CONVEX AND BILEVEL OPTIMIZATION FOR NEURO-SYMBOLIC INFERENCE AND LEARNING

Anonymous authors

Paper under double-blind review

Abstract

We address a key challenge for neuro-symbolic (NeSy) systems by leveraging convex and bilevel optimization techniques to develop a general gradient-based framework for end-to-end neural and symbolic parameter learning. The applicability of our framework is demonstrated with NeuPSL, a state-of-the-art NeSy architecture. To achieve this, we propose a smooth primal and dual formulation of NeuPSL inference and show learning gradients are functions of the optimal dual variables. Additionally, we develop a dual block coordinate descent algorithm for the new formulation that naturally exploits warm-starts. This leads to over $100 \times$ learning runtime improvements over the current best NeuPSL inference method. Finally, we provide extensive empirical evaluations across 8 datasets covering a range of tasks and demonstrate our learning framework achieves up to a 16% point prediction performance improvement over alternative learning methods.

1 INTRODUCTION

The goal of neuro-symbolic (NeSy) AI is a seamless integration of neural models for processing low-level data with symbolic frameworks to reason over high-level symbolic structures (d'Avila Garcez et al., 2002; 2009; 2019). This paper addresses an important research challenge in NeSy with the introduction of a principled and general NeSy learning framework. Further, we propose a novel inference algorithm and establish theoretical properties for a state-of-the-art NeSy system that are crucial for learning.

Our proposed learning framework builds upon NeSy energy-based models (NeSy-EBMs) (Pryor et al., 2023), a general class of NeSy systems that encompasses a variety of existing NeSy methods, including DeepProblog (Manhaeve et al., 2018; 2021), SATNet (Wang et al., 2019), logic tensor networks (Badreddine et al., 2022), and NeuPSL (Pryor et al., 2023). NeSy-EBMs use neural network outputs to parameterize an energy function and formulate an inference problem that may be non-smooth and constrained. Thus, predictions are not guaranteed to be a function of the inputs and parameters with an explicit form or to be differentiable, and traditional deep learning techniques are not directly applicable. We therefore equivalently formulate NeSy-EBM learning as a bilevel problem and, to support smooth first-order gradient-based optimization, propose a smoothing strategy that is novel to NeSy learning. Specifically, we replace the constrained NeSy energy function with its Moreau envelope. The augmented Lagrangian method for equality-constrained minimization is then applied with the new formulation.

We demonstrate the effectiveness of our proposed learning framework with NeuPSL. To ensure differentiability and provide principled forms of gradients for learning, we present a new formulation and regularization of NeuPSL inference as a quadratic program. Moreover, we introduce a dual block coordinate descent (dual BCD) inference algorithm for the quadratic program. The dual BCD algorithm is the first NeuPSL inference method that produces optimal dual variables for producing both optimal primal variables and gradients for learning. Additionally, empirical results demonstrate that dual BCD is able to effectively leverage warm starts, thus improving learning runtime.

Our key contributions are: (1) An improved formulation of the NeSy-EBM learning problem that establishes a foundation for applying smooth first-order gradient-based optimization techniques; (2) A reformulation of NeuPSL inference that is used to prove continuity properties and obtain explicit forms of gradients for learning; (3) A dual BCD algorithm for NeuPSL inference that naturally

produces statistics necessary for computing gradients for learning and that fully leverages warmstarts to improve learning runtime; (4) Two parallelization strategies for dual BCD inference; and (5) A thorough empirical evaluation demonstrating prediction performance improvements on 8 different datasets and a learning runtime speedup of up to $100 \times$.

2 RELATED WORK

NeSy AI is an active area of research that incorporates symbolic (commonly logical and arithmetic) reasoning with neural networks (Bader & Hitzler, 2005; d'Avila Garcez et al., 2009; Besold et al., 2017; De Raedt et al., 2020; Lamb et al., 2020; Giunchiglia et al., 2022). We will show that learning for a general class of NeSy systems is naturally formulated as bilevel optimization (Bracken & McGill, 1973; Colson et al., 2007; F. Bard, 2013). In other words, the NeSy learning objective is a function of predictions obtained by solving a lower-level inference problem that is symbolic reasoning. In this work, we focus on a general setting where the lower-level problem is an expressive and complex program capable of representing cyclic dependencies and ensuring the satisfaction of constraints during both learning and inference (Wang et al., 2019; Badreddine et al., 2022; Dasarth et al., 2023; Pryor et al., 2023; Cornelio et al., 2023). One prominent and tangential subgroup of such NeSy systems we would like to acknowledge enforces constraints on the structure of the symbolic model, and hence the lower-level problem, to ensure the final prediction has an explicit gradient with respect to the parameters (Xu et al., 2018; Manhaeve et al., 2021; Ahmed et al., 2022). In the deep learning community, bilevel optimization also arises in hyperparameter optimization and metalearning (Pedregosa, 2016; Franceschi et al., 2018), generative adversarial networks (Goodfellow et al., 2014), and reinforcement learning (Sutton & Barto, 2018).

Researchers typically take one of three approaches to bilevel optimization: (1) Implicit differenti*ation* methods compute or approximate the Hessian matrix at the lower-level problem solution to derive an analytic expression for the gradient of the upper-level objective called a hypergradient (Do et al., 2007; Pedregosa, 2016; Ghadimi & Wang, 2018; Rajeswaran et al., 2019; Giovannelli et al., 2022; Khanduri et al., 2023). (2) Automatic differentiation methods unroll inference into a differentiable computational graph (Stoyanov et al., 2011; Domke, 2012; Belanger et al., 2017; Ji et al., 2021). (3) Value-Function approaches reformulate the bilevel problem as a single-level constrained program using the optimal value of the lower-level objective (the *value-function*) to develop principled first-order gradient-based algorithms that do not require the calculation of Hessian matrices for the lower-level problem (V. Outrata, 1990; Liu et al., 2021; Sow et al., 2022; Liu et al., 2022; 2023; Kwon et al., 2023). Note that standard algorithms for all three approaches to bilevel optimization suggest solving the lower-level problem to derive the gradients used for optimizing the bilevel program. Principled techniques for using approximate lower-level solutions to make progress on the bilevel program is an open research direction (Pedregosa, 2016; Liu et al., 2021). Further, the lower-level problem for NeSy learning (inference) is commonly constrained. Implicit differentiation methods have been developed for bilevel optimization with lower-level constraints (Giovannelli et al., 2022; Khanduri et al., 2023). We introduce a value-function approach.

3 NeSy energy-based models

In this work, we use *NeSy energy-based models (NeSy-EBMs)* (Pryor et al., 2023) to develop a generally applicable NeSy learning framework. Here, we provide background on NeSy-EBMs and introduce a classification of losses that motivates the need for general learning algorithms.

NeSy-EBMs are a family of EBMs (LeCun et al., 2006) that use neural model predictions to define potential functions with symbolic interpretations. NeSy-EBM energy functions are parameterized by a set of neural and symbolic weights from the domains W_{nn} and W_{sy} , respectively, and quantify the compatibility of a target variable from a domain \mathcal{Y} and neural and symbolic inputs from the domains \mathcal{X}_{nn} and \mathcal{X}_{sy} : $E : \mathcal{Y} \times \mathcal{X}_{sy} \times \mathcal{X}_{nn} \times \mathcal{W}_{sy} \times \mathcal{W}_{nn} \to \mathbb{R}$. NeSy-EBM inference requires first computing the output of the neural networks, *neural inference*, and then minimizing the energy function over the targets, *symbolic inference*:

$$\arg\min_{\mathbf{w}\in\mathcal{V}} E(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{x}_{nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn}).$$
(1)

NeSy-EBM learning is finding weights to create an energy function that associates lower energies to target values near their truth in a set of training data. The training data consists of P samples that

are tuples of symbolic variables and neural network inputs: $\{S_1 = (\mathbf{y}_1, \mathbf{x}_{1,sy}, \mathbf{x}_{1,nn}), \dots, S_P = (\mathbf{y}_P, \mathbf{x}_{P,sy}, \mathbf{x}_{P,nn})\}$. Moreover, targets \mathbf{y}_i from a training sample S_i are partitioned into *labeled variables*, \mathbf{t}_i , for which there is a corresponding truth value, and *latent variables*, \mathbf{z}_i . Without loss of generality, we write $\mathbf{y}_i = (\mathbf{t}_i, \mathbf{z}_i)$. NeSy-EBM learning losses are defined using the *latent minimizer*, $\mathbf{z}_i^* \in \arg\min_{\mathbf{z}\in\mathcal{Z}} E((\mathbf{t}_i, \mathbf{z}), \mathbf{x}_{i,sy}, \mathbf{x}_{i,nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn})$, the *full minimizer*, $\mathbf{y}_i^* \in \arg\min_{\mathbf{y}\in\mathcal{Y}} E(\mathbf{y}, \mathbf{x}_{i,sy}, \mathbf{x}_{i,nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn})$, and the *latent* and *full optimal value-functions*:

$$V_{\mathbf{z}_{i}^{*}}(\mathbf{w}_{sy}, \mathbf{w}_{nn}) := E((\mathbf{t}_{i}, \mathbf{z}_{i}^{*}), \mathbf{x}_{i,sy}, \mathbf{x}_{i,nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn}),$$
(2)

$$V_{\mathbf{y}_i^*}(\mathbf{w}_{sy}, \mathbf{w}_{nn}) := E(\mathbf{y}_i^*, \mathbf{x}_{i,sy}, \mathbf{x}_{i,nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn}).$$
(3)

Note the optimal values-functions are functions of the parameters, inputs, and symbolic variables; however, to simplify notation, we only write the parameters as arguments.

Value-based learning losses depend on the model weights strictly via the optimal value-functions. Two common value-based losses for NeSy-EBMs are the latent optimal value-function (*energy loss*), and the difference between the latent and full optimal value-functions (*structured perceptron loss*) (LeCun et al., 1998; Collins, 2002):

$$L_{Energy}(E(\cdot, \cdot, \cdot, \mathbf{w}_{sy}, \mathbf{w}_{nn}), S_i) := V_{\mathbf{z}_i^*}(\mathbf{w}_{sy}, \mathbf{w}_{nn}), \tag{4}$$

$$L_{SP}(E(\cdot, \cdot, \cdot, \mathbf{w}_{sy}, \mathbf{w}_{nn}), S_i) := V_{\mathbf{z}_i^*}(\mathbf{w}_{sy}, \mathbf{w}_{nn}) - V_{\mathbf{y}_i^*}(\mathbf{w}_{sy}, \mathbf{w}_{nn}).$$
(5)

A principled first-order gradient-based method for optimizing a value-based objective only requires differentiability of the value-functions. However, performance metrics are not always aligned with value-based losses. Moreover, they are known to have degenerate solutions, e.g., weights minimizing the loss but producing a collapsed energy function (LeCun et al., 2006; Pryor et al., 2023).

Alternatively, *minimizer-based* learning losses assume the minimizer of the energy function is unique. With this assumption, energy minimization is a vector-valued function from the weight space $\mathcal{W}_{sy} \times \mathcal{W}_{nn}$ to the target space $\mathcal{Y}, \mathbf{y}_i^*(\mathbf{w}_{sy}, \mathbf{w}_{nn}) : \mathcal{W}_{sy} \times \mathcal{W}_{nn} \to \mathcal{Y}$. Then, minimizer-based losses are compositions of a differentiable supervised loss $d : \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$, and the minimizer:

$$L_d(E(\cdot, \cdot, \cdot, \mathbf{w}_{sy}, \mathbf{w}_{nn}), S_i) := d(\mathbf{y}_i^*(\mathbf{w}_{sy}, \mathbf{w}_{nn}), \mathbf{t}_i).$$
(6)

Minimizer-based losses are general and allow learning with objectives aligned with evaluation metrics. However, a direct application of a first-order gradient based method for minimizer-based learning requires the Jacobian at the minimizer. NeSy-EBM predictions are not necessarily differentiable. Even if they are differentiable, the computation of the Jacobian is often too expensive to be practical.

4 A BILEVEL NESY LEARNING FRAMEWORK

In this section, we introduce a general framework for the bilevel NeSy learning problem:

$$\underset{(\mathbf{y}_{sy},\mathbf{w}_{nn})\in\mathcal{W}_{sy}\times\mathcal{W}_{nn}}{\operatorname{arg\,min}}\sum_{i=1}^{P} \left(d(\mathbf{y}_{i},\mathbf{t}_{i}) + L_{Val}(E(\cdot,\cdot,\cdot,\mathbf{w}_{sy},\mathbf{w}_{nn}),S_{i}) \right) + \mathcal{R}(\mathbf{w}_{sy},\mathbf{w}_{nn}) \tag{7}$$

$$\underset{\mathbf{y}_{i}\in\operatorname{arg\,min}}{\operatorname{s.t.}} \mathbf{y}_{i}\in\operatorname{arg\,min}_{\mathbf{y}\in\mathcal{Y}}E(\mathbf{y},\mathbf{x}_{i,sy},\mathbf{x}_{i,nn},\mathbf{w}_{sy},\mathbf{w}_{nn}), \quad \forall i \in \{1,\cdots,P\}$$

where d and L_{Val} are a minimizer and value-based loss, respectively, and $\mathcal{R} : \mathcal{W}_{sy} \times \mathcal{W}_{nn} \to \mathbb{R}$ is a regularizer. We make the following (standard) lower-level singleton assumption.

Assumption 4.1. *E* is minimized over $\mathbf{y} \in \mathcal{Y}$ at a single point for every $(\mathbf{w}_{sy}, \mathbf{w}_{nn}) \in \mathcal{W}_{sy} \times \mathcal{W}_{nn}$.

Under Assumption 4.1, and regardless of the continuity and curvature properties of the upper and lower level objectives, (7) is equivalent to the following:

$$\arg \min_{\substack{(\mathbf{w}_{sy}, \mathbf{w}_{nn}) \in \mathcal{W}_{sy} \times \mathcal{W}_{nn} \\ (\mathbf{y}_{1}, \cdots, \mathbf{y}_{P}) \in \mathcal{Y} \times \cdots \times \mathcal{Y}}} \sum_{i=1}^{P} \left(d(\mathbf{y}_{i}, \mathbf{t}_{i}) + L_{Val}(E(\cdot, \cdot, \cdot, \mathbf{w}_{sy}, \mathbf{w}_{nn}), S_{i}) \right) + \mathcal{R}(\mathbf{w}_{sy}, \mathbf{w}_{nn})$$
s.t. $E(\mathbf{y}_{i}, \mathbf{x}_{i,sy}, \mathbf{x}_{i,nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn}) - V_{\mathbf{y}_{i}^{*}}(\mathbf{w}_{sy}, \mathbf{w}_{nn}) \le 0, \quad \forall i \in \{1, \cdots, P\}.$
(8)

The formulation in (8) is referred to as a *value-function* approach in bilevel optimization literature (V. Outrata, 1990; Liu et al., 2021; 2022; Sow et al., 2022; Kwon et al., 2023). Value-function approaches view the bilevel program as a single-level constrained optimization problem by leveraging the value-function as a tight lower bound on the lower-level objective. However, the inequality

constraints in (8) do not satisfy any of the standard *constraint qualifications* that ensure the feasible set near the optimal point is similar to its linearized approximation (Nocedal & Wright, 2006). This raises a challenge for providing theoretical convergence guarantees for constrained optimization techniques. Following a recent line of value-function approaches to bilevel programming (Liu et al., 2021; Sow et al., 2022; Liu et al., 2023), we overcome this challenge by allowing at most an $\iota > 0$ violation in each constraint in (8). With this relaxation, strictly feasible points exist and, for instance, the linear independence constraint qualification (LICQ) can hold.

Another challenge that arises from (8) is that the energy function of NeSy-EBMs is typically nondifferentiable with respect to the targets and even infinite-valued to implicitly represent constraints. As a result, penalty or augmented Lagrangian functions derived from (8) are intractable. Therefore, we substitute each instance of the energy function evaluated at the training sample i and parameterized by $(\mathbf{w}_{sy}, \mathbf{w}_{nn})$ in the constraints of (8) with the following function:

$$M_{i}(\mathbf{y}; \mathbf{w}_{sy}, \mathbf{w}_{nn}, \rho) := \inf_{\hat{\mathbf{y}} \in \mathcal{Y}} \left(E(\hat{\mathbf{y}}, \mathbf{x}_{i,sy}, \mathbf{x}_{i,nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn}) + \frac{1}{2\rho} \|\hat{\mathbf{y}} - \mathbf{y}\|_{2}^{2} \right),$$
(9)

where ρ is a positive scalar. For convex E, (9) is the Moreau envelope of the energy function Rockafellar (1970); Parikh & Boyd (2013). In general, even for non-convex energy functions, the smoothing in (9) preserves global minimizers and minimum values, i.e., $\mathbf{y}_i^*(\mathbf{w}_{sy},\mathbf{w}_{nn}) =$ $\arg\min_{\mathbf{y}} M_i(\mathbf{y}; \mathbf{w}_{sy}, \mathbf{w}_{nn}, \rho)$ and $V_{\mathbf{y}_i^*}(\mathbf{w}_{sy}, \mathbf{w}_{nn}) = \min_{\mathbf{y}} M_i(\mathbf{y}; \mathbf{w}_{sy}, \mathbf{w}_{nn}, \rho)$. Moreover, under Assumption 4.1 each M_i is finite for all $\mathbf{y} \in \mathcal{Y}$ even if the energy function is not. When the energy function is a lower semi-continuous convex function, its Moreau envelope is convex, finite, and continuously differentiable, and its gradient with respect to y is:

$$\nabla_{\mathbf{y}} M_i(\mathbf{y}; \mathbf{w}_{sy}, \mathbf{w}_{nn}, \rho) = \frac{1}{\rho} \left(\mathbf{y} - \operatorname*{arg\,min}_{\hat{\mathbf{y}} \in \mathcal{Y}} \left(\rho E(\hat{\mathbf{y}}, \mathbf{x}_{i,sy}, \mathbf{x}_{i,nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn}) + \frac{1}{2} \| \hat{\mathbf{y}} - \mathbf{y} \|_2^2 \right) \right)$$
(10)

Convexity is a sufficient but not necessary condition to ensure each M_i is differentiable with respect to the target variables. See Bonnans & Shapiro (2000) for results regarding the sensitivity of optimal value-functions to perturbations.

We propose the following relaxed and smoothed value-function approach to finding an approximate solution of (7):

$$\underset{(\mathbf{y}_{sy},\mathbf{w}_{nn})\in\mathcal{W}_{sy}\times\mathcal{W}_{nn}}{\operatorname{arg\,min}}\sum_{i=1}^{P} \left(d(\mathbf{y}_{i},\mathbf{t}_{i}) + L_{Val}(E(\cdot,\cdot,\cdot,\mathbf{w}_{sy},\mathbf{w}_{nn}),S_{i}) \right) + \mathcal{R}(\mathbf{w}_{sy},\mathbf{w}_{nn}) \tag{11}$$

s.t.
$$M_i(\mathbf{y}_i; \mathbf{w}_{sy}, \mathbf{w}_{nn}, \rho) - V_{\mathbf{y}_i^*}(\mathbf{w}_{sy}, \mathbf{w}_{nn}) \le \iota, \quad \forall i \in \{1, \cdots, P\},$$

The formulation (11) is the core of our proposed NeSy-EBM learning framework outlined in Algorithm 1. The algorithm proceeds by approximately solving instances of (11) in a sequence defined by a decreasing ι . This is a graduated approach to solving (8) with instances of (11) that are increasingly tighter approximations. Each instance of (11) is optimized using only firstorder gradients of the energy and valuefunctions with the bound-constrained augmented Lagrangian algorithm, Al-

Algorithm 1 NeSy-EBM Learning Framework

Require: Constraint Tolerance: σ^* , Movement Tolerance: ω^* , Moreau Param.: ρ , Starting points: $(\mathbf{w}_{sy}^{(0)}, \mathbf{w}_{nn}^{(0)}) \in \mathcal{W}_{sy} \times \mathcal{W}_{nn}$ 1: $\mathbf{y}_{i}^{(0)} \leftarrow (\mathbf{t}_{i}, \mathbf{z}_{i}^{*}), \forall i = 1, \cdots, P;$ 2: $\iota^{(0)} \leftarrow \max_{i \in \{1, \cdots, P\}} M_{i}(\mathbf{y}_{i}^{(0)}; \mathbf{w}_{sy}^{(0)}, \mathbf{w}_{nn}^{(0)}, \rho) - V_{\mathbf{y}_{i}^{*}}(\mathbf{w}_{sy}^{(0)}, \mathbf{w}_{nn}^{(0)});$

3: for $t = 0, 1, 2, \cdots$ do

- Find $\mathbf{w}_{sy}^{(t+1)}, \mathbf{w}_{nn}^{(t+1)}, \mathbf{y}_1^{(t+1)}, \cdots, \mathbf{y}_P^{(t+1)}$ minimizing (11) with $\iota^{(t)}$. 4: 5: if Stopping criterion satisified then Stop with: $\mathbf{w}_{sy}^{(t+1)}, \mathbf{w}_{nn}^{(t+1)}, \mathbf{y}_1^{(t+1)}, \cdots, \mathbf{y}_P^{(t+1)};$ $\iota^{(t+1)} \leftarrow \frac{1}{2} \cdot \iota^{(t)};$ 6:
- 7:

gorithm 17.4 from Nocedal & Wright (2006). Specifically, the algorithm finds approximate minimizers of the problem's augmented Lagrangian for a fixed setting of the penalty parameters using gradient descent. To simplify notation, let the equality constraints in (11) be denoted by:

$$c_i(\mathbf{y}_i, \mathbf{w}_{sy}, \mathbf{w}_{nn}; \iota) := M_i(\mathbf{y}_i; \mathbf{w}_{sy}, \mathbf{w}_{nn}, \rho) - V_{\mathbf{y}_i^*}(\mathbf{w}_{sy}, \mathbf{w}_{nn}) - \iota,$$

for each constraint indexed $i \in \{1, \dots, P\}$. Moreover, let $c(\mathbf{y}_1, \dots, \mathbf{y}_P, \mathbf{w}_{sy}, \mathbf{w}_{nn}; \iota) :=$ $[c_i(\mathbf{y}_i, \mathbf{w}_{sy}, \mathbf{w}_{nn}; \iota)]_{i=1}^P$. The augmented Lagrangian function corresponding to (11) introduces a quadratic penalty parameter μ and P linear penalty parameters $\lambda := [\lambda_i]_{i=1}^P$, as follows:

$$\mathcal{L}_{A}(\mathbf{w}_{sy}, \mathbf{w}_{nn}, \mathbf{y}_{1}, \cdots, \mathbf{y}_{p}, \mathbf{s}; \lambda, \mu, \iota) := \sum_{i=1}^{P} \left(d(\mathbf{y}_{i}, \mathbf{t}_{i}) + L_{Val}(E(\cdot, \cdot, \cdot, \mathbf{w}_{sy}, \mathbf{w}_{nn}), S_{i}) \right) \\ + \frac{\mu}{2} \sum_{i=1}^{P} \left(c_{i}(\mathbf{y}_{i}, \mathbf{w}_{sy}, \mathbf{w}_{nn}; \iota) + s_{i} \right)^{2} + \sum_{i=1}^{P} \lambda_{i} \left(c_{i}(\mathbf{y}_{i}, \mathbf{w}_{sy}, \mathbf{w}_{nn}; \iota) + s_{i} \right) + \mathcal{R}(\mathbf{w}_{sy}, \mathbf{w}_{nn}).$$
(12)

where we introduced P slack variables, $\mathbf{s} = [s_i]_{i=1}^P$, for each inequality constraint. We make the following assumption to ensure the augmented Lagrangian function is differentiable:

Assumption 4.2. Every $V_{y_i^*}$, $V_{z_i^*}$, and M_i is differentiable with respect to the weights.

We employ the bound-constrained augmented Lagrangian algorithm to solve (11) (see Appendix B for details). This method provides a principled algorithm for updating the penalty parameters and ensures fundamental convergence properties of our learning framework. Notably, we have that limit points of the iterate sequence are stationary points of $||c(y_1, \dots, y_P, w_{sy}, w_{nn}) + s||^2$ when the problem has no feasible points. When the problem is feasible and LICQ holds at the limits, they are KKT points of (11) (Theorem 17.2 in Nocedal & Wright (2006)). Convergence rates and stronger guarantees are likely possible from analyzing the structure of the energy function for specific NeSy-EBMs and is a direction for future work.

The value for ι is halved every time an approximate solution to the Lagrangian subproblem is reached. We suggest starting points for each $\mathbf{y}_i^{(0)}$ to be the latent inference minimizer and $\iota^{(0)}$ to be the maximum difference in the value-function and the smooth energy function over all $\mathbf{y}_i^{(0)}$. The outer loop of the NeSy-EBM learning framework may be stopped by either watching the progress of a training or validation evaluation metric, or by specifying a final value for ι .

5 NEUPSL AND DEEP HINGE-LOSS MARKOV RANDOM FIELDS

We demonstrate the applicability of our learning framework with Neural Probabilistic Soft Logic (NeuPSL), a general class of NeSy-EBMs designed for scalable joint reasoning (Pryor et al., 2023). In NeuPSL, relations and attributes are represented by *atoms*, and dependencies between atoms are encoded with first-order logical clauses and linear arithmetic inequalities referred to as *rules*. Atom values can be target variables, observations, or outputs from a neural network. The rules and atoms are translated into potentials measuring rule satisfaction and are aggregated to define a member of a tractable class of graphical models: *deep hinge-loss Markov random fields* (deep HL-MRF).

Definition 5.1. Let $\mathbf{g} = [g_i]_{i=1}^{n_g}$ be functions with corresponding weights $\mathbf{w}_{nn} = [\mathbf{w}_{nn,i}]_{i=1}^{n_g}$ and inputs \mathbf{x}_{nn} such that $g_i : (\mathbf{w}_{nn,i}, \mathbf{x}_{nn}) \mapsto [0, 1]$. Let $\mathbf{y} \in [0, 1]^{n_y}$ and $\mathbf{x}_{sy} \in [0, 1]^{n_x}$. A deep hinge-loss potential is a function of the form:

$$\phi(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) := (\max\{\mathbf{a}_{\phi, \mathbf{y}}^T \mathbf{y} + \mathbf{a}_{\phi, \mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{\phi, \mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}) + b_{\phi}, 0\})^p,$$

where $\mathbf{a}_{\phi,\mathbf{y}} \in \mathbb{R}^{n_y}$, $\mathbf{a}_{\phi,\mathbf{x}} \in \mathbb{R}^{n_x}$, and $\mathbf{a}_{\phi,\mathbf{g}} \in \mathbb{R}^{n_g}$ are variable coefficient vectors, $b_{\phi} \in \mathbb{R}$ is a vector of constants, and $p \in \{1, 2\}$. Let $\mathcal{T} = [\tau_i]_{i=1}^r$ denote an ordered partition of a set of m deep hingeloss potentials. Further, define $\Phi(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) := [\sum_{k \in \tau_i} \phi_k(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))]_{i=1}^r$. Let \mathbf{w}_{sy} be a vector of r non-negative symbolic weights corresponding to the partition \mathcal{T} . Then, a **deep hinge-loss energy function** is:

$$E(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{x}_{nn}, \mathbf{w}_{sy}, \mathbf{w}_{nn}) := \mathbf{w}_{sy}^T \mathbf{\Phi}(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})).$$
(13)

Let $\mathbf{a}_{c_k,\mathbf{y}} \in \mathbb{R}^{n_y}$, $\mathbf{a}_{c_k,\mathbf{x}} \in \mathbb{R}^{n_x}$, $\mathbf{a}_{c_k,\mathbf{g}} \in \mathbb{R}^{n_g}$, and $b_{c_k} \in \mathbb{R}$ for each $k \in 1, \ldots, q$ and $q \ge 0$ be vectors defining linear inequality constraints and a feasible set:

$$\mathbf{\Omega}(\mathbf{x}_{sy}, \mathbf{g}) := \left\{ \mathbf{y} \in [0, 1]^{n_y} \mid \mathbf{a}_{c_k, \mathbf{y}}^T \mathbf{y} + \mathbf{a}_{c_k, \mathbf{x}}^T \mathbf{x}_{sy} + \mathbf{a}_{c_k, \mathbf{g}}^T \mathbf{g} + b_{c_k} \le 0, \forall k = 1, \dots, q \right\}.$$

Then a **deep hinge-loss Markov random field** defines the conditional probability density:

$$P(\mathbf{y}|\mathbf{x}_{sy},\mathbf{x}_{nn}) := \begin{cases} \frac{\exp(-E(\mathbf{y},\mathbf{x}_{sy},\mathbf{x}_{nn},\mathbf{w}_{sy},\mathbf{w}_{nn}))}{\int_{\mathbf{y}\in\Omega(\cdot)}\exp(-E(\mathbf{y},\mathbf{x}_{sy},\mathbf{x}_{nn},\mathbf{w}_{sy},\mathbf{w}_{nn}))d\mathbf{y}} & \mathbf{y}\in\Omega(\mathbf{x}_{sy},\mathbf{g}(\mathbf{x}_{nn},\mathbf{w}_{nn}))\\ 0 & o.w. \end{cases}$$
(14)

NeuPSL inference is finding the MAP state of the conditional distribution defined by a deep HL-MRF, i.e., finding the minimizer of the energy function over the feasible set.

$$\min_{\mathbf{y}\in\mathbb{R}^{n_{\mathbf{y}}}} \mathbf{w}_{sy}^{T} \Phi(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) \quad \text{s.t. } \mathbf{y}\in \Omega(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})).$$
(15)

As each of the potentials are convex, (15) is a non-smooth convex linearly constrained program.

5.1 A SMOOTH FORMULATION OF INFERENCE

In this section, we introduce a primal and dual formulation of NeuPSL inference as a linearly constrained convex quadratic program (LCQP). (See Appendix C.1 for details.) In summary, m slack variables with lower bounds and $2 \cdot n_y + m$ linear constraints are defined to represent the target variable bounds and deep hinge-loss potentials. All $2 \cdot n_y + m$ variable bounds, m potentials, and $q \ge 0$ constraints are collected into a $(2 \cdot n_y + q + 2 \cdot m) \times (n_y + m)$ dimensional matrix **A** and a vector of $(2 \cdot n_y + q + 2 \cdot m)$ elements that is an affine function of the neural predictions and symbolic inputs $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$. Moreover, the slack variables and a $(n_y + m) \times (n_y + m)$ positive semi-definite diagonal matrix, $\mathbf{D}(\mathbf{w}_{sy})$, and a $(n_y + m)$ dimensional vector, $\mathbf{c}(\mathbf{w}_{sy})$, are created using the symbolic weights to define a quadratic objective. Further, we gather the original target variables and the slack variables into a vector $\nu \in \mathbb{R}^{n_y+m}$. Altogether, the regularized convex LCQP reformulation of NeuPSL inference is:

$$T(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) :=$$

$$\min_{\nu \in \mathbb{R}^{n_y + m}} \nu^T (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I}) \nu + \mathbf{c}(\mathbf{w}_{sy})^T \nu \quad \text{s.t. } \mathbf{A}\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) \le 0,$$
(16)

where $\epsilon \geq 0$ is a scalar regularization parameter added to the diagonal of **D** to ensure strong convexity (needed in the next subsection). The effect of the added regularization is empirically studied in Appendix E.3. The function $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ in (16) is the optimal value-function of the LCQP formulation of NeuPSL inference referred to in the previous section.

By Slater's constraint qualification, we have strong duality when there is a feasible solution to (16). In this case, an optimal solution to the dual problem yields an optimal solution to the primal problem. The Lagrange dual problem of (16) is:

$$\min_{\boldsymbol{\mu}\in\mathbb{R}^{2\cdot n_{\mathbf{y}}+m+q}} h(\boldsymbol{\mu}; \mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) \qquad (17)$$

$$:= \frac{1}{4} \boldsymbol{\mu}^{T} \mathbf{A} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})^{-1} \mathbf{A}^{T} \boldsymbol{\mu} + \frac{1}{2} (\mathbf{A} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})^{-1} \mathbf{c}(\mathbf{w}_{sy}) - 2\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))^{T} \boldsymbol{\mu},$$

where μ is the vector of dual variables and $h(\mu; \mathbf{w}_{sy}, \mathbf{b}(\mathbf{w}_{nn}))$ is the LCQP dual objective function. As $(\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})$ is diagonal, it is easy to invert, and thus it is practical to work in the dual space and map dual to primal variables. The dual-to-primal variable mapping is:

$$\nu \leftarrow -\frac{1}{2} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})^{-1} (\mathbf{A}^T \mu + \mathbf{c}(\mathbf{w}_{sy})).$$
(18)

On the other hand, the primal-to-dual mapping is more computationally expensive and requires calculating a pseudo-inverse of the constraint matrix A.

5.2 CONTINUITY OF INFERENCE

V

We use the LCQP formulation in (16) to establish continuity and curvature properties of the NeuPSL energy minimizer and the optimal value-function provided in the following theorem. The proof is provided in Appendix C.2.

Theorem 5.2. Suppose for any setting of $\mathbf{w}_{nn} \in \mathbb{R}^{n_g}$ there is a feasible solution to NeuPSL inference (16). Further, suppose $\epsilon > 0$, $\mathbf{w}_{sy} \in \mathbb{R}^r_+$, and $\mathbf{w}_{nn} \in \mathbb{R}^{n_g}$. Then:

- The minimizer of (16), $\mathbf{y}^*(\mathbf{w}_{sy}, \mathbf{w}_{nn})$, is a $O(1/\epsilon)$ Lipschitz continuous function of \mathbf{w}_{sy} .
- $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$, is concave over \mathbf{w}_{sy} and convex over $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$.
- $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is differentiable with respect to \mathbf{w}_{sy} . Moreover,

$$\nabla_{\mathbf{w}_{sy}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) = \mathbf{\Phi}(\mathbf{y}^*(\mathbf{w}_{sy}, \mathbf{w}_{nn}), \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})).$$

Furthermore, $\nabla_{\mathbf{w}_{sy}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is Lipschitz continuous over \mathbf{w}_{sy} .

• If there is a feasible point ν strictly satisfying the *i*'th inequality constraint of (16), i.e., $\mathbf{A}[i]\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i] < 0$, then $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is subdifferentiable with respect to the *i*'th constraint constant $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i]$. Moreover,

$$\partial_{\mathbf{b}[i]}V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) = \{\mu^*[i] \mid \mu^* \in \operatorname*{arg\,min}_{\mu \in \mathbb{R}^{2\cdot n_{\mathbf{y}}+m+q}_{\geq 0}} h(\mu; \mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))\}.$$

Furthermore, if $\mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})$ is a smooth function of \mathbf{w}_{nn} , then so is $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$, and the set of regular subgradients of $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is:

$$\hat{\partial}_{\mathbf{w}_{nn}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$$

$$\supset \nabla_{\mathbf{w}_{nn}} \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))^T \partial_{\mathbf{b}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))).$$
(19)

Theorem 5.2 provides a simple explicit form of the value-function gradient with respect to the symbolic weights and regular subgradient with respect to the neural weights. Moreover, this result is directly applicable to the Moreau envelope of the NeuPSL energy function used in Section 4 as it is a regularized value-function. Thus, Theorem 5.2 supports the principled application of Algorithm 1 for learning both the symbolic and neural weights of a NeuPSL model.

5.3 DUAL BLOCK COORDINATE DESCENT

The regular subgradients in Theorem 5.2 are functions of the optimal dual variables of the LCQP inference problem in (17). For this reason, we introduce a block coordinate descent (BCD) (Wright, 2015) algorithm for working directly with the dual LCQP formulation of inference. Details of the algorithm are provided in Appendix D. Our dual BCD algorithm is the first method specialized for the dual LCQP inference and is therefore also the first to produce optimal dual variables that directly yeild both optimal primal variables and principled gradients for learning, all without the need to compute a pseudo-inverse of the constraint matrix.

The dual BCD algorithm proceeds by successively minimizing the objective along the subgradient of a block of dual variables. For this reason, dual BCD guarantees descent at every iteration, partially explaining its effectiveness at leveraging warm-starts and improving learning runtimes. The algorithm is stopped when the primal-dual gap drops below a threshold $\delta > 0$. We suggest a practical choice of variable blocks with efficient methods for computing the objective subgradients and solving the steplength subproblems. Additionally, we develop an efficient method for identifying connected components of the factor graph defined by the deep HL-MRF yeilding a variable partition that the dual objective is additively separable over to parallelize the BCD updates. Moreover, inspired by lock-free parallelization strategies (Bertsekas & N. Tsitsiklis, 1989; Recht et al., 2011; Liu et al., 2015), we also propose a variant of the dual BCD inference algorithm that sacrifices the theoretical guaranteed descent property for significant runtime improvements. In Section 6, we show that the lock-free dual BCD algorithm consistently finds a solution satisfying the stopping criterion, and surprisingly, is still highly effective at leveraging warm starts.

6 EMPIRICAL EVALUATION

We evaluate the runtime and prediction performance of our proposed NeSy inference and parameter learning algorithms on the 8 datasets in Table 1¹. The table includes the dataset's inference task, the associated prediction performance metric, and whether the corresponding NeuPSL model has deep neural network parameters. Unless noted

Table 1: Datasets used for empirical evaluations.

Dataset	Deep	Task	Perf. Metric
CreateDebate Hasan & Ng (2013)		Stance Class.	AUROC
4Forums Walker et al. (2012)		Stance Class.	AUROC
Epinions (Richardson et al., 2003)		Link Pred.	AUROC
DDI (S. Wishart et al., 2006)		Link Pred.	AUROC
Yelp (Yelp, 2023)		Regression	MAE
Citeseer (Sen et al., 2008)	\checkmark	Node Class.	Accuracy
Cora (Sen et al., 2008)	\checkmark	Node Class.	Accuracy
MNIST-Add.(Manhaeve et al., 2018)	 ✓ 	Image Class.	Accuracy

otherwise, all experiments are run on 5 splits and the average and standard deviation of times and performance metric values are reported. Details on the datasets, hardware specifications, hyperparameter searches, and model architectures are provided in Appendix E.

For learning experiments in Section 6.2 and Section 6.3, NeuPSL models with weights trained using value-based learning losses, e.g., energy and structured perceptron (SP), use mirror descent (Kivinen & Warmuth, 1997; Shalev-Shwartz, 2012) on the symbolic weights constrained to the unit simplex and Adam (P. Kingma & Lei Ba, 2017) for the neural weights. NeuPSL models with weights trained using minimizer-based losses, e.g., mean squared error (MSE) and binary cross entropy (BCE), use our proposed NeSy learning framework in Algorithm 1 with a scaled energy loss term added to the objective as in (7). Moreover, optimization of the augmented Lagrangian, line 4 of Algorithm 1, is performed using the bound constrained augmented Lagrangian algorithm (Appendix B) with mirror descent on the symbolic weights and Adam for the neural weights.

¹All code and data is available at https://github.com/convexbilevelnesylearning.

6.1 INFERENCE RUNTIME

We begin by examining the runtime of symbolic inference. We evaluate the alternating direction method of multipliers (ADMM) Boyd et al. (2010), the current state-of-the-art inference algorithm for NeuPSL, and our proposed inference algorithms: connected component parallel dual BCD (CC D-BCD) and lock-free parallel dual BCD (LF D-BCD). We also evaluate the performance of Gurobi, a leading off-the-shelf optimizer, and subgradient descent (GD) in Appendix E.4. All inference algorithms have access to the same computing resources . We run a hyperparameter search,

Table 2: Time in seconds for inference us-
ing ADMM and our proposed CC D-BCD
and LF D-BCD algorithms on each dataset.

	ADMM	CC D-BCD	LF D-BCD
CreateDebate 4Forums	9.98 ± 1.13 15.17 ± 0.74	$\begin{array}{c c} 0.05 \pm 0.02 \\ 0.11 \pm 0.02 \end{array}$	$\begin{array}{c} 0.05 \pm 0.03 \\ 0.05 \pm 0.01 \end{array}$
Epinions	0.36 ± 0.041	1.84 ± 0.4	0.26 ± 0.04
Citeseer Cora	0.63 ± 0.07 0.71 ± 0.07	1.36 ± 0.24 6.46 ± 3.5	0.49 ± 0.08 0.79 ± 0.19
DDI	7.85 ± 0.28	31.47 ± 0.17	1.76 ± 0.17
Yelp MNIST-Add1	$\begin{array}{c} {\bf 6.37 \pm 1.19} \\ {11.45 \pm 1.32} \\ \\ {205 \pm 0.0} \end{array}$	$48.44 \pm 3.82 \\ 10.23 \pm 1.04 \\ 20.00 \pm 0.00 \\ 10.00 \pm 0.00 \\ 10.0$	7.58 ± 0.48 115 ± 45
MNIST-Add2	285 ± 66	29.09 ± 8.00	$1,189 \pm 16$

detailed in Appendix E.4, for each algorithm, and the configuration yielding a prediction performance that is within a standard deviation of the best and completed with the lowest runtime is reported. All algorithms are stopped when the L_{∞} norm of the primal variable change between iterates is less than 0.001.

The total average inference runtime in seconds for each algorithm and model is provided in Table 2. Surprisingly, despite the potential for an inexact solution to the BCD steplength subproblem, LF D-BCD is faster than CC D-BCD in the first 7 datasets and demonstrates up to $6 \times$ speedup over CC D-BCD in Yelp. However, in MNIST-Add datasets, CC D-BCD is up to $10 \times$ faster than LF D-BCD as there is a high number of tightly connected components, one for each addition instance. This behavior highlights the complementary strengths of the two parallelization strategies. LF D-BCD should be applied to problems with larger factor graph representations that are connected while CC D-BCD is effective when there are many similarly sized connected components.

6.2 LEARNING RUNTIME

Next, we study how the algorithms applied to solve inference affect the learning runtime with the SP and MSE losses. Specifically, we examine the cumulative time required for ADMM and D-BCD inference to complete 500 weight updates on the first 7 datasets in Table 1 and 100 weight updates on MNIST-Add datasets. Hyperparameters used for SP and MSE learning are reported in Appendix E.5. For inference, we apply the Table 3: Cumulative time in seconds for ADMM and D-BCD inference during learning with SP and MSE losses.

	\$	P	MSE		
	ADMM	D-BCD	ADMM	D-BCD	
CreateDebate	10.68 ± 8.63	0.34 ± 0.36	49.00 ± 31.23	0.62 ± 0.09	
4Forums	11.87 ± 12.81	0.65 ± 0.05	67.09 ± 13.79	1.11 ± 0.16	
Epinions	12.54 ± 0.37	1.33 ± 0.06	17.48 ± 0.62	2.27 ± 0.98	
Citeseer	167 ± 37	41.57 ± 6.39	225 ± 32	70.01 ± 5.86	
Cora	183 ± 26	48.16 ± 5.82	241 ± 37	$\textbf{79.62} \pm \textbf{13.77}$	
DDI	$4,554 \pm 13$	19.65 ± 0.30	$7,652 \pm 218$	52.78 ± 4.23	
Yelp	$1,835 \pm 47$	114 ± 4	$2,250 \pm 100$	170 ± 12	
MNIST-Add1	$1,624 \pm 34$	232 ± 44	$2,942 \pm 109$	$2,738 \pm 93$	
MNIST-Add2	TIME-OUT	804 ± 106	TIME-OUT	$4,291 \pm 114$	

same hyperparameters used in the previous section and the fastest parallelization method for D-BCD.

Table 3 shows that the D-BCD algorithm consistently results in the lowest total inference runtime, validating it's ability to leverage warm starts to improve learning runtimes. Notably, on the DDI dataset, D-BCD achieves roughly a $100 \times$ speedup over ADMM. Moreover, on MNIST-Add2, ADMM timed out with over 6 hours of inference time for SP and MSE learning, while D-BCD accumulated less than 0.5 and 1.2 hours of inference runtime on average for SP and MSE, respectively.

6.3 LEARNING PREDICTION PERFORMANCE

In our final experiment, we analyze the prediction performance of NeuPSL models trained with our NeSy-EBM learning framework. A hyperparameter search (detailed in Appendix E.6) is performed over learning steplengths, regularizations, and parameters for Algorithm 1.

HL-MRF learning We first evaluate the prediction performance on non-deep vari-

Table 4: Prediction performance of HL-MRF modelstrained on value and minimizer-based losses.

		Energy	SP		MSE	BCE
CreateDebate		64.76 ± 9.54	$ 64.68 \pm 11.05$	5	65.33 ± 11.98	64.83 ± 9.70
4Forums		62.96 ± 6.11	63.15 ± 6.40		64.22 ± 6.41	64.85 ± 6.01
Epinions		78.96 ± 2.29	79.85 ± 1.62		81.18 ± 2.21	80.89 ± 2.32
Citeseer		70.29 ± 1.54	70.92 ± 1.33		71.22 ± 1.56	71.94 ± 1.17
Cora		54.30 ± 1.74	74.16 ± 2.32		81.05 ± 1.41	81.07 ± 1.31
DDI		94.54 ± 0.00	94.61 ± 0.00		94.70 ± 0.00	95.08 ± 0.00
Yelp	Τ	18.11 ± 0.34	18.57 ± 0.66		18.14 ± 0.36	17.93 ± 0.50

ants of NeuPSL models for the first 7 datasets, i.e., only symbolic weights are learned. Table 4 shows that across all 7 datasets, NeuPSL models trained with Algorithm 1 obtain a better average prediction performance than those trained using a valued-based loss. On the Cora dataset, the NeuPSL model fit with the BCE loss achieves over a 6% point improvement over SP, the higher-performing value-based loss.

Deep HL-MRF learning Next, we evaluate the prediction performance of deep NeuPSL models. Here, we study the standard low-data setting for Citeseer and Cora. Specifically, results are averaged over 10 randomly sampled splits using 5% of the nodes for training, 5% of the nodes for validation, and 1,000 nodes for testing. We also report the prediction performance of the same strong baseline models used in Pryor et al. (2023) for this task: DeepStochLog (Winters et al., 2022), and a Graph Convolutional Network (GCN) (Kipf & Welling, 2017). Additionally, we investigate performance on MNIST-Addition, a widely used NeSy evaluation task first introduced by Manhaeve et al. (2018). In MNIST-Addition, models must determine the sum of two lists of MNIST images, for example, ([3] + [sr] = 8). The challenge stems from the lack of labels for the MNIST images; only the final sum of the equation is provided during training, 8 in this example. Implementation details for the neural and symbolic components of the NeuPSL models for both citation network and MNIST-Add experiments are provided in Appendix E.6.

Table 5: Accuracy of DeepStochlog, GCN, and NeuPSL on Citeseer and Cora.

				Net	ıP	SL	
	DeepStochlog	GCN	Energy	SP		MSE	BCE
Citeseer Cora	$ \begin{array}{c} 62.68 \pm 3.84 \\ 71.28 \pm 1.98 \end{array} $	$\left \begin{array}{c} 67.42 \pm 0.66 \\ 80.32 \pm 1.11 \end{array}\right.$	$\begin{array}{c} 69.63 \pm 1.33 \\ 80.41 \pm 1.81 \end{array}$	$ \begin{vmatrix} 69.78 \pm 1.42 \\ 78.59 \pm 4.93 \end{vmatrix} $		$\begin{array}{c} 69.62 \pm 1.27 \\ 81.48 \pm 1.45 \end{array}$	$\begin{array}{c} 69.64 \pm 1.33 \\ 81.28 \pm 1.45 \end{array}$

Table 6: Accuracy of CNN, LTN, DeepProblog and NeuPSL on MNIST-Addition.

		Nei	uPSL			
	Additions	CNN	LTN	DeepProblog	Energy	BCE
	300	17.16 ± 00.62	69.23 ± 15.68	85.61 ± 01.28	87.96 ± 01.58	88.84 ± 02.07
MNIST-Add1	3,000	78.99 ± 01.14	93.90 ± 00.51	92.59 ± 01.40	95.60 ± 0.91	95.70 ± 0.84
		01.31 ± 00.23		71.37 ± 03.90		76.00 ± 2.61
MNIST-Add2	1,500	01.69 ± 00.27	71.79 ± 27.76	87.44 ± 02.15	90.56 ± 0.61	93.04 ± 2.26

Table 5 shows that fitting the neural network weights of a NeuPSL model with our NeSy-EBM learning framework is effective. NeuPSL models fit with the MSE and BCE losses consistently outperform both DeepStochlog and the GCN baseline. Moreover, Table 6 demonstrates NeuPSL models trained with Algorithm 1 and a BCE loss can achieve up to a 16% point performance improvement over those trained with a value-based loss.

7 LIMITATIONS

Our learning framework is limited to NeSy-EBMs satisfying the two assumptions made in Section 4. While we advance the theory for NeuPSL to show it meets the assumptions, we do not know how to support NeSy-EBMs with non-differentiable value-functions. One approach is to substitute the inference program with a principled approximation. Lastly, although the idea to leverage inference algorithms such as BCD that effectively use warm-starts and improve learning runtimes is general, the inference algorithms were implemented for a NeSy system with an LCQP structure.

8 CONCLUSIONS AND FUTURE WORK

We introduced a general learning framework for NeSy-EBMs and demonstrated its applicability with NeuPSL. Additionally, we proposed a novel NeuPSL inference formulation and algorithm with practical and theoretical advantages. A promising direction for future work is to extend the learning framework to support approximate inference solutions for estimating the objective gradient to further improve learning runtimes. In addition, the empirical results presented in this work motivate generalizing and applying our learning framework to additional NeSy systems and tasks.

REFERENCES

- Kareem Ahmed, Stefano Teso, Kai-Wei Chang, Guy Van den Broeck, and Antonio Vergari. Semantic probabilistic layers for neuro-symbolic learning. In *NeurIPS*, 2022.
- Stephen Bach, Matthias Broecheler, Bert Huang, and Lise Getoor. Hinge-loss Markov random fields and probabilistic soft logic. *Journal of Machine Learning Research (JMLR)*, 18(1):1–67, 2017.
- Sebastian Bader and Pascal Hitzler. Dimensions of neural-symbolic integration A structured survey. *arXiv*, 2005.
- Samy Badreddine, Artur d'Avila Garcez, Luciano Serafini, and Michael Spranger. Logic tensor networks. *AI*, 303(4):103649, 2022.
- David Belanger, Bishan Yang, and Andrew McCallum. End-to-end learning for structure prediction energy networks. In *ICML*, 2017.
- Dimitri Bertsekas. Control of Uncertain Systems with a Set-Membership Description of Uncertainty. PhD thesis, MIT, 1971.
- Dimitri Bertsekas. Convex Optimization Theory. Athena Scientific, 2009.
- Dimitri Bertsekas and John N. Tsitsiklis. Parallel and Distributed Computation: Numerical Methods. Prentice Hall, 1989.
- Tarek R. Besold, Artur S. d'Avila Garcez, Sebastian Bader, Howard Bowman, Pedro M. Domingos, Pascal Hitzler, Kai-Uwe Kühnberger, Luís C. Lamb, Daniel Lowd, Priscila Machado Vieira Lima, Leo de Penning, Gadi Pinkas, Hoifung Poon, and Gerson Zaverucha. Neural-symbolic learning and reasoning: A survey and interpretation. arXiv, 2017.
- Joseph Bonnans and Alexander Shapiro. Optimization problems with perturbations: A guided tour. *SIAM Review*, 40(2):228–264, 1998.
- Joseph Bonnans and Alexander Shapiro. *Perturbation Analysis of Optimization Problems*. Springer, 2000.
- Stephen Boyd and Lieven Vandenberghe. Convex Optimization. Cambridge University Press, 2004.
- Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning (FTML)*, 3(1):1–122, 2010.
- Jerome Bracken and James T. McGill. Mathematical programs with optimization problems in the constraints. *Operations Research*, 21(1):37–44, 1973.
- Michael Collins. Discriminative training methods for hidden Markov models: Theory and experiments with perceptron algorithms. In *EMNLP*, 2002.
- Benoît Colson, Patrice Marcotte, and Gilles Savard. An overview of bilevel optimization. *Annals of Operations Research*, 153(1):235–256, 2007.
- Thomas H. Cormen, Charles E. Leiserson, Ronald L. Rivest, and Clifford Stein. *Introduction to Algorithms, Third Edition.* The MIT Press, 3rd edition, 2009.
- Cristina Cornelio, Jan Stuehmer, Shell Xu Hu, and Timothy Hospedales. Learning where and when to reason in neuro-symbolic inference. In *ICLR*, 2023.
- John Danskin. The theory of max-min, with applications. *SIAM Journal on Applied Mathematics*, 14(4):641–664, 1966.
- Sridhar Dasarth, Sai Akhil Puranam, Karmvir Aingh Phogat, Sunil Reddy Tiyyagura, and Nigel Duffy. Deeppsl: End-to-end perception and reasoning. In *IJCAI*, 2023.
- Artur d'Avila Garcez, Marco Gori, Luís C. Lamb, Luciano Serafini, Michael Spranger, and Son N. Tran. Neural-symbolic computing: An effective methodology for principled integration of machine learning and reasoning. *Journal of Applied Logics*, 6(4):611–632, 2019.

- Artur S. d'Avila Garcez, Krysia Broda, and Dov M. Gabbay. *Neural-Symbolic Learning Systems: Foundations and Applications*. Springer, 2002.
- Artur S. d'Avila Garcez, Luís C. Lamb, and Dov M. Gabbay. *Neural-Symbolic Cognitive Reasoning*. Springer, 2009.
- Luc De Raedt, Sebastijan Dumančić, Robin Manhaeve, and Giuseppe Marra. From statistical relational to neuro-symbolic artificial intelligence. In *IJCAI*, 2020.
- Chuong Do, Chuan-Sheng Foo, and Andrew Ng. Efficient multiple hyperparameter learning for log-linear models. In *NeurIPS*, 2007.
- Justin Domke. Generic methods for optimization-based modeling. In AISTATS, 2012.
- Jonathan F. Bard. *Practical Bilevel Optimization: Algorithms and Applications*. Springer Science & Business Media, 2013.
- Luca Franceschi, Paolo Frasconi, Saverio Salzo, Riccardo Grazzi, and Massimiliano Pontil. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, 2018.
- Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. Arxiv, 2018.
- Tommaso Giovannelli, Griffin Kent, and Luis Nune Vicente. Inexact bilevel stochastic gradient methods for constrained and unconstrained lower-level problems. *Arxiv*, 2022.
- Eleonora Giunchiglia, Mihaela Catalina Stoian, and Thomas Lukasiewicz. Deep learning with logical constraints. In *IJCAI*, 2022.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, 2014.
- Kazi Saidul Hasan and Vincent Ng. Stance classification of ideological debates: Data, models, features, and constraints. In *IJCNLP*, 2013.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. In *ICLR*, 2017.
- Kaiyi Ji, Junjie Yang, and Yingbin Liang. Bilevel optimization: Convergence analysis and enhanced design. In *ICML*, 2021.
- Prashant Khanduri, Ioannis Tsaknakis, Yihua Zhang, Jia Liu, Sijia Liu, Jiawei Zhang, and Mingyi Hong. Linearly constrained bilevel optimization: A smoothed implicit gradient approach. In *ICML*, 2023.
- Thomas Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *ICLR*, 2017.
- Jyrki Kivinen and Manfred K. Warmuth. Exponentiated gradient versus gradient descent for linear predictors. *Information and Computation*, 132(1):1–63, 1997.
- Pigi Kouki, Shobeir Fakhraei, James Foulds, Magdalini Eirinaki, and Lise Getoor. Hyper: A flexible and extensible probabilistic framework for hybrid recommender systems. In *RecSys*, 2015.
- Jeongyeol Kwon, Dohyun Kwon, Stephen Wright, and Robert Nowak. A fully first-order method for stochastic bilevel optimization. In *ICML*, 2023.
- Luís C. Lamb, Artur d'Avila Garcez, Marco Gori, Marcelo O. R. Prates, Pedro H. C. Avelar, and Moshe Y. Vardi. Graph neural networks meet neural-symbolic computing: A survey and perspective. In *IJCAI*, 2020.
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

- Yann LeCun, Sumit Chopra, Raia Hadsell, Marc'Aurelio Ranzato, and Fu Jie Huang. A tutorial on energy-based learning. *Predicting Structured Data*, 1(0), 2006.
- Bo Liu, Mao Ye, Stephen Wright, Peter Stone, and Qiang Liu. Bome! bilevel optimization made easy: A simple first-order approach. In *NeurIPS*, 2022.
- Ji Liu, Stephen J. Wright, Christopher Rè, Victor Bittorf, and Srikrishna Sridhar. An asynchronous parallel stochastic coordinate descent aglorithm. *Journal of Machine Learning Research (JMLR)*, 16:285–322, 2015.
- Risheng Liu, Xuan Liu, Xiaoming Yuan, Shangzhi Zeng, and Jin Zhang. A value-function-based interior-point method for non-convex bi-level optimization. In *ICML*, 2021.
- Risheng Liu, Xuan Liu, Shangzhi Zeng, Jin Zhang, and Yixuan Zhang. Value-function-based sequential minimization for bi-level optimization. *Arxiv*, 2023.
- Robin Manhaeve, Sebastijan Dumancic, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. DeepProbLog: Neural probabilistic logic programming. In *NeurIPS*, 2018.
- Robin Manhaeve, Sebastijan Dumančić, Angelika Kimmig, Thomas Demeester, and Luc De Raedt. Neural probabilistic logic programming in DeepProbLog. *Artificial Intelligence (AI)*, 298:103504, 2021.

Jorge Nocedal and Stephen J. Wright. Numerical Optimization. Springer, 2006.

- Diedrik P. Kingma and Jimmy Lei Ba. Adam: A method for stochastic optimization. In ICLR, 2017.
- Neal Parikh and Stephen Boyd. Proximal algorithms. *Foundations and Trends in Machine Learning* (*FTML*), 3(1):123–231, 2013.
- Fabian Pedregosa. Hyperparameter optimization with approximate gradient. In *ICML*, 2016.
- Connor Pryor, Charles Dickens, Eriq Augustine, Alon Albalak, William Yang Wang, and Lise Getoor. Neupsl: Neural probabilistic soft logic. In *IJCAI*, 2023.
- Aravind Rajeswaran, Chelsea Finn, Sham M. Kakade, and Sergey Levine. Meta-learning with implicit gradients. In *NeurIPS*, 2019.
- Benjamin Recht, Christopher Re, Stephen Wright, and Feng Niu. Hogwild!: A lock-free approach to prallelizing stochastic gradient descent. In *NeurIPS*, 2011.
- Matthew Richardson, Rakesh Agrawal, and Pedro Domingos. Trust management for the semantic web. In *ISWC*, 2003.
- R. Tyrrell Rockafellar. Convex Analysis. Princeton University Press, 1970.
- R. Tyrrell Rockafellar. Conjugate duality and optimization. In *Regional Conference Series in Applied Mathematics*, 1974.
- R. Tyrrell Rockafellar and Roger J-B Wets. Variational Analysis. Springer, 1997.
- David S. Wishart, Craig Knox, An Chi Guo, Savita Shrivastava, Murtaza Hassanali, Paul Stothard, Zhan Chang, and Jennifer Woolsey. Drugbank: A comprehensive resource for in silico drug discovery and exploration. *Nucleic Acids Research*, 34:D668–D672, 2006.
- Prithviraj Sen, Galileo Mark Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina Eliassi-Rad. Collective classification in network data. *AI Magazine*, 29(3):93–106, 2008.
- Shai Shalev-Shwartz. Online learning and online convex optimization. Foundations and Trends in Machine Learning (FTML), 4(2):107–194, 2012.
- Daouda Sow, Kaiyi Ji, Ziwei Guan, and Yingbin Liang. A primal-dual approach to bilevel optimization with multiple inner minima. *Arxiv*, 2022.

- Dhanya Sridhar, James Foulds, Marilyn Walker, Bert Huang, and Lise Getoor. Joint models of disagreement and stance in online debate. In *Annual Meeting of the Association for Computational Linguistics (ACL)*, 2015.
- Dhanya Sridhar, Shobeir Fakhraei, and Lise Getoor. A probabilistic approach for collective similarity-based drug-drug interaction prediction. *Bioinformatics*, 32(20):3175–3182, 2016.
- Veselin Stoyanov, Alexander Ropson, and Jason Eisner. Empirical risk minimization of graphical model parameters given approximate inference, decoding, and model structure. In AISTATS, 2011.
- Richard Sutton and Andrew Barto. Reinforcement Learning: An Introduction. The MIT Press, 2018.
- Jivrí V. Outrata. On the numerical solution of a class of stackelberg problems. Methods and Models of Operations Research, 34(4):255–277, 1990.
- Marilyn Walker, Jean Fox Tree, Pranav Anand, Rob Abbott, and Joseph King. A corpus for research on deliberation and debate. In *LREC*, 2012.
- Po-Wei Wang, Priya Donti, Bryan Wilder, and Zico Kolter. Satnet: Bridging deep learning and logical reasoning using a differentiable satisfiability solver. In *ICML*, 2019.
- Thomas Winters, Giuseppe Marra, Robin Manhaeve, and Luc De Raedt. DeepStochLog: Neural stochastic logic programming. In AAAI, 2022.
- Stephen J. Wright. Coordinate descent algorithms. Mathematical Programming, 151:3–34, 2015.
- Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang, and Guy Van den Broeck. A semantic loss function for deep learning with symbolic knowledge. In *ICML*, 2018.
- Yelp. Yelp open dataset, 2023. URL www.yelp.com/dataset.

A APPENDIX

This appendix includes the following sections: extended bilevel NeSy Learning framework, extended NeuPSL and deep hinge-loss Markov random fields, extended dual block coordinate descent, and extended empirical evaluation.

Code for running the experiments along with all data, models, and hyperparameters are available at: https://github.com/convexbilevelnesylearning/experimentscripts. Code for the NeuPSL implementation of our proposed learning framework and inference algorithms is available at: https://github.com/convexbilevelnesylearning/psl.

B EXTENDED BILEVEL NESY LEARNING FRAMEWORK

In this section we provide the further details on our proposed NeSy learning framework. A complete version of Algorithm 1 is provided in Algorithm 2.

Algorithm	2 Fu	ll NeS	v-EBM	Learning	Framework

Require: Constraint Tolerance: σ^* , Movement Tolerance: ω^* , Moreau Param.: ρ Starting points: $\mu^{(0)} > 1, \lambda_1^{(0)}, \dots, \lambda_P^{(0)}, (\mathbf{w}_{sy}^{(0)}, \mathbf{w}_{nn}^{(0)}) \in \mathcal{W}_{sy} \times \mathcal{W}_{nn}$ 1: $\mathbf{y}_{i}^{(0)} \leftarrow (\mathbf{t}_{i}, \mathbf{z}_{i}^{*}), \forall i = 1, \cdots, P;$ 2: $\iota^{(0)} \leftarrow \max_{i \in \{1, \cdots, P\}} M_{i}(\mathbf{y}_{i}^{(0)}; \mathbf{w}_{sy}^{(0)}, \mathbf{w}_{nn}^{(0)}, \rho) - V_{\mathbf{y}_{i}^{*}}(\mathbf{w}_{sy}^{(0)}, \mathbf{w}_{nn}^{(0)});$ 3: for $t = 0, 1, 2, \cdots$ do 4: Set $\omega^{(0)} = \frac{1}{\mu^{(0)}}$, and $\sigma^{(0)} = \frac{1}{(\mu^{(0)})^{0.1}}$ 5: for $k = 0, 1, 2, \cdots$ do Find $(\mathbf{w}_{sy}^{(k)}, \mathbf{w}_{nn}^{(k)}) \in \mathcal{W}_{sy} \times \mathcal{W}_{nn}, (\mathbf{y}_1^{(k)}, \cdots, \mathbf{y}_P^{(k)}) \in \mathcal{Y} \times \cdots \times \mathcal{Y}, \text{ and } \mathbf{s}^{(k)} \in \mathbb{R}_{\geq 0}^P \text{ s.t.}$ 6: $\delta^{(k)} \leftarrow \delta(\mathbf{w}_{su}^{(k)}, \mathbf{w}_{nn}^{(k)}, \mathbf{y}_{1}^{(k)}, \cdots, \mathbf{y}_{p}^{(k)}, \mathbf{s}^{(k)}; \lambda^{(k)}, \mu^{(k)}, \iota^{(k)}) \leq \omega^{(k)};$ if $\left(\sum_{i=1}^{P} c_i(\mathbf{y}_i^{(k)}, \mathbf{w}_{sy}^{(k)}, \mathbf{w}_{nn}^{(k)}, \iota^{(k)}) + s_i\right) < \sigma^{(k)}$ then 7: if $\left(\sum_{i=1}^{P} c_i(\mathbf{y}_i^{(k)}, \mathbf{w}_{sy}^{(k)}, \mathbf{w}_{nn}^{(k)}, \iota^{(k)}) + s_i\right) < \sigma^*$ and $\delta^{(k)} \leq \omega^*$ then 8: Break with the approximate solution: $\mathbf{w}_{sy}^{(k)}, \mathbf{w}_{nn}^{(k)}, \mathbf{y}_{1}^{(k)}, \cdots, \mathbf{y}_{P}^{(k)}, \mathbf{s}^{(k)};$ $\lambda_{i}^{(k+1)} \leftarrow \lambda_{i}^{(k)} + \mu^{(k)} \left(c_{i}(\mathbf{y}_{i}^{(k)}, \mathbf{w}_{sy}^{(k)}, \mathbf{w}_{nn}^{(k)}, \iota^{(k)}) + s_{i} \right), \quad \forall i = 1, \cdots, P;$ $\mu^{(k+1)} \leftarrow \mu^{(k)}; \quad \sigma^{(k+1)} \leftarrow \frac{\sigma^{(k)}}{(\mu^{(k+1)})^{0.9}}; \quad \omega^{(k+1)} \leftarrow \frac{\omega^{(k)}}{\mu^{(k+1)}};$ 9: 10: $\begin{array}{l} \textbf{else} \\ \mu^{(k+1)} \leftarrow 2 \cdot \mu^{(k)}; \ \lambda_i^{(k+1)} \leftarrow \lambda_i^{(k)}, \ \forall i = 1, \cdots, P; \\ \sigma^{(k+1)} \leftarrow \frac{1}{(\mu^{(k+1)})^{0.1}}; \ \omega^{(k+1)} \leftarrow \frac{1}{\mu^{(k+1)}}; \end{array}$ 11: 12: 13: 14: if Stopping criterion satisified then Stop with the approximate solution: $\mathbf{w}_{sy}^{(k)}, \mathbf{w}_{nn}^{(k)}, \mathbf{y}_1^{(k)}, \cdots, \mathbf{y}_P^{(k)}, \mathbf{s}^{(k)};$ 15: $\mu^{(0)} \leftarrow \mu^{(k)}; \ \lambda_i^{(0)} \leftarrow \lambda_i^{(k)}, \ \forall i = 1, \cdots, P;$ $\iota^{(t+1)} \leftarrow \frac{1}{2} \cdot \iota^{(t)};$ 16: 17:

As stated in the main paper, each instance of (11) is optimized using the bound constrained augmented Lagrangian algorithm, Algorithm 17.4 from Nocedal & Wright (2006). This algorithm is applied in lines 4 through 16 in Algorithm 2. The algorithm iteratively finds approximate minimizers of the problem's augmented Lagrangian, (12), for a fixed setting of the penalty parameters using randomized incremental gradient descent, line 6 in Algorithm 2. Specifically, gradient descent is applied to find an approximate minimizer of (12) satisfying the following stopping criterion:

$$\delta(\mathbf{w}_{sy}, \mathbf{w}_{nn}, \mathbf{y}_{1}, \cdots, \mathbf{y}_{p}, \mathbf{s}; \lambda, \mu, \iota) := \\ \|\mathbf{w}_{sy} - \Pi \left(\mathbf{w}_{sy} - \nabla_{\mathbf{w}_{sy}} \mathcal{L}_{A} \right) \| + \|\mathbf{w}_{nn} - \Pi \left(\mathbf{w}_{nn} - \nabla_{\mathbf{w}_{nn}} \mathcal{L}_{A} \right) \| \\ + \sum_{i=1}^{P} \|\mathbf{y}_{i} - \Pi \left(\mathbf{y}_{i} - \nabla_{\mathbf{y}_{i}} \mathcal{L}_{A} \right) \| + \|\mathbf{s} - \Pi \left(\mathbf{s} - \nabla_{\mathbf{s}} \mathcal{L}_{A} \right) \| \le \omega,$$
(20)

where $\omega > 0$ is a positive tolerance that is updated with the Lagrange variables. Further, note the Lagrangian gradients are evaluated at the iterate specified as arguments of δ . Practically, the parameter movement between an epoch of incremental gradient descent is used to approximate δ .

As stated in the main paper, employing the bound constrained augmented Lagrangian algorithm to solve the instances of (11) ensures fundamental convergence properties of our learning framework. Specifically, theorem 17.2 in Nocedal & Wright (2006) is applicable to Algorithm 2. This theorem states that limit points of the iterate sequence are stationary points of $||c(\mathbf{y}_1, \dots, \mathbf{y}_P, \mathbf{w}_{sy}, \mathbf{w}_{nn})| + s||^2$ when they are infeasible or, when the LICQ holds and the iterates are feasible, are KKT points of (11).

C EXTENDED NEUPSL AND DEEP HINGE-LOSS MARKOV RANDOM FIELDS

In this section, we expand on the smooth formulation of NeuPSL inference and provide proofs for the continuity results presented in Section 5.2.

C.1 EXTENDED SMOOTH FORMULATION OF INFERENCE

Recall the primal formulation of NeuPSL inference restated below:

$$\underset{\mathbf{y}\in\mathbb{R}^{n_{\mathbf{y}}}}{\arg\min} \mathbf{w}_{sy}^{T} \boldsymbol{\Phi}(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) \quad \text{s.t. } \mathbf{y}\in \boldsymbol{\Omega}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})).$$
(21)

Importantly, note the structure of the deep hinge-loss potentials defining Φ :

$$\phi_k(\mathbf{y}, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) := (\max\{\mathbf{a}_{\phi_k, y}^T \mathbf{y} + \mathbf{a}_{\phi_k, \mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{\phi_k, \mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}) + b_{\phi_k}, 0\})^{p_k}.$$
(22)

The LCQP NeuPSL inference formulation is defined using ordered index sets: I_S for the partitions of squared hinge potentials (indices k which for all $j \in t_k$ the exponent term $p_j = 2$) and I_L for the partitions of linear hinge potentials (indices k which for all $j \in t_k$ the exponent term $p_j = 1$). With the index sets, we define

$$\mathbf{W}_{S} := \begin{bmatrix} w_{\mathbf{I}_{S}[1]}\mathbf{I} & 0 & \cdots & 0\\ 0 & w_{\mathbf{I}_{S}[2]}\mathbf{I} & & \\ \vdots & & \ddots \end{bmatrix} \quad \text{and} \quad \mathbf{w}_{L} := \begin{bmatrix} w_{\mathbf{I}_{L}[1]}\mathbf{1} \\ w_{\mathbf{I}_{L}[2]}\mathbf{1} \\ \vdots \end{bmatrix}$$
(23)

Let $m_S := |\cup_{\mathbf{I}_S} t_k|$ and $m_L := |\cup_{\mathbf{I}_L} t_k|$, be the total number of squared and linear hinge potentials, respectively, and define slack variables $\mathbf{s}_S := [s_j]_{j=1}^{m_S}$ and $\mathbf{s}_L := [s_j]_{j=1}^{m_L}$ for each of the squared and linear hinge potentials, respectively. NeuPSL inference is equivalent to the following LCQP:

$$\min_{\mathbf{y}\in[0,1]^{n_y},\,\mathbf{s}_S\in\mathbb{R}^{m_S},\,\mathbf{s_H}\in\mathbb{R}_+^{m_L}}\,\mathbf{s}_S^T\mathbf{W}_S\mathbf{s}_S+\mathbf{w}_L^T\mathbf{s}_L\tag{24a}$$

s.t.
$$\mathbf{a}_{c_i,\mathbf{y}}^T \mathbf{y} + \mathbf{a}_{c_i,\mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{c_i,\mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn},\mathbf{w}_{nn}) + b_{c_i} \le 0 \quad \forall i = 1,\ldots,q,$$
 (24b)

$$\mathbf{a}_{\phi_j,\mathbf{y}}^T \mathbf{y} + \mathbf{a}_{\phi_j,\mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{\phi_j,\mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn},\mathbf{w}_{nn}) + b_{\phi_j} - s_j \le 0 \quad \forall j \in I_S \cup I_L.$$
(24c)

We ensure strong convexity by adding a square regularization with parameter ϵ to the objective. Let the bound constraints on y and s_L and linear inequalities in the LCQP be captured by the $(2 \cdot n_y + q + m_S + 2 \cdot m_L) \times (n_y + m_S + m_L)$ matrix A and $(2 \cdot n_y + q + m_S + m_L)$ dimensional vector $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$. More formally, $\mathbf{A} := [a_{ij}]$ where a_{ij} is the coefficient of a decision variable in the implicit and explicit constraints in the formulation above:

$$a_{i,j} := \begin{cases} 0 & (i \leq q) \land (j \leq m_S + m_L) \\ \mathbf{a}_{c_i,\mathbf{y}}[j - (m_S + m_L)] & (i \leq q) \land (j > m_S + m_L) \\ 0 & (q < i \leq q + m_S + m_L) \land (j \leq m_S + m_L) \land (j \neq i - q) \\ -1 & (q < i \leq q + m_S + m_L) \land (j \leq m_S + m_L) \land (j = i - q) \\ \mathbf{a}_{\phi_{i-q},\mathbf{y}}[j - (m_S + m_L)] & (q < i \leq q + m_S + m_L) \land (j > m_S + m_L) \\ 0 & (q + m_S + m_L < i \leq q + m_S + 2 \cdot m_L + n_y) \\ \land (j \neq i - (q + m_L)) & \land (j = i - (q + m_L)) \\ -1 & (q + m_S + m_L < i \leq q + m_S + 2 \cdot m_L + n_y) \\ \land (j = i - (q + m_L)) & (q + m_S + 2 \cdot m_L + n_y) \\ \land (j \neq i - (q + m_S + m_L)) \\ 1 & (q + m_S + 2 \cdot m_L + n_y < i \leq q + m_S + 2 \cdot m_L + 2 \cdot n_y) \\ \land (j = i - (q + m_S + m_L)) \\ 1 & (q + m_S + 2 \cdot m_L + n_y < i \leq q + m_S + 2 \cdot m_L + 2 \cdot n_y) \\ \land (j = i - (q + m_S + m_L)) \end{cases}$$

$$(25)$$

Furthermore, $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) = [b_i(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))]$ is the vector of constants corresponding to each constraint in the formulation above:

$$b_{i}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$$
(26)
$$:= \begin{cases} \mathbf{a}_{c_{i}, \mathbf{x}_{sy}}^{T} \mathbf{x}_{sy} + \mathbf{a}_{c_{i}, \mathbf{g}}^{T} \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}) + b_{c_{i}} & i \leq q \\ \mathbf{a}_{\phi_{i-q}, \mathbf{x}_{sy}}^{T} \mathbf{x}_{sy} + \mathbf{a}_{\phi_{i-q}, \mathbf{g}}^{T} \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}) + b_{\phi_{i-q}} & q < i \leq q + m_{S} + m_{L} \\ 0 & q + m_{S} + m_{L} < i \\ \leq q + m_{S} + 2 \cdot m_{L} + n_{y} \\ -1 & q + m_{S} + 2 \cdot m_{L} + n_{y} < i \\ \leq q + m_{S} + 2 \cdot m_{L} + n_{y} < i \\ \leq q + m_{S} + 2 \cdot m_{L} + 2 \cdot n_{y} \end{cases}$$
(27)

Note that $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$ is a linear function of the neural network outputs, hence, if $\mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})$ is a smooth function of the neural parameters, then $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$ is also smooth.

With this notation, the regularized inference problem is:

$$V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) := \min_{\mathbf{y}, \mathbf{s}_{S}, \mathbf{s}_{H}} \begin{bmatrix} \mathbf{s}_{S} \\ \mathbf{s}_{L} \\ \mathbf{y} \end{bmatrix}^{T} \begin{bmatrix} \mathbf{W}_{S} + \epsilon I & 0 & 0 \\ 0 & \epsilon I & 0 \\ 0 & 0 & \epsilon I \end{bmatrix} \begin{bmatrix} \mathbf{s}_{S} \\ \mathbf{s}_{L} \\ \mathbf{y} \end{bmatrix} + \begin{bmatrix} 0 \\ \mathbf{w}_{L} \\ 0 \end{bmatrix}^{T} \begin{bmatrix} \mathbf{s}_{S} \\ \mathbf{s}_{L} \\ \mathbf{y} \end{bmatrix}$$

s.t. $\mathbf{A} \begin{bmatrix} \mathbf{s}_{S} \\ \mathbf{s}_{L} \\ \mathbf{y} \end{bmatrix} + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) \le 0.$ (28)

For ease of notation, let

$$D(\mathbf{w}_{sy}) := \begin{bmatrix} \mathbf{W}_S & 0 & 0\\ 0 & 0 & 0\\ 0 & 0 & 0 \end{bmatrix}, \ \mathbf{c}(\mathbf{w}_{sy}) := \begin{bmatrix} 0\\ \mathbf{w}_L\\ 0 \end{bmatrix}, \ \nu := \begin{bmatrix} \mathbf{s}_S\\ \mathbf{s}_L\\ \mathbf{y} \end{bmatrix}.$$
(29)

Then the regularized primal LCQP MAP inference problem is concisely expressed as

$$\min_{\nu \in \mathbb{R}^{n_{\mathbf{y}}+m_{S}+m_{L}}} \nu^{T} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I}) \nu + \mathbf{c}(\mathbf{w}_{sy})^{T} \nu$$
s.t. $\mathbf{A}\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) \le 0.$
(30)

By Slater's constraint qualification, we have strong-duality when there is a feasible solution. In this case, an optimal solution to the dual problem yields an optimal solution to the primal problem. The

Lagrange dual problem of (30) is

$$\arg\max_{\mu\geq 0} \min_{\nu\in\mathbb{R}^{n_{\mathbf{y}}+m_{S}+m_{L}}} \nu^{T}(\mathbf{D}(\mathbf{w}_{sy})+\epsilon\mathbf{I})\nu + \mathbf{c}(\mathbf{w}_{sy})^{T}\nu + \mu^{T}(\mathbf{A}\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$$

$$= \arg\max_{\mu\geq 0} -\frac{1}{4}\mu^{T}\mathbf{A}(\mathbf{D}(\mathbf{w}_{sy})+\epsilon\mathbf{I})^{-1}\mathbf{A}^{T}\mu$$

$$-\frac{1}{2}(\mathbf{A}(\mathbf{D}(\mathbf{w}_{sy})+\epsilon\mathbf{I})^{-1}\mathbf{c}(\mathbf{w}_{sy}) - 2\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))^{T}\mu$$
(31)

where $\mu = [\mu_i]_{i=1}^{n_{\mu}}$ are the Lagrange dual variables. For later reference, denote the negative of the Lagrange dual function of MAP inference as:

$$h(\mu; \mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$$
(32)
$$:= \frac{1}{4} \mu^T \mathbf{A} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})^{-1} \mathbf{A}^T \mu + \frac{1}{2} (\mathbf{A} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})^{-1} \mathbf{c}(\mathbf{w}_{sy}) - 2\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))^T \mu.$$

The dual LCQP has more decision variables but is only over non-negativity constraints rather than the complex polyhedron feasible set. The dual-to-primal variable translation is:

$$\nu = -\frac{1}{2} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})^{-1} (\mathbf{A}^T \boldsymbol{\mu} + \mathbf{c}(\mathbf{w}_{sy}))$$
(33)

As $(\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})$ is diagonal, it is easy to invert and hence it is practical to work in the dual space to obtain a solution to the primal problem.

C.2 EXTENDED CONTINUITY OF INFERENCE

We now provide background on sensitivity analysis that we then apply in our proofs on the continuity properties of NeuPSL inference.

C.2.1 BACKGROUND

Theorem C.1 (Boyd & Vandenberghe (2004) p. 81). *If for each* $\mathbf{y} \in A$, $f(\mathbf{x}, \mathbf{y})$ *is convex in* \mathbf{x} *then the function*

$$g(\mathbf{x}) := \sup_{\mathbf{y} \in \mathcal{A}} f(\mathbf{x}, \mathbf{y}) \tag{34}$$

is convex in \mathbf{x} .

Theorem C.2 (Boyd & Vandenberghe (2004) p. 81). *If for each* $\mathbf{y} \in A$, $f(\mathbf{x}, \mathbf{y})$ *is concave in* \mathbf{x} *then the function*

$$g(\mathbf{x}) := \inf_{\mathbf{y} \in \mathcal{A}} f(\mathbf{x}, \mathbf{y})$$
(35)

is concave in x.

Definition C.3 (Convex Subgradient: Boyd & Vandenberghe (2004) and Shalev-Shwartz (2012)). Consider a convex function $f : \mathbb{R}^n \to [-\infty, \infty]$ and a point $\overline{\mathbf{x}}$ with $f(\overline{\mathbf{x}})$ finite. For a vector $\mathbf{v} \in \mathbf{R}^n$, one says that \mathbf{v} is a (convex) subgradient of f at $\overline{\mathbf{x}}$, written $\mathbf{v} \in \partial f(\overline{\mathbf{x}})$, iff

$$f(\mathbf{x}) \ge f(\overline{\mathbf{x}}) + \langle \mathbf{v}, \mathbf{x} - \overline{\mathbf{x}} \rangle, \quad \forall \mathbf{x} \in \mathbf{R}^n.$$
 (36)

Definition C.4 (Closedness: Bertsekas (2009)). If the epigraph of a function $f : \mathbb{R}^n \to [-\infty, \infty]$ is a closed set, we say that f is a closed function.

Definition C.5 (Lower Semicontinuity: Bertsekas (2009)). The function $f : \mathbb{R}^n \to [-\infty, \infty]$ is *lower semicontinuous* (lsc) at a point $\overline{\mathbf{x}} \in \mathbb{R}^n$ if

$$f(\overline{\mathbf{x}}) \le \liminf_{k \to \infty} f(\mathbf{x}_k),\tag{37}$$

for every sequence $\{\mathbf{x}_k\} \subset \mathbb{R}^n$ with $\mathbf{x}_k \to \overline{\mathbf{x}}$. We say f is *lsc* if it is lsc at each $\overline{\mathbf{x}}$ in its domain.

Theorem C.6 (Closedness and Semicontinuity: Bertsekas (2009) Proposition 1.1.2.). For a function $f : \mathbb{R}^n \to [-\infty, \infty]$, the following are equivalent:

- 1. The level set $V_{\gamma} = \{\mathbf{x} \mid f(\mathbf{x}) \leq \gamma\}$ is closed for every scalar γ .
- 2. *f* is lsc.
- 3. f is closed.

The following definition and theorem are from Rockafellar & Wets (1997) and they generalize the notion of subgradients to non-convex functions and the chain rule of differentiation, respectively. For complete statements see Rockafellar & Wets (1997) Rockafellar & Wets (1997).

Definition C.7 (Regular Subgradient: Rockafellar & Wets (1997) Definition 8.3). Consider a function $f : \mathbb{R}^n \to [-\infty, \infty]$ and a point $\overline{\mathbf{x}}$ with $f(\overline{\mathbf{x}})$ finite. For a vector $\mathbf{v} \in \mathbf{R}^n$, one says that \mathbf{v} is a regular subgradient of f at $\overline{\mathbf{x}}$, written $\mathbf{v} \in \hat{\partial} f(\overline{\mathbf{x}})$, iff

$$f(\mathbf{x}) \ge f(\overline{\mathbf{x}}) + \langle \mathbf{v}, \mathbf{x} - \overline{\mathbf{x}} \rangle + o(\mathbf{x} - \overline{\mathbf{x}}), \quad \forall \mathbf{x} \in \mathbf{R}^n,$$
(38)

where the o(t) notation indicates a term with the property that

$$\lim_{t \to 0} \frac{\mathbf{o}(t)}{t} = 0.$$
(39)

The relation of the regular subgradient defined above and the more familiar convex subgradient is the addition of the $o(\mathbf{x} - \overline{\mathbf{x}})$ term. Evidently, a convex subgradient is a regular subgradient.

Theorem C.8 (Chain Rule for Regular Subgradients: Rockafellar & Wets (1997) Theorem 10.6). Suppose $f(\mathbf{x}) = g(F(\mathbf{x}))$ for a proper, lsc function $g : \mathbb{R}^m \to [-\infty, \infty]$ and a smooth mapping $F : \mathbb{R}^n \to \mathbb{R}^m$. Then at any point $\overline{\mathbf{x}} \in dom f = F^{-1}(dom g)$ one has

$$\hat{\partial}f(\overline{\mathbf{x}}) \supset \nabla F(\overline{\mathbf{x}})^T \hat{\partial}g(F(\overline{\mathbf{x}})),\tag{40}$$

where $\nabla F(\overline{\mathbf{x}})^T$ is the Jacobian of F at $\overline{\mathbf{x}}$.

Theorem C.9 (Danskin's Theorem: Danskin (1966) and Bertsekas (1971) Proposition A.22). Suppose $\mathcal{Z} \subseteq \mathbb{R}^m$ is a compact set and $g(\mathbf{x}, \mathbf{z}) : \mathbb{R}^n \times \mathcal{Z} \to (-\infty, \infty]$ is a function. Suppose $g(\cdot, \mathbf{z}) : \mathbb{R}^n \to \mathbb{R}$ is closed proper convex function for every $\mathbf{z} \in \mathcal{Z}$. Further, define the function $f : \mathbb{R}^n \to \mathbb{R}$ such that

$$f(\mathbf{x}) := \max_{\mathbf{z} \in \mathcal{Z}} g(\mathbf{x}, \mathbf{z}).$$

Suppose f is finite somewhere. Moreover, let $\mathcal{X} := int(dom f)$, i.e., the interior of the set of points in \mathbb{R}^n such that f is finite. Suppose g is continuous on $\mathcal{X} \times \mathcal{Z}$. Further, define the set of maximizing points of $g(\mathbf{x}, \cdot)$ for each \mathbf{x}

$$Z(\mathbf{x}) = \operatorname*{arg\,max}_{\mathbf{z}\in\mathcal{Z}} g(\mathbf{x}, \mathbf{z})$$

Then the following properties of f hold.

- *1. The function* $f(\mathbf{x})$ *is a closed proper convex function.*
- 2. For every $\mathbf{x} \in \mathcal{X}$,

$$\partial f(\mathbf{x}) = conv \left\{ \partial_{\mathbf{x}} g(\mathbf{x}, \mathbf{z}) \, | \, \mathbf{z} \in Z(\mathbf{x}) \right\}. \tag{41}$$

Corollary C.10. Assume the conditions for Danskin's Theorem above hold. For every $\mathbf{x} \in \mathcal{X}$, if $Z(\mathbf{x})$ consists of a unique point, call it \mathbf{z}^* , and $g(\cdot, \mathbf{z}^*)$ is differentiable at \mathbf{x} , then $f(\cdot)$ is differentiable at \mathbf{x} , and

$$\nabla f(\mathbf{x}) := \nabla_{\mathbf{x}} g(\mathbf{x}, \mathbf{z}^*). \tag{42}$$

Theorem C.11 (Bonnans & Shapiro (1998) Theorem 4.2, Rockafellar (1974) p. 41). Let X and U be Banach spaces. Let K be a closed convex cone in the Banach space U. Let $G : X \to U$ be a convex mapping with respect to the cone $\mathbf{C} := -\mathbf{K}$ and $f : X \to (-\infty, \infty]$ be a (possibly infinite-valued) convex function. Consider the following convex program and its optimal value function:

$$v_P(\mathbf{u}) := \min_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x})$$
(P)
s.t. $G(\mathbf{x}) + \mathbf{u} \in \mathbf{K}.$

Moreover, consider the (Lagrangian) dual of the program:

$$v_D(\mathbf{u}) := \max_{\lambda \in \mathbf{K}^-} \min_{\mathbf{x} \in \mathbf{X}} f(\mathbf{x}) + \lambda^T (G(\mathbf{x}) + \mathbf{u})$$
(D)

Suppose $v_P(\mathbf{0})$ is finite. Further, suppose the feasible set of the program is nonempty for all \mathbf{u} in a neighborhood of $\mathbf{0}$, i.e.,

$$\mathbf{0} \in int\{G(\mathbf{X}) - \mathbf{K}\}. \tag{43}$$

Then,

- 1. There is no primal dual gap at u = 0, i.e., $v_P(0) = v_D(0)$.
- 2. The set, Λ_0 , of optimal solutions to the dual problem with $\mathbf{u} = 0$ is non-empty and bounded.
- 3. The optimal value function $v_P(\mathbf{u})$ is continuous at $\mathbf{u} = 0$ and $\partial v_P(\mathbf{0}) = \Lambda_0$.

Theorem C.12 (Bonnans & Shapiro (2000) Proposition 4.3.2). *Consider two optimization problems over a non-empty feasible set* Ω *:*

$$\min_{\mathbf{x}\in\Omega} f_1(\mathbf{x}) \quad and \quad \min_{\mathbf{x}\in\Omega} f_2(\mathbf{x}) \tag{44}$$

where $f_1, f_2 : \mathcal{X} \to \mathbb{R}$. Suppose f_1 has a non-empty set \mathbf{S} of optimal solutions over Ω . Suppose the second order growth condition holds for \mathbf{S} , i.e., there exists a neighborhood \mathcal{N} of \mathbf{S} and a constant $\alpha > 0$ such that

$$f_1(\mathbf{x}) \ge f_1(\mathbf{S}) + \alpha (dist(\mathbf{x}, \mathbf{S}))^2, \quad \forall \mathbf{x} \in \mathbf{\Omega} \cap \mathcal{N},$$
(45)

where $f_1(\mathbf{S}) := inf_{\mathbf{x} \in \mathbf{\Omega}} f_1(\mathbf{x})$. Define the difference function:

$$\Delta(\mathbf{x}) := f_2(\mathbf{x}) - f_1(\mathbf{x}). \tag{46}$$

Suppose $\Delta(\mathbf{x})$ is L-Lipschitz continuous on $\mathbf{\Omega} \cap \mathcal{N}$. Let $\mathbf{x}^* \in \mathcal{N}$ be an δ -solution to the problem of minimizing $f_2(\mathbf{x})$ over $\mathbf{\Omega}$. Then

$$dist(\mathbf{x}^*, \mathbf{S}) \le \frac{L}{\alpha} + \sqrt{\frac{\delta}{\alpha}}.$$
 (47)

C.2.2 PROOFS

We provide proofs of theorems presented in the main paper and restated here for completeness. *Theorem 5.2.* Suppose for any setting of $\mathbf{w}_{nn} \in \mathbb{R}^{n_g}$ there is a feasible solution to NeuPSL inference (16). Further, suppose $\epsilon > 0$, $\mathbf{w}_{sy} \in \mathbb{R}^r_+$, and $\mathbf{w}_{nn} \in \mathbb{R}^{n_g}$. Then:

- The minimizer of (16), $\mathbf{y}^*(\mathbf{w}_{sy}, \mathbf{w}_{nn})$, is a $O(1/\epsilon)$ Lipschitz continuous function of \mathbf{w}_{sy} .
- $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))))$, is concave over \mathbf{w}_{sy} and convex over $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$.
- $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is differentiable with respect to \mathbf{w}_{sy} . Moreover,

$$\nabla_{\mathbf{w}_{sy}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) = \mathbf{\Phi}(\mathbf{y}^*(\mathbf{w}_{sy}, \mathbf{w}_{nn}), \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})).$$

Furthermore, $\nabla_{\mathbf{w}_{sy}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is Lipschitz continuous over \mathbf{w}_{sy} .

• If there is a feasible point ν strictly satisfying the *i'th* inequality constraint of (16), i.e., $\mathbf{A}[i]\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i] < 0$, then $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is subdifferentiable with respect to the *i'th* constraint constant $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i]$. Moreover,

$$\partial_{\mathbf{b}[i]} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) = \{\mu^*[i] \mid \mu^* \in \operatorname*{arg\,min}_{\mu \in \mathbb{R}^{2 \cdot n_{\mathbf{y}} + m + q}_{> 0}} h(\mu; \mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))\}.$$

Furthermore, if $\mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})$ is a smooth function of \mathbf{w}_{nn} , then so is $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$, and the set of regular subgradients of $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is:

$$\hat{\partial}_{\mathbf{w}_{nn}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$$

$$\supset \nabla_{\mathbf{w}_{nn}} \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))^T \partial_{\mathbf{b}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))).$$
(48)

Proof of Theorem 5.2. We first show the minimizer of the LCQP formulation of NeuPSL inference, ν^* , with $\epsilon > 0$, $\mathbf{w}_{sy} \in \mathbb{R}^r_+$, and $\mathbf{w}_{nn} \in \mathbb{R}^{n_g}$ is a Lipschitz continuous function of \mathbf{w}_{sy} . Suppose $\epsilon > 0$. To show continuity over $\mathbf{w}_{sy} \in \mathbb{R}^r_+$, first note the matrix $(\mathbf{D} + \epsilon \mathbf{I})$ is positive definite and the primal inference problem (17) is an ϵ -strongly convex LCQP with a unique minimizer denoted by $\nu^*(\mathbf{w}_{sy}, \mathbf{w}_{nn})$. We leverage the Lipschitz stability result for optimal values of constrained problems from Bonnans & Shapiro (2000) and presented here in Theorem C.12. Define the primal objective as an explicit function of the weights:

$$f(\nu, \mathbf{w}_{sy}, \mathbf{w}_{nn}) := \nu^T (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I})\nu + \mathbf{c}^T (\mathbf{w}_{sy})\nu$$
(49)

Note that the solution $\nu^* = \begin{bmatrix} \mathbf{s}_S^* \\ \mathbf{s}_L^* \\ \mathbf{y}^* \end{bmatrix}$ will always be bounded, since from (24c) in LCQP we always

have for all $j \in I_S \cup I_L$,

$$0 \le s_j^* = \max(\mathbf{a}_{\phi_k,y}^T \mathbf{y}^* + \mathbf{a}_{\phi_k,\mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{\phi_k,\mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn},\mathbf{w}_{nn}) + b_{\phi_k}, 0)$$
(50)

$$\leq \|\mathbf{a}_{\phi_k,y}\| + |\mathbf{a}_{\phi_k,\mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{\phi_k,\mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn},\mathbf{w}_{nn}) + b_{\phi_k}|.$$
(51)

Thus, setting these trivial upper bounds for s_j will not change the solution of the problem. We can henceforth consider the problem in a bounded domain $\|\nu\| \leq C$ where C does not depend on w's.

Let $\mathbf{w}_{1,sy}, \mathbf{w}_{2,sy} \in \mathbb{R}^r_+$ and $\mathbf{w}_{nn} \in \mathcal{W}_{nn}$ be arbitrary. As $\epsilon > 0$, $f(\nu, \mathbf{w}_{1,sy}, \mathbf{w}_{nn})$ is strongly convex in ν and it therefore satisfies the second-order growth condition in ν . Define the difference function:

$$\Delta_{\mathbf{w}_{sy}}(\nu) := f(\nu, \mathbf{w}_{2,sy}, \mathbf{w}_{nn}) - f(\nu, \mathbf{w}_{1,sy}, \mathbf{w}_{nn})$$
(52)

$$=\nu^{T}(\mathbf{D}(\mathbf{w}_{2,sy})+\epsilon\mathbf{I})\nu+\mathbf{c}^{T}(\mathbf{w}_{2,sy})\nu-\left(\nu^{T}(\mathbf{D}(\mathbf{w}_{1,sy})+\epsilon\mathbf{I})\nu+\mathbf{c}^{T}(\mathbf{w}_{1,sy})\nu\right)$$
(53)

$$= \nu^{T} (\mathbf{D}(\mathbf{w}_{2,sy}) - \mathbf{D}(\mathbf{w}_{1,sy}))\nu + (\mathbf{c}(\mathbf{w}_{2,sy}) - \mathbf{c}(\mathbf{w}_{1,sy}))^{T}\nu.$$
(54)

The difference function $\Delta_{\mathbf{w}_{sy}}(\nu)$ over $\mathcal N$ has a finitely bounded gradient:

$$\|\nabla \Delta_{\mathbf{w}_{sy}}(\nu)\|_{2} = \left\|2(\mathbf{D}(\mathbf{w}_{2,sy}) - \mathbf{D}(\mathbf{w}_{1,sy}))\nu + \mathbf{c}(\mathbf{w}_{2,sy}) - \mathbf{c}(\mathbf{w}_{1,sy})\right\|_{2}$$
(55)

$$\leq \|\mathbf{c}(\mathbf{w}_{2,sy}) - \mathbf{c}(\mathbf{w}_{1,sy})\|_2 + 2\|(\mathbf{D}(\mathbf{w}_{2,sy}) - \mathbf{D}(\mathbf{w}_{1,sy}))\nu\|_2$$
(56)

$$\leq \|\mathbf{w}_{2,sy} - \mathbf{w}_{1,sy}\|_{2} + 2\|\mathbf{w}_{2,sy} - \mathbf{w}_{1,sy}\|_{2} \|\nu\|_{2}$$
(57)

$$\leq \|\mathbf{w}_{2,sy} - \mathbf{w}_{1,sy}\|_2 (1+2C) =: L_{\mathcal{N}}(\mathbf{w}_{1,sy}, \mathbf{w}_{2,sy}).$$
(58)

Thus, the distance function, $\Delta_{\mathbf{w}_{sy}}(\nu)$ is $L_{\mathcal{N}}(\mathbf{w}_{1,sy}, \mathbf{w}_{2,sy})$ -Lipschitz continuous over \mathcal{N} . Therefore, by Bonnans & Shapiro (2000) (Theorem C.12), the distance between $\nu^*(\mathbf{w}_{1,sy}, \mathbf{w}_{nn})$ and $\nu^*(\mathbf{w}_{2,sy}, \mathbf{w}_{nn})$ is bounded above:

$$\|\nu^{*}(\mathbf{w}_{2,sy},\mathbf{w}_{nn}) - \nu^{*}(\mathbf{w}_{1,sy},\mathbf{w}_{nn})\|_{2} \leq \frac{L_{\mathcal{N}}(\mathbf{w}_{1,sy},\mathbf{w}_{2,sy})}{\epsilon} = \frac{(1+2C)}{\epsilon} \|\mathbf{w}_{2,sy} - \mathbf{w}_{1,sy}\|_{2}.$$
(59)

Therefore, the function $\nu^*(\mathbf{w}_{sy}, \mathbf{w}_{nn})$ is $O(1/\epsilon)$ -Lipschitz continuous in \mathbf{w}_{sy} for any \mathbf{w}_{nn} .

Next, we prove curvature properties of the value-function with respect to the weights. Observe NeuPSL inference is an infimum over a set of functions that are concave (affine) in \mathbf{w}_{sy} . Therefore, by Theorem C.2, we have that $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is concave in \mathbf{w}_{sy} .

We use a similar argument to show $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is convex in the constraint constants, $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$. Assuming for any setting of the neural weights, $\mathbf{w}_{nn} \in \mathbb{R}^{n_g}$, there is a feasible solution to the NeuPSL inference problem, then (16) satisfies the conditions for Slater's constraint qualification. Therefore, strong duality holds, i.e., $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is equal to the optimal value of the dual inference problem (31). Observe that the dual NeuPSL inference problem is a supremum over a set of functions convex (affine) in $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$. Therefore, by Theorem C.1, we have that $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$ is convex in $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$.

We can additionally prove convexity in b from first principles. For simplicity we fix other parameters, and write the objective and the value function as $Q(\nu)$ and $V(\mathbf{b})$. Given two values \mathbf{b}_1 and

b₂, let corresponding optimal values of (33) be ν_1 and ν_2 . Take any $\alpha \in [0, 1]$, note that when $\mathbf{b} = \alpha \mathbf{b}_1 + (1 - \alpha) \mathbf{b}_2$, then $\alpha \nu_1 + (1 - \alpha) \nu_2$ is feasible for this **b**. Because we take the inf over all ν s, the optimal ν for this **b** might be even smaller. Thus, we have (for convex quadratic objective Q) that

$$V(\alpha b_{1} + (1 - \alpha)b_{2}) \leq Q(\alpha \nu_{1} + (1 - \alpha)\nu_{2})$$

$$\leq \alpha Q(\nu_{1}) + (1 - \alpha)Q(\nu_{2})$$

$$= \alpha V(b_{1}) + (1 - \alpha)V(b_{2}),$$
(60)

which shows that V is convex in **b**.

Next, we prove (sub)differentiability properties of the value-function. Suppose $\epsilon > 0$. First, we show the optimal value function, $V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$, is differentiable with respect to the symbolic weights. Then we show subdifferentiability properties of the optimal value function with respect to the constraint constants. Finally, we apply the Lipschitz continuity of the minimzer result to show the gradient of the optimal value function is Lipschitz continuous with respect to \mathbf{w}_{sy} .

Starting with differentiability with respect to the symbolic weights, \mathbf{w}_{sy} , note, the optimal value function of the regularized LCQP formulation of NeuPSL inference, (16), is equivalently expressed as the following maximization over a continuous function in the primal target variables, \mathbf{y} , the slack variables, \mathbf{s}_S and \mathbf{s}_L , and the symbolic weights, \mathbf{w}_{sy} :

$$V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))$$

$$= -\left(\max_{\mathbf{y}, \mathbf{s}_{\mathbf{H}}, \mathbf{s}_{\mathbf{L}}} - \left(\begin{bmatrix}\mathbf{s}_{S}\\\mathbf{s}_{L}\\\mathbf{y}\end{bmatrix}^{T}\begin{bmatrix}\mathbf{W}_{S} + \epsilon I & 0 & 0\\ 0 & \epsilon I & 0\\ 0 & 0 & \epsilon I\end{bmatrix}\begin{bmatrix}\mathbf{s}_{S}\\\mathbf{s}_{L}\\\mathbf{y}\end{bmatrix} + \begin{bmatrix}\mathbf{0}\\\mathbf{w}_{L}\\\mathbf{0}\end{bmatrix}^{T}\begin{bmatrix}\mathbf{s}_{S}\\\mathbf{s}_{L}\\\mathbf{y}\end{bmatrix}\right)\right)$$

$$s.t. \quad \mathbf{A}\begin{bmatrix}\mathbf{s}_{S}\\\mathbf{s}_{L}\\\mathbf{y}\end{bmatrix} + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}) \le 0,$$

$$(61)$$

where the matrix \mathbf{W}_s and vector \mathbf{w}_L are functions of the symbolic parameters \mathbf{w}_{sy} as defined in (23). Moreover, the objective above is and convex (affine) in \mathbf{w}_{sy} . Additionally, note that the decision variables can be constrained to a compact domain without breaking the equivalence of the formulation. Specifically, the target variables are constrained to the box $[0, 1]^{\mathbf{n}_y}$, while the slack variables are nonnegative and have a trivial upper bound derived from (24c):,

$$0 \leq s_j^* = \max(\mathbf{a}_{\phi_k,y}^T \mathbf{y}^* + \mathbf{a}_{\phi_k,\mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{\phi_k,\mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}) + b_{\phi_k}, 0)$$

$$\leq \|\mathbf{a}_{\phi_k,y}\| + |\mathbf{a}_{\phi_k,\mathbf{x}_{sy}}^T \mathbf{x}_{sy} + \mathbf{a}_{\phi_k,\mathbf{g}}^T \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}) + b_{\phi_k}|,$$
(62)

for all $j \in I_S \cup I_L$. Therefore, the negative optimal value function satisfies the conditions for Danskin's theorem Danskin (1966) (stated in Appendix C.2.1). Moreover, as there is a single unique solution to the inference problem when $\epsilon > 0$, and the quadratic objective in (16) is differentiable for all $\mathbf{w}_{sy} \in \mathbb{R}^r_+$, we can apply Corollary C.10. The optimal value function is therefore concave and differentiable with respect to the symbolic weights with

$$\nabla_{\mathbf{w}_{sy}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) = \mathbf{\Phi}(\mathbf{y}^*, \mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})).$$
(63)

Next, we show subdifferentiability of the optimal value-function with respect to the constraint constants, $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$. Suppose at a setting of the neural weights $\mathbf{w}_{nn} \in \mathbb{R}^{n_g}$ there is a feasible point ν for the NeuPSL inference problem. Moreover, suppose ν strictly satisfies the i'thinequality constraint of (16), i.e., $\mathbf{A}[i]\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i] < 0$. Observe that the following strongly convex conic program is equivalent to the LCQP formulation of NeuPSL inference, (16):

$$\min_{\nu \in \mathbb{R}^{n_{\mathbf{y}}+m_{S}+m_{L}}} \nu^{T} (\mathbf{D}(\mathbf{w}_{sy}) + \epsilon \mathbf{I}) \nu + \mathbf{c}(\mathbf{w}_{sy})^{T} \nu + P_{\Omega \setminus i}(\nu)$$
s.t. $\mathbf{A}[i] \nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i] \in \mathbb{R}_{\leq 0},$
(64)

where $P_{\Omega \setminus i}(\nu) : \mathbb{R}^{n_y+m_s+m_L} \to \{0, \infty\}$ is the indicator function identifying feasibility w.r.t. all the constraints of the LCQP formulation of NeuPSL inference in (16) except the *i'th* constraint: $\mathbf{A}[i]\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i] \leq 0$. In other words, in the conic formulation above only the *i'th* constraint is explicit. Note that $\mathbb{R}_{\leq 0}$ is a closed convex cone in \mathbb{R} . Moreover, both the objective in the program and the mapping $G(\nu) := \mathbf{A}[i]\nu + \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i]$ are convex. Lastly, note the constraint qualification (43) is similar to Slater's condition in the case of (64) which is satisfied by the supposition there exists a feasible ν that strictly satisfies the *i'th* inequality constraint of (16). Therefore, (64) satisfies the conditions of Theorem C.11. Thus, the value function is continuous in the constraint constant $\mathbf{b}(\mathbf{x}_{su}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))[i]$ at \mathbf{w}_{nn} and

$$\partial_{\mathbf{b}[i]} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))) = \{\mu^*[i] \mid \mu^* \in \operatorname*{arg\,min}_{\mu \in \mathbb{R}^{2 \cdot n_{\mathbf{y}} + m + q}_{\geq 0}} h(\mu; \mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})))\}$$

Moreover, when $\mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn}))$ is a smooth function of the neural weights \mathbf{w}_{nn} , then we can apply the chain rule for regular subgradients, Theorem C.8, to get

$$\hat{\partial}_{\mathbf{w}_{nn}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})) \supset \nabla \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})^T \partial_{\mathbf{b}} V(\mathbf{w}_{sy}, \mathbf{b}(\mathbf{x}_{sy}, \mathbf{g}(\mathbf{x}_{nn}, \mathbf{w}_{nn})).$$
(66)

To prove the optimal value function is Lipschitz smooth over \mathbf{w}_{sy} , it is equivalent to show it is continuously differentiable and that all gradients have bounded magnitude. To show the value function is continuously differentiable, we first apply the result asserting the minimizer is unique and a continuous function of the symbolic parameters \mathbf{w}_{sy} . Therefore, the optimal value function gradient is a composition of continuous functions, hence continuous in \mathbf{w}_{sy} . The fact that the value function has a bounded gradient magnitude follows from the fact that the decision variables \mathbf{y} have a compact domain over which the gradient is finite; hence a trivial and finite upper bound exists on the gradient magnitude.

D EXTENDED DUAL BLOCK COORDINATE DESCENT

We introduce a novel block coordinate descent (BCD) algorithm for the dual LCQP formulation of NeuPSL inference in (31), a bound-constrained, strongly convex quadratic program. In this section, we omit the symbolic and neural weights from the function arguments to simplify notation. We define U_i , i = 1, 2, ..., p to be a cover of the dual variable components $\{1, 2, ..., n_y + m + q\}$. In practice, blocks are defined as a single dual variable corresponding to a constraint from the feasible set or a deep hinge-loss function, along with the dual variables corresponding to the bounds of the primal variables in the constraint or hinge-loss.

We will deal with a slightly more general objective,

$$h(\mu) := \frac{1}{2} \mu^T \mathbf{A} \tilde{\mathbf{D}} \mathbf{A}^T \mu + \tilde{\mathbf{c}}^T \mu,$$
(67)

(65)

from which we can recover (32) by replacing $\tilde{\mathbf{D}} \leftarrow (\mathbf{D} + \epsilon \mathbf{I})^{-1}$ and $\tilde{\mathbf{c}} \leftarrow \mathbf{A}(\mathbf{D} + \epsilon \mathbf{I})^{-1}\mathbf{c} - 2\mathbf{b}$.

We will use the superscript $\cdot^{(l)}$ to denote values in the *l*-th iteration and subscript $\cdot_{[i]}$ for the values corresponding to the block U_i . The row submatrix of **A** that corresponds to block *i* is denote by $\mathbf{A}_{[i]}$.

At each iteration l, we choose one block $i \in \{1, 2, ..., p\}$ at random and compute the subvector of $\nabla h(\mu^{[l]})$ that corresponds to this block,

$$\mathbf{d}_{[i]}^{(l)} := \nabla_{[i]} h(\mu^{(l)}) = (\mathbf{A} \tilde{\mathbf{D}} \mathbf{A}^T \mu^{(l)} + \tilde{\mathbf{c}})_{[i]}.$$
(68)

Defining $\mathbf{d}^{(l)}$ to be the vector in \mathbb{R}^N whose *i*th block is $\mathbf{d}_{[i]}^{(l)}$ with zeros elsewhere, we perform a line search along the negative of this direction. Note that

$$h(\mu^{(l)} - \alpha \mathbf{d}^{(l)}) = \frac{1}{2} \alpha^2 \mathbf{d}^{(l)T} \mathbf{A} \tilde{\mathbf{D}} \mathbf{A}^T \mathbf{d}^{(l)} - \alpha \mathbf{d}^{(l)T} (\mathbf{A} \tilde{\mathbf{D}} \mathbf{A}^T \mu^{(l)} + \tilde{\mathbf{c}}) + \mathbf{constant}$$
(69)

$$= \frac{1}{2} \alpha^2 \mathbf{d}_{[i]}^{(l)T} \mathbf{A}_{[i]} \tilde{\mathbf{D}} \mathbf{A}_{[i]}^T \mathbf{d}_{[i]}^{(l)} - \alpha \mathbf{d}_{[i]}^{(l)T} \mathbf{d}_{[i]}^{(l)} + \mathbf{constant}.$$
 (70)

The unconstrained minimizer of this expression is

$$\alpha_l^* = \frac{\mathbf{d}_{[i]}^{(l)T} \mathbf{d}_{[i]}^{(l)}}{\mathbf{d}_{[i]}^{(l)T} \mathbf{A}_{[i]} \tilde{\mathbf{D}} \mathbf{A}_{[i]}^T \mathbf{d}_{[i]}^{(l)}}.$$
(71)

Given the nonnegativity constraints, we also need to ensure that $\mu_{[i]}^{(l)} - \alpha \mathbf{d}_{[i]}^{(l)} \ge 0$. Therefore, our choice of steplength is

$$\alpha_l = \min\left\{\alpha_l^*, \min_{j \in U_i : \mathbf{d}_j^{(l)} > 0} \frac{\mu_j^{(l)}}{\mathbf{d}_j^{(l)}}\right\}.$$
(72)

To save some computation, we introduce intermediate variables $\mathbf{f}^{(l)} := \mathbf{A}^T \mathbf{d}^{(l)} = \mathbf{A}_{[i]}^T \mathbf{d}_{[i]}^{(l)}$, and $\mathbf{m}^{(l)} := \mathbf{A}^T \mu^{(l)}$. With the intermediate variables, the updates of the BCD algorithm are:

$$\mathbf{d}_{[i]}^{(l)} \leftarrow \mathbf{A}_{[i]} \tilde{\mathbf{D}} \mathbf{m}^{(l)} + \tilde{\mathbf{c}}_{[i]}, \, \mathbf{f}^{(l)} \leftarrow \mathbf{A}_{[i]}^T \mathbf{d}_{[i]}^{(l)}$$
(73)

$$\mathbf{m}^{(l+1)} \leftarrow \mathbf{A}^T (\mu^{(l)} - \alpha_l \mathbf{d}^{(l)}) = \mathbf{m}^{(l)} - \alpha_l \mathbf{f}^{(l)}.$$
(74)

With the steplength suggested by (72), descent is guaranteed at each iteration. This property is partially why the dual BCD algorithm is effective at leveraging warmstarts which is valuable for improving the runtime of learning algorithms, as is demonstrated in Section 6.2.

Algorithm 3 Dual LCQP Block Coordinate Descent

- 1: Set l = 0 and compute an initial feasible point $\mu^{(0)}$;
- 2: Compute $\mathbf{m}^{(0)} = \mathbf{A}^T \mu^{(0)};$
- 3: while Stopping Criterion Not Satisfied do
- 4:
- 5:
- $S_k \leftarrow \text{Permutation}([1, 2, ..., p]);$ **for all** $i \in S_k$ (in order) **do** Compute $\mathbf{d}_{[i]}^{(l)} \leftarrow \mathbf{A}_{[i]} \tilde{\mathbf{D}} \mathbf{m}^{(l)} + \tilde{\mathbf{c}}_{[i]}; \quad \mathbf{f}^{(l)} \leftarrow \mathbf{A}_{[i]}^T \mathbf{d}_{[i]}^{(l)};$ 6:

7: Compute
$$\alpha_l \leftarrow \min \left\{ \frac{\mathbf{d}_{[i]}^{(l)T} \mathbf{d}_{[i]}^{(l)}}{\mathbf{f}^{(l)T} \tilde{\mathbf{D}} \mathbf{f}^{(l)}}, \min_{j \in U_i: \mathbf{d}_j^{(l)} > 0} \frac{\mu_j^{(l)}}{\mathbf{d}_j^{(l)}} \right\}$$

8:
$$\mu_{[i]}^{(l+1)} \leftarrow \mu_{[i]}^{(l)} - \alpha_l \mathbf{d}_{[i]}^{(l)}; \quad \mu_{[j]}^{(l+1)} \leftarrow \mu_{[j]}^{(l)} \text{ for all } j \neq i$$

9:
$$\mathbf{m}^{(l+1)} \leftarrow \mathbf{m}^{(l)} - \alpha_l \mathbf{f}^{(l)};$$

10: $l \leftarrow l + 1;$

11: $k \leftarrow k + 1;$

As strong duality holds for the LCQP formulation of deep HL-MRF inference, stopping when the primal-dual gap is below a given threshold $\delta > 0$, is a principled stopping criterion. Formally, at any iteration Algorithm 3 applied to (31), we recover an estimate of the primal variable v from (33)and terminate when the gap between the primal and the dual objective falls below δ . The stopping criterion is checked after every permutation block has been completely iterated over.

Connected Component Parallel D-BCD Oftentimes, the NeuPSL dual inference objective is additively separable over partitions of the variables. In this case, the dual BCD algorithm is parallelizable over the partitions. We propose identifying the separable components via the primal objective and constraints. More formally, prior to the primal problem instantiation, a disjoint-set data structure (Cormen et al., 2009) is initialized such that every primal variable belongs to a single unique disjoint set. Then, during instantiaion, the disjoint-set data structure is maintained to preserve the property that two primal variables exist in the same set if and only if they occur together with a non-zero coefficient in a constraint or a potential. This is achieved by merging the sets of variables in every generated constraint or potential. This process is made extremely efficient with a path compression strategy implemented to optimize finding set representatives. This parallelization strategy is empirically studied in Section 6 where we refer to it as CC D-BCD.

Lock Free Parallel D-BCD In general, there may only be a few connected components in the factor graph of the inference problem. In this case, D-BCD cannot fully leverage computational resources using the CC D-BCD parallelization strategy. One solution to overcome this issue and preserve the guaranteed descent property is to lock access and updates to dual variables. In other words, processes checkout locks on the dual variables to access and update its value and corresponding statistics. Unfortunately, in practice there is too much overlap in the blocks for this form of synchronization to see runtime improvements. For this reason, we additionally propose a method that sacrifices the theoretical guaranteed descent property of the dual BCD algorithm for significant runtime improvements. Our approach is inspired by lock free parallelization strategies in optimization literature (Bertsekas & N. Tsitsiklis, 1989; Recht et al., 2011; Liu et al., 2015). Specifically, rather than having processes checkout locks on dual variables for the entire iteration, we only assume dual and intermediate variable updates are atomic. This assumption ensures the dual variables and intermediate variables are synchronized across processes. However, the steplength subproblem solution and the gradient may be incorrect. Despite this, in Section 6.1 we show this distributed variant of the dual BCD algorithm consistently finds a solution satisfying the stopping criterion and realizes significant runtime improvements over the CC D-BCD algorithm in some datasets.

E EXTENDED EMPIRICAL EVALUATION

In this section, we provide additional details on the datasets and NeuPSL models used in experiments, hardware used to run experiments, an additional evaluation on the effect of the LCQP regularization on the prediction performance of NeuPSL, more inference runtime experiments, and the hyperparameter details for all the experiments in the main paper.

E.1 DATASETS AND NEUPSL MODELS

In this section, we provide additional information on all five evaluation datasets and corresponding NeuPSL models.

E.1.1 4FORUMS AND CREATEDEBATE

Stance-4Forums and Stance-CreateDebate are two datasets containing dialogues from online debate websites: 4forums.com and createdebate.com, respectively. In this paper, we study stance classification, i.e., the task of identifying the stance of a speaker in a debate as being for or against.

The 5 data splits and the NeuPSL model we evaluate in this paper originated from Sridhar et al. (2015). The data and NeuPSL models are available at: https://github.com/linqs/psl-examples/tree/main/stance-4forums and https://github.com/linqs/psl-examples/tree/main/stance-createdebate.

E.1.2 Epinions

Epinions is a trust network with 2,000 individuals connected by 8,675 directed edges representing whether they know each other and whether they trust each other Richardson et al. (2003). We study link prediction, i.e., we predict if two individuals trust each other.

In each of the 5 data splits, the entire network is available, and the prediction performance is measured on $\frac{1}{8}$ of the trust labels. The remaining set of labels are available for training. We use The NeuPSL model from Bach et al. (2017). The data and NeuPSL model are available at https://github.com/lings/psl-examples/tree/main/epinions.

E.1.3 CITESEER AND CORA

Citeseer and Cora are citation networks introduced by Sen et al. (2008). For Citeseer, 3, 312 documents are connected by 4, 732 edges representing citation links. For Cora, 2, 708 documents are connected by 5, 429 edges representing citation links. We study node classification, i.e., we classify the documents into one of 6 topics for Citeseer and 7 topics for Cora.

We study two different data settings for evaluations. For the inference and learning runtime experiments and the HL-MRF learning prediction performance experiments, Section 6.1, Section 6.2, and Section 6.3, respectively, the data is split following Bach et al. (2017). Specifically, for each of the 5 folds, 1/2 of the nodes are sampled and specify a graph for training, and the remaining 1/2 of the nodes define the graph for testing. 1/2 of the node labels are observed for both the training and test graphs. For the deep HL-MRF learning prediction performance setting, Section 6.3, for each of the 10 folds, we randomly sample 5% of the node labels for training 5% of the node labels for validation and 1,000 for testing.

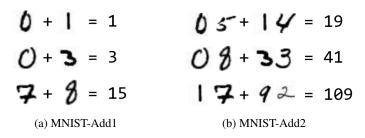


Figure 1: Example of MNIST-Add1 and MNIST-Add2.

Moreover, we use three different NeuPSL models for this dataset. The inference and learning runtime experiment models are from Bach et al. (2017) Bach et al. (2017). The data and NeuPSL models for these experiments are available at: https://github.com/ lings/psl-examples/tree/main/citeseer and https://github.com/lings/ psl-examples/tree/main/cora for Citeseer and Cora, respectively. The models for HL-MRF learning prediction performance experiments are extended versions of those in the inference and learning runtime experiments. Specifically, a copy of each rule is made that is specialized for the topic. Moreover, topic propagation across citation links is considered for papers with differing topics. For instance, the possibility of a citation from a paper with topic 'A' could imply a paper is more or less likely to be topic 'B'. The extended models are available at https://github.com/convexbilevelnesylearning/ experimentscripts/hlmrf_learning/psl-extended-examples. The models for deep HL-MRF learning prediction performance experiments are from Pryor et al. (2023). The data and models are available at: https://github.com/lings/neupsl-ijcai23.

E.1.4 DDI

Drug-drug interaction (DDI) is a network of 315 drugs and 4,293 interactions derived from the DrugBank database (S. Wishart et al., 2006). The edges in the drug network represent interactions and seven different similarity metrics. In this paper, we perform link prediction, i.e., we infer unknown drug-drug interactions.

The 5 data splits and the NeuPSL model we evaluate in this paper originated from Sridhar et al. (2016). The data and NeuPSL models are available at: https://github.com/lings/psl-examples/tree/main/drug-drug-interaction.

E.1.5 Yelp

Yelp is a network of 34, 454 users and 3, 605 items connected by 99, 049 edges representing ratings. The task is to predict missing ratings, i.e., regression, which could be used in a recommendation system.

In each of the 5 folds, 80% of the ratings are randomly sampled and available for training, and the remaining 20% is held out for testing. We use The NeuPSL model from Kouki et al. (2015). The data and NeuPSL model are available at: https://github.com/linqs/psl-examples/tree/main/yelp.

E.1.6 MNIST-ADDITION

MNIST Addition is a canonical NeSy image classification dataset first introduced by Manhaeve et al. (2018). In MNIST-Addition, models must determine the sum of two lists of MNIST images, for example, [3] + [sr] = 8. The challenge stems from the lack of labels for the MNIST images; only the final sum of the equation is provided during training, 8 in this example. 5 MNIST-Addition train splits are generated by randomly sampling, without replacement, $n \in \{600, 6, 000, 50, 000\}$ unique MNIST images from the original MNIST dataset and converted to MNIST additions. Specifically, additions are created by creating n/2 non-overlapping pairs of digits from the sample for MNIST-Add2. Then the MNIST image labels are then added together, as shown in Fig. 1, to define the addition label used in the task.

Order	Layer	Parameter	Value
1	ResNet18 (He et al., 2016)		
2	Fully Connected	Input Shape Output Shape Activation	128 64 ReLU
3	Fully Connected	Input Shape Output Shape Activation	64 10 Gumbel Softmax (Jang et al., 2017)

Table 7: Neural architecture used in NeuPSL MNIST-Add models.

This process is repeated to create five corresponding validation and test splits, with 1,0000 MNIST examples being sampled per test split from the original MNIST dataset.

MNIST-Add1 The MNIST-Add1 NeuPSL model integrates the neural component summarized in Table 7 with the symbolic model summarized in Fig. 2. The symbolic model contains the following predicates:

- NEURAL (Img, X) The NEURAL predicate is the class probability for each image as inferred by the neural network. Img is MNIST image identifier and X is a digit class that the image may represent.
- **DIGITSUMONESPLACE** ($\mathbf{W}, \mathbf{X}, \mathbf{Y}, \mathbf{Z}$) The DIGITSUMONESPLACE predicate represents whether the ones place of the sum of the digits W, X, and Y is Z. For example, substituting 0, 1, 2, and 3 for W, X, Y and Z, the predicate value would be 1, but substituting 1, 1, 2, and 3 for W, X, Y and Z would be 0 since $1 + 1 + 2 + 3 \neq 3$.
- **DIGITSUMTENSPLACE**($\mathbf{W}, \mathbf{X}, \mathbf{Y}, \mathbf{Z}$) The DIGITSUMTENSPLACE predicate represents whether the tens place of the sum of the digits W, X, and Y is Z. For example, substituting 0, 1, 2, and 0 for W, X, Y and Z, the predicate value would be 1, but substituting 0, 1, 9, and 0 for W, X, Y and Z would be 0 since 1 + 1 + 2 + 3 = 10, i.e., the tens place digit of the sum is 1 not 0.
- SUMPLACE(Img1, Img2, place, Z) The SUMPLACE predicate is the probability that the digits represented in the images identified by arguments Img1 and Img2 add up to a number with a place's place of Z.
- CARRY(Img1, Img2, W) The CARRY predicate represents is the probability that the digits represented in the images identified by arguments Img1 and Img2 add up to a number with a carry value of W. For example, images representing digits 1 and 2 do not have a carry. These variables are considered latent in the NeuPSL model as there are no truth labels for carries.
- **POSSIBLEDIGIT**(\mathbf{X} , \mathbf{Z}) The POSSIBLEDIGITS predicate represents whether a digit (X) can be included in a sum that equals a number (Z). For example, POSSIBLEDIGITS(9,0) would return 0 as no positive digit when added to 9 will equal 0. Conversely, POSSIBLEDIGITS(9,17) would return 1 as 8 added to 9 equals 17.
- IMAGESUM(Img1, Img2, Z) The IMAGESUM predicate is the probability that the digits represented by the images specified by Img1 and Img2 sum up to the number indicated by the argument Z. This predicate instantiates decision variables, i.e., variables from this predicate are not fixed during inference and learning as described in the NeSy EBM, NeuPSL, and Inference and Learning sections.
- **PLACEDREPRESENTATION**(**Z**₁₀, **Z**₁, **Z**) The PLACEDREPRESENTATION predicate represents whether the number Z has tens place digit Z₁₀ and ones place digit Z₁.

MNIST-Add2 The MNIST-Add2 NeuPSL model integrates the neural component summarized in Table 7 with the symbolic model summarized in Fig. 3. The symbolic model contains the following predicates:

• SUMPLACE(Img1, Img2, Img3, Img4, place, Z) The SUMPLACE predicate is the probability that the digits represented in the images identified by arguments Img1, Img2, Img3, and Img4 add up to a number with a place's place of Z.

$$\begin{split} w_1 : \text{DigitSumOnesPlace}('0', \textbf{X}, \textbf{Y}, \textbf{Z}) \land \text{Neural}(\texttt{Img1}, \textbf{X}) \land \text{Neural}(\texttt{Img2}, \textbf{Y}) \rightarrow \text{SumPlace}(\texttt{Img1}, \texttt{Img2}, '1', \textbf{Z}) \\ w_2 : \text{DigitSumTensPlace}('0', \textbf{X}, \textbf{Y}, \textbf{Z}) \land \text{Neural}(\texttt{Img1}, \textbf{X}) \land \text{Neural}(\texttt{Img2}, \textbf{Y}) \rightarrow \text{SumPlace}(\texttt{Img1}, \texttt{Img2}, '10', \textbf{Z}) \end{split}$$

```
w_3: DIGITSUMONESPLACE(W, 0', 0', 2) \land CARRY(Img1, Img2, W) \rightarrow SUMPLACE(Img1, Img2, '10', Z)
```

```
\begin{split} w_4: \texttt{SUMPLACE}(\texttt{ImageId1},\texttt{ImageId2},\texttt{'1'},\texttt{Z}_1) \land \texttt{SUMPLACE}(\texttt{ImageId1},\texttt{ImageId2},\texttt{'10'},\texttt{Z}_{10}) \\ & \land \texttt{PLACEDREPRESENTATION}(\texttt{Z}_{10},\texttt{Z}_1,\texttt{Z}) \rightarrow \texttt{IMAGESUM}(\texttt{ImageId1},\texttt{ImageId2},\texttt{Z}) \end{split}
```

Figure 2: Summarized NeuPSL MNIST-Add1 Symbolic Model. The full model is available at: https://github.com/convexbilevelnesylearning/experimentscripts/mnist_addition/neupsl_models.

```
 \begin{split} & w_1: \text{DIGITSUMONESPLACE}('0', X, Y, Z) \land \text{NEURAL}(\text{Img2}, X) \land \text{NEURAL}(\text{Img4}, Y) \rightarrow \text{SUMPLACE}(\text{Img1}, \text{Img2}, \text{Img3}, \text{Img4}'1', Z) \\ & w_2: \text{DIGITSUMTENSPLACE}('0', X, Y, Z) \land \text{NEURAL}(\text{Img2}, X) \land \text{NEURAL}(\text{Img4}, Y) \rightarrow \text{CARRY}(\text{Img2}, \text{Img4}, Z) \\ & w_3: \text{DIGITSUMONESPLACE}(W, X, Y, Z) \land \text{NEURAL}(\text{Img1}, X) \land \text{NEURAL}(\text{Img3}, Y) \land \text{CARRY}(\text{Img2}, \text{Img4}, W) \\ & \rightarrow \text{SUMPLACE}(\text{Img1}, \text{Img2}, \text{Img3}, \text{Img4}'10', Z) \\ & w_4: \text{DIGITSUMTENSPLACE}(W, X, Y, Z) \land \text{NEURAL}(\text{Img1}, X) \land \text{NEURAL}(\text{Img3}, Y) \land \text{CARRY}(\text{Img2}, \text{Img4}, W) \\ & \rightarrow \text{SUMPLACE}(\text{Img1}, \text{Img2}, \text{Img3}, \text{Img4}'100', Z) \\ & w_5: \text{DIGITSUMONESPLACE}(W, '0', '0', Z) \land \text{CARRY}(\text{Img1}, \text{Img3}, W) \rightarrow \text{SUMPLACE}(\text{Img1}, \text{Img2}, \text{Img4}'100', Z) \\ & w_6: \text{SUMPLACE}(\text{ImageId1}, \text{ImageId2}, \text{ImageId3}, \text{ImageId4}, '1', Z_{10}) \\ & \land \text{SUMPLACE}(\text{ImageId1}, \text{ImageId2}, \text{ImageId3}, \text{ImageId4}, '100', Z_{100}) \land \text{PLACEDREPRESENTATION}(Z_{100}, Z_{10}, Z_{10}
```

 $\land SUMPLACE(ImageId1, ImageId2, ImageId3, ImageId4, '100', Z_{100}) \land PLACEDREPRESENTATION(Z_{100}, Z_{10}, Z_{1}, Z) \rightarrow IMAGESUM(ImageId1, ImageId2, ImageId3, ImageId4, Z)$

Figure 3: Summarized NeuPSL MNIST-Add2 Symbolic Model. The full model is available at: https://github.com/convexbilevelnesylearning/experimentscripts/mnist_addition/neupsl_models.

- **POSSIBLEONESDIGIT**(**X**, **Z**) The POSSIBLEONESDIGIT predicate represents whether a digit (X) can be included in a sum as a ones place digit that equals a number (Z). For example, POSSIBLEDIGITS(9,0) would return 0 as no positive digit when added to 9 will equal 0. Conversely, POSSIBLEDIGITS(9,17) would return 1 as 8 added to 9 equals 17.
- **POSSIBLETENSDIGIT**(**X**, **Z**) The POSSIBLETENSDIGITS predicate represents whether a digit (X) can be included in a sum as a tens place digit that equals a number (Z). For example, POSSIBLEDIGITS(9,0) would return 0 as no positive digit when added to 90, 91, \cdots , 99 will equal 0. Conversely, POSSIBLEDIGITS(9,97) would return 1 as 7 added to 90 equals 97, for instance.
- IMAGESUM(Img1, Img2, Img3, Img4, Z) The IMAGESUM predicate is the probability that the digits represented by the images specified by Img1, Img2, Img3, and Img4 sum up to the number indicated by the argument Z. This predicate instantiates decision variables, i.e., variables from this predicate are not fixed during inference and learning as described in the NeSy EBM, NeuPSL, and Inference and Learning sections.
- PLACEDREPRESENTATION(**Z**₁₀₀, **Z**₁₀, **Z**₁, **Z**) The PLACEDREPRESENTATION predicate represents whether the number Z has hundereds place digit Z₁₀₀, tens place digit Z₁₀ and ones place digit Z₁.

E.2 HARDWARE

All timing experiments were performed on an Ubuntu 22.04.1 Linux machine with Intel Xeon Processor E5-2630 v4 at 3.10GHz and 128 GB of RAM.

E.3 DUAL BCD AND REGULARIZATION

The regularization parameter added to the LCQP formulation of NeuPSL inference in (16) ensures strong convexity of the optimal value of the energy function. However, adding regularization makes the new formulation an approximation. In this section, the runtime and prediction performance of the D-BCD inference algorithm is evaluated at varying levels of regularization to understand its effect on NeuPSL inference. The regularization parameter varies in the range $\epsilon \in \{100, 10, 1, 0.1, 0.01\}$. The D-BCD algorithm is stopped when the primal-dual gap drops below $\delta = 0.1$ Inference time is provided in seconds, and the performance metric is consistent with Table 1. Results are provided in Table 8.

Table 8: D-BCD Inference time in seconds and prediction performance with varying values for the LCQP regularization parameter ϵ .

Dataset	ϵ	Time (sec)	Perf.
	100	0.02 ± 0.01	64.77 ± 10.61
	10	0.02 ± 0.01	64.83 ± 10.53
CreateDebate (AUROC)	1	0.02 ± 0.01	64.74 ± 10.67
	0.1	0.05 ± 0.02	65.39 ± 9.07
	0.01	0.42 ± 0.51	66.01 ± 9.35
	100	0.11 ± 0.02	61.31 ± 6.17
	10	0.10 ± 0.03	61.26 ± 6.16
4Forums (AUROC)	1	0.09 ± 0.01	61.12 ± 6.18
	0.1	0.43 ± 0.11	62.73 ± 5.46
	0.01	7.11 ± 3.05	62.31 ± 5.47
	100	0.33 ± 0.05	72.59 ± 2.27
	10	0.28 ± 0.04	72.69 ± 2.21
Epinions (AUROC)	1	0.33 ± 0.05	74.24 ± 1.95
	0.1	1.08 ± 0.16	77.05 ± 1.06
	0.01	5.21 ± 0.37	77.45 ± 0.70
	100	0.95 ± 0.14	71.28 ± 1.31
	10	1.00 ± 0.12	71.28 ± 1.30
Citeseer (Accuracy)	1	1.48 ± 0.29	71.59 ± 1.01
	0.1	7.01 ± 1.57	71.75 ± 1.10
	0.01	62.41 ± 14.67	71.92 ± 1.09
	100		
	10		
Cora (Accuracy)	1	7.36 ± 4.19	81.48 ± 1.70
	0.1	42.24 ± 25.06	81.88 ± 1.82
	0.01	269.45 ± 49.50	81.79 ± 1.72
	100	24.56 ± 0.25	94.85 ± 0.00
	10	29.23 ± 0.59	94.85 ± 0.00
DDI (AUROC)	1	47.15 ± 0.95	94.82 ± 0.00
	0.1	280.62 ± 5.19	94.80 ± 0.00
	0.01	266.07 ± 42.68	94.81 ± 0.00
	100	105.60 ± 5.03	0.23 ± 0.01
	10	$3,239 \pm 81$	0.22 ± 0.01
Yelp (MAE)	1	$3,227 \pm 54$	0.19 ± 0.01
	0.1	421 ± 202	0.18 ± 0.00
	0.01	$2,472 \pm 297$	0.18 ± 0.00

Table 8 shows there is a consistent correlation between the LCQP regularization parameter and the runtime and performance of inference. As ϵ increases, there is a significant decrease in the runtime performance as the D-BCD algorithm can find a solution with a gradient meeting the stopping criterion in fewer iterations. Notably, for the Citeseer inference problem, the D-BCD algorithm realizes a roughly $45 \times$ speedup. On the other hand, while the runtime performance improves with increasing ϵ , the prediction performance can sometimes decay. There is a tradeoff between runtime and prediction performance when setting the ϵ regularization parameter.

E.4 EXTENDED INFERENCE RUNTIME

Table 9: Inference time in seconds for each inference optimization technique.

	Gurobi	GD	ADMM	CC D-BCD	LF D-BCD
Citeseer		$\begin{array}{c} 34.63 \pm 0.33 \\ 47.17 \pm 0.61 \\ 48.66 \pm 1.24 \\ 6.061 + 46 \end{array}$	0.63 ± 0.07	1.36 ± 0.24	$ \begin{vmatrix} 0.26 \pm 0.04 \\ 0.49 \pm 0.08 \\ 0.79 \pm 0.19 \\ 7.58 \pm 0.48 \end{vmatrix} $

Dataset	Parameter	Range	Final Value
Course Dalasta	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	1.0
CreateDebate	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01\}$	0.1
4Forums	ADMM Step Length	{10.0, 1.0, 0.1, 0.01}	1.0
41 01 01115	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01\}$	0.1
	GD Step Length	$\{10.0, 1.0, 0.1, 0.01, 0.001\}$	0.01
Epinions	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	0.1
	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01\}$	0.1
	GD Step Length	$\{10.0, 1.0, 0.1, 0.01, 0.001\}$	0.1
Citeseer	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	10.0
	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01\}$	10.0
	GD Step Length	$\{10.0, 1.0, 0.1, 0.01, 0.001\}$	0.1
Cora	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	10.0
	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01\}$	10.0
DDI	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	1.0
221	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01\}$	10.0
	GD Step Length	$\{10.0, 1.0, 0.1, 0.01, 0.001\}$	0.001
Yelp	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	1.0
	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01\}$	0.1
MNIST-Add1	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	1.0
	LCQP Regularization	{100, 10, 1, 0.1, 0.01, 0.001}	0.001
MNIST-Add2	ADMM Step Length	$\{10.0, 1.0, 0.1, 0.01\}$	1.0
	LCQP Regularization	$\{100, 10, 1, 0.1, 0.01, 0.001\}$	0.001

Table 10: Hyperparameter ranges and final values for the inference runtime experiments.

This section details the hyperparameter settings and search process for the inference runtime experiments in Section 6.1. The GD, ADMM, and D-BCD algorithms are stopped when the L_{∞} norm of the primal variable change between iterates is less than 0.001. For the D-BCD algorithms, the regularization parameter from Appendix E.3 resulting in the fastest runtime and yielding a prediction performance within a standard error of the best is used. The default Gurobi optimizer hyperparameters are used. Table 10 reports the range of hyperparameters searched over and the final values. Furthermore, for the MNIST-Add1 and MNIST-Add2 models, the highest performing trained neural models for each split from the performance experiments in Section 6.3 are used.

Table 9 reports the average and standard deviation of the inference runtime for Gurobi, GD, ADMM, and D-BCD algorithms on 4 of the datasets from Table 1. As in the main paper, we see the D-BCD algorithms are competitive with ADMM, the current state of the art optimizer for NeuPSL inference. Moreover, here we see the LF D-BCD algorithm is also competitive with Gurobi for a single round of inference.

E.5 EXTENDED LEARNING RUNTIME

This section provides details of the hyperparameter settings for the learning runtime experiments in Section 6.2. For both learning losses, a negative log regularization with coefficient 1.0e - 3 on the symbolic weights is added to the learning loss as suggested by Pryor et al. (2023). For ADMM inference on both learning losses, the same steplength from the inference runtime experiment is used for the first 7 datasets in Table 1. Similarly, for D-BCD inference on both learning losses, the same regularization parameter from the inference runtime experiment is used for the first 7 datasets in Table 1. For the MNIST-Add experiments, we use the regularization parameter $\epsilon = 1.0e - 3$ and ADMM steplength 1.0 as the values were found to achieve the highest final validation prediction performance.

Mirror descent is applied to learn the symbolic weights for both SP and MSE losses. The mirror descent steplength is set to a default value of 1.0e - 3 for the first 7 datasets in Table 1. For the MNIST-Add datasets the mirror descent steplength is set to 1.0e - 14 as in this problem we only need to learn the neural weights. The Adam steplength for the neural component of the MNIST-Add models is set to a default value of 1.0e - 3.

Our learning framework, Algorithm 1, is used to fit the MSE learning loss. We set the initial squared penalty parameter to a default value of 2.0 for all datasets. Moreover, for the first 7 datasets in Table 1 we set the Moreau parameter to 0.01, the energy loss coefficient to 0.1, and the steplength on the target variables y to 0.01. For the MNIST-Add datasets we set the Moreau parameter to 1.0e - 3, the energy loss coefficient to 1.0.0, and the steplength on the target variables y to 1.0e - 3.

E.6 EXTENDED LEARNING PREDICTION PERFORMANCE

This section details the hyperparameter settings and search process for the prediction performance experiments in Section 6.3. For all learning losses, a negative log regularization with coefficient 1.0e - 3 on the symbolic weights is added to the learning loss as suggested by Pryor et al. (2023). The remaining hyperparameter search and setting details are described separately for the HL-MRF learning and deep HL-MRF learning experiments.

HL-MRF Learning The LF D-BCD algorithm is used for inference in all experiments. Moreover, the D-BCD algorithm is stopped when the primal-dual gap drops below $\delta = 1.0e - 2$ CreateDebate, 4Forums, Epinions, Citeseer, Cora, and DDI while the primal-dual threshold is set to $\delta = 1.0e - 1$ to adjust to the larger scale of the dataset. For all learning losses, the learning algorithm is stopped when the training evaluation metric stops improving after 50 epochs. For the MSE and BCE losses trained with Algorithm 1, the final objective difference tolerance was set to 0.1 for the smaller CreateDebate, 4Forums, and Epinions datasets and 1 for Citeseer, Cora, DDI, and Yelp. Moreover, the initial squared penalty coefficient is set to 2 for all datasets. The remaining hyperparameters are searched over the ranges specified in Table 11. The hyperparameter value with the best performance metric on the first fold is selected.

Deep HL-MRF Learning

For deep HL-MRF learning in the citation network and MNIST-Add evaluations reported in Table 5 and Table 6, respectively the validation set is used to determine when to stop the learning algorithms and what weights to use for final evaluations. Specifically, after every learning step the model performance is measured on the validation data, and when 50 consecutive steps finish without improvement, the learning algorithm is stopped. For citation network datasets, the model obtaining the best validation metric averaged across all splits are used for final test evaluation. For MNIST-Add datasets, the model obtaining the best validation metric on the first split is used for final test evaluation across all splits. Table 12 and Table 13 report the range of hyperparameters searched over and the final values resulting in the highest validation prediction performance for citation network datasets and MNIST-Add datasets, respectively.

Dataset	Learning Loss	Parameter	Range	Final Value
	Energy	Mirror Descent Step Length LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 2
	SP	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		LCQP Regularization Mirror Descent Step Length		1.0e - 3 1.0e - 2
		y Step Length	$\{1.0e - 3, 1.0e - 2\}\$ $\{1.0e - 2, 1.0e - 1\}$	1.0e - 2 1.0e - 1
CreateDebate	MSE	Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1 1.0e - 2
		Energy Loss Coefficient LCQP Regularization	$\{0, 1.0e - 1, 1, 10\}\$ $\{1.0e - 3, 1.0e - 2\}$	$\begin{array}{c} 0.1 \\ 1.0e - 3 \end{array}$
		Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
	BCE	y Step Length Moreau Parameter	$\{1.0e - 2, 1.0e - 1\}$ $\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1 1.0e - 2
		Energy Loss Coefficient	$\{0, 1.0e - 1, 1, 10\}$	10
		LCQP Regularization Mirror Descent Step Length	${1.0e - 3, 1.0e - 2}$ ${1.0e - 3, 1.0e - 2}$	1.0e - 2 1.0e - 3
	Energy	LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
	SP	Mirror Descent Step Length LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 3 1.0e - 3
4Forums		Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
	MSE	y Step Length Moreau Parameter	$\{1.0e - 2, 1.0e - 1\}$ $\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 2 1.0e - 3
	MOL	Energy Loss Coefficient	$\{0, 1.0e - 1, 1, 10\}$	0
		LCQP Regularization Mirror Descent Step Length	${1.0e - 3, 1.0e - 2}$ ${1.0e - 3, 1.0e - 2}$	1.0e - 3 1.0e - 3
		y Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 2, 1.0e - 1\}$	1.0e - 3 1.0e - 2
	BCE	Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 3
		Energy Loss Coefficient LCQP Regularization	$\{0, 1.0e - 1, 1, 10\}\$ $\{1.0e - 3, 1.0e - 2\}$	$\begin{array}{c} 0 \\ 1.0e - 3 \end{array}$
	Energy	Mirror Descent Step Length	$ \begin{array}{c} \{1.0e-3, 1.0e-2\} \\ \{1.0e-3, 1.0e-2\} \\ \{1.0e-3, 1.0e-2\} \\ \{1.0e-3, 1.0e-2\} \end{array} $	1.0e - 3 1.0e - 3
		LCQP Regularization Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 3
	SP	LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
		Mirror Descent Step Length y Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 2, 1.0e - 1\}$	1.0e - 2 1.0e - 1
Epinions	MSE	Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1 1.0e - 2
		Energy Loss Coefficient LCOP Regularization	$\{0, 1.0e - 1, 1, 10\}\$ $\{1.0e - 3, 1.0e - 2\}$	0.1 1.0e - 2
	BCE	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		y Step Length Moreau Parameter	$\{1.0e - 2, 1.0e - 1\}$ $\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1 1.0e - 2
		Energy Loss Coefficient	$\{0, 1.0e - 1, 1, 10\}$	1
		LCQP Regularization Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 2 1.0e - 3
	Energy	LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
	SP	Mirror Descent Step Length LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 3
Citeseer		Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 2, 1.0e - 2\}$	1.0e - 3 1.0e - 3
	MSE	y Step Length Moreau Parameter	$\{1.0e - 2, 1.0e - 1\}$ $\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 2 1.0e - 2
		Energy Loss Coefficient	$\{0, 1.0e - 1, 1, 10\}$	1
		LCQP Regularization Mirror Descent Step Length	${1.0e - 3, 1.0e - 2}$ ${1.0e - 3, 1.0e - 2}$	1.0e - 2 1.0e - 2
	BCE	y Step Length Moreau Parameter	$\{1.0e - 2, 1.0e - 1\}$ $\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1 1.0e - 3
		Energy Loss Coefficient	$\{0, 1.0e - 1, 1, 10\}$	0
	-	LCQP Regularization Mirror Descent Step Length	${1.0e - 3, 1.0e - 2}$ ${1.0e - 3, 1.0e - 2}$	1.0e - 3 1.0e - 3
	Energy	LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2 1.0e - 3 1.0e - 3
	SP	Mirror Descent Step Length LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 3
		Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
6	MSE	y Step Length Moreau Parameter	$\{1.0e - 2, 1.0e - 1\}$ $\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1 1.0e - 2
Cora		Energy Loss Coefficient	$\{0, 1.0e - 1, 1, 10\}$	0.1
		LCQP Regularization Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 2
	ner	y Step Length	$\{1.0e - 2, 1.0e - 1\}$	1.0e - 1
	BCE	Moreau Parameter Energy Loss Coefficient	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$ $\{0, 1.0e - 1, 1, 10\}$	1.0e - 2 0.1
		LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
	Energy	Mirror Descent Step Length LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 2
	SP	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
		LCQP Regularization Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 2 1.0e - 3
	MSE	y Step Length	$\{1.0e - 2, 1.0e - 1\}$	1.0e - 1
DDI		Moreau Parameter Energy Loss Coefficient	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$ $\{0, 1.0e - 1, 1, 10\}$	1.0e - 3 0.1
		LCQP Regularization Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 2 1.0e - 2
	BCE	y Step Length	$\{1.0e - 2, 1.0e - 1\}$	1.0e - 2
		Moreau Parameter Energy Loss Coefficient	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$ $\{0, 1.0e - 1, 1, 10\}$	1.0e - 2 0.1
		LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
	Energy	Mirror Descent Step Length LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3 1.0e - 2
		Mirror Descent Step Length	${1.0e - 3, 1.0e - 2}$ ${1.0e - 3, 1.0e - 2}$ ${1.0e - 3, 1.0e - 2}$	1.0e - 2 1.0e - 3
	SP			1.0e - 2
	SP	LCQP Regularization	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
		LCQP Regularization Mirror Descent Step Length y Step Length	$\{1.0e - 3, 1.0e - 2\}$ $\{1.0e - 2, 1.0e - 1\}$	1.0e - 3 1.0e - 2
Yelp	SP MSE	LCQP Regularization Mirror Descent Step Length y Step Length Moreau Parameter	$ \{1.0e - 3, 1.0e - 2\} \\ \{1.0e - 2, 1.0e - 1\} \\ \{1.0e - 3, 1.0e - 2, 1.0e - 1\} \\ \{0, 1.0e - 1, 1, 10\} $	1.0e - 1
Yelp		LCQP Regularization Mirror Descent Step Length y Step Length Moreau Parameter Energy Loss Coefficient LCQP Regularization	$ \{1.0e - 3, 1.0e - 2\} \\ \{1.0e - 2, 1.0e - 1\} \\ \{1.0e - 3, 1.0e - 2, 1.0e - 1\} \\ \{0, 1.0e - 1, 1, 10\} $	1.0e - 1 10 1.0e - 2
Yelp		LCQP Regularization Mirror Descent Step Length y Step Length Moreau Parameter Energy Loss Coefficient LCQP Regularization Mirror Descent Step Length	$ \{1.0e - 3, 1.0e - 2\} \\ \{1.0e - 2, 1.0e - 1\} \\ \{1.0e - 3, 1.0e - 2, 1.0e - 1\} \\ \{0, 1.0e - 1, 1, 10\} $	1.0e - 1 10 1.0e - 2
Yelp		LCQP Regularization Mirror Descent Step Length y Step Length Moreau Parameter Energy Loss Coefficient LCQP Regularization	$ \{ 1.0e - 3, 1.0e - 2 \} \\ \{ 1.0e - 2, 1.0e - 1 \} \\ \{ 1.0e - 3, 1.0e - 2, 1.0e - 1 \} $	1.0e - 1 10

Table 11: Hyperparameter ranges and final values for the HL-MRF learning prediction performance experiments in Table 4.

Dataset	Loss	Parameter	Range	Final Value
Citeseer	Energy	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
	SP	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
		Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		y Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
	MSE	Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1
		Energy Loss Coefficient	$\{1.0e - 1, 1, 10\}$	1.0e - 1
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
		Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
	BCE	y Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
		Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1
		Energy Loss Coefficient	$\{1.0e - 1, 1, 10\}$	1.0e - 1
		LCQP Regularization	$\{1.0e-3\}$	1.0e - 3
·	Energy	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
	SP	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
	MSE	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		y Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
Cora		Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1
Cora		Energy Loss Coefficient	$\{1.0e - 1, 1, 10\}$	1
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
	BCE	Mirror Descent Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		y Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 2
		Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 1
		Energy Loss Coefficient	$\{1.0e - 1, 1, 10\}$	1.0e - 1
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3

Table 12: Hyperparameter ranges and final values for the deep HL-MRF learning prediction performance experiments on Citeseer and Cora.

Table 13: Hyperparameter ranges and final values for the deep HL-MRF learning prediction performance experiments on MNIST-Add datasets.

Dataset	Loss	Parameter	Range	Final Value
	Enongy	Mirror Descent Step Length	$\{1.0e - 14\}$	1.0e - 14
	Energy	Adam Step Length	$\{1.0e - 4, 1.0e - 3\}$	1.0e - 3
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
MNIST-Add1		Mirror Descent Step Length	$\{1.0e - 14\}$	1.0e - 14
		Adam Step Length	$\{1.0e - 4, 1.0e - 3\}$	1.0e - 3
	BCE	y Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
		Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 2
		Energy Loss Coefficient	$\{1.0e - 1, 1, 10\}$	10
		LCQP Regularization	$\{1.0e-3\}$	1.0e - 3
MNIST-Add2	Energy	Mirror Descent Step Length	$\{1.0e - 14\}$	1.0e - 14
		Adam Step Length	$\{1.0e - 4, 1.0e - 3\}$	1.0e - 3
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3
	BCE	Mirror Descent Step Length	$\{1.0e - 14\}$	1.0e - 14
		Adam Step Length	$\{1.0e - 4, 1.0e - 3\}$	1.0e - 4
		y Step Length	$\{1.0e - 3, 1.0e - 2\}$	1.0e - 3
		Moreau Parameter	$\{1.0e - 3, 1.0e - 2, 1.0e - 1\}$	1.0e - 3
		Energy Loss Coefficient	$\{1.0e - 1, 1, 10\}$	10
		LCQP Regularization	$\{1.0e - 3\}$	1.0e - 3