quad \tilde{\Psi}_{\tilde{c}}[\mathbf{u}](\mathbf{z}_{\tilde{c}} \mid \tilde{\mathbf{w}}_{\tilde{c}}) = \begin{cases} \prod_{c_i} \Psi_{c_i}[\mathbf{u}](\mathbf{z}_{c_i} \mid \mathbf{w}_{c_i}) \Phi_{\tilde{c}}^{(ii)}[\mathbf{u}](\mathbf{z}_{\tilde{c}} \mid \mathbf{w}_{c}, \tilde{\mathbf{w}}_{c}) &, \text{if } \tilde{c} \in \mathcal{B}_{\cup}, \\ \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c}) \Phi_{\tilde{c}}^{(iii)}[\mathbf{u}]((\mathbf{z}_{\tilde{c}} \mid \mathbf{w}_{c}, \tilde{\mathbf{w}}_{c}) &, \text{if } \tilde{c} \in \mathcal{B}_{\subset}, \\ \tilde{\Psi}_{\tilde{c}}[\mathbf{u}](\mathbf{z}_{\tilde{c}} \mid \tilde{\mathbf{w}}_{\tilde{c}}) &, \text{if } \tilde{c} \in \mathcal{B}_{\text{new}} \\ \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c}) &, \text{if } \tilde{c} \in \mathcal{C}_{\tilde{\mathcal{G}}} \end{cases}$$

We call  $\Phi_{\tilde{c}}^{(ii)}$ ,  $\Phi_{\tilde{c}}^{(iii)} > 0$   $\tilde{q}$ -excess factors of type II relative to q on  $\tilde{c}$  defined on the variables of  $\tilde{c}$ . In fact, the two functions play the role of the discrepancy at a distribution level when new maximal clique  $\tilde{c}$  has been created by connecting existing maximal cliques  $c_i$  and by enlarging an existing maximal clique. When  $\tilde{c} \in \mathcal{B}_{\text{new}}$ , there is no need to express the clique potential through the potentials of  $q(\cdot \mid \mathbf{w})$ . For simplicity, we assume that clique potentials on common maximal cliques between  $\mathcal{G}$  and  $\tilde{\mathcal{G}}$  do not change. However, one can consider different potentials and in that case, a term  $\Phi$  should be introduced similar to (ii) and (iii). For convenience, we establish

one last unifying terminology. We call

(3.7) 
$$\Phi_{\mathbf{u}}^{\mathbf{I}}(\mathbf{Z}) := \prod_{c \in \mathcal{B}} \Phi_{c}[\mathbf{u}](\mathbf{Z}_{c})$$

$$(3.7) \qquad \Phi_{\mathbf{u}}^{\mathrm{I}}(\mathbf{Z}) := \prod_{c \in \mathcal{B}} \Phi_{c}[\mathbf{u}](\mathbf{Z}_{c})$$

$$(3.8) \qquad \Phi_{\mathbf{u}}^{\mathrm{II}}(\mathbf{Z}) := \prod_{\tilde{c} \in \mathcal{B}_{\mathrm{new}}} \tilde{\Psi}_{\tilde{c}}[\mathbf{u}](\mathbf{Z}_{\tilde{c}} \mid \tilde{\mathbf{w}}_{\tilde{c}}) \prod_{\tilde{c} \in \mathcal{B}_{\cup}} \Phi_{\tilde{c}}^{(ii)}[\mathbf{u}](\mathbf{Z}_{\tilde{c}}) \prod_{\tilde{c} \in \mathcal{B}_{\subseteq}} \Phi_{\tilde{c}}^{(iii)}[\mathbf{u}]((\mathbf{Z}_{\tilde{c}}) \mid \tilde{\mathbf{w}}_{\tilde{c}})$$

total  $\tilde{q}$ -excess factor of type I and II relative to q respectively. The total  $\tilde{q}$ -excess factor of type I relative to q captures all the parameters changes while the total  $\tilde{q}$ excess factor of type II relative to q captures all the structural discrepancies. In the case of MRF, we drop the context  $\mathbf{u}$  from (3.7) and (3.8) and  $\mathbf{Z}$  is replaced by Y. Equations (3.3)-(3.8) are explicitly specified in medical diagnostics application in Section 5 and its detailed analysis in Appendix D, and in statistical mechanics, Section 6. In Type III, there exists the total  $\tilde{q}$ -excess factor of type III relative to q. However, due to the high degree of discrepancies, we cannot interrelate maximal cliques of Type III model with q, and by extension each  $\tilde{q}$ -excess factor cannot be determined. The next results are straightforward but essential in our calculations. To avoid heavy notation, we remind that  $q(\cdot) = q(\cdot \mid \mathbf{w})$  and  $\tilde{q}(\cdot) = \tilde{q}(\cdot \mid \tilde{\mathbf{w}})$ .

Partition function of alternative models. Based on the above description of alternative models, the partition function of  $\tilde{q}$  is given in the next lemma.

LEMMA 3.2. Let  $(\mathbf{Z}, q)$  be a rMRF. Then for any alternative rMRF  $(\mathbf{Z}, \tilde{q})$  of Type i with i = I, II its partition function is expressed as:

(3.9) 
$$\tilde{Z}_{\mathbf{u}}(\tilde{\mathbf{w}}) = E_q[\Phi_{\mathbf{u}}^{\mathbf{i}}] Z_{\mathbf{u}}(\mathbf{w})$$

where  $\Phi_{\mathbf{u}}^{i}$  are given by (3.7) and (3.8).

*Proof.* The proof is based on the method of interrelating the distribution q and  $\tilde{q}$ , utilizing the total  $\tilde{q}$ -excess factors relative to q given by (3.7) and (3.8). The explicit computation is provided in Appendix B.1.

**Likelihood ratio.** The following lemma provides the likelihood ratio between  $\tilde{q}$  and q and constitutes the key ingredient for the simplification of (4.2) and the UQ bounds provided in (4.1).

LEMMA 3.3. Let  $(\mathbf{Z}, q)$  be a rMRF. Then for any alternative rMRF  $(\mathbf{Z}, \tilde{q})$  of Type i with i = I, II, the corresponding likelihood ratio satisfies:

(3.10) 
$$\frac{d\tilde{q}}{dq} = \frac{\Phi_{\mathbf{u}}^{i}}{E_{q}[\Phi_{\mathbf{u}}^{i}]}$$

where  $\Phi_{\mathbf{u}}^{i}$  is given by (3.7) and (3.8).

The proof is omitted as the lemma is a direct consequence of the method of interrelating two distributions discussed above and Lemma 3.2. Note that both results hold for MRFs denoted by  $(\mathbf{Y}, p)$  and  $(\mathbf{Y}, \tilde{p})$ , dropping the context **u** from  $\Phi_{\mathbf{u}}^{\mathbf{i}}$  in (3.10).

**3.2.** KL divergence. As we see in Section 4, our UQ methods rely on the KL divergence as means to measure "distance" between baseline and alternative MRFs. The fact that it scales correctly with the dimension of the baseline model [23] as well as the commonalities in parameters and structure between baseline and alternative models combined with the Hammersley-Clifford Theorem allows the KL divergence to be expressed in a simplified and informative form. In particular, we show that KL divergence (which is finite due to the positive probabilities q and  $\tilde{q}$ ) depends only on the total  $\tilde{q}$ -excess factor relative to q given by (3.7) and (3.8). To simplify the notation, we omit the dependence of  $\mathbf{Z}$  from  $\kappa_i$ , f and  $\Phi_n^i$ .

LEMMA 3.4. Let  $(\mathbf{Y}, p^{\mathbf{w}})$ ,  $(\mathbf{Y}, \tilde{p}^{\tilde{\mathbf{w}}})$  be two MRFs defined over graphs  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and  $\tilde{\mathcal{G}} = (\mathcal{V}, \tilde{\mathcal{E}})$  respectively. Let  $\mathbf{u}$  be a context and  $\mathcal{M} \subset \mathcal{V}$ . We consider the corresponding rMRFs  $(\mathbf{Z}, q)$ ,  $(\mathbf{Z}, \tilde{q})$ .

a. If  $\tilde{q}$  is Type i, with i = I or II, then the KL divergence is given

$$R(\tilde{q}||q) = E_{\tilde{q}} \left[ \log \frac{\tilde{q}}{q} \right] = E_q \left[ \frac{\tilde{q}}{q} \log \frac{\tilde{q}}{q} \right]$$

$$= E_{\tilde{q}} [\log \Phi_{\mathbf{u}}^{i}] - \log E_q [\Phi_{\mathbf{u}}^{i}] = \frac{1}{E_q [\Phi_{\mathbf{u}}^{i}]} E_q \left[ \Phi_{\mathbf{u}}^{i} \log \Phi_{\mathbf{u}}^{i} \right] - \log E_q [\Phi_{\mathbf{u}}^{i}]$$
(3.11)

where  $\Phi_{\mathbf{u}}^{i}$  is defined in (3.7) and (3.8) accordingly.

b. If  $\tilde{q}$  is Type i, with i=I or II, then for any f satisfying (4.3), the KL divergence is given by

$$(3.12) \qquad \mathcal{R}(\tilde{q}\|q) = C_{\mathbf{i}} E_{\tilde{q}}[f] + \frac{E_{q}\left[\kappa_{\mathbf{i}}\Phi_{\mathbf{u}}^{\mathbf{i}}\right]}{E_{q}\left[\Phi_{\mathbf{u}}^{\mathbf{i}}\right]} - \log E_{q}\left[\Phi_{\mathbf{u}}^{\mathbf{i}}\right], \quad \Phi_{\mathbf{u}}^{\mathbf{i}}(\mathbf{Z}) = e^{C_{\mathbf{i}}f(\mathbf{Z}) + \kappa_{\mathbf{i}}(\mathbf{Z})}$$

*Proof.* a. We express the KL divergence as follows

$$R(\tilde{q}||q) = E_{\tilde{q}} \left[ \log \frac{\tilde{q}}{q} \right] = E_q \left[ \frac{\tilde{q}}{q} \log \frac{\tilde{q}}{q} \right]$$

Then, we use Theorem 3.3 and we obtain (3.11). For b., we additionally recall (4.3).

Remark 3.5. As mentioned in Theorem 3.3, the result holds for MRFs denoted by  $(\mathbf{Y}, p)$  and  $(\mathbf{Y}, \tilde{p})$ , dropping the context  $\mathbf{u}$  from  $\Phi_{\mathbf{u}}^{\mathbf{i}}$ .

4. Main Results. In this section, we present a information-based UQ method on the predictions for QoIs for general MRFs/rMRFs by quantifying the model uncertainty for MRFs/rMRFs arising from statistical learning of graph models or from physical modeling. Our starting point is the Donsker-Varadhan variational principle [22], which in turn implies the *Gibbs Variational principle* for the KL divergence (see [15, 23]):

$$(4.1) \qquad \sup_{\lambda > 0} \left\{ \frac{-\Lambda_p^f(-\lambda) - R(\tilde{q}\|q)]}{\lambda} \right\} \le E_{\tilde{q}}[f] \le \inf_{\lambda > 0} \left\{ \frac{\Lambda_q^f(\lambda) + R(\tilde{q}\|q)]}{\lambda} \right\}$$

As mentioned earlier, we focus on KL divergence as it scales correctly with the dimension of the baseline model [23]. In the above inequality, q is the baseline rMRF and  $\tilde{q}$  is an alternative model in the ambiguity set defined in (3.1). We note that at a MRF point of view, (4.1) holds as well. Moreover,  $\Lambda_q^f(\lambda)$  is the cumulant generating function (CGF) computed with respect to p given by

(4.2) 
$$\Lambda_q^f(\lambda) := \log E_q[e^{\lambda f}]$$

while f is a QoI. The class of QoI that we examine here as discussed in the next subsection.

We take advantage of the total  $\tilde{q}$ -excess factors relative to q, likelihood ratio and an explicit formula for KL divergence on MRFs/rMRFs (see Lemma 3.4) in Section 3

as well as in handling of the inherent high-dimensionality of such graphical models and we obtain tight and scalable, information-based bounds on the predictions for QoIs. Finally, we prove tightness of the UQ bounds, i.e. we prove that the bounds are attainable by MRFs/rMRFs, we compute their probability distribution and we develop a strategy to determine their associated graph structures.

**4.1. Quantities of Interest.** We primarily consider two classes of QoIs  $f(\mathbf{Z})$ . The first has QoIs that are expressed as a characteristic function on events of interest such as (5.1) in the medical diagnostics example presented in Section 5. The second class consists of QoIs that are sufficient statistics for the models q and  $\tilde{q}$  and are also present in the total  $\tilde{q}$ -excess factor of type I and II relative to q, i.e. we consider  $f(\mathbf{Z})$  that satisfies

$$f(\mathbf{Z}) = \frac{1}{C_{\rm i}} \left( \log \Phi_{\mathbf{u}}^{\rm i}(\mathbf{Z}) + \kappa_{\rm i}(\mathbf{Z}) \right), \ \ {\rm i} = {\rm I, II}. \label{eq:force_fit}$$

for some non-zero constant  $C_i \equiv C_i(\mathbf{w}, \tilde{\mathbf{w}}, \mathbf{u}) < 1$  and a function  $\kappa_i(\cdot) \equiv \kappa_i(\cdot \mid \mathbf{w}, \tilde{\mathbf{w}}, \mathbf{u})$  that may depend on  $\mathbf{w}, \tilde{\mathbf{w}}, \mathbf{u}$ , see also (3.12). Such a class covers observables involved in finite size effects and phase diagrams for statistical mechanics models examined later (e.g. averages of spins given by (6.12)). The CGF given by (4.2) is computable for QoIs in both classes.

**4.2.** UQ bounds. The next theorem is an UQ result on rMRFs that is obtained by consolidating the total  $\tilde{q}$ -excess factors relative to q, likelihood ratio, KL divergence and QoIs. Part (a) provides the UQ bounds for a general QoI and hence we use such bounds for QoIs examined in the medical diagnostics application in Section 5. Part (b) is particularly applicable for QoIs satisfy (4.3), so they are exploited by the statistical mechanics section.

THEOREM 4.1. Let  $(\mathbf{Y}, p)$ ,  $(\mathbf{Y}, \tilde{p})$  be two MRFs defined over graphs  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  and  $\tilde{\mathcal{G}} = (\mathcal{V}, \tilde{\mathcal{E}})$  respectively. Let  $\mathbf{u}$  be a context and  $\mathcal{M} \subset \mathcal{V}$ . We consider the corresponding rMRFs  $(\mathbf{Z}, q)$ ,  $(\mathbf{Z}, \tilde{q})$ . If  $\tilde{q}$  is of Type i, with i = I or II, then (a) for any QoI  $f(\mathbf{Z})$ , the following bounds hold

$$(4.4) \quad \pm E_{\tilde{q}}[f] \leq \inf_{\lambda > 0} \frac{1}{\lambda} \left\{ \log E_q[e^{\pm \lambda f}] + \frac{1}{E_q[\Phi^{\mathbf{i}}_{\mathbf{u}}]} E_q\left[\Phi^{\mathbf{i}}_{\mathbf{u}} \log \Phi^{\mathbf{i}}_{\mathbf{u}}\right] - \log E_q[\Phi^{\mathbf{i}}_{\mathbf{u}}] \right\}$$

(b) for any QoI  $f(\mathbf{Z})$  that satisfies (4.3), the following bounds hold:

$$(4.5) \qquad \pm E_{\tilde{q}}[f] \leq \frac{1}{1 - C_{\mathbf{i}}} \inf_{\lambda > 0} \frac{1}{\lambda} \left\{ \log E_{q}[e^{\pm \lambda f}] - \log E_{q} \left[\Phi_{\mathbf{u}}^{\mathbf{i}}\right] + \frac{E_{q} \left[\kappa_{i} \Phi_{\mathbf{u}}^{\mathbf{i}}\right]}{E_{q} \left[\Phi_{\mathbf{u}}^{\mathbf{i}}\right]} \right\}$$

where  $\Phi_{\mathbf{u}}^{i}$  is the total  $\tilde{q}$ -excess factor relative to q given by (3.7) and (3.8),  $\kappa_{i}$  and  $C_{i}$  are defined in (4.3). Note that when  $\tilde{q}$  is of Type I,  $\tilde{Z}_{\mathbf{u}}(\tilde{\mathbf{w}}) = Z_{\mathbf{u}}(\tilde{\mathbf{w}})$ .

The proof given in Appendix B.2 is based on Lemma 3.4 and the characterization of the exponential integrals. An application to a single parameter exponential family is given in Appendix B.2.

**4.3. Tightness of UQ bounds for MRFs/rMRFs.** Here, we prove that the inequalities (4.4) and (4.5) are tight i.e. they become an equality for a suitable model  $\tilde{q} \in \mathcal{Q}^{\eta}$  given by (3.1) standing for the worst case scenarios. The practical interpretation of the tightness of UQ bounds is that these distributions are reasonable as they belong to the ambiguity set in (3.1).

THEOREM 4.2. Let  $(\mathbf{Z}, q)$  be a rMRF defined in subsection 2.2.1 and  $f(\mathbf{Z})$  be a QoI with finite MGF  $E_q[e^{\lambda f(\mathbf{Z})}]$  in a neighborhood of the origin. Then there exist  $0 < \eta_{\pm} \leq \infty$  such that for any  $\eta \leq \eta_{\pm}$  there exist probability measures  $q^{\pm} = q^{\pm}(\eta) \in \mathcal{Q}_{\eta}$ , where  $\mathcal{Q}_{\eta}$  is given in (3.1), such that (4.4) and (4.5) become an equality. Furthermore,  $q^{\pm} = q^{\lambda_{\pm}}$  with

$$(4.6) dq^{\lambda \pm} = \frac{e^{\lambda \pm f}}{E_q[e^{\lambda \pm f}]} dq$$

and  $\lambda_{\pm}$  being the unique solutions of  $R(q^{\lambda_{\pm}}||q) = \eta$ . In particular, the total  $q^{\pm}$ -excess factor relative to q denoted by  $\Phi_{\mathbf{u}}^{\pm}$ , satisfies

$$\Phi_{\mathbf{n}}^{\pm} = e^{\lambda_{\pm}f}$$
 and  $C_{i} = \lambda_{\pm}$ ,  $\kappa_{i} = 0$  respectively.

*Proof.* See Appendix B.3.

The result holds also for MRFs. The corresponding quantities involved in the theorem are denoted by  $p, p^{\lambda_{\pm}}$  and  $\Phi^{\pm}$ .

Remark 4.3. For convenience we use its MRF version. Given a baseline MRF  $(\mathbf{Y}, p)$ , its associated graph  $\mathcal{G}$  and a QoI f, Theorem 4.2 guarantees the existence of probability distributions  $p^{\lambda_{\pm}}$  such that (4.4) and (4.5) become an equality (this is not an unlikely extreme case) and also specifies the distributions explicitly. However, it does not imply how different the associated graphs of  $p^{\pm}$  are, compared to the graph associated to p or grossly speaking, if they are Type I or II. Depending on f, there are cases where this can be determined. In fact, by recalling the Hammersley-Clifford Theorem, we express

$$(4.7) dp^{\lambda \pm} = \frac{e^{\lambda \pm f}}{E_p[e^{\lambda \pm f}]} \frac{1}{Z(\mathbf{w})} \prod_{c \in \mathcal{C}_{\mathcal{G}}} \Psi_c(\mathbf{y}_c \mid \mathbf{w}_c) = \frac{1}{Z^{\pm}(\lambda^{\pm}, \mathbf{w})} \prod_{c \in \mathcal{C}_{\mathcal{G}}} e^{\lambda \pm f} \Psi_c(\mathbf{y}_c \mid \mathbf{w}_c)$$

where  $Z^{\pm}(\lambda^{\pm}, \mathbf{w}) = E_p[e^{\lambda_{\pm} f}]Z(\mathbf{w})$  is the partition function of  $p^{\lambda_{\pm}}$ .

We turn our attention to the product in (4.7). Each factor is defined on a maximal clique of  $\mathcal{G}$  apart from  $e^{\lambda_{\pm}f}$ . We focus on f; Suppose that f is a QoI with domain  $\mathrm{Dom}(f)$  and cannot be written as a sum of more than two functions e.g. sample average. If there is a maximal clique  $c_0$  such that  $\mathrm{Dom}(f) \subseteq c_0$ , then it turns out that all clique potentials of  $p^{\lambda_{\pm}}$  and p are equal except  $\tilde{\Psi}_{c_0} = e^{\lambda_{\pm}f}\Psi_{c_0}$ , and hence

(4.8) 
$$dp^{\lambda_{\pm}} = \frac{1}{E_p[e^{\lambda_{\pm}f}]Z(\mathbf{w})} \underbrace{e^{\lambda_{\pm}f}\Psi_{c_0}}_{\tilde{\Psi}_{c_0}} \prod_{c \neq c_0} \Psi_c(\mathbf{y}_c \mid \mathbf{w}_c)$$

The associated graphs of  $p^{\lambda_{\pm}}$  are apparently of Type I as no change on maximal cliques occurs. If  $\mathrm{Dom}(f) \cap c \neq \emptyset$  for more than two maximal cliques c, then the associated graphs to  $p^{\lambda_{\pm}}$  have been changed and thus are Type II. An example is discussed in subsection 5.1. On the other hand, if f can be expressed as a sum of some functions  $f = \sum_i f_i$ , then we may have more than one candidate graphs associated to  $p^{\lambda_{\pm}}$  either Type I or II. In fact, the exponential can be factorized further (e.g.  $e^{\lambda_{\pm}f} = \prod_i e^{\lambda_{\pm}f_i}$ ), giving rise to more than one ways of matching the clique potentials in the sense of (4.8).

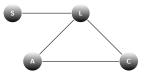
Remark 4.4. The parameter  $\eta$  in Theorem 4.2 is also called misspecification parameter, and can be thought of as a non-parametric "stress test" for the rMRF, and

can be tuned by hand so one can explore how the level of uncertainty affects QoIs. Alternatively,  $\eta$  can be computed as the KL divergence from the available data (e.g data used to construct the baseline model in Medical Diagnostics, Section 5) in the form of a histogram or a KDE and thus subs for the distance of the baseline model from the unknown true model, [30].

5. UQ for Medical Diagnostics. Let us introduce a simple example from medical diagnostics. We exploit its simplicity and low dimensionality to demonstrate MRF modeling with parameters and structure learned from data as well as the types of uncertainties that arise naturally in MRF modeling.

**Setup.** Consider the problem of investigating interdependence (structure) and its strength (parameters) between Smoking (S), Asthma (A), Lung cancer (L), and Cough (C), [20]. It is assumed there are prior expert knowledge and data encoded by a probabilistic model (distribution)  $p^*$  defined on  $\{S, C, L, A\}$ . Due to limitations in expert knowledge and data, the true distribution  $p^*$  itself may be altogether unknown. This, in turn, forces us to build a surrogate baseline model p, which therefore is uncertain in ways we will specify next.

Baseline MRF. Let  $\mathcal{D} = \{\mathbf{d}[1], \dots, \mathbf{d}[N]\}$  be a large collection of patient records sampled from  $p^*$ . Using a structure-learning algorithm on the data  $\mathcal{D}$  (for instance, greedy score-based structure search algorithm for log-linear models [49, 38]), a model with the structure of  $\mathcal{G}$  illustrated in Figure 2, (Left) is built, [20]. We assume that the graph is undirected as the directionality associated with the variable dependencies is not known (or is not expected). Subsequently, by parameter learning (for instance, using maximum likelihood estimation [49]) the weights  $\mathbf{w}$  become specified from the available data. From now on the resulting model ( $\mathcal{G}, \mathbf{w}, p$ ) is called the baseline model.



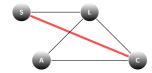


Fig. 2. (Left) MRF structure  $(\mathbf{Y},p) = (\{S,C,L,A\},p)$  over  $\mathcal{G}$  with joint probability distribution p.  $S \in \{s_0,s_1\}$ ,  $L \in \{l_0,l_1\}$ ,  $A \in \{a_0,a_1\}$  and  $C \in \{c_0,c_1\}$ . For example, the values  $s_0$  and  $s_1$  can be thought as smoking and non-smoking respectively, and so forth. The random variables  $\mathbf{Y} = \{Y_1,Y_2,Y_3,Y_4\} = \{S,L,A,C\}$  are accordingly attached to the nodes in  $\mathcal{V} = \{1,\ldots,4\}$  with edges in  $\mathcal{E} = \{1-2,2-3,2-4,3-4\}$ . The class of maximal cliques is  $\mathcal{C}_{\mathcal{G}} = \{\{1,2\},\{2,3,4\}\}$ . (Right) A Type II model  $(\tilde{\mathcal{G}},\tilde{\mathbf{w}},\tilde{p})$  over  $\mathbf{Y} = \{S,C,L,A\}$  with joint probability distribution  $\tilde{p}$ . The associated graph is demonstrated by  $\tilde{\mathcal{G}} = (\mathcal{V},\tilde{\mathcal{E}})$  with  $\tilde{\mathcal{E}} = \mathcal{E} \cup \{1-4\}$ . The new edge is shown in red color.

As in [20], the joint probability distribution could be a log-linear model ([49], Section 4.4) and thanks to Hammersley-Clifford Theorem, is factorized over the maximal cliques with clique potentials  $\Psi_c(\mathbf{y}_c \mid \mathbf{w}_c) = e^{w_c f_c(\mathbf{y}_c)}$ ,  $\mathbf{w} = \{\mathbf{w}_c\}_{c \in \mathcal{C}_{\mathcal{G}}}$ , where  $f_c$  is often called feature.

Alternative models. Both learning steps can induce uncertainties in structure and/or parameters on the baseline. Next, we model and quantify such uncertainties by considering alternative models to the baseline of Type I and II: we focus on graphical models that may have been obtained by learning structure and parameters from either a different data set  $\tilde{\mathcal{D}} = \{\tilde{\mathbf{d}}[1], \dots, \tilde{\mathbf{d}}[\tilde{N}]\}$  or the same data set  $\mathcal{D}$  but with different prior (expert) knowledge. We denote the corresponding alternative models  $(\tilde{\mathcal{G}}, \tilde{\mathbf{w}}, \tilde{p})$ 

and assume they can be also represented by a MRF with  $\tilde{p} > 0$  in the class of loglinear models with clique potentials being given by  $\tilde{\Psi}_c(\mathbf{y}_c) = e^{\tilde{w}_c \tilde{f}_c(\mathbf{y}_c)}$  be the clique potential. We consider the following QoIs defined as:

(5.1) 
$$g(\mathbf{Y}) = \mathbf{1}_A$$
, for any event of interest  $A \subset \Omega$ .

For instance,  $A = \{\text{patient is smoker with asthma}\} = \{\omega = (\omega_1, \omega_2, \omega_3, \omega_4) : \omega_1 = s_0, \omega_3 = a_0\}.$ 

**Type I.** We consider the class of log-linear models  $\tilde{p}$  over  $\mathcal{G}$  with weight change in one maximal clique after learning weights from  $\tilde{\mathcal{D}}$ . Let c be the maximal clique that a weight change occurred. If  $p_{\rm I} \equiv p(B_c)$  and  $a \in [-1,1]$  (depends on  $\tilde{p}$ ), then for any event of interest A, the following holds:

$$(5.2) \pm \tilde{p}(A) \le \inf_{\lambda > 0} \frac{1}{\lambda} \left\{ \log \left( \frac{p(A)e^{\pm \lambda} + 1 - p(A)}{e^{aw_c}p_{\rm I} + 1 - p_{\rm I}} \right) - \frac{aw_c e^{aw_c}p_{\rm I}}{e^{aw_c}p_{\rm I} + 1 - p_{\rm I}} \right\}$$

where  $a \in [-1,1]$  stands for the model uncertainty of alternative models of Type I and  $w_c$  is the weight on c of p. The derivation of the UQ bounds in (5.2) is given in Appendix D, while their demonstration as functions of the uncertainty parameter a for any event of interest A with p(A) = 0.3 and when  $p_I = 0.2$  is given in Figure 3.

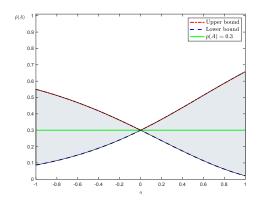
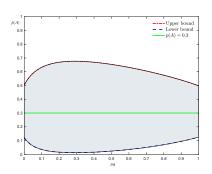


Fig. 3. For any event of interest, A with p(A) = 0.3, the red dashed-dot and the blue dashed curve are the upper bound and lower bound for  $\tilde{p}(A)$  provided in (5.2), computed as functions of the weight change a.

**Type II.** We consider the class of log-linear models  $\tilde{p}$  over  $\tilde{\mathcal{G}}$  with  $\tilde{\mathcal{V}} = \mathcal{V}$ ,  $\tilde{\mathcal{E}} = \mathcal{E} \cup e$ , where the new edge e (e.g see Figure 2, (Right)) enlarges an already existing maximal clique  $\tilde{c}$  in the sense of the analysis in subsection 3.1 after structure-learning from  $\tilde{\mathcal{D}}$ . The model uncertainties lie in the binary function  $\tilde{f}_{\tilde{c}}$  defined on  $\tilde{c}$  and the new weight  $\tilde{\mathbf{w}}_{\tilde{c}}$ . The binary function  $f_{\tilde{c}}$  induces a set  $B_{\tilde{c}} = \{(\omega_1, \omega_2, \omega_3, \omega_4) : \tilde{f}_{\tilde{c}}(\omega_{\tilde{c}}) = 1\}$ . The set  $B_{\tilde{c}}$  satisfies one of the following:  $B_{\tilde{c}} \cap B_c = \emptyset$  or  $B_{\tilde{c}} \cap B_c \neq \emptyset$ . For  $B_{\tilde{c}} \cap B_c = \emptyset$ , if  $p_{\rm I} \equiv p(B_c)$ ,  $p_{\rm II} \equiv p(B_{\tilde{c}})$  and  $a \in \mathbb{R}$ , then for any event of interest A, the following holds:

$$\pm \tilde{p}(A) \leq \inf_{\lambda > 0} \frac{1}{\lambda} \left\{ \log \left( \frac{p(A)e^{\pm \lambda} + 1 - p(A)}{1 - (1 - e^{(1+a)w_c})p_{\text{II}} - (1 - e^{-w_c})p_{\text{I}}} \right) - \frac{w_c e^{-w_c} p_{\text{I}} - (1 + a)w_c e^{(1+a)w_c} p_{\text{II}}}{1 - (1 - e^{(1+a)w_c})p_{\text{II}} - (1 - e^{-w_c})p_{\text{I}}} \right\}$$
(5.3)

The derivation of the UQ bounds in (5.3) is given in Appendix D while their demonstration for any event A with p(A) = 0.3 as functions of the uncertainty parameters a (when  $p_{\rm I} = 0.2$ ,  $w_c = 1.5$  and  $p_{\rm II} = 0.7$ ) and  $p_{\rm II}$  (when  $p_{\rm I} = 0.2$ ,  $w_c = 1.5$ , a = -0.2) is given in Figure 4. Note that the case where  $B_{\tilde{c}} \cap B_c \neq \emptyset$  is more complicated. However, the KL divergence is still explicitly computable (see Remark D.1).



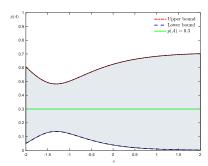


Fig. 4. A is an event of interest with p(A) = 0.3. (Left) For  $p_I = 0.2$ ,  $w_c = 1.5$  and a = -0.2, the red dash-dot and the blue dashed curves are the upper bound and lower bound for  $\tilde{p}(A)$  provided in (5.3), computed as functions of  $p_{II}$ . (Right) For  $p_I = 0.2$ ,  $w_c = 1.5$ ,  $p_{II} = 0.7$ , the red curve and the blue are the upper bound and lower bound for  $\tilde{p}(A)$ , computed as functions of the weight change  $a \in [-2, 2]$ .

- **5.1. Tightness.** Let g be the QoI given by (5.1). By applying Theorem 4.2, there exist probability measures  $p^{\pm} = p^{\pm}(\eta) \in \mathcal{Q}^{\eta}$ , where  $\mathcal{Q}^{\eta}$  is given in (3.1), such that (4.4) becomes an equality and  $p^{\pm} = p^{\lambda_{\pm}}$  are given by  $dp^{\lambda_{\pm}} = \frac{e^{\lambda_{\pm} 1_A}}{p(A)e^{\lambda_{\pm}} + 1 p(A)} dp$  and  $\lambda_{\pm}$  being the unique solutions of  $R(p^{\lambda_{\pm}} || p) = \eta$ . Depending on the event of interest A, we can determine the graph associated with  $p^{\lambda_{\pm}}$ . Specifically, if  $A = \bigcap_i A_i$  where all  $A_i$  are defined on the same maximal clique of  $\mathcal{G}$  given in Figure 2, then the graph associated with  $p^{\lambda_{\pm}}$  is  $\mathcal{G}$  and hence both models are Type I. If at least two  $A_i, A_j$  are defined on different maximal cliques, the associated graphs are different than  $\mathcal{G}$ , e.g. let  $A = \{\text{patient}$  is smoker with asthma $\} = \{\omega = (\omega_1, \omega_2, \omega_3, \omega_4) : \omega_1 = s_0, \omega_3 = a_0\} = \{\omega : \omega_1 = s_0\} \cap \{\omega : \omega_3 = a_0\}$ . Since the total  $p^{\pm}$ -excess factor relative to  $p \Phi^{\pm} = e^{\lambda_{\pm} 1_A}$  cannot be further factorized, the new graph has the same set of nodes with an extra edge 1 3, that is  $\tilde{\mathcal{E}} = \mathcal{E} \cup \{1 3\}$ . In that case, both models are Type II.
- 6. UQ for Statistical Mechanics. Large-scale physical systems of interacting particles such as gases, liquids, and solids, are at the core of statistical mechanics and in particular of equilibrium statistical mechanics. The macroscopic properties of a system can be understood through its underlying microscopic description which fundamentally requires the microscopic states and an interaction between microscopic constituents. Statistical mechanics models such as the Ising model are fundamental in ML, especially energy-based probabilistic models (generally defined as (6.6)) such as Boltzmann machines [38]. Furthermore, methods from equilibrium statistical mechanics combined with information theory can provide first insights into profound cornerstones of deep learning. For example, although we use the KL divergence defined in Lemma 3.4 for UQ, KL between an energy-based model and available data equals to the difference between Gibbs and Helmholtz free energy and is a natural "distance" to use for statistical learning. Note that both UQ and statistical learning can be considered as dual concepts, [9]. A more extensive analysis on these ideas,

and generally on the intersection between statistical mechanics—also including non-equilibrium statistical mechanics—and deep learning have been reviewed in [3].

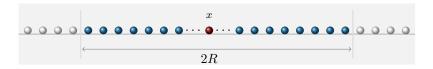


Fig. 5. One-dimensional Ising spin lattice on  $\Delta$  (light gray area with blue, red, and white particles). The spin located at  $x \in \Delta$  (red particle) interacts only with spins located at y in  $B_x(R)$  (blue particles) with strength of interaction J(x,y). The red spin does not interact with the white ones as they are located at distance greater than R from x.

**6.1.** Ising Model. An illustrative example is the Ising model, where the space of all microstates is the collection of all spin configurations on a bounded region  $\Delta \subseteq \mathbb{Z}^d$ :

$$\Omega := \{\pm 1\}^{\Delta} = \left\{ \sigma_{\Delta} = \{\sigma_{\Delta}(x)\}_{x \in \Delta} : \sigma_{\Delta}(x) \in \{+1, -1\} \right\}$$

as in Figure 5, [57, 49]. An interaction between spins can be short, long range or a combination (such as Lennard-Jones potential, [58]), positive (ferromagnetism), etc, [57, 31, 34]. Here we consider a d-dimensional Ising spin system on  $\Delta$  with a generic interaction  $\mathbf{J} = \{J(x,y) : x,y \in \Delta\}$  satisfying three properties: for all  $x,y \in \Delta$  and  $z \in \mathbb{R}^d$ 

(6.1) 
$$J(x+z, y+z) = J(x, y)$$
 (translational invariance)

(6.2) 
$$J(x,y) = J(y,x)$$
 (symmetry)

(6.3) 
$$\sum_{x \neq 0} |J(0,x)| < \infty \qquad \text{(summability)}$$

and an external field,  $h \in \mathbb{R}$ . Let R > 0 be the length of the range of interaction. For  $x \in \mathbb{Z}^d$ ,  $B_x(R) = \{y \in \mathbb{Z}^d : ||x - y||_d \le R\}$  is the set of all spins that the spin located at x interacts with and  $||x - y||_d := \sqrt{\sum_{i=1}^d |x_i - y_i|^2}$ . For convenience, we denote  $B_{x,R}^{\neq} := B_x(R) \setminus x$ .

**6.1.1.** Boundary conditions. Boundary conditions are a fundamental concept in statistical mechanics, [62]. For simplicity, let us assume that  $\Delta$  is a hypercube. We consider a system where particles not only interact with particles in  $\Delta$ , but also with particles "outside" of  $\Delta$ . Let  $\bar{\sigma}_{\Delta^c}$  be a given fixed configuration of spins on the complement of  $\Delta$  denoted by  $\Delta^c$ , see Figure 10. The Hamiltonian energy of the system is given by:

(6.4) 
$$H^{\mathbf{J},h}(\sigma_{\Delta}|\bar{\sigma}_{\Delta^c}) = H^{\mathbf{J},h}(\sigma_{\Delta}) - \sum_{x \in \Delta} \sum_{y \in \Delta^c} J(x,y)\sigma_{\Delta}(x)\sigma_{\Delta}(y)$$

where

(6.5) 
$$H^{\mathbf{J},h}(\sigma_{\Delta}) = -\frac{1}{2} \sum_{x \in \Delta} \sum_{y \in \Delta} J(x,y) \sigma_{\Delta}(x) \sigma_{\Delta}(y) - h \sum_{x \in \Delta} \sigma_{\Delta}(x)$$

The Gibbs measure with boundary condition  $\bar{\sigma}_{\Delta^c}$  is defined as

(6.6) 
$$\mu_{\mathbf{J},\beta,h}^{\Delta}(\sigma_{\Delta} \mid \bar{\sigma}_{\Delta^{c}}) = \frac{1}{Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J},\beta,h)} e^{-\beta H^{\mathbf{J},h}(\sigma_{\Delta}|\bar{\sigma}_{\Delta^{c}})}.$$

where  $Z_{\bar{\sigma}_{\Delta^c}}(\mathbf{J}, \beta, h) = \sum_{\sigma_{\Delta}} e^{-\beta H^{\mathbf{J}, h}(\sigma_{\Delta}|\bar{\sigma}_{\Delta^c})}$  is the partition function.

**6.1.2. rMRF** formulation. A system with configuration as boundary conditions does not admit a MRF description. So, we describe the system using rMRFs. The set of nodes is  $\mathbb{Z}^d$ , the set of edges can be constructed by looking at all (x,y) such that  $||x-y||_d \leq R$  and the context is  $\mathbf{u} = \bar{\sigma}_{\Delta^c}$ , which corresponds to a fixed boundary condition. Then  $(\sigma_{\Delta}, \mu_{\mathbf{J},\beta,h}^{\Delta}(\cdot \mid \bar{\sigma}_{\Delta^c}))$  is a rMRF with maximal cliques  $c_x = \{y \in \Delta : y \in B_{x,R}^{\neq}\}$  (spins in  $c_x$  interact with all spins in  $c_x$ ). Let  $\mathbf{w} = \{\mathbf{w}_{c_x}\}_{x \in \Delta}$  with  $\mathbf{w}_{c_x} = \{J_{c_x}, \beta, h\}$  and  $\mathbf{J}_{c_x} = \{J_{(x,y)} : y \in c_x\}$ . We express each clique potential as

(6.7) 
$$\Psi_{c_x} = \exp \left\{ \beta \sigma_{\Delta}(x) \left( h + \frac{1}{2} \sum_{\substack{y \in \Delta \\ y \in B_{x,R}^{\neq}}} J(x,y) \sigma_{\Delta}(y) + \sum_{\substack{y \in \Delta^c \\ y \in B_{x,R}^{\neq}}} J(x,y) \bar{\sigma}_{\Delta^c}(y) \right) \right\}$$

Note that we may resume the full notation when we needed, that is  $\Psi_{c_x} \equiv \Psi_{c_x}[\bar{\sigma}_{\Delta^c}](\sigma_{c_x} \mid \mathbf{w}_{c_x})$  where  $\sigma_{c_x}$  is the Ising spin configuration defined on all  $y \in c_x$ .

### 6.2. UQ Formulation.

**6.2.1.** Alternative models. We consider models on a lattice with perturbed interaction in the strength (Type I) and/or range (Type II) such as truncated or long range interaction. Given **J** as in subsection 6.1, an interaction F(x,y) satisfying (6.1)-(6.3) with length of range  $R_F$ , we say that  $\tilde{\mathbf{J}}^{\mathbf{F}} = {\tilde{J}^F(x,y) : x,y \in \mathbb{Z}^d}$  is a perturbed interaction if

(6.8) 
$$\tilde{J}^F(x,y) = J(x,y)\mathbf{1}_{\|x-y\|_d \le R} + F(x,y)\mathbf{1}_{\|x-y\|_d \le R_F} + F(x,y)\mathbf{1}_{\|x-y\|_d > R_F}$$

We say that a perturbed interaction is Type I iff

(6.9) 
$$R = R_F \text{ and supp}(F) = \{(x, y) : ||x - y||_d \le R_F\}.$$

We say that a perturbed interaction is Type II iff

(6.10) 
$$R = R_F \text{ and supp}(F) = \{(x, y) : ||x - y||_d > R_F\}.$$

The rMRF formulation of the system with  $\tilde{\mathbf{J}}^{\mathbf{F}}$  goes similarly as in subsection 6.1.2. Note that the graph representation simplifies a possible complexity of J, F and  $\tilde{J}^F$  as we connect nodes x, y according to the range of J, F and  $\tilde{J}^F$  and assign the corresponding strengths J(x, y), F(x, y) and  $\tilde{J}^F(x, y)$ .

## **6.2.2.** Total $\tilde{q}_{\Delta}$ -excess factor relative to $q_{\Delta}$ .

LEMMA 6.1. Let  $\tilde{\mathbf{J}}^{\mathbf{F}}$  be defined in subsection 6.2.1 with support given by (6.9) or (6.10), and  $q_{\Delta}(\cdot) := \mu_{\mathbf{J},\beta,h}^{\Delta}(\cdot \mid \bar{\sigma}_{\Delta^c})$ ,  $\tilde{q}_{\Delta}(\cdot) := \mu_{\tilde{\mathbf{J}}^{\mathbf{F}},\beta,\tilde{h}}^{\Delta}(\cdot \mid \bar{\sigma}_{\Delta^c})$  be the corresponding Gibbs measures defined in (6.6). The total  $q_{\Delta}$ -excess factor for i = I, II is given by

$$\Phi^{\mathbf{i}}_{\bar{\sigma}_{\Delta^{c}}}(\sigma_{\Delta}) = \exp\left\{\beta \sum_{x \in \Delta} \sigma_{\Delta}(x) \left( (\tilde{h} - h) + \frac{1}{2} \sum_{y \in A^{\mathbf{i}}_{x} \cap \Delta} F(x, y) \sigma_{\Delta}(y) + \sum_{y \in A^{\mathbf{i}}_{x} \cap \Delta^{c}} F(x, y) \bar{\sigma}_{\Delta^{c}}(y) \right) \right\}$$
(6.11)

where for each  $x \in \Delta$ ,  $A_x^{\mathrm{I}} = B_x(R)$  and  $A_x^{\mathrm{II}} = B_x(R)^c$ , with  $B_x(R)^c$  being the complement of  $B_x(R)$ .

The proof is straightforward (see Appendix E.2). Both  $(\tilde{h} - h)$  and F(x, y) in the total  $\tilde{q}_{\Delta}$ -excess factor relative to  $q_{\Delta}$  point out how different the external fields and interactions are respectively, as the latter satisfies  $F(x, y) = \tilde{J}^F(x, y) - J(x, y)$ .

**6.2.3.** Quantities of Interest. The use of phase diagrams is central in physics and material science. A phase diagram is defined as a graphical representation of equilibrium states under different thermodynamic parameters such as external field h, temperature T and pressure P. It is typically computed in the thermodynamic limit (i.e a limiting process with  $\Delta \nearrow \mathbb{Z}^d$  such that the ratio between inter-atomic distances and macroscopic lengths vanishes), [57]. Equilibrium states are characterized by order parameters such as magnetization. For that, we consider the following observable

(6.12) 
$$m(\sigma_{\Delta}) := \frac{1}{|\Delta|} \sum_{x \in \Delta} \sigma_{\Delta}(x)$$

where  $|\Delta|$  stands for the volume of a hypercube  $\Delta \subset \mathbb{Z}^d$ . As  $\Delta$  invades the whole  $\mathbb{Z}^d$ , the expectation of  $m(\sigma_{\Delta})$  yields the magnetization. Other QoIs could also be considered e.g. correlation functions  $v(\sigma_{\Delta}) = \frac{1}{|\Delta|^2} \sum_{x \in \Delta} \sum_{y \in \Delta} \sigma_{\Delta}(x) \sigma_{\Delta}(y)$ .

**6.2.4.** Cumulant Generating Function. Let  $\Delta$  be a hypercube in  $\mathbb{Z}^d$ . Given a configuration  $\bar{\sigma}_{\Delta^c}$ , the baseline model is an Ising model with interaction **J** defined in subsection 6.1. We compute the cumulant generating function defined by (4.2) w.r.t the baseline model  $q_{\Delta}$  (the computation is given in (E.3)):

$$(6.13) \qquad \Lambda_{q_{\Delta};|\Delta|m(\sigma_{\Delta})}(\pm\lambda) = \beta|\Delta| \left( P_{h\pm\frac{\lambda}{\beta},\beta,\mathbf{J}}^{\Delta}(\bar{\sigma}_{\Delta^{c}}) - P_{h,\beta,\mathbf{J}}^{\Delta}(\bar{\sigma}_{\Delta^{c}}) \right)$$

where  $P_{h,\beta,\mathbf{J}}^{\Delta}$  stands for the thermodynamic pressure, [57], defined as

$$P^{\Delta}_{h,\beta,\mathbf{J}}(\bar{\sigma}_{\Delta^c}) := \frac{Z(\mathbf{J},\beta,h,\bar{\sigma}_{\Delta^c})}{\beta|\Delta|}.$$

**6.2.5.** KL Divergence. Here we utilize Lemma 3.4 and specify the KL divergence in terms of  $\kappa_i$  and  $\Phi_{\mathbf{u}}$  as involved in (3.12) when the alternative models are Ising models with a perturbed interaction  $\tilde{\mathbf{J}}^{\mathbf{F}}$  defined in subsection 6.2.1. Then we bound it by using Lemma 6.3. Before that, we use a well-established tool in statistical mechanics referred to as norm- $\|\cdot\|_1$ , [62] to alternatively bound the KL divergence. After all, we conclude that our UQ approach gives a narrower area (i.e the area between the upper and lower UQ bound) provided by Theorem 4.1 and thus smaller uncertainty, see Figure 6.

uncertainty, see Figure 6. Norm- $\|\cdot\|_1$ : Let  $\Phi_{\Delta,\sigma_{\Delta^c}}^{h,\beta,\mathbf{J}}(\sigma_X)$  be the following quantity: (6.14)

$$\Phi_{\Delta,\bar{\sigma}_{\Delta^{c}}}^{h,\beta,\mathbf{J}}(\sigma_{X}) = \begin{cases} -\frac{1}{2}\beta J(x,y)\sigma_{\Delta}(x)\sigma_{\Delta}(y) &, X = \{x,y\}, \ x \neq y, \\ -\beta\sigma_{\Delta}(x)\left(h + \sum_{y \in B_{x,R}^{\neq} \cap \Delta^{c}} J(x,y)\bar{\sigma}_{\Delta^{c}}(y)\right) &, X = \{x\} \\ 0 &, \text{otherwise} \end{cases}$$

and similarly we define  $\Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{\tilde{h},\beta,\tilde{\mathbf{J}}^{\mathbf{F}}}(\sigma_X)$ . Then,

(6.15) 
$$\beta H^{\mathbf{J},h}(\sigma_{\Delta}|\bar{\sigma}_{\Delta^c}) = \sum_{X:X\cap\Delta\neq\emptyset} \Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{h,\beta,\mathbf{J}}(\sigma_X)$$

Also,  $\beta H^{\tilde{\mathbf{J}}^{\mathbf{F}},\tilde{h}}(\sigma_{\Delta}|\bar{\sigma}_{\Delta^c})$  is defined similarly. Then the norm- $||\cdot||_1$  of  $\Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{h,\beta,\mathbf{J}} - \Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{\tilde{h},\beta,\tilde{\mathbf{J}}^{\mathbf{F}}}$  is defined as

where  $\|\Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{h,\beta,\mathbf{J}} - \Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{\tilde{h},\beta,\tilde{\mathbf{J}}^{\mathbf{F}}}\|_{\infty} = \sup_{\sigma_X} |\Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{h,\beta,\mathbf{J}}(\sigma_X) - \Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{\tilde{h},\beta,\tilde{\mathbf{J}}^{\mathbf{F}}}(\sigma_X)| \text{ for } X \subset \mathbb{Z}^d.$ 

LEMMA 6.2. Let F be an interaction satisfying (6.1)-(6.3) with support given by (6.9) or (6.10), then

$$R(\tilde{q}_{\Delta} || q_{\Delta}) \leq 2|\Delta| \|\Phi_{\Delta,\tilde{\sigma}_{\Delta^c}}^{h,\beta,\mathbf{J}} - \Phi_{\Delta,\tilde{\sigma}_{\Delta^c}}^{\tilde{h},\beta,\tilde{\mathbf{J}}^{\mathbf{F}}}\|_{1} \leq 2\beta|\Delta| \left( |\tilde{h} - h| + \sum_{x \neq 0} |F(0,x)| \right).$$

Let us turn to our approach developed in Section 4. We recall the total  $\tilde{q}_{\Delta}$ -excess factor relative to  $q_{\Delta}$  from subsection 6.2.2 as well as the quantities from Section 4.1, and we express  $\log \Phi_{\bar{\sigma}_{\Delta^c}}^{i}(\sigma_{\Delta}) = C_i |\Delta| m(\sigma_{\Delta}) + \kappa_i(\sigma_{\Delta})$  with (6.17)

$$C_{\mathrm{I}} = \beta(\tilde{h} - h) < 1, \quad \kappa_{\mathrm{I}}(\sigma_{\Delta}) = \beta \sum_{x \in \Delta} \sigma_{\Delta}(x) \left( \frac{1}{2} \sum_{y \in A^{1} \cap \Delta} F(x, y) \sigma_{\Delta}(y) \right) + \beta F(\Delta | \bar{\sigma}_{\Delta^{c}})$$

where  $F(\Delta|\bar{\sigma}_{\Delta^c}) = \sum_{x \in \Delta} \sum_{y \in A^i_{\sigma} \cap \Delta^c} F(x,y) \bar{\sigma}_{\Delta^c}(y)$ . We bound  $\kappa_{\rm I}(\sigma_{\Delta})$  as

(6.18) 
$$0 \le \kappa_{\mathrm{I}}(\sigma_{\Delta}) \le \beta |\Delta| \left(\frac{1}{2} + 2R \frac{|\partial \Delta|}{|\Delta|}\right) \sum_{x \ne 0} |F(x, y)|$$

where we use the next lemma.

LEMMA 6.3. Let L and  $\partial \Delta$  be the side and the boundary of the hypercube  $\Delta$  respectively with  $L >> R_F$ . Then, for any interaction  $\mathbf{F} = \{F(x,y) : x,y \in \mathbb{Z}^d\}$  satisfying (6.1)-(6.3) and range  $R_F$ , the following holds:

(i) If the support of F is given by (6.9), then

$$\sum_{x \in \Delta} \sum_{\substack{y \in \Delta^c \\ y \in B_{x,R_F}^{\neq}}} F(x,y) \le R_F |\partial \Delta| \sum_{x \ne 0} |F(0,x)|.$$

(ii) If the support of F is given by (6.10), then

$$\sum_{x \in \Delta} \sum_{y \in \Delta^c} F(x, y) \le R_F |\Delta| \sum_{x \ne 0} |F(0, x)|$$

*Proof.* The bounds are straightforward once we split the sum as follows:

$$\sum_{x \in \Delta} \sum_{\substack{y \in \Delta^c \\ y \in B_{\mathcal{F}_{R,\Gamma}}^{\mathcal{F}}}} F(x,y) = \sum_{\substack{x \in \Delta \\ dist(x,\Delta^c) \leq R_F}} \sum_{\substack{y \in \Delta^c \\ y \in B_{\mathcal{F}_{R,\Gamma}}^{\mathcal{F}}}} F(x,y) + \sum_{\substack{x \in \Delta \\ dist(x,\Delta^c) > R_F}} \sum_{\substack{y \in \Delta^c \\ dist(x,\Delta^c) > R_F}} F(x,y) \leq R_F |\Delta|.$$

where  $dist(x, \Delta^c) = \inf\{\|x - y\| : y \in \Delta^c\}$ . Note that when  $L \ll R_F$ , both (i) and (ii) are bounded by  $R_F |\Delta| \sum_{x \neq 0} |F(0, x)|$ .

6.3. UQ for finite-size effects and boundary conditions. Having computed all the ingredients needed for the analysis of subsections 3.2, 4.1 and 4.2 under the above statistical mechanics formulation through rMRFs, we capture the behavior of  $m(\sigma_{\Delta})$  given in (6.12) with respect to the perturbed model  $\tilde{q}_{\Delta}$ . The analysis from now on refers to models of Type I. Although Type II models can be worked on similarly, one example of Type II is discussed in Appendix F. To get the UQ bounds for  $E_{\tilde{q}_{\Delta}}[m(\sigma_{\Delta})]$ , for  $f(\mathbf{Z}) = |\Delta|m(\sigma_{\Delta})$  we can either apply (4.1) using the crude bound in Lemma 6.2:

$$(6.19) \pm E_{\tilde{q}_{\Delta}}[m(\sigma_{\Delta})] \leq \inf_{\lambda > 0} \left\{ \frac{P_{h \pm \frac{\lambda}{\beta}, \beta, \mathbf{J}}^{\Delta} - P_{h, \beta, \mathbf{J}}^{\Delta}}{\lambda/\beta} + 2\frac{\beta}{\lambda}(|\tilde{h} - h| + \mathcal{F}) \right\}$$

or Theorem 4.1:

$$(6.20) \pm E_{\tilde{q}_{\Delta}}[m(\sigma_{\Delta})] \leq \frac{1}{1 - \beta(\tilde{h} - h)} \inf_{\lambda > 0} \left\{ \frac{P_{h \pm \frac{\lambda}{\beta}, \beta, \mathbf{J}}^{\Delta} - P_{h, \beta, \mathbf{J}}^{\Delta}}{\lambda/\beta} + \frac{\beta}{\lambda} \mathcal{F} \left( 1 + R_F \frac{|\partial \Delta|}{|\Delta|} \right) \right\}$$

with  $\partial \Delta$  being the boundary of the hypercube  $\Delta$  and  $\mathcal{F} := \sum_{x \neq 0} |F(0,x)|$  which is bounded due to the property (6.3) and  $R_F = R$ .

Furthermore, inequality (6.20) implies a new UQ formula for systems with a fixed configuration outside of the domain that here is considered as a Dirichlet-type boundary condition. In particular it allows us to quantify the effect of the boundary conditions on  $\partial \Delta$  on the QoIs, as can be seen more clearly when  $\tilde{h} = h$ . Note, the term  $\frac{|\partial \Delta|}{|\Delta|}$  in (6.20) comes from a more careful bound on the KL divergence using Lemma 6.3 while this term has been eliminated in (6.19) due to the relative crudeness of the bound of KL in Lemma 6.2, see also Fig. 6.

**6.4.** UQ for Phase Diagrams. Here we capture the phase diagram of the perturbed model  $\tilde{q}_{\Delta}$  looking at the magnetization defined in subsection 6.2.3. We study the limit of the bounds obtained in subsection 6.3. The high-dimensionality of statistical mechanics models requires scalable bounds at the thermodynamic limit. In fact, the MGF and the KL divergence scale correctly with the size of the system  $|\Delta|$  (all are multiplied by  $|\Delta|$  see (6.13), Lemma 6.2 and (6.18)). Let  $M(\tilde{\mathbf{J}}^{\mathbf{F}}, \tilde{\beta}, \tilde{h})$  be the limit as  $\Delta \nearrow \mathbb{Z}^d$  of  $E_{\tilde{q}_{\Delta}}[m(\sigma_{\Delta})]$ . Then the limit  $\Delta \nearrow \mathbb{Z}^d$  of (6.20):

(6.21) 
$$\pm M(\tilde{\mathbf{J}}^F, \beta, \tilde{h}) \le \frac{1}{1 - \beta(\tilde{h} - h)} \inf_{\lambda > 0} \left\{ \frac{\left(P_{h \pm \frac{\lambda}{\beta}, \beta, \mathbf{J}} - P_{h, \beta, \mathbf{J}}\right)}{\lambda/\beta} + \frac{\beta}{\lambda} \mathcal{F} \right\}$$

with  $\lim_{\Delta \nearrow \mathbb{Z}^d} P_{h,\beta,\mathbf{J}}^{\Delta} = P_{h,\beta,\mathbf{J}}$  by Theorem 2.3.3.1 in [57] and  $\lim_{\Delta \nearrow \mathbb{Z}^d} \frac{|\partial \Delta|}{|\Delta|} = 0$ , while in the limit of (6.19) the thermodynamic pressure is only replaced by its limit  $P_{h,\beta,\mathbf{J}}$ . The bounds for the  $\tilde{\beta} \neq \beta$  can be adjusted similarly.

**6.5.** Ising-Kac Model. Here we consider an Ising-spin model with a Kac-type interaction as a baseline model. Such a model combines sufficient complexity—since it is not a mean field model—but it is still analytically fairly tractable to serve as a good benchmark problem for high-dimensional rMRF. We illustrate the uncertainty area of the phase diagram for both (6.21) and the limit of (6.19). when the alternative models are a Kac perturbation and a truncated Kac interaction.

An Ising-spin model with a Kac-type interaction behaves like a mean field (or Van der Waals model in gas lattice) in the limit with the convexity of free energy emerging naturally in the limit, contrary to mean field or Curie-Weiss models where Maxwell's equal area law is required to refine the non-convex free energy (double well shape), [57]. Such a discrepancy comes from the fact that each spin interacts with all particles in the same way and independently. The idea of Kac was to keep such a picture on large regions but relatively small compared to the range of interaction. Then, the thermodynamically incorrect of the free energy (i.e. the non-convex free energy) on these large regions looks refined at the scale of interaction. Therefore, the system contains a two-scale behavior that was carried out by introducing a small parameter  $\gamma > 0$  known as  $Kac\ scaling$ . As we suppose that an Ising spin model is endowed by such an interaction, the model has overall three scales: the lattice spacing is 1, the range of interaction is  $\gamma^{-1}$  while the size of the system is much larger than  $\gamma^{-1}$  and all are well-separated, contrary to the mean field model where the range of interaction is the same as the size of the system. Next, we formally introduce the model.

**6.6.** Mathematical Background of Ising-Kac Model. A Kac-type interaction is defined as

(6.22) 
$$J_{\gamma}(x,y) = \gamma^{d} J(\gamma x, \gamma y), \quad x, y \in \mathbb{Z}^{d}$$

where  $\gamma$  is a positive parameter sufficiently small and J is a non-negative (ferromagnetic interaction), even, symmetric function (i.e J(r,r')=J(r',r) for every  $r,r'\in\mathbb{R}^d$ ), translational invariant (i.e J(r,r')=J(r'+a,r+a) for every  $r,r'\in\mathbb{R}^d$  and  $a\in\mathbb{R}^d$ ) function such that J(r)=0 for all |r|>1,  $\int_{\mathbb{R}^d}J(r)dr=\mathcal{J}$  and  $J\in C^2(\mathbb{R}^d)$ . The use of  $\mathbf{J}_\gamma$  stands for the collection of  $J_\gamma(x,y)$ , that is  $\mathbf{J}_\gamma=\{J_\gamma(x,y)\}_{\mathbb{Z}^d\times\mathbb{Z}^d}$ . As  $\gamma$  becomes smaller, more particles are included in a spin neighborhood with  $\gamma^{-1}$  diameter and while the strength of the interactions becomes weaker.

Let  $\Delta$  be a bounded,  $\mathcal{P}_{\mathbb{R}^d}^{(l)}$ -measurable region, with  $L >> \gamma^{-1}$  (see Appendix C.1),  $\beta > 0$  be the inverse temperature,  $h \in \mathbb{R}$  be the external magnetic field and  $\bar{\sigma}_{\Delta^c}$  be a given configuration on its complement (see Figure 5 with  $R = \gamma^{-1}$ ).

**Hamiltonian energy.** The Hamiltonian energy of a spin configuration  $\sigma_{\Delta}$  given  $\bar{\sigma}_{\Delta^c}$ :

$$H_{\gamma}^{\mathbf{J},h}(\sigma_{\Delta} \mid \bar{\sigma}_{\Delta^{c}}) = -\frac{1}{2} \sum_{x \neq y \in \Delta} J_{\gamma}(x,y) \sigma_{\Delta}(x) \sigma_{\Delta}(y) - \sum_{\substack{x \in \Delta, \\ y \in \Delta^{c}}} J_{\gamma}(x,y) \sigma_{\Delta}(x) \bar{\sigma}_{\Delta^{c}}(y)$$

$$(6.23) \qquad -h \sum_{x \in \Delta} \sigma_{\Delta}(x), \qquad \text{Hamiltonian enery}.$$

Finite volume Gibbs measure. The Gibbs measure given a fixed boundary condition  $\bar{\sigma}_{\Delta^c}$  is defined as follows:

$$\mu_{\mathbf{J},\beta,h}^{\Delta,\gamma}(\cdot\mid\bar{\sigma}_{\Delta^{c}}) = \frac{1}{Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J},\beta,h)}e^{-\beta H_{\gamma}^{\mathbf{J},h}(\sigma_{\Delta};\bar{\sigma}_{\Delta^{c}})}, \qquad \text{finite volume Gibbs measure}$$

where  $Z_{\bar{\sigma}_{\Delta^c}}(\mathbf{J}, \beta, h)$  is the normalization (partition function). To simplify the notation, we shall often drop  $\gamma$  and the given configuration in the complement of  $\Delta$  from the Gibbs measure, resuming the full notation when needed, and therefore we write  $\mu_{\beta,\Delta}^{J,h} \equiv \mu_{\beta,\Delta,\gamma}^{\bar{\sigma}_{\Delta^c},J,h}$ .

**Thermodynamic pressure.** The thermodynamic pressure for the Ising-Kac model denoted by  $P_{\mathbf{J},\beta,h}^{\Delta,\gamma}$  is defined as

(6.25) 
$$P_{\mathbf{J},\beta,h}^{\Delta,\gamma}(\bar{\sigma}_{\Delta^{c}}) := \frac{\log Z_{\bar{\sigma}_{\mathbf{I}^{c}}}(\mathbf{J},\beta,h)}{\beta|\Delta|}$$

Its Lebowitz-Penrose (LP) limit (i.e  $\lim_{\gamma \to 0} \lim_{\Delta \nearrow \mathbb{Z}^d}$ )  $p_{\mathbf{J},\beta,h}$  is given by (6.26)

$$p_{\mathbf{J},\beta,h} := -\inf_{m \in [-1,1]} \{-hm + \phi_{\mathbf{J},\beta,0}(m)\}, \qquad \phi_{\mathbf{J},\beta,h}(m) := \left\{-\frac{\mathcal{J}}{2}m^2 - hm\right\} - \frac{1}{\beta}I(m)$$

(see also Appendix C.4 for further discussion). The rMRF formulation of such a model and its perturbations considered next is structured analogously to the ones in subsection 6.1.2 and for that we omit it.

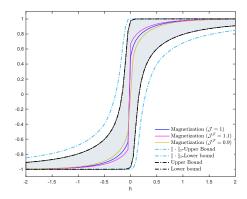


Fig. 6. The curves in blue, magenta and dark yellow color are the magnetizations of the Ising model with Kac interaction at inverse temperature  $\beta = \tilde{\beta} = 1.1$ ,  $\tilde{h} = h$ , and total strength  $\mathcal{J} = 1$ ,  $\tilde{\mathcal{J}}^F = 1.1$  (a = 0.1) and 0.9 (a = -0.1) (validation) respectively. The black dashed-dot curves are the UQ upper and lower bounds provided by Corollary 6.5 and viewed as functions of  $h \in [-2,2]$ . The gray area depicts the size of the uncertainty region. The light blue dashed-dot curves are the UQ upper and lower bounds obtained using norm- $\|\cdot\|_1$ . The uncertainty area of the phase diagram in grey color is significantly better than the uncertainty area between the light blue dashed-dot curves. This comes from the fact that the difference between the limit of (6.19) and (6.21) lies on the term  $\frac{\beta}{\lambda}\mathcal{F}$  which is multiplied by 2.

# **6.6.1. Phase Diagram of Perturbed Kac model.** Let define a perturbation of a Kac potential.

Definition 6.4. Let  $F_{\gamma}$  be an even function satisfying (6.1)-(6.3) and (6.22) with length of range  $\gamma^{-1}$  and  $\mathcal{F} := \int_{\mathbb{R}^d} F(r) dr$ . We define

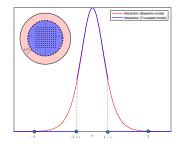
(6.27) 
$$\tilde{J}_{\gamma}^{F}(x,y) = J_{\gamma}(x,y) + F_{\gamma}(x,y), \text{ such that } \mathcal{F} = a\mathcal{J}, \quad a \in [-1,1]$$

The parameter a represents the percentage of increase or decrease of the total strength of interaction  $\tilde{\mathcal{J}}^F := \int_{\mathbb{R}^d} \tilde{J}^F(r) dr = (1+a)\mathcal{J}$ .

COROLLARY 6.5. Let  $\tilde{J}^F$  be the interaction given in Definition 6.4. Then, for  $\gamma > 0$  small enough, the UQ bounds (6.19) and (6.20) hold for  $R_F = R = \gamma^{-1}$  and  $\mathcal{F} = |a|\mathcal{J}$ . The thermodynamic pressure  $P_{\mathbf{J},\beta,h}^{\Delta,\gamma}$  is given in (6.25). Let  $M(\tilde{\mathbf{J}}^F,\beta,\tilde{h})$  be the LP-limit of  $E_{\tilde{q}_{\Delta}}[m(\sigma_{\Delta})]$ . Then, the UQ bounds (6.21) and LP-limit of (6.19) hold with the LP-limit of  $P_{\mathbf{J},\beta,h}^{\Delta,\gamma}$  being  $p_{\mathbf{J},\beta,h}$  given in (6.26).

Remark 6.6. (6.19) represents crude bounds as norm- $\|\cdot\|_1$  (subsection 6.2.5) has been used, while (6.20) obtained by Theorem 4.1, includes more detail. The difference is illustrated in Figure 6. Furthermore, even if there is a  $\gamma^{-1}$  in the term  $2\gamma^{-1}\frac{|\partial\Delta|}{|\Delta|}$  in (6.19), the order of the LP-limit makes it vanish as  $L \to \infty$ .

Validation. Given  $\beta$ , h,  $\mathcal{J}$  and a tolerance  $\eta > 0$ , we can construct with the use of norm- $\|\cdot\|_1$  and Lemma 6.2 a class of models such that  $\mathcal{Q}^{\mathrm{I}}_{\eta} := \{\tilde{q}_{\Delta} : 2\beta a \mathcal{J} \leq \eta\}$ . This is subclass of  $\mathcal{Q}^{\eta}$  defined in (1.1) with the KL divergence in place of d. In Figure 6,  $\beta = 1.1$  and  $\mathcal{J} = 1$  while the external field h varies from -2 and 2. The positive parameter  $\eta = 0.1$  and the perturbed model with 10% decrease (a = -0.1) of the total strength (magnetization in magenta color) is in  $\mathcal{Q}^{\mathrm{I}}_{0.1}$  as demonstrated in dark yellow color.



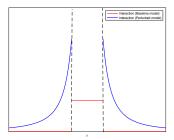


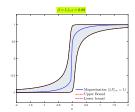
Fig. 7. (Left) The red curve is a Kac interaction and the blue curve is a truncation of it. The two curves coincide at all r with  $|r| \leq 1 - \epsilon$ . The embedded picture demonstrates the two interactions at the microscopic level. The red particle located at the site  $x \in \Delta \subset \mathbb{Z}^2$  interacts with the particles in the blue and the light red through  $J_{\gamma}$ . The particle interacts only with the particles in the blue area through  $\tilde{J}_{\gamma}^{-J}$  with range  $\gamma^{-1}(1-\epsilon)$ . (Right) The red curve is an example of Kac interaction (piecewise constant) with  $J(r) = \mathbf{1}_{r \leq \frac{1}{2}}(r)$  and the blue curve is a perturbation given by  $G(r) = \frac{a}{r^2} \mathbf{1}_{r > \frac{1}{2}}(r)$ , for some a > 0.

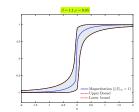
**6.6.2.** Phase diagram of Truncated Potential. From a computational point of view, macroscopic properties of high dimensional systems can be studied through simulation models where one can consider an appropriate truncated interaction which can reduce the computational overhead associated with the interaction [68, Chapter 3]. In our context, a truncated interaction can be thought of as: The support of the interaction J is [-1,1] as in Fig. 7. J is cut off at  $1-\epsilon$  and  $-1+\epsilon$  for some parameter  $\epsilon \in [-1,1]$ . Then the resulting interaction is called truncated interaction of J and its support is  $[-1+\epsilon,1-\epsilon]$  of length  $2\epsilon$ . The introduced parameter  $\epsilon$  quantifies the impact of the truncation of the interaction J. Moreover, Fig 8 quantifies how the uncertainty area becomes smaller as  $\epsilon$  becomes smaller (and hence the truncated interaction tends to be the original J). We mathematically define such an interaction as follows:

Definition 6.7. Let  $0 < \epsilon < 1$ . We define the truncated interaction as

(6.28) 
$$\tilde{J}^{-J}(0,r) = \begin{cases} J(0,r) &, |r| \le 1 - \epsilon \\ 0 &, otherwise \end{cases}$$

The truncated model can be viewed as Type II. However, to be consistent with the assumption  $\mathcal{E} \subset \tilde{\mathcal{E}}$  in Definition 3.1, we view it as perturbed interaction of Type I arising from the subtraction of J (also explains the notation  $\tilde{J}^{-J}$  in (6.28)) on regions of radius greater than  $1 - \epsilon$  as illustrated in Figure 7.





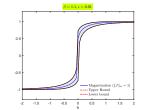


Fig. 8. The three graphs demonstrate the uncertainty area in gray color for different values of  $\epsilon$ . In all graphs, the blue solid line is the magnetization of d-sing model with Kac interaction at inverse temperature  $\beta=1.1$ ,  $\|J\|_{\infty}=1$  and  $\tilde{h}=h$ . The black dashed-dot curves are the upper and lower bound of magnetization of the truncated interaction  $\tilde{J}^{-J}$ , viewed as functions of h. (Left)  $\epsilon=0.09$ . (Center)  $\epsilon=0.05$ . (Right)  $\epsilon=0.01$ .

COROLLARY 6.8. Let  $\tilde{J}^{-J}$  be the interaction given in Definition 6.7. Then, for  $0 < \epsilon < 1$  and  $\gamma > 0$  small enough, the UQ bounds (6.19) and (6.20) hold for  $R_{-J} = \gamma^{-1}$  and  $\mathcal{F} \le \epsilon \|J\|_{\infty}$ . The thermodynamic pressure  $P_{\mathbf{J},\beta,h}^{\Delta,\gamma}$  is given in (6.25). Let  $M(\tilde{\mathbf{J}}^F,\beta,\tilde{h})$  be the LP-limit of  $E_{\tilde{q}_{\Delta}}[m(\sigma_{\Delta})]$ . Then, the UQ bounds (6.21) and LP-limit of (6.19) hold with the limit of  $P_{\mathbf{J},\beta,h}^{\Delta,\gamma}$  being  $p_{\mathbf{J},\beta,h}$  given by (6.26).

Remark 6.9. Given  $\beta$ ,  $||J||_{\infty}$ , we can choose  $\epsilon \equiv \epsilon(\beta, ||J||_{\infty})$  sufficiently small. Consequently, the phase diagram of the two models are close to each other as the uncertainty area is very small (Figure 8). The parameter  $\epsilon$  quantifies the length of the area that one cuts off the initial interaction.

The same methods are applicable to other perturbations, e.g. the very long range in Appendix F and perturbations in "contexts"/configuration as boundary conditions.

Conclusion and future work. In this article, we developed an information-based uncertainty quantification method for Markov Random Fields/rMRFs. We considered a surrogate (baseline) MRF/rMRF constructed by physical modeling or by learning structure and parameters from data, and we quantify uncertainties inherited from data, modeling choices, or numerical approximations, that are also propagated in predictions for QoIs. Our UQ method quantifies uncertainties not only in parameters but also in structure as well as is capable in handling of the inherent high-dimensionality of systems that admit a MRF/rMRF formulation. This was achieved by obtaining tight and scalable, information-based bounds on the predictions for QoIs.

We demonstrated our UQ method in an example from medical diagnostics as well as several high dimensional equilibrium statistical mechanics models defined on bounded domains with suitable boundary conditions. We aim to extend the developed approach to non-equilibrium statistical mechanics systems [58] also arising in ML [3]. Furthermore, motivated by [30] we plan to develop robust uncertainty quantification for Bayesian networks defined on Directed Acyclical Graphs.

**Acknowledgments:** The research of M. K. was partially supported by the NSF HDR TRIPODS CISE-1934846. The research of P. B. and M. K., was partially supported by the Air Force Office of Scientific Research (AFOSR) under the grant FA-9550-18-1-0214.

Appendix A. Reduced Markov Random Fields (rMRFs). Let  $\mathbf{Y} = \{Y_i\}_{i \in \mathcal{V}}$  be a MRF indexed by a set of nodes  $\mathcal{V}$  (finite or infinite) of a graph  $\mathcal{G}$ . Let us consider  $\mathcal{M} \subset \mathcal{V}$ . Let also  $\mathbf{U} = \{Y_i\}_{i \in \mathcal{M}}$  and  $\mathbf{u}$  be an assignment to them, namely  $\mathbf{U} = \mathbf{u}$ . If  $\mathbf{Z} := \{Y_i\}_{i \in \mathcal{V} \setminus \mathcal{M}}$ , how does the underlying graph corresponding to  $\mathbf{Z} \mid \mathbf{U} = \mathbf{u}$  look like? Can the conditional probability  $p(\mathbf{z} \mid \mathbf{U} = \mathbf{u})$  still keep a

product structure/factorization as the joint distribution given in (2.1)? To answer the questions, we need a special class of MRF which is called *reduced Markov Random Fields* (rMRFs).

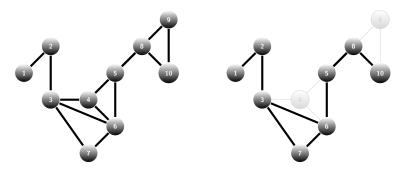


FIG. 9. The set of nodes is  $\mathcal{V} = \{1, \cdots, 10\}$  and  $\mathcal{M} = \{4, 9\}$ . Left:  $\mathbf{Y} = \{Y_i\}_{i=1}^{10}$  with joint distribution p is a MRF over  $\mathcal{G}$ . The set of maximal cliques is given by  $\mathcal{C}_{\mathcal{G}} = \{\{1, 2\}, \{2, 3\}, \{3, 4, 6, \}, \{3, 6, 7\}, \{4, 5, 6\}, \{5, 8\}, \{8, 9, 10\}\}$ . Right:  $\mathbf{Z} = \{Y_i\}_{i \in \mathcal{V} \setminus \mathcal{M}}$  with joint distribution q is the corresponding rMRF over  $\mathcal{G}'$  with  $\mathbf{U} = \{Y_4, Y_9\}$  and  $\mathbf{u} = \{u_4, u_9\}$ . The rMRF is demonstrated by removing the node 4 and 9 (faded nodes) from the graph  $\mathcal{G}$ .  $\mathcal{C}_{\mathbf{U}} = \{\{3, 4, 6, 7\}, \{4, 5, 6\}, \{8, 9, 10\}\}$  while  $\mathcal{C}_{\emptyset} = \{\{1, 2\}, \{2, 3\}, \{5, 8\}\}$ .

DEFINITION A.1. Let  $\mathbf{Y} = \{Y_i\}_{i \in \mathcal{V}}$  be a collection of random variables indexed by a set of nodes  $\mathcal{V}$  (finite or infinite) of a graph  $\mathcal{G}$ . If  $(\mathbf{Y}, p)$  is a MRF,  $\mathbf{u}$  a context,  $\mathcal{M} \subset \mathcal{V}$  and  $\mathbf{U} = \{Y_i\}_{i \in \mathcal{M}}$ , we define as **reduced Markov Random Field**, a MRF  $\mathbf{Z} = \{Y_i\}_{i \in \mathcal{V} \setminus \mathcal{M}}$  indexed by the set of nodes  $\mathcal{V} \setminus \mathcal{M}$  of the subgraph  $\mathcal{G}[\mathcal{V} \setminus \mathcal{M}]$  with joint distribution  $\mathbb{Q}$  such that

(A.1) 
$$q(\mathbf{z}) \equiv \mathbb{Q}(\mathbf{Z} = \mathbf{z}) := p(\mathbf{z} \mid \mathbf{U} = \mathbf{u}).$$

Therefore,  $\mathbf{Z} \mid \mathbf{U} = \mathbf{u}$  could be thought as a induced subgraph of  $\mathcal{G}$  with set of nodes  $\mathcal{V} \setminus \mathcal{M}$ , that is eliminating any node corresponding to random variables  $\mathbf{U}$  and any edge adjacent to them. Furthermore, according to Definition A.1,  $\mathbf{Z}$  is clearly MRF and therefore the conditional probability  $p(\mathbf{z} \mid \mathbf{U} = \mathbf{u})$  is expected to have a product structure. All the above are summarized in the following proposition:

PROPOSITION A.2. Let  $\mathbf{Y}$  be a MRF with probability distribution p > 0 parametrized by some parameters  $\mathbf{w} = \{\mathbf{w}_c\}_{c \in \mathcal{C}_{\mathcal{G}}}$  given in (2.1) and let  $\mathbf{U}, \mathbf{Z}$  be defined as in the beginning of the subsection. Then, q parametrized by  $\mathbf{w}$  is expressed as

(A.2) 
$$q^{\mathbf{w}}(\mathbf{z}) \equiv p(\mathbf{z} \mid \mathbf{U} = \mathbf{u}, \mathbf{w}) = \frac{1}{Z_{\mathbf{u}}(\mathbf{w})} \prod_{c \in \mathcal{C}_{C}} \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c})$$

where for every  $c \in C_{\mathcal{G}}$ 

(A.3) 
$$\Psi_c[\mathbf{u}](\mathbf{z}_c \mid \mathbf{w}_c) := \Psi_c(\mathbf{z}_c, \mathbf{u}_c \mid \mathbf{w}_c)$$

Moreover,  $Z_{\mathbf{u}}(\mathbf{w})$  is given by

(A.4) 
$$Z_{\mathbf{u}}(\mathbf{w}) = \sum_{\mathbf{Y}} \prod_{c \in \mathcal{C}_{\mathcal{C}}} \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c})$$

We refer to [49] and [55] for further discussion about MRFs, rMRFs and the proof of the Hammersley-Clifford Theorem and Propostion A.2.

**A.1. Partition of the class of maximal cliques.** We further investigate the structure of the class of all maximal cliques. Precisely, we collect  $c \in \mathcal{C}_{\mathcal{G}}$  such that  $\mathbf{U} \cap \mathbf{Y}_c \neq \emptyset$ . This leads to partition the set of maximal cliques  $\mathcal{C}_{\mathcal{G}} = \mathcal{C}_{\mathbf{U}} \sqcup \mathcal{C}_{\emptyset}$  with

(A.5) 
$$C_{\mathbf{U}} = \{c : \mathbf{U} \cap \mathbf{Y}_c \neq \emptyset\} \text{ and } C_{\emptyset} = \{c : \mathbf{U} \cap \mathbf{Y}_c = \emptyset\}.$$

(see example shown in Figure 9). On top of that, the partition of  $\mathcal{C}_{\mathcal{G}}$  makes the joint distributions q take the form

(A.6) 
$$q(\mathbf{z}) = P_{\Psi}^{\mathbf{w}}[\mathbf{u}](\mathbf{z}) = \frac{1}{Z_{\mathbf{u}}(\mathbf{w})} \prod_{c \in \mathcal{C}_{\mathbf{u}}} \Psi_{c}(\mathbf{y}_{c} \mid \mathbf{w}_{c}) \prod_{c \in \mathcal{C}_{\mathbf{u}}} \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c})$$

Appendix B. Proofs of the main results.

**B.1. Proof of Lemma 3.2.** In the following computation we use either (3.2) for type I or (3.6) for type II:

$$\begin{split} \tilde{Z}_{\mathbf{u}}(\tilde{\mathbf{w}}) &= \sum_{\mathbf{z}} \prod_{\tilde{c}} \tilde{\Psi}_{\tilde{c}}[\mathbf{u}](\mathbf{z}_{\tilde{c}} \mid \tilde{\mathbf{w}}_{\tilde{c}}) \\ &= \sum_{\mathbf{z}} \prod_{c} \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c}) \Phi_{\mathbf{u}}^{i}(\mathbf{z}) \\ &= \sum_{\mathbf{z}} \Phi_{\mathbf{u}}^{i}(\mathbf{z}) \prod_{c} \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c}) \\ &= Z_{\mathbf{u}}(\mathbf{w}) \sum_{\mathbf{z}} \Phi_{\mathbf{u}}^{i}(\mathbf{z}) \prod_{c} \Psi_{c}[\mathbf{u}](\mathbf{z}_{c} \mid \mathbf{w}_{c}) \frac{1}{Z_{\mathbf{u}}(\mathbf{w})} \\ &= Z_{\mathbf{u}}(\mathbf{w}) E_{q}[\Phi_{\mathbf{u}}^{i}(\mathbf{z})] \end{split}$$

**B.2. Proof of Theorem 4.1.** We are mostly based on the proof of the characterization of the exponential integrals (see, e.g. [22]). Let the probability measure R be defined by

$$dR/dq = e^{f(\mathbf{Z})}/E_q[f(\mathbf{Z})].$$

Note that  $\mathcal{R}(\tilde{q}||q) < \infty$ , since  $q, \tilde{q} > 0$ . Thus,

(B.1) 
$$-\mathcal{R}(\tilde{q}||q) + E_{\tilde{q}}[f(\mathbf{Z})] = -\mathcal{R}(\tilde{q}||R) + \log E_q[e^{f(\mathbf{Z})}] \le \log E_q[e^{f(\mathbf{Z})}].$$

where for the last inequality we use that  $\mathcal{R}(\tilde{q}||R) \geq 0$  and  $\mathcal{R}(\tilde{q}||R) = 0$  iff  $\tilde{q} = R$  [22, Lemma 1.4.1]. For part a., we combine (3.11) of Lemma 3.4 and (B.1) and we get

$$E_{\tilde{q}}[f(\mathbf{Z})] \le \log E_q[e^{f(\mathbf{Z})}] + \frac{1}{E_q[\Phi_{\mathbf{u}}^i]} E_q\left[\Phi_{\mathbf{u}}^i \log \Phi_{\mathbf{u}}^i\right] - \log E_q[\Phi_{\mathbf{u}}^i]$$

By replacing  $f(\mathbf{Z})$  to  $\pm \lambda f(\mathbf{Z})$ , we obtain

$$\pm E_{\tilde{q}}[f(\mathbf{Z})] \le \frac{1}{\lambda} \left\{ \log E_q[e^{\pm \lambda f(\mathbf{Z})}] + \frac{1}{E_q[\Phi_{\mathbf{u}}^i]} E_q\left[\Phi_{\mathbf{u}}^i \log \Phi_{\mathbf{u}}^i\right] - \log E_q[\Phi_{\mathbf{u}}^i] \right\}$$

By optimizing over  $\lambda > 0$  (see [15] and [53]), the following tight estimates are obtained:

$$\pm E_{\tilde{q}}[f(\mathbf{Z})] \leq \inf_{\lambda > 0} \frac{1}{\lambda} \left\{ \log E_q[e^{\pm \lambda f(\mathbf{Z})}] + \frac{1}{E_q[\Phi_{\mathbf{u}}^i]} E_q\left[\Phi_{\mathbf{u}}^i \log \Phi_{\mathbf{u}}^i\right] - \log E_q[\Phi_{\mathbf{u}}^i] \right\}$$

Part b. is proved similarly, utilizing (3.12) instead of (3.11).

Example B.1. (Single-parameter exponential families) This is a straightforward example and a simple illustration of the ideas in the proof of part b., Theorem 4.1, giving us insights on how well the ideas work together with a rearranging argument. The simplicity of this example arises from the fact that the exponential family is single parametric and therefore the structural part is not present. The probability density function of a random variable X with range R(X), is given by

$$p^{\theta}(x) = P^{\theta}(X = x) = e^{\theta\phi(x) - F(\theta)}$$

taken with respect to some measure  $d\nu$  where  $F(\theta) = \log \int_x e^{\theta \phi(x)} \nu(dx)$  and  $\phi(x)$  is a real-valued function also known as sufficient statistic. Suppose a second probability density function of the same single-parameter exponential family associated with  $\phi$ 

$$p^{\theta+\zeta}(x) = P^{\theta+\zeta}(X=x) = e^{(\theta+\zeta)\phi(x)-F(\theta+\zeta)}$$

for some  $\zeta < 1$ . One may want to investigate how sensitive the model is in such a change in  $\theta$  by  $\zeta$  with respect to  $\phi(X)$  as means to bound  $E_{P^{\theta+\zeta}}[\phi(X)]$  or to find the error in replacing the first distribution by the "perturbed" one and phrased as bound  $E_{P^{\theta+\zeta}}[\phi(X)] - E_{P^{\theta}}[\phi(X)]$ . The second exponential family is apparently a perturbation on parameters by  $\zeta$ , so we can think of the model as Type I. In addition, after employing UQ bounds, the cumulant generating function and KL divergence are the two main ingredients to compute: for any  $\lambda > 0$ ,

$$\Lambda_{P^{\theta}}^{\phi}(\lambda) = \log E_{P^{\theta}}[e^{\lambda\phi(X)}] = F(\theta + \lambda) - F(\theta)$$

$$R(P^{\theta+\zeta}||P^{\theta}) = \zeta E_{P^{\theta+\zeta}}[\phi(X)] - \log E_{P^{\theta}}[e^{\zeta\phi(X)}]$$

The above expression for KL divergence comes from the calculation of expressing  $F(\theta + \lambda)$  in terms of  $F(\theta)$  and for that every term is computed with respect to  $P^{\theta}$ . By substituting the quantities to the UQ bounds and by doing a delicate rearrangement of terms that is feasible because the QoI is a sufficient statistic for the model, we get

$$\pm E_{P^{\theta+\zeta}}[\phi(X)] \leq \frac{1}{1-\zeta} \inf_{\lambda>0} \left\{ \frac{F(\theta+\lambda) - F(\theta)}{\lambda} + \frac{1}{\lambda} \log E_{P^{\theta}}[e^{\zeta\phi(X)}] \right\}$$

**B.3. Proof of Theorem 4.2.** The existence and the explicit form of the distribution  $q^{\pm}$  relies on [39], Theorem 2. Consequently, given a QoI f, we identify the total  $\tilde{q}$ -excess factor relative to q explicitly, that is  $\Phi_{\mathbf{u}}^{\pm} = e^{\lambda_{\pm}f}$ . However, the new element is that by utilizing the Hammersley-Clifford Theorem,  $q^{\pm}$  defined on  $\mathbf{Z}$  are rMRFs, lie in the class  $\mathcal{Q}_{\mathcal{P}}^{\eta}$  and the total  $\tilde{q}$ -excess factor relative to q is explicitly determined.

Appendix C. Coarse-Graining, Kac and Hamiltonian Estimates.



Fig. 10. One-dimensional Ising spin lattice on  $\Delta$  (white spins) with configuration boundary conditions on the complement of  $\Delta$  denoted as  $\bar{\sigma}_{\Delta^c}$  (black spins).

**C.1. Coarse-graining.** We divide  $\mathbb{R}^d$  into cubes of side  $l = \gamma^{-1/2}$ . We denote by  $\mathcal{P}^{(l)}_{\mathbb{R}^d}$  the partition of  $\mathbb{R}^d$ . Namely, for every  $i \in l\mathbb{Z}^d$  we set

(C.1) 
$$I_{\gamma,i} = \{ r \in \mathbb{R}^d : i_k \le r_k \le i_k + l, k = 1, \dots, d \}$$

 $(r_k \text{ and } i_k \text{ being the } k\text{-th coordinate of } r \text{ and } i)$ . Then we call

(C.2) 
$$\mathcal{P}_{\mathbb{R}^d}^{(l)} = \{ I_{\gamma,i} : i \in l\mathbb{Z}^d \},$$

the collection of all the above cubes.

DEFINITION C.1 ([57]). (1) A function f(r) is  $\mathcal{P}_{\mathbb{R}^d}^{(l)}$ -measurable, if it is constant in each cube  $I_{\gamma,i}$ ,  $i \in l\mathbb{Z}^d$ .

- (2) A region  $\Delta \subset \mathbb{R}^d$  is  $\mathcal{P}^{(l)}_{\mathbb{R}^d}$ -measurable, if it can be written as a union of cubes of  $\mathcal{P}^{(l)}_{\mathbb{R}^d}$  (or its characteristic is  $\mathcal{P}^{(l)}_{\mathbb{R}^d}$ -measurable).
  - (3) Any  $\Delta \subset \mathbb{Z}^d$  can be identified as a union of cubes with length 1.
  - (4) The size of each cube is given by

(C.3) 
$$|I_{\gamma,i}| = |I| = l^d = \gamma^{-d/2}$$

for every  $i \in l\mathbb{Z}^d$ . For notational simplicity, we drop  $\gamma$  from  $I_{\gamma,i}$ .

For any bounded region  $\Delta$   $\mathcal{P}_{\mathbb{R}^d}^{(l)}$ -measurable, we denote  $\Delta := \Delta \cap \mathbb{Z}^d$ . Hence,  $\mathbf{I}_i = I_i \cap \mathbb{Z}^d$ .

C.2. Coarse-grained Interaction. We introduce a new interaction  $\bar{J}_{\gamma}$  which describes the interaction between cubes. More precisely, for every  $i, j \in l\mathbb{Z}^d$  with  $i \neq j$ , we consider

(C.4) 
$$\bar{J}_{\gamma}(i,j) = \frac{1}{|I|^2} \sum_{x \in \mathbf{I}_i} \sum_{y \in \mathbf{I}_i} J_{\gamma}(x,y),$$

and for i = j, we define

(C.5) 
$$\bar{J}_{\gamma}(i,i) = \frac{1}{|I|(|I|-1)} \sum_{x \in \mathbf{I}_i} \sum_{\substack{x \in \mathbf{I}_i, \\ y \neq x}} J_{\gamma}(x,y)$$

LEMMA C.2. For fixed and small  $\gamma > 0$ , for any  $x \in \mathbf{I}_i$  and any  $y \in \mathbf{I}_j$ ,  $i, j \in l\mathbb{Z}^d$  with  $i \neq j$ , we have

(C.6) 
$$|J_{\gamma}(x,y) - \bar{J}_{\gamma}(i,j)| \le \gamma^{d+\frac{1}{2}} ||DJ||_{\infty} \mathbf{1}_{|x-y| < 2\gamma^{-1}}.$$

Also, for any  $i \in l\mathbb{Z}^d$  and any  $x, y \in \mathbf{I}_i$ , we have

(C.7) 
$$|J_{\gamma}(x,y) - \bar{J}_{\gamma}(i,i)| \le \gamma^d ||J||_{\infty}$$

*Proof.* Let  $x \in \mathbf{I}_i$  and any  $y \in \mathbf{I}_j$ ,  $i, j \in l\mathbb{Z}^d$  with  $i \neq j$ , we have

$$\begin{split} |J_{\gamma}(x,y) - \tilde{J}_{\gamma}(x,y)| &= |J_{\gamma}(x,y) - \frac{1}{|I|^2} \sum_{z \in \mathbf{I}_i} \sum_{w \in \mathbf{I}_j} J_{\gamma}(z,w)| \\ &\leq \frac{1}{|I|^2} \sum_{z \in I_i} \sum_{w \in I_j} |J_{\gamma}(x,y) - J_{\gamma}(z,w)| \\ &\leq \frac{1}{|I|^2} \sum_{z \in I_i} \sum_{w \in I_j} \gamma^d \|DJ\|_{\infty} \gamma |x - y - z + w| \mathbf{1}_{|x - y| \leq \gamma^{-1}} \\ &\leq \frac{1}{|I|^2} |I|^2 \gamma^d \|DJ\|_{\infty} \gamma \gamma^{-1/2} \mathbf{1}_{|x - y| \leq \gamma^{-1}} \\ &= \gamma^{d + \frac{1}{2}} \|DJ\|_{\infty} \mathbf{1}_{|x - y| \leq \gamma^{-1}} \end{split}$$

We prove (C.7) similarly.

**C.3.** Coarse-grained Hamiltonian Energy. In this section we analyze the Hamiltonian energy by using the new interaction defined in (C.4) and the estimates in Lemma C.2. We start by introducing some notation: for any  $r \in \mathbb{R}^d$ , we define the following quantity as block spin configuration:

(C.8) 
$$\sigma^{(\gamma^{-1/2})}(r) := \frac{1}{|I|} \sum_{x \in \mathbf{I}_r} \sigma_{I_i}(x)$$

so that

$$\sigma^{(\gamma^{-1/2})}(r) = \frac{1}{|I|} \int_{I} \sigma^{(1)}(r') dr'$$

Let  $\Delta \subset \mathbb{R}^d$  be  $\mathcal{P}^{(l)}_{\mathbb{R}^d}$ -measurable region. We denote by  $\mathcal{M}^{(\gamma^{-1/2})}_{\Delta}$  all  $\mathcal{P}^{(l)}_{\mathbb{R}^d}$ -measurable functions on  $\Delta$  with values in

(C.9) 
$$M^{(\gamma^{-1/2})} := \{-1, -1 + \frac{1}{\gamma^{-d/2}}, \dots, 1 - \frac{1}{\gamma^{-d/2}}, 1\}$$

For any bounded  $\mathcal{P}_{\mathbb{R}^d}^{(l)}$ -measurable region  $\Delta$  and  $m_{\Delta} \in \mathcal{M}_{\Delta}^{(\gamma^{-1/2})}$ , we define as coarse-grained Hamiltonian energy

$$\bar{H}_{\gamma,h}^{\bar{\mathbf{J}}}(m_{\Delta}; m_{\Delta^{c}}) := \int_{\Delta} \phi_{\beta,h}(m_{\Delta}(r))dr + \frac{1}{4} \int_{\Delta} \int_{\Delta} J_{\gamma}(r, r')[m_{\Delta}(r) - m_{\Delta}(r')]^{2} dr dr' 
+ \frac{1}{2} \int_{\Delta} \int_{\Delta^{c}} J_{\gamma}(r, r')[m_{\Delta}(r) - m_{\Delta^{c}}(r')]^{2} dr dr' 
- \frac{1}{2} \int_{\Delta} \int_{\Delta^{c}} J_{\gamma}(r, r')m_{\Delta^{c}}(r')^{2} dr dr' 
+ \frac{1}{\beta} \int_{\Delta} I(m_{\Delta}(r)) dr$$
(C.10)

where

(C.11) 
$$I(m) := -\frac{1-m}{2} \log \frac{1-m}{2} - \frac{1+m}{2} \log \frac{1+m}{2}$$

with  $\phi_{\mathbf{J},\beta,h}(m)$  being given in (6.26). We recall that  $\mathcal{J} = \int_{\mathbb{R}^d} J(r) dr$ .

LEMMA C.3. Let  $\Delta$  be any bounded  $\mathcal{P}_{\mathbb{R}^d}^{(l)}$ -measurable region  $\Delta$ , then there exists a constant C > 0 such that the following estimate holds:

$$\left| H_{\gamma,h}^{\mathbf{J}}(\sigma_{\Delta}; \bar{\sigma}_{\Delta^c}) - \bar{H}_{\gamma,h}^{\bar{\mathbf{J}}}(\sigma_{\Delta}^{(\gamma^{-1/2})}; \bar{\sigma}_{\Delta^c}^{(\gamma^{-1/2})}) \right| \leq C|\Delta|\gamma^{1/2},$$

where  $\sigma_{\Delta}^{(\gamma^{-1/2})}$  and  $\bar{\sigma}_{\Delta^c}^{(\gamma^{-1/2})}$  are defined in (C.8).

C.4. Estimates for the thermodynamic pressure of an Ising-Kac model. We recall that

$$P_{\mathbf{J},\beta,h}^{\Delta,\gamma}(\bar{\sigma}_{\Delta^c}) := \frac{\log Z_{\bar{\sigma}_{\mathbf{I}^c}}(\mathbf{J},\beta,h)}{\beta|\Delta|}$$

and

$$p_{\mathbf{J},\beta,h} := -\inf_{m \in [-1,1]} \{-hm + \phi_{\mathbf{J},\beta,0}(m)\}$$

If  $\epsilon(\gamma) = \gamma^{1/2} + \gamma^{d/2} \log \gamma^{-1}$ , then the following bounds hold: there exist constants c, c' > 0 such that

$$(\mathrm{C.13}) \hspace{1cm} P_{\mathbf{J},\beta,h}^{\Delta,\gamma}(\bar{\sigma}_{\Delta^c}) \leq p_{\mathbf{J},\beta,h} + \left(c\frac{\gamma^{-1}}{L} + c\epsilon(\gamma)\right), \hspace{1cm} \mathbf{Upper\ Bound}$$

Let  $m^*$  be the minimizer of  $\phi_{\mathbf{J},\beta,h}$ , then  $p_{\mathbf{J},\beta,h} = -\phi_{\mathbf{J},\beta,h}(m^*)$ , then (C.14)

$$P_{\mathbf{J},\beta,h}^{\Delta,\gamma}(\bar{\sigma}_{\Delta^c}) \geq p_{\mathbf{J},\beta,h} - |\phi_{\mathbf{J},\beta,h}([m^*]_{\gamma}) - \phi_{\mathbf{J},\beta,h}(m^*)| - c\epsilon(\gamma) - c'\frac{\gamma^{-1}}{L}, \qquad \mathbf{Lower\ Bound}$$

where  $[m^*]_{\gamma}$  is the value in (C.9) closest to  $m^*$ .

C.5. Limit as  $\Delta \nearrow \mathbb{Z}^d$  and then  $\gamma \to 0$ . By using the estimates for the hamiltonian energy given in (C.12), (C.13) and (C.14) we can prove that

(C.15) 
$$\limsup_{\gamma \to 0} \lim_{\Delta \nearrow \mathbb{Z}^d} P_{\mathbf{J},\beta,h}^{\Delta,\gamma}(\bar{\sigma}_{\Delta^c}) \le p_{\mathbf{J},\beta,h}$$

(C.16) 
$$\liminf_{\gamma \to 0} \lim_{\Delta \nearrow \mathbb{Z}_d} P_{\mathbf{J},\beta,h}^{\Delta,\gamma}(\bar{\sigma}_{\Delta^c}) \ge p_{\mathbf{J},\beta,h}$$

and therefore if  $P_{\mathbf{J},\beta,h}^{\gamma} := \lim_{\Delta \to \mathbb{Z}^d} P_{\mathbf{J},\beta,h}^{\Delta,\gamma}$ , then

(C.17) 
$$\lim_{\gamma \to 0} P_{\mathbf{J},\beta,h}^{\gamma} = p_{\mathbf{J},\beta,h} = -\inf_{m \in [-1,1]} \{-hm + \phi_{\mathbf{J},\beta,0}(m)\}.$$

Hence, the thermodynamic pressure converges to the mean field pressure at the LPlimit, namely

$$\lim_{\gamma \to 0} \lim_{\Delta \nearrow \mathbb{Z}^d} P_{\mathbf{J},\beta,h}^{\Delta,\gamma} = p_{\mathbf{J},\beta,h}$$

where  $p_{\mathbf{J},\beta,h}$  is defined in (6.26). The convexity properties are provided by the limit as  $\Delta \nearrow \mathbb{Z}^d$  and then preserved by  $\gamma \to 0$ .

**C.6.** Thermodynamics of an Ising-spin model with a Kac potential. It is shown that when  $\gamma > 0$  is sufficiently small, the phase diagram of an Ising-spin model with a Kac potential is close to the phase diagram of a mean field model. Precisely, in [13, 11] (see also [57]) it is proved that for  $d \geq 2$ , if  $h \neq 0$  then there exists a unique DLR measure, [58]. If h = 0 there is a critical value of inverse temperature  $\beta_c(\gamma) > 0$  such that for any  $\beta < \beta_c(\gamma)$ , there exists one DLR measure while for  $\beta > \beta_c(\gamma)$  there are at least two distinct DLR measures  $\mu_{\beta,\gamma}^{\pm}$ . Finally, there is an absence of phase transition when  $\gamma$  is kept small (for more details see [57, 58] and references therein).

### Appendix D. Detailed Analysis of Medical Diagnostics.

**D.1. Baseline model.** Let us consider the undirected graph in Figure 2, [20] denoted by  $\mathcal{G}$ . The class of maximal cliques is  $\mathcal{C}_{\mathcal{G}} = \{\{1,2\},\{2,3,4\}\}$ . The distribution defined over the graph is a log-linear model with clique potentials given by  $\Psi_c(\mathbf{y}_c \mid \mathbf{w}_c) = e^{w_c f_c(\mathbf{y}_c)}$ , where all the weights  $\mathbf{w}_c$ , and the binary functions  $f_c$  are known. For example, for  $c = \{1,2\}$ ,  $\mathbf{w}_{\{1,2\}} = 1.5$  and

$$f_{\{1,2\}}(\mathbf{y}_{\{1,2\}}) = \begin{cases} 1 & , \mathbf{y}_{\{1,2\}} \in \{(s_1, l_1), (s_1, l_0), (s_0, l_0)\} \\ 0 & , \mathbf{y}_{\{1,2\}} \in \{(s_0, l_1)\} \end{cases}$$

Each binary function  $f_c$  induces a set  $B_c = \{(\omega_1, \omega_2, \omega_3, \omega_4) : f_c(\omega_c) = 1\}$ . For example,  $B_{\{1,2\}} := \{\omega : \omega_{\{1,2\}} \in \{(s_1, l_1), (s_1, l_0), (s_0, l_0)\}\}$ . We compare predictions between the baseline and alternatives of Type I and II (see Section D.2) for the following QoIs:

$$g(\mathbf{Y}) = \mathbf{1}_A$$
, for any event of interest  $A \subset \Omega$ .

For instance,  $A = \{\text{patient is smoker with asthma}\} = \{\omega = (\omega_1, \omega_2, \omega_3, \omega_4) : \omega_1 = s_0, \omega_3 = a_0\}.$ 

#### D.2. Alternative models.

**D.2.1.** Type I. First, we consider the class of log-linear models  $\tilde{p}$  over  $\mathcal{G}$  with weight change in one maximal clique. Let c be the maximal clique that a weight change occurred. Then the clique potential is given by

$$\tilde{\Psi}_c(\mathbf{y}_c) = e^{\tilde{w}_c f_c(\mathbf{y}_c)}$$

The weight after increasing or decreasing by 100a% equals to  $\tilde{w}_c = (1+a)w_c$ , where  $a \in [-1,1]$  stands for the model uncertainty of alternative models of Type I and  $w_c$  is the weight on c of the baseline model p. For example, for  $\{1,2\}$ , the corresponding clique potential is expressed as

$$\begin{split} \tilde{\Psi}_{\{1,2\}}(\mathbf{y}_{\{1,2\}} \mid \tilde{\mathbf{w}}_{\{1,2\}}) &= e^{\tilde{w}_{\{1,2\}}\tilde{f}_{\{1,2\}}(\mathbf{y}_{\{1,2\}})} \\ &= \Psi_{\{1,2\}}(\mathbf{y}_{\{1,2\}} \mid \mathbf{w}_{\{1,2\}}) \Phi_{\{1,2\}}(\mathbf{y}_{\{1,2\}} \mid \tilde{\mathbf{w}}_{\{1,2\}}) \end{split}$$

with

$$\Phi_{\{1,2\}}(\mathbf{y}_{\{1,2\}}\mid \tilde{\mathbf{w}}_{\{1,2\}}) = e^{-0.2w_{\{1,2\}}f_{\{1,2\}}(\mathbf{y}_{\{1,2\}})}$$

since we consider the simplest case where  $\tilde{f}_{\{1,2\}}(\mathbf{y}_{\{1,2\}}) = f_{\{1,2\}}(\mathbf{y}_{\{1,2\}})$  as well as the fact that  $\tilde{w}_{\{1,2\}} - w_{\{1,2\}} = -0.2w_{\{1,2\}}$ . Note that  $\mathcal{B} = \{c\}$ , where  $\mathcal{B}$  defined in subsection 3.1.

**Derivation of** (5.2): We compute all the quantities involved in (4.4) explicitly. Let us start with the cumulant generating function:

$$\Lambda_p^f(\lambda) = \log E_p[e^{\lambda g}] = \log \left( \sum_{\mathbf{y} \in A} e^{\lambda g} p(\mathbf{y}) + \sum_{\mathbf{y} \notin A} e^{\lambda g} p(\mathbf{y}) \right)$$
$$= \log \left( e^{\lambda} p(A) + 1 - p(A) \right)$$

It is straightforward to see that

(D.1) 
$$\frac{d\tilde{p}}{dp} = \frac{\Phi^{\mathrm{I}}}{E_{p}[\Phi^{\mathrm{I}}]} = \frac{e^{aw_{c}f_{c}}}{e^{aw_{c}}p_{\mathrm{I}} + 1 - p_{\mathrm{I}}}.$$

and we now go through the computation of  $E_p[\Phi^{\rm I}]$ :

$$E_{p}[\Phi^{\mathbf{I}}] = \sum_{\mathbf{y}} \Phi^{\mathbf{I}}(\mathbf{y}) p(\mathbf{y}) = \sum_{\mathbf{y}} e^{aw_{c}f_{c}(\mathbf{y}_{c})} p(\mathbf{y}_{c})$$

$$= \sum_{\mathbf{y} \in B_{c}} e^{aw_{c}f_{c}(\mathbf{y}_{c})} p(\mathbf{y}) + \sum_{\mathbf{y} \notin B_{c}} e^{aw_{c}f_{c}(\mathbf{y}_{c})} p(\mathbf{y})$$

$$= e^{aw_{c}} p_{\mathbf{I}} + 1 - p_{\mathbf{I}}.$$

Similarly, we prove that

(D.2) 
$$E_p[\Phi^i \log \Phi^i] = aw_c e^{aw_c} p_I$$

Overall, by recalling (3.11) the KL divergence equals to

$$R(\tilde{p}||p) = \frac{aw_c e^{aw_c} p_{\rm I}}{e^{aw_c} p_{\rm I} + 1 - p_{\rm I}} - \log\left(e^{aw_c} p_{\rm I} + 1 - p_{\rm I}\right)$$

**D.2.2.** Type II. We consider the class of log-linear models  $\tilde{p}$  over  $\tilde{\mathcal{G}}$  with  $\tilde{\mathcal{V}} = \mathcal{V}$ ,  $\tilde{\mathcal{E}} = \mathcal{E} \cup e$ , where e is a new edge (for example, see Figure 2, (Right)). We assume that the edge e enlarges an already existing maximal clique in the sense of the analysis in subsection 3.1. The model uncertainties arising from structure-learning from either a new data set  $\tilde{\mathcal{D}}$  and/or different prior knowledge; see for example Figure 2 (Right) lie in the binary function  $\tilde{f}_{\tilde{c}}$  defined on  $\tilde{c}$  and the new weight  $\tilde{\mathbf{w}}_{\tilde{c}}$ , where  $\tilde{c}$  is the enlargement of an existing maximal clique c. The weight  $\tilde{\mathbf{w}}_{\tilde{c}}$  can also be expressed with respect to  $w_c$ :  $\tilde{w}_{\tilde{c}} = (1+a)w_c$ . This time  $a \in \mathbb{R}$ , not necessarily in [-1,1] as before (e.g  $w_c = 1.5$  and  $\tilde{w}_{\tilde{c}} = 5$ ). Then the corresponding clique potential is given by

$$\tilde{\Psi}_{\tilde{c}}(\mathbf{y}_{\tilde{c}}) = e^{\tilde{w}_{\tilde{c}}\tilde{f}_{\tilde{c}}(\mathbf{y}_{\tilde{c}})} = e^{(1+a)w_c\tilde{f}_{\tilde{c}}(\mathbf{y}_{\tilde{c}})}$$

The binary function  $f_{\tilde{c}}$  induces a set  $B_{\tilde{c}} = \{(\omega_1, \omega_2, \omega_3, \omega_4) : \tilde{f}_{\tilde{c}}(\omega_{\tilde{c}}) = 1\}$ . For example, Let  $\tilde{\mathcal{G}} \neq \mathcal{G}$  (also  $\mathcal{C}_{\mathcal{G}} \neq \mathcal{C}_{\tilde{\mathcal{G}}}$ ) and  $\mathbf{w} \neq \tilde{\mathbf{w}}$ . Intuitively, a change on the set of edges can be thought of as structure-learning from either a new data set  $\tilde{\mathcal{D}}$  and/or different prior knowledge; see for example Figure 2, (Right) where only one new edge has been added.

The set  $B_{\tilde{c}}$  satisfies one of the following:  $B_{\tilde{c}} \cap B_c = \emptyset$  or  $B_{\tilde{c}} \cap B_c \neq \emptyset$ . Note that  $\mathcal{B}_{\subset} = \{\tilde{c}\}$  and  $\mathcal{B}_{\cup} = \mathcal{B}_{new} = \emptyset$  with  $\mathcal{B}_{\subset}, \mathcal{B}_{\cup}$  and  $\mathcal{B}_{new}$  are defined in subsection 3.1.

**Derivation of** (5.3): The cumulant generating function is the same as in the derivation of (5.2). Let us compute the expected value of the total  $\tilde{p}$ -excess factor of type

II relative to p with respect to p:

$$E_{p}[\Phi^{\mathrm{II}}] = \sum_{\mathbf{y}} \Phi^{\mathrm{II}}(\mathbf{y}) p(\mathbf{y}) = \sum_{\mathbf{y}} e^{(1+a)w_{c}\tilde{f}_{\tilde{c}} - w_{c}f_{c}} p(\mathbf{y})$$

$$= \sum_{\mathbf{y} \in B_{c}} e^{aw_{c}f_{c}(\mathbf{y}_{c})} p(\mathbf{y}) + \sum_{\mathbf{y} \in B_{\tilde{c}}} e^{(1+a)w_{c}\tilde{f}_{\tilde{c}} - w_{c}f_{c}} p(\mathbf{y}) + \sum_{\mathbf{y} \notin B_{c} \cup B_{\tilde{c}}} e^{(1+a)w_{c}\tilde{f}_{\tilde{c}} - w_{c}f_{c}} p(\mathbf{y})$$

$$(D.3) = e^{(1+a)w_{c}} p_{\mathrm{II}} + e^{-w_{c}} p_{\mathrm{I}} + 1 - p_{\mathrm{I}} - p_{\mathrm{II}}.$$

We split the sum into the three sums since  $B_c \cap B_{\tilde{c}} = \emptyset$ . Similarly, we prove that

(D.4) 
$$E_p[\Phi^{II} \log \Phi^{II}] = -w_c e^{-w_c} p_I + (1+a)w_c e^{(1+a)w_c} p_{II}$$

Overall, by recalling (3.11) the KL divergence equals to

$$R(\tilde{p}||p) = \frac{-w_c e^{-w_c} p_{\rm I} + (1+a) w_c e^{(1+a)w_c} p_{\rm II}}{e^{(1+a)w_c} p_{\rm II} + e^{-w_c} p_{\rm I} + 1 - p_{\rm I} - p_{\rm II}} - \log\left(-w_c e^{-w_c} p_{\rm I} + (1+a) w_c e^{(1+a)w_c} p_{\rm II}\right)$$

Remark D.1. If  $B_c \cap B_{\tilde{c}} \neq \emptyset$ , then we need to split the sum of (D.3) as follows: Let  $U \equiv B_c \cap B_{\tilde{c}}$ , then

$$\begin{split} E_{p}[\Phi^{\mathrm{II}}] &= \sum_{\mathbf{y}} \Phi^{\mathrm{II}}(\mathbf{y}) p(\mathbf{y}) = \sum_{\mathbf{y}} e^{(1+a)w_{c}\tilde{f}_{\tilde{c}} - w_{c}f_{c}} p(\mathbf{y}) \\ &= \sum_{\mathbf{y} \in B_{c} \setminus U} e^{aw_{c}f_{c}(\mathbf{y}_{c})} p(\mathbf{y}) + \sum_{\mathbf{y} \in B_{\tilde{c}} \setminus U} e^{(1+a)w_{c}\tilde{f}_{\tilde{c}} - w_{c}f_{c}} p(\mathbf{y}) + \sum_{\mathbf{y} \in U} e^{(1+a)w_{c}\tilde{f}_{\tilde{c}} - w_{c}f_{c}} p(\mathbf{y}) \\ &+ \sum_{\mathbf{y} \notin B_{c} \cup B_{\tilde{c}}} e^{(1+a)w_{c}\tilde{f}_{\tilde{c}} - w_{c}f_{c}} p(\mathbf{y}) \\ &= e^{(1+a)w_{c}}(p_{\mathrm{II}} - p(U)) + e^{-w_{c}}(p_{\mathrm{I}} - p(U)) + e^{aw_{c}} p(U) + 1 - p_{\mathrm{I}} - p_{\mathrm{II}} + p(U). \end{split}$$

Note that  $p_{I}, p_{II}$  and p(U) are computable as p is known.

### Appendix E. Analysis of UQ for Statistical Mechanics.

**E.1. Proof of Lemma 6.2.** It is not difficult to show (see also Proposition II.1.2 and Lemma II.2.2C in [62]) that

$$|\log Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J}, \beta, h) - \log Z_{\bar{\sigma}_{\Delta^{c}}}(\tilde{\mathbf{J}}^{\mathbf{F}}, \beta, h)| \leq \beta \|H^{\mathbf{J}, h}(\sigma_{\Delta}|\bar{\sigma}_{\Delta^{c}}) - H^{\tilde{\mathbf{J}}^{\mathbf{F}}, h}(\sigma_{\Delta}|\bar{\sigma}_{\Delta^{c}})\|_{\infty}$$

$$\leq |\Delta| \|\Phi_{\Delta, \bar{\sigma}_{\Delta^{c}}}^{h, \beta, \mathbf{J}} - \Phi_{\Delta, \bar{\sigma}_{\Delta^{c}}}^{h, \beta, \tilde{\mathbf{J}}^{\mathbf{F}}}\|_{1}$$
(E.1)

which in turn gives

(E.2) 
$$R(\tilde{q}_{\Delta} || q_{\Delta}) \le 2|\Delta| \|\Phi_{\Delta, \bar{\sigma}_{\Delta} c}^{h, \beta, \mathbf{J}} - \Phi_{\Delta, \bar{\sigma}_{\Delta} c}^{h, \beta, \tilde{\mathbf{J}}^{\mathbf{F}}} \|_{\mathbf{I}}$$

since

$$R(\tilde{q}_{\Delta} || q_{\Delta}) = \beta \left( E_{\tilde{q}_{\Delta}} [H^{\mathbf{J},h}(\sigma_{\Delta} | \bar{\sigma}_{\Delta^{c}})] - E_{q_{\Delta}} [H^{\tilde{\mathbf{J}}^{\mathbf{F}},\tilde{h}}(\sigma_{\Delta} | \bar{\sigma}_{\Delta^{c}})] \right) + \log Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J},\beta,h) - \log Z_{\bar{\sigma}_{\Delta^{c}}}(\tilde{\mathbf{J}}^{\mathbf{F}},\beta,\tilde{h})$$

A straightforward bound yields that

$$\|\Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{h,\beta,\mathbf{J}} - \Phi_{\Delta,\bar{\sigma}_{\Delta^c}}^{h,\beta,\tilde{\mathbf{J}}^{\mathbf{F}}}\|_1 \le \beta \left( |\tilde{h} - h| + \sum_{x \ne 0} |F(0,x)| \right).$$

**E.2. Proof of Lemma 6.2.2.** It is a straightforward computation after subtracting the hamiltonian energies with interaction J and

$$\tilde{J}^F(x,y) = J(x,y) \mathbf{1}_{\|x-y\|_d \le R} + F(x,y) \mathbf{1}_{\|x-y\|_d \le R}$$
, Type I,

and

$$\tilde{J}^F(x,y) = J(x,y) \mathbf{1}_{\|x-y\|_d \le R} + F(x,y) \mathbf{1}_{\|x-y\|_d \ge R}, \text{ Type II}$$

**E.2.1.** Cumulant generating function for  $f(\mathbf{Z}) = |\Delta| m(\sigma_{\Delta})$ .

$$\Lambda_{q_{\Delta};|\Delta|m(\sigma_{\Delta})}(\pm\lambda) = \log E_{q_{\Delta}}[e^{\lambda|\Delta|\frac{1}{|\Delta|}\sum_{x\in\Delta}\sigma_{\Delta}(x)}]$$

$$= \log \left(\frac{1}{Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J},\beta,h)}\sum_{\sigma_{\Delta}}e^{\lambda\sum_{x\in\Delta}\sigma_{\Delta}(x)}e^{-\beta H^{\mathbf{J},h}(\sigma_{\Delta}|\sigma_{\Delta^{c}})}\right)$$

$$= \log \left(e^{\lambda\sum_{x\in\Delta}\sigma_{\Delta}(x)-\beta H^{\mathbf{J},h}(\sigma_{\Delta}|\sigma_{\Delta^{c}})}\right) - \log Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J},\beta,h)$$
(E.3)
$$:= \log Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J},\beta,h\pm\frac{\lambda}{\beta}) - \log Z_{\bar{\sigma}_{\Delta^{c}}}(\mathbf{J},\beta,h)$$

Then by using the definition of the thermodynamic pressure in (6.25), we get:

(E.4) 
$$\frac{1}{|\Delta|} \Lambda_{q_{\Delta}; |\Delta| m(\sigma_{\Delta})}(\pm \lambda) = \beta \left( P_{h \pm \frac{\lambda}{\beta}, \beta, \mathbf{J}}^{\Delta, \gamma} - P_{h, \beta, \mathbf{J}}^{\Delta, \gamma} \right)$$

Appendix F. Phase diagram of a long range perturbation.

- **F.1.** Thermodynamics of a long range perturbation of 1-dimensional Kac model. There is a significant number of works in the literature studying the phase diagram of one-dimensional ferromagnetic Ising model with long range interactions of the form  $1/r^k$  with k indicating the decay of interaction and  $k \leq 2$ . For k < 2, the occurrence of phase transition has been proved (see [26, 27, 28]). For k = 2, the existence of a spontaneous magnetization at low temperature is proved in [33]. The establishment of the existence of phase transition, proving the discontinuity of the magnetization at a critical point, also known as *Thouless effect*, was proved by Aizenman et al in [1]. In [12], the authors study the phase diagram of the system with interaction defined in (F.1) with F given in Definition F.1 as illustrated in the right graph of Figure 7. Precisely, they have shown that there is a critical value of the inverse temperature depending on a and  $\gamma$  sufficiently small such that the system exhibits phase transition.
- **F.1.1. Phase diagram of a long range perturbation.** We consider a one dimensional ferromagnetic Ising spin system with interactions that correspond to a  $1/r^2$  long range perturbation of the usual Kac model, see the right picture of Figure 7.

DEFINITION F.1. Let  $J_{\gamma}^{\text{pwc}}(x,y) = \gamma^d \mathbf{1}_{|x-y| \leq \frac{\gamma-1}{2}}$  (i.e. a special case of Kac-type interaction where in fact  $J_{\gamma}^{\text{pwc}}(x,y)$  is piecewise constant interaction). Then we define

(F.1) 
$$\tilde{J}_{\gamma}^{F}(x,y) = \begin{cases} J_{\gamma}^{\text{pwc}} &, 0 \le |x-y| \le (2\gamma)^{-1} \\ F(x,y) &, |x-y| > (2\gamma)^{-1}, \end{cases}$$

with  $F(x,y) = \frac{a}{|x-y|^2}$  for some number  $a \in (0,\infty)$ , Figure 7 (right).

The range of the perturbation F is clearly Type II. We derive the UQ bounds as follows:

(F.2) 
$$\log \Phi_{\bar{\sigma}_{\Delta^{c}}}^{i}(\sigma_{\Delta}) = \beta \sum_{x \in \Delta} \sigma_{\Delta}(x) \Big( \tilde{h} - h + \frac{1}{2} \sum_{y \in A_{x}^{\Pi} \cap \Delta} F(x, y) \sigma_{\Delta}(y) + \sum_{y \in A_{x}^{\Pi} \cap \Delta^{c}} F(x, y) \bar{\sigma}_{\Delta^{c}}(y) \Big)$$

then  $C^{\text{II}} := \beta(\tilde{h} - h)$  and

$$\kappa_{\mathrm{II}} := \beta \sum_{x \in \Delta} \sigma_{\Delta}(x) \left( \frac{1}{2} \sum_{y \in A_{x}^{\mathrm{II}} \cap \Delta} F(x, y) \sigma_{\Delta}(y) + \sum_{y \in A_{x}^{\mathrm{II}} \cap \Delta^{c}} F(x, y) \bar{\sigma}_{\Delta^{c}}(y) \right)$$

We bound  $\kappa_{\rm II}$  based on the following:

$$\begin{split} \sum_{x \in \Delta} \sum_{y \in A_x^{\mathrm{II}} \cap \Delta} F(x,y) &\leq |\Delta| \sum_{y \in A_x^{\mathrm{II}}} F(0,y) = |\Delta| \sum_{y \in A_0^{\mathrm{II}}} \frac{a}{y^2} \\ &= \gamma |\Delta| \sum_{y \in A_0^{\mathrm{II}}} \frac{\gamma a}{(\gamma y)^2} \leq C \gamma |\Delta| \end{split}$$
 (F.3)

for some constant C arises from  $\sum_{y \in A_0^{\text{II}}} \frac{a}{y^2} < \infty$ . Then  $\kappa_{\text{II}} \leq 2C\gamma |\Delta|$  and the UQ bounds for long range perturbation with  $\beta(\tilde{h} - h) < 1$  are

$$(F.4) \qquad \pm E_{\tilde{q}_{\Delta}}[m(\sigma_{\Delta})] \leq \frac{1}{1 - \beta(\tilde{h} - h)} \inf_{\lambda > 0} \left\{ \frac{P_{h \pm \frac{\lambda}{\beta}, \beta, \mathbf{J}}^{\Delta, \gamma} - P_{h, \beta, \mathbf{J}}^{\Delta, \gamma}}{\lambda/\beta} + \frac{\beta}{\lambda} 2C\gamma \right\}$$

In the LP-limit we get

$$(F.5) \pm M(\tilde{\mathbf{J}}^F, \beta, \tilde{h}) \leq \frac{1}{1 - \beta(\tilde{h} - h)} \inf_{\lambda > 0} \left\{ \frac{p_{h \pm \frac{\lambda}{\beta}, \beta, \mathbf{J}} - p_{h, \beta, \mathbf{J}}}{\lambda/\beta} \right\}$$

### REFERENCES

- [1] M. AIZENMAN, J. CHAYES, L. CHAYES, C. NEWMAN, Discontinuity of the magnetization in one dimensional  $\frac{1}{|x-y|^2}$  percolation, Ising and Potts models, J. Stat. Phys., 50 (1988), pp. 1–40.
- [2] R. Atar, K. Chowdhary, P. Dupuis, Robust bounds on risk-sensitive functionals via Rényi divergence, SIAM/ASA J. Uncertain. Quantif., 3 (1) (2015), pp. 18-33.
- [3] Y. Bahri, J. Kadmon, J. Pennington, S. S. Schoenholz, J. Sohl-Dickstein, S. Ganguli Statistical Mechanics of Deep Learning, Annual Review of Condensed Matter Physics 2020 11:1, 501–528
- [4] U. BASU, M. KRÜGER, A. LAZARESCU, C. MAES, Frenetic aspects of second order response, Phys. Chem. Chem. Phys. 17 (9) (2015), pp. 6653–6666.
- [5] R. J. Baxter, Exactly Solved Models in Statistical Mechanics, Courier Corporation, 2007.
- [6] J. BIRRELL AND L. REY-BELLET, Uncertainty quantification for Markov processes via variational principles and functional inequalities, arXiv:1812.05174 (2018).
- [7] J. BIRRELL AND L. REY-BELLET Concentration Inequalities and Performance Guarantees for Hypocoercive MCMC Samplers, arXiv:1907.11973 (2019).
- [8] J. BIRRELL, P. DUPUIS, M. A. KATSOULAKIS, L. REY-BELLET, J. WANG, Distributional Robustness and Uncertainty Quantification for Rare Events, arXiv:1911.09580 (2019).
- [9] J. BIRRELL, M. A. KATSOULAKIS, Y. PANTAZIS, Optimizing variational representations of divergences and accelerating their statistical estimation, arXiv e-prints, (2020), arXiv:2006.08781, https://arxiv.org/abs/2006.08781.

- [10] A. BOVIER, Statistical Mechanics, Extreme Values, and Disordered Systems, University Lecture Notes.
- [11] A. BOVIER, M. ZAHRADNÍK, The low-temperature phase of Kac-Ising models, J. Stat. Phys., 87 (1997), pp. 311–332.
- [12] M. CASSANDRO, I. MEROLA, M. E. VARES, Study of a Long Range Perturbation of a One-Dimensional Kac Model, J. of Stat Phys., 142 (2011), pp. 487–509.
- [13] M. CASSANDRO, E. PRESUTTI, Phase transitions in Ising systems with long but finite range interactions. Markov Process. Related Fields, 2 (1996), pp. 241–262.
- [14] H. CHAN AND A. DARWICHE. Sensitivity analysis in Markov networks. In International Joint Conference on Artificial Intelligence, (2005).
- [15] K. CHOWDHARY AND P. DUPUIS, Distinguishing and integrating aleatoric and epistemic variation inuncertainty quantification, ESAIM Math. Model. Numer. Anal., 47 (2013), pp. 635– 662.
- [16] DARWICHE A. Modeling and reasoning with Bayesian networks. New York: Cambridge University Press ELT, 2009.
- [17] J. DE BOCK, A. ANTONUCCI, AND C. P. DE CAMPOS. Global sensitivity analysis for MAP inference in graphical models. In Neural Information Processing Systems, 27 (2014), pp. 2690-2698
- [18] S. Della Pietra, V. Della Pietra, J. Lafferty. Inducing features of random fields. IEEE Transactions on Pattern Analysis and Machine Intelligence, 19 (1997), pp. :380–392,
- [19] G. DIEZEMANN, Nonlinear response theory for Markov processes: simple models for glassy relaxation, Phys. Rev. E, 85 (2012), pp. 051502.
- [20] P. DOMINGOS, STATISTICAL RELATIONAL LEARNING, Tutorial in ICML, Oregon State University-Corvallis, OR, USA, https://icml.cc/Conferences/2007/tutorials.html, (2007).
- [21] P. Domingos, D. Lowd, Markov Logic: An Interface Layer for Artificial Intelligence, Morgan & Claypool (2009).
- [22] P. Dupuis, R. Ellis, A Weak Convergence Approach to the Theory of Large Deviations, Wiley Series in Probability and Statistics, 1997.
- [23] P. Dupuis, M. A. Katsoulakis, Y. Pantazis, P. Plecháč, Path-space information bounds for uncertainty quantification and sensitivity analysis of stochastic dynamics. SIAM/ASA J. Uncertain. Quantif., 4(1) (2016), pp. 80–111.
- [24] P. DUPUIS, M. A. KATSOULAKIS, Y. PANTAZIS, L. REY-BELLET, Sensitivity Analysis for Rare Events based on Rényi Divergence, to appear Ann. Appl. Probab.
- [25] P. Dupuis, Y. Mao, Formulation and properties of a divergence used to compare probability measures without absolute continuity, arXiv:1911.07422 (2019).
- [26] F.J. DYSON, Existence of phase transition in a one-dimensional Ising ferromagnetic, Comm. Math. Phys., 12 (1969), pp. 91–107.
- [27] F.J. DYSON, Non-existence of spontaneous magnetization in a one-dimensional Ising ferromagnet, Math. Phys., 12 (1969), pp. 212–215.
- [28] F.J. DYSON, An Ising ferromagnet with discontinuous long-range order, Comm. Math. Phys., 21 (1971), pp. 269–283.
- [29] M. S. Eldred, B. M. Adams, D. M. Gay, L. P. Swiler, K. Haskell, W. J. Bohnhoff, J. P. Eddy, W. E. Hart, J. Paul Watson, P. D. Hough, and T. G. Kolda. Dakota, A multilevel parallel object-oriented framework for design optimization, parameter estimation, uncertainty quantification, and sensitivity analysis: Version 5.0 user's manual, sandia. Technical report, Sandia, (2009).
- [30] J. Feng, J. L. Lansford, M. A. Katsoulakis, D. G. Vlachos, Explainable and trustworthy artificial intelligence for correctable modeling in chemical sciences, Sci. Adv., 6 (2020), eabc3204.
- [31] S. FRIEDLI, Y. VELENIK, Statistical Mechanics of Lattice Systems: A Concrete Mathematical Introduction, Cambridge University Press, 2017.
- [32] N. FRIEDMAN, I. NACHMAN, D. PEER, Learning bayesian network structure from massive datasets: The sparse candidate algorithm, UAI' 99: Proceedings of the Fifteenth conference on Uncertainty in artificial intelligence, July 1999, pp. 206–215.
- [33] J. FRÖHLICH, T. SPENCER, The phase transition in the one-dimensional Ising model with \(\frac{1}{r^2}\) interaction energy. Commun. Math. Phys., 84 (1982), pp. 87–101.
- [34] G. GALLAVOTTI, Statistical Mechanics, a Short Treatise, Text and Monographs in Physics, Springer, Berlin, 1999.
- [35] L. EL GHAOUI, M. OKS, F. OUSTRY, Worst-case value-at-risk and robust portfolio optimization: A conic programming approach. Oper. Res., 51(4) (2003,), pp. 543–556.
- [36] Z. GHAHRAMANI, Probabilistic Machine Learning and Artificial Intelligence, Nature, 521(7553) 2015, pp. 452-459.

- [37] P. GLASSERMAN, X. Xu, 2014. Robust Risk Measurement and Model Risk, Quant. Finance 14, 1 (2014), pp. 29–58.
- [38] I. GOODFELLOW, Y. BENGIO, A. COURVILLE, Deep Learning, MIT press, 2016.
- [39] K. GOURGOULIAS, M. A. KATSOULAKIS, L. REY-BELLET, J. WANG, How Biased Is Your Model? Concentration Inequalities, Information and Model Bias, IEEE Trans. Inform. Theory, 66 (2020), pp. 3079-3097.
- [40] G. R. GRIMMETT, Probability on graphs: random processes on graphs and lattices, Institute of Mathematical Statistics Textbooks 1, Cambridge University Press, 2010.
- [41] E. J. Hall, S. Taverniers, M. A. Katsoulakis, D. M. Tartakovsky, GINNs: Graph-Informed Neural Networks for Multiscale Physics, arXiv:2006.14807 (2020).
- [42] J. M. Hammersley, P. Clifford, Markov fields on finite graphs and lattices, 1971.
- [43] AJ. Hartemink, DK. Gifford, TS. Jaakkola, RA. Young, Using graphical models and genomic expression data to statistically validate models of genetic regulatory networks, Pacific Symposium on Biocomputing, Hawaii, 2001.
- [44] T. HASTIE, R. TIBSHIRANI, J. FRIEDMAN, The Elements of Statistical Learning: Data Mining, Inference, and Prediction, Springer Series in Statistics. Springer Verlag, 2001.
- [45] X. HE, R. S. ZEMEL, M. A. CARREIRA-PERPIÑIAN, Multiscale conditional random fields for image labeling. Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004., Washington, DC, USA, 2004, pp. II-II, doi: 10.1109/CVPR.2004.1315232.
- [46] M. A. KATSOULAKIS, L. REY-BELLET, J. WANG, Scalable information inequalities for uncertainty quantification, J. Comp. Phys., 336 (2017), pp. 513–545.
- [47] DIEDERIK P. KINGMA, M. WELLING, Auto-Encoding Variational Bayes. In The 2nd International Conference on Learning Representations (ICLR), 2013.
- [48] C. Kipnis, C. Landim, Scaling limits of interacting particle systems, Springer-Verlag Vol. 320, 1999.
- [49] D. KOLLER, N. FRIEDMAN, Probabilistic Graphical Models: Principles and Techniques, MIT Press, 2009.
- [50] R. G. KRISHNAN, U. SHALIT, D. SONTAG, Structured Inference Networks for Nonlinear State Space Models, arXiv:1609.09869, 2016.
- [51] L. LANDAU, E. LIFSHITZ, Perspectives in Theoretical Physics (ed L. P.Pitaevski) pp. 287–297, Pergamon, 1992
- [52] S. Lauritzen, Graphical Models. Oxford University Press, 1996. ISBN: 0-19-852219-3.
- [53] J. LI AND D. XIU, Computation of failure probability subject to epistemic uncertainty, SIAM J. Sci. Comput., 34 (2012), pp. A2946–A2964.
- [54] J. MOUSSOURIS, Gibbs and Markov Random Systems with Constraints, J. Stat. Phys., 10 (1974), pp. 11–33. issn: 0022-4715. doi: 10.1007/BF01011714.
- [55] K. P. Murphy, Machine Learning: A Probabilistic Perspective, MIT Press, 2012.
- [56] FUCHUN PENGAND, A. McCallum, Accurate Information Extraction from Research Papers using Conditional Random Fields, HLT-NAACL, 2004.
- [57] E. Presutti, Scaling Limits in Statistical Mechanics and Microstructures in Continuum Mechanics, Springer, 2000.
- [58] E. Presutti, From equilibrium to nonequilibrium statistical mechanics. Phase transitions and the Fourier law, Braz. J. Probab. Stat., 29, Number 2 (2015), pp. 211-281.
- [59] D. RUELLE, A review of linear response theory for general differentiable dynamical systems.
   Nonlinearity, 22 (2009), pp. 855–870, doi:10.1088/0951-7715/22/4/009.
- [60] A. SALTELLI, M. RATTO, T. ANDRES, F. CAMPOLONGO, J. CARIBONI, D. GATELLI, M. SAISANA, S. TARANTOLA, Global sensitivity analysis: the primer, John Wiley & Sons, 2008.
- [61] K. Sato, Y. Sakakibara, RNA secondary structural alignment with conditional random fields, Bioinformatics, 2005 Sep 1;21 Suppl 2:ii237–42. doi: 10.1093/bioinformatics/bti1139. PMID: 16204111.
- [62] B. Simon, The Statistical Mechanics of Lattice Gases, Vol. 1, Princeton University Press, 2014.
- [63] R. C. SMITH, Uncertainty Quantification: Theory, Implementation, and Applications, SIAM Computational Science & Engineering Series: Philadelphia, PA, USA, 2014, pp. 382.
- [64] B. TASKAR, P. ABBEEL, D. KOLLER, Discriminative probabilistic models for relational data, In Eighteenth Conference on Uncertainty in Artificial Intelligence (UAI02), (2002), pp. 485– 494, Edmonton, Canada.
- [65] S. TAVERNIERS, F. J. ALEXANDER, D. M. TARTAKOVSKY, Noise propagation in hybrid models of nonlinear systems: The Ginzburg-Landau equation. J. Comput.Phys., 262 (2014), pp. 313–324, doi: https://doi.org/10.1016/j.jcp.2014.01.015.
- [66] S. TAVERNIERS, E. J. HALL, M. A. KATSOULAKIS, D. M. TARTAKOVSKY, Mutual Information for Explainable Deep Learning of Multiscale Systems, arXiv:2009.04570, 2020.

- [67] A. B. TSYBAKOV, Introduction to Nonparametric Estimation, Springer Science & Business, Media, 2008.
- [68] M. Tuckerman, Statistical mechanics: theory and molecular simulation, Oxford university
- press, 2010.
  [69] K. UM, E. J. HALL, M. A. KATSOULAKIS, D. M. TARTAKOVSKY, Causality and Bayesian Network PDEs for multiscale representations of porous media, J. Comput. Phys., 394 (2019),
- [70] K. Um, E. J. Hall, M. A. Katsoulakis, D. M. Tartakovsky, Causality and Bayesian Network PDEs for multiscale representations of porous media, J. Comput. Phys., 394 (2019), pp. 658–678.
- [71] M. Wainwright, M. Jordan, Graphical models, exponential families, and variational inference. Technical Report 649, Department of Statistics, University of California, Berkeley (2003)