

---

# Is Parameter Learning via Weighted Model Integration Tractable?

---

Zhe Zeng\*<sup>1</sup> Paolo Morettin\*<sup>2</sup> Fanqi Yan<sup>3</sup> Antonio Vergari<sup>1</sup> Andrea Passerini<sup>4</sup> Guy Van den Broeck<sup>1</sup>

<sup>1</sup>University of California, Los Angeles, {zhezeng, aver, guyvdb}@cs.ucla.edu

<sup>2</sup>KU Leuven, paolo.morettin@kuleuven.be

<sup>3</sup>University of Texas, Austin, fanqi\_yan@lsec.cc.ac.cn

<sup>4</sup>University of Trento, Italy, andrea.passerini@unitn.it

## Abstract

Weighted Model Integration (WMI) is a recent and general formalism for reasoning over hybrid continuous/discrete probabilistic models with logical and algebraic constraints. While many works have focused on inference in WMI models, the challenges of learning them from data have received much less attention. Our contribution is twofold. First, we provide novel theoretical insights on the problem of estimating the parameters of these models from data in a tractable way, generalizing previous results on maximum-likelihood estimation (MLE) to the broader family of log-linear WMI models. Second, we show how our results on WMI can characterize the tractability of inference and MLE for another widely used class of probabilistic models, Hinge Loss Markov Random Fields (HL-MRFs). Specifically, we bridge these two areas of research by reducing marginal inference in HL-MRFs to WMI inference, and thus we open up new interesting applications for both model classes.

## 1 INTRODUCTION

Solving many complex AI tasks require both logical and numerical reasoning in uncertain environments. Probabilistic inference in hybrid continuous/discrete settings poses significant challenges, in particular when logical and algebraic constraints are considered, such as laws of physics or safety requirements. In many discrete probabilistic models, reducing marginal inference to Weighted Model Counting (WMC) [Sang et al., 2005] is a state-of-the-art technique, thanks to its support for arbitrary propositional logic constraint. Recently, Weighted Model Integration (WMI) [Belle et al., 2015] emerged as a generalization of WMC to hybrid domains. By extending propositional logic with algebraic relations over continuous variables and generalizing

the weight function to piecewise densities, WMI unifies marginal inference in a large class of probabilistic models under the same formalism.

WMI-based inference has been receiving increasing interest in recent years, with advancements both on practical algorithms [Morettin et al., 2021] and on theoretical aspects [Zeng et al., 2020a]. While its theoretical groundwork was initially developed to address inference and parameter learning in hybrid Markov Random Fields, over the last five years it has been shown that a much larger class of problems can be reduced to WMI, such as marginal inference in Mixed Sum-Product Networks [Morettin et al., 2020], or verification of fairness properties in probabilistic programs [Albarghouthi et al., 2017]. Yet, many questions related to learning this class of distributions from data are still open.

In this paper, we investigate weight learning for WMI models in a principled way. First, we show under which conditions weight learning can be done in a tractable way and by doing this we greatly extend the class of WMI models amenable to efficient learning. Additionally, we show how a popular class of probabilistic models, namely Hinge Loss Markov Random Fields (HL-MRFs) [Bach et al., 2017], fall into this class by reducing marginal inference in HL-MRFs to WMI. Leveraging recent results on the tractability of WMI-based inference, we then characterize a subclass of HL-MRFs that admit tractable maximum likelihood estimation of their parameters and marginal inference.

## 2 BACKGROUND

**Notation.** Uppercase letters denote random variables ( $X, B$ ) and lowercase letters denote their assignments ( $x, b$ ). We use bold for sets of variables ( $\mathbf{X}, \mathbf{B}$ ), and their joint assignments ( $\mathbf{x}, \mathbf{b}$ ). We use capital Greek letters for logical formulas ( $\Gamma, \Delta$ ). Literals are atomic formulas or their negation, and are denoted using either  $\ell$  or lowercase Greek letters ( $\gamma, \delta$ ). We let  $\mathbf{x} \models \Delta$  denote the satisfaction of a

formula  $\Delta$  by an assignment  $\mathbf{x}$ . Its corresponding indicator function is  $\llbracket \mathbf{x} \models \Delta \rrbracket$ .

A state-of-the-art approach for answering probabilistic queries in many discrete models reduces the problem to *weighted model counting* (WMC), i.e. the task of computing the weighted sum of the models (solutions) of a propositional formula  $\Delta$ :

$$\text{WMC}(\Delta, \mathcal{W}) = \sum_{\mu \models \Delta} \mathcal{W}(\mu) \quad (1)$$

Typically, it is assumed that the weight of a model factorizes as the product of non-negative constant weights associated with each literal (an atom or its negation) in the solution, i.e.  $\mathcal{W}(\mu) = \prod_{\ell \in \mathcal{L}} w_\ell^{\llbracket \mu \models \ell \rrbracket}$ , where  $\llbracket \cdot \rrbracket : \mathbb{B} \rightarrow \{0, 1\}$  is the indicator function,  $\mathcal{L} = \text{Atoms}(\Delta) \cup \{\neg A \mid A \in \text{Atoms}(\Delta)\}$  denotes the set of all the literals and  $w_\ell \in \mathbb{R}^+$  is the positive weight of a literal  $\ell$  (unweighted literals are assumed to have weight 1).

**Example 1** Consider the formula  $\Delta = (A \rightarrow B)$  and weights  $w_A = 2, w_{\neg A} = 3, w_B = 5$ , then:

$$\text{WMC}(\Delta, \mathcal{W}) = \underbrace{\mathcal{W}(A \wedge B)}_{2 \cdot 5} + \underbrace{\mathcal{W}(\neg A \wedge B)}_{3 \cdot 5} + \underbrace{\mathcal{W}(\neg A \wedge \neg B)}_{3 \cdot 1} = 28.$$

The pair  $\langle \Delta, \mathcal{W} \rangle$  encodes an unnormalized joint distributions over the propositional variables in  $\Delta$ , thus the normalized probability of a formula  $\Delta_Q$  given evidence expressed as another formula  $\Delta_E$  is computed as:

$$\text{Pr}_{\langle \Delta, \mathcal{W} \rangle}(\Delta_Q \mid \Delta_E) = \frac{\text{WMC}(\Delta \wedge \Delta_Q \wedge \Delta_E, \mathcal{W})}{\text{WMC}(\Delta \wedge \Delta_E, \mathcal{W})}. \quad (2)$$

Algorithmic advances in both exact and approximate WMC motivated its generalization to hybrid continuous/discrete settings [Belle et al., 2015].

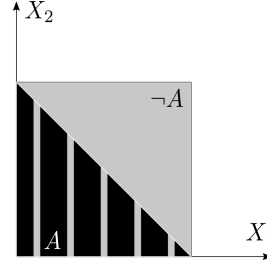
Given a set of continuous random variables  $\mathbf{X}$  and a set of Boolean variables  $\mathbf{A}$ , the support of the joint probability  $\text{Pr}(\mathbf{X}, \mathbf{A})$  is encoded by means of a *satisfiability modulo theories* (SMT) [Barrett et al., 2010] formula, i.e. a (typically quantifier-free) expression containing both propositional and *theory* atoms connected with the usual logical connectives. Specifically, the theory atoms encode algebraic constraints over  $\mathbf{X}$ , often restricted to the theory of *linear algebra over reals* ( $\mathcal{LRA}$ ), where atoms have form  $(\mathbf{c}^T \mathbf{X} \leq b)$ .

**Example 2** The following SMT- $\mathcal{LRA}$  formula:

$$\Delta = (0 \leq X_1) \wedge (X_1 \leq 1) \wedge (0 < X_2) \wedge (X_2 \leq 1) \\ \wedge (A \rightarrow (X_1 + X_2 \leq 1))$$

has 3 satisfying truth assignments ( $X_1, X_2 \in [0, 1]$  are always true):

$$\left\{ \overbrace{A \wedge (X_1 + X_2 \leq 1) \wedge \dots}^{\mu_1}, \overbrace{\neg A \wedge (X_1 + X_2 \leq 1) \wedge \dots}^{\mu_2}, \right. \\ \left. \overbrace{\neg A \wedge \neg (X_1 + X_2 \leq 1) \wedge \dots}^{\mu_3} \right\}$$



for a total (unweighted) volume of  $\frac{3}{2}$ .

In addition to the usual constant weights associated with propositional literals, algebraic literals  $\ell$  on variables  $\mathbf{X}_\ell \subseteq \mathbf{X}$  can be mapped to functions over  $\mathbf{X}_\ell$ . Thus, the weight function  $\mathcal{W}(\mu, \mathbf{x}) = \prod_{\ell \in \mathcal{L}} w_\ell(\mathbf{x}_\ell)^{\llbracket \mu, \mathbf{x} \models \ell \rrbracket}$  is a piecewise function, associating an

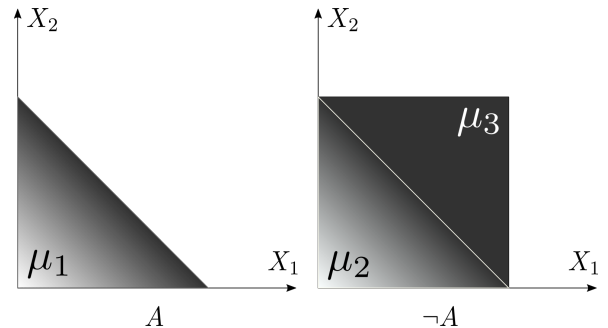
unnormalized density to each solution  $\mu \models \Delta$ . In this context  $\mathcal{L}$  denotes the set of weighted literals, which does not necessarily have to appear in  $\Delta$ . Intuitively, a truth assignment to the literals  $\mathcal{L}$  univocally defines a region of the piecewise distribution. Similarly to WMC, the weighted model integral is computed by summing over the solutions of the SMT formula and, additionally, by integrating the density function in the continuous domain:

$$\text{WMI}(\Delta, \mathcal{W}) = \sum_{\mu \models \Delta} \int_{\mathbf{x} \models \mu} \mathcal{W}(\mu, \mathbf{x}) d\mathbf{x} \quad (3)$$

Piecewise polynomials, which admit closed-form solutions when integrated over arbitrary polytopes, are the most investigated class of joint density in the WMI literature.

**Example 3** Consider the formula in Example 2 and a single weighted literal  $\ell = (X_1 + X_2 \leq 1)$  with  $w_\ell(X_1, X_2) = X_1 + X_2$ , i.e.  $\mathcal{W}(X_1, X_2) = (X_1 + X_2)^{\llbracket X_1 + X_2 \leq 1 \rrbracket}$ . Then, the weighted model integral is:

$$\text{WMI}(\Delta, \mathcal{W}) = \int_{\mu_1} X_1 + X_2 dX_1 dX_2 \\ + \int_{\mu_2} X_1 + X_2 dX_1 dX_2 \\ + \int_{\mu_3} 1 dX_1 dX_2 = 2 \cdot \frac{1}{3} + \frac{1}{2}$$



W.l.o.g., for notational convenience from here on we will focus on WMI problems on continuous variables only. This is possible since any WMI problem on continuous and Boolean variables can be reduced in polytime to a new WMI problem on continuous variables only, without changing its inference complexity class, by properly introducing

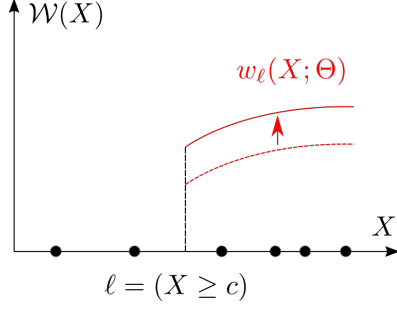


Figure 1: Graphical representation of the parameter estimation task considered in this work.

auxiliary variables to account for Boolean variables [Zeng and Van den Broeck, 2019]. We remark that the results on parameter learning reported in the following sections hold for hybrid distributions too. In fact, weighted propositional variables would contribute with piecewise constant factors only, a subcase of our generalized theory that was previously investigated by Belle et al. [2015]. We define the normalized density of any data point  $\mathbf{x}$  as:

$$p(\mathbf{x}) = \frac{\prod_{\ell \in \mathcal{L}} w_{\ell}(\mathbf{x})^{\llbracket \mathbf{x} \models \ell \rrbracket}}{\text{WMI}(\Delta, \mathcal{W})}, \quad \mathbf{x} \models \Delta \quad (4)$$

Characterizing the dependency structure among the variables in a model is crucial for characterizing the tractability of probabilistic inference. This structure can be captured with the notion of *primal graph*.

**Definition 2.1.** (*Primal Graph*) Given a model  $\langle \Delta, \mathcal{W} \rangle$ , its primal graph  $\mathcal{G}_{\langle \Delta, \mathcal{W} \rangle} = (\mathcal{V}, \mathcal{E})$  is the undirected graph whose vertex set  $\mathcal{V}$  is the index set of variables in formula  $\Delta$  and  $\mathcal{L}$ , and whose edge set  $\mathcal{E}$  has edge  $i - j$  iff variable  $X_i$  and variable  $X_j$  appear together in one clause  $\Gamma \in \Delta$  or in one weighted literal  $\ell \in \mathcal{L}$ .

### 3 MAXIMUM LIKELIHOOD ESTIMATION OF WEIGHTED MODEL INTEGRATION PARAMETERS

We now consider the problem of estimating the parameters of a structured, piecewise distribution from data. Specifically, we are concerned with learning the parameters in the literal weights  $w_{\ell}$  for  $\ell \in \mathcal{L}$  as illustrated in Figure 1, assuming that the SMT formula  $\Delta$  and the literal set  $\mathcal{L}$  is fixed and that a dataset consisting of i.i.d. fully observed samples is provided.

We adopt the canonical parameter estimation approach, maximum likelihood estimation (MLE), which gives estimation for parameters by solving the optimization problem below,

$$\Theta^* = \arg \max_{\Theta} L(\Theta) \quad (5)$$

where  $L(\Theta)$  is the log-likelihood given the dataset, i.e.,

$$L(\Theta) = \log \prod_{\mathbf{x} \in \mathcal{D}} p(\mathbf{x}; \Theta).$$

Plugging in the point-wise density  $p$  in Eq. 4 we get:

$$L(\Theta) = \sum_{\mathbf{x} \in \mathcal{D}} \sum_{\ell \in \mathcal{L}} \log w_{\ell}(\mathbf{x})^{\llbracket \mathbf{x} \models \ell \rrbracket} - |\mathcal{D}| \cdot \log \text{WMI}(\Delta, \mathcal{W}; \Theta). \quad (6)$$

#### 3.1 CONVEXITY OF MLE

In general, the optimization problem in Equation 5 does not have a closed form solution. For this reason, we consider optimizing our objective via iterative methods. Before that, we are interested in analyzing the convexity property of the log-likelihood  $\frac{\partial L(\Theta)}{\partial \Theta}$  for characterizing convergence of the optimization problem. Given the SMT formula  $\Delta$  and the set of weighted literals  $\mathcal{L}$ , we explore under what conditions the log-likelihood objective is concave, which guarantees convergence to the global optimum when iterative methods such as gradient ascent are used [Nesterov and Nemirovskii, 1994].

As reported in the original WMI paper [Belle et al., 2015], MLE of the parameters in weight functions can be solved optimally using standard convex optimization tools for constant weight functions, that is, when the function  $w_{\ell}(\mathbf{x})$  is constant for any  $\ell \in \mathcal{L}$ . We generalize this result to weight functions that are *log-linear with respect to their parameters*. Specifically, the weight functions take their form as

$$w_{\ell}(\mathbf{x}_{\ell}; \Theta_{\ell}) = \exp\{\Theta_{\ell}^{\top} \mathbf{f}_{\ell}(\mathbf{x}_{\ell})\}$$

where  $\mathbf{f}_{\ell}(\cdot) = (f_{\ell}^1(\cdot), \dots, f_{\ell}^K(\cdot))^{\top}$  is a vector of features defined over variable subset  $\mathbf{X}_{\ell}$ . Log-linear functions are an expressive weight family that generalizes many functions adopted in previous works. Some example choices for feature functions  $\mathbf{f}(\cdot)$  could be 1) constant functions, in which case the WMI problem has constant weights, 2) monomials, in which case the WMI problem has exponentiated polynomial weights, 3) log-polynomials, in which case the WMI problem has polynomial weights.

Notice that when the weight functions are log-linear, the first term in the log-likelihood in Equation 6 is a linear combination of parameters, i.e.,

$$\sum_{\mathbf{x} \in \mathcal{D}} \sum_{\ell \in \mathcal{L}} \log w_{\ell}(\mathbf{x})^{\llbracket \mathbf{x} \models \ell \rrbracket} = \sum_{\mathbf{x} \in \mathcal{D}} \sum_{\ell \in \mathcal{L}} \llbracket \mathbf{x} \models \ell \rrbracket \cdot \Theta_{\ell}^{\top} \mathbf{f}_{\ell}(\mathbf{x}_{\ell}).$$

We thus investigate the convexity of the log-partition function term  $Z(\Theta) = \text{WMI}(\Delta, \mathcal{W}; \Theta)$ .

**Proposition 3.1.** *The log-partition function  $\log Z(\Theta)$  is convex if the weight functions are log-linear in their parameters.*

**Sketch of Proof** Denote the  $i$ -th parameter in weight function  $w_\ell$  as  $\Theta_i^\ell$ . The first-order derivative of the log-partition function  $\log Z(\Theta)$  with respect to a parameter  $\Theta_i^\ell$  is as follows.

$$\frac{\partial}{\partial \Theta_i^\ell} \log Z(\Theta) = \mathbb{E}_{\mathbf{x} \sim p} \left[ \frac{\partial}{\partial \Theta_i^\ell} \log w_\ell(\mathbf{x}) \mathbb{1}^{\{\mathbf{x} \models \ell\}} \right]$$

For brevity, denote  $\frac{\partial}{\partial \Theta_i^\ell} \log w_\ell(\mathbf{x}) \mathbb{1}^{\{\mathbf{x} \models \ell\}}$  by notation  $a_i^\ell$ . Further, the second-order derivative of the log-partition function with respect to parameters  $\Theta_i^\ell$  and  $\Theta_j^{\ell'}$  is as follows.

$$\frac{\partial^2}{\partial \Theta_i^\ell \partial \Theta_j^{\ell'}} \log Z(\Theta) = \mathbb{E}_{\mathbf{x} \sim p} [a_i^\ell a_j^{\ell'}]$$

Consider  $a_i^\ell$  as random variables distributed according to distributions defined by parameter  $\Theta$ . Then the Hessian matrix of  $\log Z(\Theta)$  is the co-variance matrix of random variables  $a_i^\ell$  which is positive definite. Therefore the log-partition function  $\log Z(\Theta)$  is a convex function in parameters  $\Theta$ .  $\square$

**Proposition 3.2.** *The log-likelihood function  $L(\Theta; \mathcal{D})$  is concave if the weight functions are log-linear in their parameters.*

*Proof.* This result directly follows from the linearity of the first term and the convexity of  $\log Z(\Theta)$  (Prop. 3.1).  $\square$

The concavity of the log-likelihood function allows the problem of computing the maximum likelihood to be formulated as a convex minimization problem and optimally solved with gradient descent methods.

### 3.2 TRACTABILITY OF GRADIENTS

The next question that naturally follows is about the tractability of the gradient computation. The partial derivative of the log-likelihood with respect to a parameter  $\Theta_\ell^k$  is:

$$\begin{aligned} \frac{\partial}{\partial \Theta_\ell^k} L(\Theta) &= \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\mathbb{1}^{\{\mathbf{x} \models \ell\}} \cdot f_\ell^k(\mathbf{x})] \\ &\quad - \mathbb{E}_{\mathbf{x} \sim p(\mathbf{X}; \Theta)} [\mathbb{1}^{\{\mathbf{x} \models \ell\}} \cdot f_\ell^k(\mathbf{x})] \end{aligned} \quad (7)$$

Notice that the first expectation is a weighted count of the dataset  $\mathcal{D}$  and thus can be computed in time complexity  $\Theta(|\mathcal{D}|)$ . The second expectation requires the computation of the partition function, i.e., a WMI problem, which is intractable in general, and therefore, approximation algorithms are needed to compute the gradient updates. Recently, Zeng et al. [2020a] characterized the tractable class of WMI problems in terms conditions on both their structure and in the form of the weight functions [Zeng et al., 2020b]. A family of weight functions satisfies *tractable weight conditions* (TWC) iff:

- i) it is closed under product;
- ii) it admits efficient computation of antiderivatives;
- iii) it is closed under definite integration over each variable.

Some examples when the weight function families satisfy TWCs include the cases where the features functions are constants, linear functions or log-polynomial functions. Next we present the previous result on the tractability of WMI problems which is necessary for deriving the tractability of the partial derivatives in MLE.

**Proposition 3.3.** [Zeng et al., 2020a] *Let  $\mathcal{WMJ}(\Omega, \log(n), tr)$  be the class of WMI problems with primal graph  $\mathcal{G}$  with diameter of size  $\Theta(\log(n))$  and treewidth  $tr$ . When  $\mathcal{WMJ}(\Omega, \log(n), tr)$  has a parametric weight function family  $\Omega$  that satisfies the TWCs,  $\mathcal{WMJ}(\Omega, \log(n), tr)$  is a tractable WMI class for inference if-and-only-if treewidth  $tr = 1$ .*

**Theorem 3.4.** *If the WMI problem  $\text{WMI}(\Delta, \mathcal{W})$  is in the tractable WMI problem class  $\mathcal{WMJ}(\Omega, \log(n), 1)$  and the feature functions  $f_\ell^k(\mathbf{x})$  are in a function family that satisfies the TWCs, then the partial derivative in Equation 7 can be tractably computed.*

*Proof.* Given that the first term in Equation 7 can be computed linearly in  $|\mathcal{D}|$ , it suffices to show that the second term is a ratio of two tractable WMI problems:

$$\mathbb{E}_{\mathbf{x} \sim p(\mathbf{X}; \Theta)} [\mathbb{1}^{\{\mathbf{x} \models \ell\}} \cdot f_\ell^k(\mathbf{x})] = \frac{\text{WMI}(\Delta \wedge \ell, \mathcal{W}')}{\text{WMI}(\Delta, \mathcal{W})},$$

where for all  $\ell^*, w_{\ell^*}' \in \mathcal{W}'$ ,  $w_{\ell^*} \in \mathcal{W}$ , if  $\ell^* \neq \ell$ , weight function  $w_{\ell^*}' = w_{\ell^*}$ ; otherwise, the weight function is defined as  $w_{\ell^*}'(\mathbf{x}) = w_{\ell^*}(\mathbf{x}) f_\ell^k(\mathbf{x})$ . The tractability of  $\text{WMI}(\Delta \wedge \ell, \mathcal{W}')$  follows from the tractability of  $\text{WMI}(\Delta, \mathcal{W})$ , given that they have the same primal graph and that  $w_{\ell^*}'(\mathbf{x})$  is the product of two functions that satisfies the TWC. Therefore, the partial derivative in Equation 7 can be efficiently computed.  $\square$

With the results above, we characterized which families of distributions that are amenable to WMI-based inference admit both tractable and exact MLE of their parameters. In the following, we present a novel reduction of marginal inference in HL-MRFs to WMI and generalize these results to this popular class of models.

## 4 MARGINAL INFERENCE FOR HL-MRFs VIA WMI

Hinge-loss Markov Random fields (HL-MRFs) [Bach et al., 2017] are a recently proposed statistical relational learning framework used to model highly structured data. This family of models allows tractable and exact MAP inference when

weighted constraints are defined using Łukasiewicz logic, a fuzzy relaxation of Boolean logic. However, exact marginal inference for HL-MRFs is intractable in general [Embar et al., 2019]. Therefore, weight learning via MLE for HL-MRFs approximates expectations with MAP states. We investigate the connection between HL-MRFs and WMI problems and propose a reduction from HL-MRFs to WMI models. This reduction is profound in two senses. Firstly, it makes the HL-MRFs amenable to the MLE-based parameter learning approach as proposed in Section 3. Secondly, it allows the characterization of tractability of marginal inference for HL-MRFs with the conditions reported in Prop. 3.3. To the best of our knowledge, *this is the first class of tractable HL-MRFs for exact marginal inference.*

#### 4.1 HINGE-LOSS MARKOV RANDOM FIELDS

Before describing the reduction to WMI, we provide the main definitions in HL-MRFs. W.l.o.g, in the following we assume that all the random variables are unobserved.

**Definition 4.1.** (*potential*) Let  $\mathbf{X} = (X_1, \dots, X_n)$  be a vector of  $n$  variables with joint domain  $\mathcal{D} = [0, 1]^n$ . Let  $\phi = (\phi_1, \dots, \phi_m)$  be a vector of  $m$  continuous **potentials** of the form

$$\phi_j(\mathbf{x}) = (\max\{\ell_j(\mathbf{x}), 0\})^{p_j} \quad (8)$$

where  $\ell_j$  is a linear function of  $\mathbf{X}$  and  $p_j \in \{1, 2\}$ .

Typically,  $\ell_j$  are interpreted as clauses in Łukasiewicz logic, with  $\phi_j(\mathbf{x})$  representing the (possibly squared) distance to the satisfaction of the linear constraint  $\ell_j \leq 0$ .

**Definition 4.2.** (*constrained hinge-loss energy function*) For  $\mathbf{X} \in \mathcal{D}$ , given a vector of  $m$  non-negative free parameters, i.e., weights,  $\alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}^m$ , a **constrained hinge-loss energy function**  $f_\alpha$  is defined as follows.

$$f_\alpha(\mathbf{x}) = \sum_{j=1}^m \alpha_j \phi_j(\mathbf{x}). \quad (9)$$

This energy function quantifies the (weighted) satisfaction of the linear constraints above. These soft constraints are complemented with a set of hard constraints (i.e. with infinite weights) that are defined separately for convenience.

**Definition 4.3.** (*feasible set*) Let  $\mathbf{c} = (c_1, \dots, c_r)$  be a vector of  $r$  linear constraint functions associated with index sets denoting equality constraints  $\mathcal{E}$  and inequality constraints  $\mathcal{I}$ , which define the **feasible set**

$$\tilde{\mathcal{D}} = \left\{ (\mathbf{x}) \in \mathcal{D} \mid \begin{array}{l} \forall k \in \mathcal{E}, c_k(\mathbf{x}) = 0 \\ \forall k \in \mathcal{I}, c_k(\mathbf{x}) \leq 0 \end{array} \right\}. \quad (10)$$

Then, the probability density function is defined over the feasible set and is inversely proportional to the constrained hinge-loss energy function.

**Definition 4.4.** (*HL-MRF*) A **hinge-loss Markov random field**  $p$  over random variables  $\mathbf{X}$  is a probability density defined as follows: if  $\mathbf{x} \notin \tilde{\mathcal{D}}$ , then  $p(\mathbf{x}) = 0$ ; if  $\mathbf{x} \in \tilde{\mathcal{D}}$ , then

$$p(\mathbf{x}; \alpha) = \frac{1}{Z(\alpha)} \exp(-f_\alpha(\mathbf{x})) \quad (11)$$

where  $Z(\alpha)$  is the partition function, i.e.,

$$Z(\alpha) = \int_{\mathbf{x} | \mathbf{x} \in \tilde{\mathcal{D}}} \exp(-f_\alpha(\mathbf{x})) d\mathbf{x}. \quad (12)$$

The *templated HL-MRFs* can be similarly defined, where the potentials are defined by the templates denoted by  $\mathcal{T} = (t_1, \dots, t_S)$  where each  $t_s$  is the set of indices of the potentials defined by the template. The templates are associated with weights  $\mathbf{W} = (W_1, \dots, W_S)$ . Then the sum of potentials defined by a template  $t_s$  is  $\Phi_s(\mathbf{x}) = \sum_{i \in t_s} \phi_i(\mathbf{x})$ . For each potential  $\phi_i(\mathbf{x})$  with its index  $i$  in  $t_s$ , they are assigned weights to be the weight of the template, that is,  $\alpha_i = W_s$ . Together, the constrained hinge-loss energy function of a templated HL-MRF is as follows.

$$f_{\mathbf{W}}(\mathbf{x}) = \mathbf{W}^\top \Phi(\mathbf{x}) \quad (13)$$

where  $\Phi(\mathbf{x}) = (\Phi_1(\mathbf{x}), \dots, \Phi_S(\mathbf{x}))$ . The templated HL-MRFs are often used in the parameter learning when the weights are considered to be shared among potentials [Bach et al., 2017].

#### 4.2 FROM HL-MRFs TO WMI

This section presents the reduction from HL-MRFs to WMI models that bridge the two frameworks. It allows the inference on HL-MRFs to be amenable to the tractability results on WMI problems as well as WMI solvers.

**Theorem 4.5.** For any HL-MRF  $p$  over random variables  $\mathbf{X}$ , there exists a WMI problem  $\text{WMI}(\Delta, \mathcal{W})$  with per-literal weights whose WMI density denoted by  $p_\Delta$  equals to the HL-MRF  $p$ .

*Proof.* Given an HL-MRF  $p$ , construct an  $\text{SMT}(\mathcal{L}\mathcal{R}\mathcal{A})$  formula  $\Delta$  in CNF form as follows.

$$\Delta = \underbrace{\bigwedge_{i \in [n]} (0 \leq x_i) \wedge (x_i \leq 1)}_{\text{joint domain } \Delta_{\mathcal{D}}} \underbrace{\bigwedge_{k \in \mathcal{E}} (c_k(\mathbf{x}) = 0) \bigwedge_{k \in \mathcal{I}} (c_k(\mathbf{x}) \leq 0)}_{\text{hard constraints } \Delta_{\tilde{\mathcal{D}}}} \quad (14)$$

Then construct a set of per-literal weights  $\mathcal{W}$  with the set of literals  $\mathcal{L} = \{\ell_j(\mathbf{x}) \geq 0\}_{j \in [m]}$  as follows. For brevity, here the per-literal weight is denoted by  $\mathcal{W}_j$  with subscript

being the index of the  $\mathcal{LR}\mathcal{A}$  atom  $\ell_j$  in its associated literal  $\ell_j(\mathbf{x}) \geq 0$ .

$$\mathcal{W}_j(\mathbf{x}) = \exp(-\alpha_j \ell_j^{p_j}(\mathbf{x})) \quad (15)$$

With the above notations, the WMI density defined by Equation 4 is as follows.

$$p_\Delta(\mathbf{x}) = \frac{1}{Z_\Delta(\mathcal{W})} \prod_{j \in [m]} \mathcal{W}_j(\mathbf{x})^{\llbracket \mathbf{x} \models \ell_j(\mathbf{x}) \geq 0 \rrbracket}, \quad \mathbf{x} \models \Delta. \quad (16)$$

with

$$Z_\Delta(\mathcal{W}) = \int_{\mathbf{x} \models \Delta} \prod_{j \in [m]} w_j(\mathbf{x})^{\llbracket \mathbf{x} \models \ell_j(\mathbf{x}) \geq 0 \rrbracket} d\mathbf{x}. \quad (17)$$

Moreover, the HL-MRF  $P$  in Definition 4.4 has an equivalent formulation by using indicator function.

$$p(\mathbf{x}) = \frac{1}{Z(\boldsymbol{\alpha})} \llbracket \mathbf{x} \in \mathcal{D} \rrbracket \cdot \llbracket \mathbf{x} \in \tilde{\mathcal{D}} \rrbracket \cdot \exp(-f_\alpha(\mathbf{x})) \quad (18)$$

Next we show that the un-normalized HL-MRF  $p$  is proportional to the WMI density  $p_\Delta$ . Notice that  $\llbracket \mathbf{x} \in \mathcal{D} \rrbracket = \llbracket \mathbf{x} \models \Delta_{\mathcal{D}} \rrbracket$  and  $\llbracket \mathbf{x} \in \tilde{\mathcal{D}} \rrbracket = \llbracket \mathbf{x} \models \Delta_{\tilde{\mathcal{D}}} \rrbracket$  hold. Moreover, since sub-formula  $\Delta_S$  is always SAT, it holds that

$$\llbracket \mathbf{x} \in \mathcal{D} \rrbracket \cdot \llbracket \mathbf{x} \in \tilde{\mathcal{D}} \rrbracket = \llbracket \mathbf{x} \models \Delta_{\mathcal{D}} \rrbracket \cdot \llbracket \mathbf{x} \models \Delta_{\tilde{\mathcal{D}}} \rrbracket = \llbracket \mathbf{x} \models \Delta \rrbracket. \quad (19)$$

Further since it holds that  $\max\{\ell_j(\mathbf{x}), 0\} = \llbracket \ell_j(\mathbf{x}) \geq 0 \rrbracket \cdot \ell_j(\mathbf{x})$ , we could rewrite the potentials in HL-MRF using indicator functions.

$$\phi_j(\mathbf{x}) = \llbracket \mathbf{x} \models \ell_j(\mathbf{x}) \geq 0 \rrbracket \cdot \ell_j^{p_j}(\mathbf{x}) \quad (20)$$

By Equations 19 and Equation 20, the following equality on un-normalized HL-MRF  $p$  holds for any  $\mathbf{x} \models \Delta$ .

$$\begin{aligned} Z(\boldsymbol{\alpha})p(\mathbf{x}) &= \prod_{j \in [m]} \exp(-\alpha_j \phi_j(\mathbf{x})) \\ &= \prod_{j \in [m]} \exp(-\alpha_j \ell_j^{p_j}(\mathbf{x}))^{\llbracket \mathbf{x} \models \ell_j(\mathbf{x}) \geq 0 \rrbracket} \\ &= \prod_{j \in [m]} \mathcal{W}_j(\mathbf{x})^{\llbracket \mathbf{x} \models \ell_j(\mathbf{x}) \geq 0 \rrbracket} = Z_\Delta(\mathcal{W})p_\Delta(\mathbf{x}) \end{aligned}$$

Thus it holds that  $p(\mathbf{x}) \propto p_\Delta(\mathbf{x})$ . By the same arguments, it can be shown that HL-MRF  $p$  and WMI density  $p_\Delta$  have the same partition functions, which finishes the proof for  $p = p_\Delta$ .  $\square$

What immediately follows the reduction is the tractability of marginal inference on HL-MRFs. When the potentials in the HL-MRF model is linear, the resulting WMI problem has its weight function to be in the exponentiated linear function family that satisfies the TWCs [Zeng et al., 2020a]. This gives the first non-trivial class of tractable HL-MRFs for exact marginal inference so far. Specifically, the tractability of the HL-MRF model is characterized by model structures in the form of primal graphs.

**Definition 4.6.** (*Primal Graph of HL-MRFs*) Given an HL-MRF model  $p$ , its primal graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  is the undirected graph, whose vertex set  $\mathcal{V}$  is the index set of variables over which  $p$  is defined, and whose edge set  $\mathcal{E}$  has edge  $i - j$  iff variables  $X_i$  and variable  $X_j$  appear together either in a potential or in a constraint in the feasible set  $\tilde{\mathcal{D}}$ .

**Corollary 4.7.** Given a HL-MRF  $p$  whose potentials are as defined in Definition 4.1 with  $p_i = 1$ , if it has its primal graph with tree-width one and with diameter of  $\Theta(\log(n))$  with  $n = |\mathcal{X}|$ , then the marginal inference of  $p$  is tractable.

### 4.3 MLE FOR HL-MRFs

Besides the tractable marginal inference for HL-MRFs, what follows the reduction is that the parameter learning of WMI problems can be leveraged for the parameter learning of the HL-MRFs, and our analysis on the optimality and the tractability of MLE approach naturally generalize from WMI to HL-MRFs. To derive MLE-based parameter learning, we first compute the partial derivative of the log-likelihood for the templated HL-MRFs in Equation 13 as follows.

$$\begin{aligned} \frac{\partial}{\partial W_s} L(\boldsymbol{\Theta}) &= -\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} [\llbracket \mathbf{x} \models \wedge_{i \in t_s} \ell_i \rrbracket \cdot \Phi_s(\mathbf{x})] \\ &\quad + \mathbb{E}_{\mathbf{x} \sim p(\mathbf{X}; \boldsymbol{\Theta})} [\llbracket \mathbf{x} \models \wedge_{i \in t_s} \ell_i \rrbracket \cdot \Phi_s(\mathbf{x})] \end{aligned} \quad (21)$$

In previous work, since the expectation over the templated HL-MRFs is intractable in general, the MAP state is used to approximate the expectation [Bach et al., 2017]. However, our HL-MRFs-to-WMI reduction provides an alternative way to compute the expectation by using WMI solvers on the reduced WMI problem. Specifically, the expectation can be computed exactly and tractably when an exact WMI solver is applicable [Moretting et al., 2021]. Moreover, if the reduced WMI problem satisfies the assumptions in Theorem 3.4, the partial derivative for the templated HL-MRFs above can be computed exactly in tractable time. Below we show one of such cases.

**Proposition 4.8.** Given a HL-MRF  $p$  with the same assumptions as in Corollary 4.7, then the partial derivative in Equation 21 can be tractably computed.

## 5 RELATED WORK

Although most of the research output has focused on purely continuous or purely discrete distributions, in recent years probabilistic modelling for hybrid domains has received increasing attention.

Undirected hybrid models with potentials in the exponential family have been investigated for pairwise distributions [Lee and Hastie, 2015] and more general cases [Yang et al., 2014]. Density Estimation Trees (DETs) [Ram and Gray,

2011] address the problem of learning a distribution with axis-aligned piecewise constant components. Mixed Sum-Product Networks (MSPNs) [Molina et al., 2018, Vergari et al., 2019] learn mixtures of categorical and univariate piecewise polynomial distributions. In contrast with our setting, this line of research focuses on learning probabilistic relationships and the proposed learning and inference techniques disregard the algebraic and logical structures in the distributions.

The maximum-likelihood estimation of the parameters of piecewise distributions with logical and algebraic constraints, here investigated in the log-linear case, was initially addressed for piecewise constants by Belle et al. [2015]. Constraint learning in SMT, which is closely related to the problem of estimating the support of a structured hybrid distribution from data, was addressed in both supervised [Kolb et al., 2018] and unsupervised Morettin et al. [2020] settings. The latter work additionally propose LARIAT, a full pipeline for learning hybrid structured distributions from data by renormalizing DETs or MSPNs with the learned constraints. This approach heavily relies on greedy procedures for learning both the structure and the parameters of the model, while we focus on a principled way of estimating the parameters of these distributions via MLE.

While HL-MRFs were initially developed to efficiently answer MAP queries, approximate marginal inference has been investigated by Embar et al. [2019]. Besides approximating MLE with MAP states, parameter learning in these models has been addressed using maximum pseudolikelihood, large-margin estimation [Bach et al., 2017] and, more recently, with Bayesian optimization techniques [Srinivasan et al., 2020].

## 6 CONCLUSION

In this work we investigated in a rigorous and principled manner the problem of estimating the parameters of hybrid models that account for logical and algebraic relationships among variables. We showed that MLE can be solved optimally for piecewise distributions that are log-linear in their parameters and, leveraging recent theoretical insights on inference in this setting, is also tractable when the structure of the model satisfy certain conditions.

We extended the above results to HL-MRFs by presenting a novel reduction of marginal inference in HL-MRFs to WMI, thus identifying a subclass of models that admit tractable and exact marginal inference and MLE-based parameter learning.

For future work, we aim at developing and evaluating practical algorithms for HL-MRFs based on the presented reduction. While we showed that log-linear weight functions admit optimal MLE, whether it is a necessary condition is still an open problem. In this work we focused on esti-

imating the parameters in the weight function, but we plan to investigate the problem of estimating the coefficients in the weighted literals, i.e. finding the optimal piecewise decomposition of the joint probability.

## Acknowledgements

This work is supported in part by NSF grants #CCF-1837129, #IIS-1956441, #IIS-1943641, DARPA grant #N66001-17-2-4032, a Sloan Fellowship, and gifts from Intel and Facebook Research. ZZ is supported by a NEC Student Research Fellowship.

## References

- Aws Albarghouthi, Loris D’Antoni, Samuel Drews, and Aditya V. Nori. Fairsquare: Probabilistic verification of program fairness. *Proc. ACM Program. Lang.*, (OOPSLA):80:1–80:30, 2017.
- Stephen H Bach, Matthias Broecheler, Bert Huang, and Lise Getoor. Hinge-loss markov random fields and probabilistic soft logic. *Journal of Machine Learning Research*, 18: 1–67, 2017.
- Clark Barrett, Leonardo de Moura, Silvio Ranise, Aaron Stump, and Cesare Tinelli. The SMT-LIB initiative and the rise of SMT. In *Proceedings of the 6th international conference on Hardware and software: verification and testing*, pages 3–3. Springer-Verlag, 2010.
- Vaishak Belle, Andrea Passerini, and Guy Van den Broeck. Probabilistic inference in hybrid domains by weighted model integration. In *Proceedings of IJCAI*, pages 2770–2776, 2015.
- Varun Embar, Sriram Srinivasan, and Lise Getoor. Tractable marginal inference for hinge-loss markov random fields. In *Third ICML workshop on Tractable Probabilistic Modeling*, 2019.
- Samuel Kolb, Stefano Teso, Andrea Passerini, and Luc De Raedt. Learning smt (Ira) constraints using smt solvers. In *IJCAI*, volume 18, pages 2333–2340, 2018.
- Jason D Lee and Trevor J Hastie. Learning the structure of mixed graphical models. *Journal of Computational and Graphical Statistics*, 24(1):230–253, 2015.
- Alejandro Molina, Antonio Vergari, Nicola Di Mauro, Sri-  
raam Natarajan, Floriana Esposito, and Kristian Kersting. Mixed sum-product networks: A deep architecture for hybrid domains. In *Thirty-second AAAI conference on artificial intelligence*, 2018.
- Paolo Morettin, Samuel Kolb, Stefano Teso, and Andrea Passerini. Learning weighted model integration distributions. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5224–5231, 2020.
- Paolo Morettin, Pedro Zuidberg Dos Martires, Samuel Kolb, and Andrea Passerini. Hybrid probabilistic inference with logical and algebraic constraints: a survey. In *Proceedings of the 30th International Joint Conference on Artificial Intelligence*, 2021.
- Yurii Nesterov and Arkadii Nemirovskii. *Interior-point polynomial algorithms in convex programming*. SIAM, 1994.
- Parikshit Ram and Alexander G Gray. Density estimation trees. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 627–635, 2011.
- Tian Sang, Paul Beame, and Henry A Kautz. Performing bayesian inference by weighted model counting. In *AAAI*, volume 5, pages 475–481, 2005.
- Sriram Srinivasan, Golnoosh Farnadi, and Lise Getoor. Bowl: Bayesian optimization for weight learning in probabilistic soft logic. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10267–10275, 2020.
- Antonio Vergari, Alejandro Molina, Robert Peharz, Zoubin Ghahramani, Kristian Kersting, and Isabel Valera. Automatic bayesian density analysis. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 5207–5215, 2019.
- Eunho Yang, Yulia Baker, Pradeep Ravikumar, Genevera Allen, and Zhandong Liu. Mixed graphical models via exponential families. In *Artificial Intelligence and Statistics*, pages 1042–1050, 2014.
- Zhe Zeng and Guy Van den Broeck. Efficient search-based weighted model integration. *Proceedings of UAI*, 2019.
- Zhe Zeng, Paolo Morettin, Fanqi Yan, Antonio Vergari, and Guy Van den Broeck. Probabilistic inference with algebraic constraints: Theoretical limits and practical approximations. *Advances in Neural Information Processing Systems*, 33, 2020a.
- Zhe Zeng, Paolo Morettin, Fanqi Yan, Antonio Vergari, and Guy Van den Broeck. Scaling up hybrid probabilistic inference with logical and arithmetic constraints via message passing. In *Proceedings of the International Conference of Machine Learning (ICML)*, 2020b.



## A PROOFS

**Proof of Proposition 3.1.** Denote the  $i$ -th parameter in weight function  $w_{\ell'}$  as  $\Theta_i^{\ell'}$ . The first-order derivative of the log-partition function  $\log Z(\Theta)$  with respect to a parameter  $\Theta_i^{\ell'}$  can be computed as

$$\begin{aligned} \frac{\partial}{\partial \Theta_i^{\ell'}} \log Z(\Theta) &= \frac{1}{Z(\Theta)} \cdot \frac{\partial}{\partial \Theta_i^{\ell'}} Z(\Theta) \\ &= \frac{1}{Z(\Theta)} \cdot \left[ \int_x \prod_{\ell \neq \ell'} w_{\ell}(x)^{\mathbb{1}[x=\ell]} \cdot \frac{\partial}{\partial \Theta_i^{\ell'}} w_{\ell'}(x)^{\mathbb{1}[x=\ell']} dx \right] \\ &= \frac{1}{Z(\Theta)} \cdot \left[ \int_x \prod_{\ell} w_{\ell}(x)^{\mathbb{1}[x=\ell]} \cdot \frac{1}{w_{\ell'}(x)^{\mathbb{1}[x=\ell']}} \cdot \frac{\partial}{\partial \Theta_i^{\ell'}} w_{\ell'}(x)^{\mathbb{1}[x=\ell']} dx \right] \\ &= \frac{1}{Z(\Theta)} \cdot \left[ \int_x \prod_{\ell} w_{\ell}(x)^{\mathbb{1}[x=\ell]} \cdot \left[ \frac{\partial}{\partial \Theta_i^{\ell'}} \log w_{\ell'}(x)^{\mathbb{1}[x=\ell']} \right] dx \right] \end{aligned}$$

Notice that the last equation follows from the fact that  $\frac{1}{w_{\ell'}(x)^{\mathbb{1}[x=\ell']}} \cdot \frac{\partial}{\partial \Theta_i^{\ell'}} w_{\ell'}(x)^{\mathbb{1}[x=\ell']} = \frac{\partial}{\partial \Theta_i^{\ell'}} \log w_{\ell'}(x)^{\mathbb{1}[x=\ell']}$ . Given the WMI density of  $x$  in Equation 4, we have that

$$\frac{\partial}{\partial \Theta_i^{\ell'}} \log Z(\Theta) = \mathbb{E}_{x \sim p} \left[ \frac{\partial}{\partial \Theta_i^{\ell'}} \log w_{\ell'}(x)^{\mathbb{1}[x=\ell']} \right]$$

For brevity, denote  $\frac{\partial}{\partial \Theta_i^{\ell}} \log w_{\ell}(x)^{\mathbb{1}[x=\ell]}$  by notation  $a_i^{\ell}$ . Further, the second-order derivative of the log-partition function with respect to parameters  $\Theta_i^{\ell}$  and  $\Theta_j^{\ell'}$  is as follows.

$$\begin{aligned} \frac{\partial^2}{\partial \Theta_i^{\ell'} \partial \Theta_j^{\ell''}} \log Z(\Theta) &= \frac{\partial}{\partial \Theta_j^{\ell''}} \frac{\partial}{\partial \Theta_i^{\ell'}} \log Z(\Theta) \\ &= \frac{\partial}{\partial \Theta_j^{\ell''}} \frac{1}{Z(\Theta)} \cdot \frac{\partial}{\partial \Theta_i^{\ell'}} Z(\Theta) + \frac{1}{Z(\Theta)} \cdot \frac{\partial^2}{\partial \Theta_i^{\ell'} \partial \Theta_j^{\ell''}} Z(\Theta) \end{aligned}$$

where the first term can be rewritten as

$$\begin{aligned} &\frac{\partial}{\partial \Theta_j^{\ell''}} \frac{1}{Z(\Theta)} \cdot \frac{\partial}{\partial \Theta_i^{\ell'}} Z(\Theta) \\ &= -\frac{1}{Z^2(\Theta)} \frac{\partial}{\partial \Theta_j^{\ell''}} Z(\Theta) \cdot \frac{\partial}{\partial \Theta_i^{\ell'}} Z(\Theta) \\ &= -\mathbb{E}_{x \sim p} [a_j^{\ell''}] \cdot \mathbb{E}_{x \sim p} [a_i^{\ell'}] \end{aligned}$$

and the second term, given that the weight functions are log-linear in the parameters, can be rewritten as

$$\begin{aligned} &\frac{1}{Z(\Theta)} \cdot \frac{\partial^2}{\partial \Theta_i^{\ell'} \partial \Theta_j^{\ell''}} Z(\Theta) \\ &= \frac{1}{Z(\Theta)} \cdot \frac{\partial}{\partial \Theta_j^{\ell''}} \left[ \frac{\partial}{\partial \Theta_i^{\ell'}} Z(\Theta) \right] \\ &= \mathbb{E}_{x \sim p} [a_j^{\ell''} a_i^{\ell'}] \end{aligned}$$

We can view  $a_i^{\ell}$  as random variables distributed according to distribution defined by parameters  $\Theta$ . Then the Hessian

matrix of  $\log Z(\Theta)$  is the covariance matrix of random variables  $a_i^{\ell}$  which is positive semidefinite. Therefore, the log-partition function  $\log Z(\Theta)$  is a convex function in  $\Theta$ , which finishes our proof.  $\square$