Learning with Local Search MCMC Layers

Anonymous Author(s)

Affiliation Address email

Abstract

Integrating combinatorial optimization layers into neural networks has recently attracted significant research interest. However, many existing approaches lack theoretical guarantees or fail to perform adequately when relying on inexact solvers. This is a critical limitation, as many operations research problems are NP-hard, often necessitating the use of neighborhood-based local search heuristics. These heuristics iteratively generate and evaluate candidate solutions based on an acceptance rule. In this paper, we introduce a theoretically-principled approach for learning with such inexact combinatorial solvers. Inspired by the connection between simulated annealing and Metropolis-Hastings, we propose to transform problem-specific neighborhood systems used in local search heuristics into proposal distributions, implementing MCMC on the combinatorial space of feasible solutions. This allows us to construct differentiable combinatorial layers and associated loss functions. Replacing an exact solver by a local search strongly reduces the computational burden of learning on many applications. We demonstrate our approach on a large-scale dynamic vehicle routing problem with time windows.

1 Introduction

2

5

6

7

8

9 10

11

12

13

14

15

16

Models that combine neural networks and combinatorial optimization have recently attracted significant attention [14, 39, 8, 6, 50, 5, 34, 43, 7]. Such models enable the transformation of learned continuous latent representations into structured discrete outputs that satisfy complex constraints. They enrich combinatorial optimization algorithms by providing them with context-dependent features, making decisions more resilient to uncertainty. An important subset of this line of research involves integrating, within a neural network, a linear programming layer of the form:

$$\theta \mapsto \underset{\boldsymbol{y} \in \mathcal{Y}}{\operatorname{argmax}} \langle \theta, \boldsymbol{y} \rangle \subseteq \underset{\boldsymbol{y} \in \operatorname{conv}(\mathcal{Y})}{\operatorname{argmax}} \langle \theta, \boldsymbol{y} \rangle,$$
 (1)

this is often called the maximum a posteriori (MAP) problem [51]. The main challenge in using 24 such layers lies in their end-to-end model training. Indeed, as piecewise-constant, discontinuous 25 functions, such layers break the differentiable programming computational graph, and prevent one 26 from backpropagating meaningful gradients from the final output of the model to its parameters. 27 28 Many approaches have been proposed to derive relaxations and loss functions for this setting; see Appendix A for a detailed review and overview of relevant related work. To methodologically 29 position our work, Table 1 provides a high-level overview of foundational approaches, contrasting 30 them based on the type of oracle they assume access to. Some rely on an oracle for solving a 31 regularized version of Eq. (1), such as a quadratic or entropy-regularized program. They typically 32 perform a single oracle call per data point. Some other approaches assume access to an oracle 33

where \mathcal{Y} is a finite set of feasible outputs. In the graphical model and structured prediction literature,

35 36

34

for solving the original linear program (i.e., a MAP oracle), but perform multiple oracle calls, for

smoothing reasons. Their theoretical guarantees usually assume an oracle returning exact solutions.

Table 1: The proposed approach leverages the neighborhood systems used by local search heuristics (inexact solvers) to obtain a differentiable combinatorial layer when usual oracles are not available.

	Regularization	Oracle	Approach
Differentiable DP (2009, 2018)	Entropy	Exact marginal	DP
SparseMAP (2018)	Quadratic	Exact MAP	Frank-Wolfe
Barrier FW (2015)	TRW Entropy	Exact MAP	Frank-Wolfe
IntOpt (2020)	Log barrier	Interior point solver	Primal-Dual
Perturbed optimizers (2020)	Implicit via noise	Exact MAP	Monte-Carlo
Blackbox solvers (2020)	None	Exact MAP	Interpolation
Contrastive divergences (2000)	Entropy	Gibbs / Langevin sampler	MCMC
Proposed	Entropy	Local search	MCMC

Unfortunately, many problems in operations research are NP-hard in nature (e.g., routing, scheduling, network design), making access to an exact oracle difficult. In contrast, operations research applications often rely on local search heuristics, such as simulated annealing. These heuristics iteratively generate a neighbor of the current solution, and either accept it or reject it based on an acceptance rule. The aim of this work is to provide a theoretically-principled approach for learning with such inexact solvers. To do so, we propose to leverage unexploited links between neighborhood-based, local search heuristics used to approximately solve combinatorial problems, and Markov chain Monte-Carlo (MCMC) methods used to perform approximate marginal inference in graphical models.

Contributions. (i) We enable the integration of local search heuristics as layers into ML models, by converting their neighborhood systems into proposal distributions for a discrete MCMC sampler over the combinatorial set of solutions. (ii) We extend our framework to handle local search heuristics that leverage a diversity of neighborhood systems, enabling this powerful class of solvers to be used as a unified MCMC sampler. (iii) We show that there exist Fenchel-Young losses [8] whose stochastic gradients are given by the proposed layer (even with a single MCMC iteration), leading to principled learning algorithms in both supervised and unsupervised settings, for which we provide a convergence analysis. (iv) We demonstrate our approach on the EURO Meets NeurIPS 2022 challenge [27], a large-scale, ML-enriched dynamic vehicle routing problem with time windows (DVRPTW), which involves an intractable combinatorial optimization problem. In Appendix B, we also empirically validate the quality of the proposed gradient estimators through abundant experiments.

Problem setup. In this paper, our goal is to learn models with an optimization layer of the form

$$\widehat{\boldsymbol{y}}: \boldsymbol{\theta} \mapsto \operatorname*{argmax}_{\boldsymbol{y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}), \tag{2}$$

 $\widehat{\pmb{y}}: \pmb{\theta} \mapsto \operatorname*{argmax}_{\pmb{y} \in \mathcal{Y}} \langle \pmb{\theta}, \, \pmb{y} \rangle + \varphi(\pmb{y}), \tag{2}$ where $\pmb{\theta} \in \mathbb{R}^d$ and $\mathcal{Y} \subset \mathbb{R}^d$ is a finite but combinatorially-large set. This formulation is a generalization of the standard linear objective in Eq. (1). The function $\varphi: \mathcal{Y} \to \mathbb{R}$ is an integral part of the problem definition, capturing any structural costs or preferences (e.g., routing distances, fixed costs) that are independent of θ . We focus on settings where this optimization problem is intractable and only heuristic algorithms are available to obtain an approximate solution. We distinguish between two settings. In the unsupervised setting, our goal will be to learn $\theta \in \mathbb{R}^d$ from observations $y_1, \dots, y_N \in \mathbb{R}^d$. In the supervised setting, we will assume that $\theta = q_W(x)$ and our goal will be to learn the parameters W from observation pairs $(x_1, y_1), \ldots, (x_N, y_N)$.

Local search based MCMC layers

In this section, we show how to design principled combinatorial layers without relying on exact MAP solvers, by transforming local search heuristics into MCMC algorithms.

2.1 From local search to MCMC 69

39

41

42

43

45

49

50

51

52

53

60

61

62

63

70

71

72

73

Local search and neighborhood systems. Local search heuristics [19] iteratively generate a neighbor $y' \in \mathcal{N}(y^{(k)})$ of the current solution $y^{(k)}$, and either accept it or reject it based on an acceptance rule. In this context, a neighborhood system \mathcal{N} defines, for each feasible solution $\boldsymbol{u} \in \mathcal{Y}$, a set of neighbors $\mathcal{N}(y) \subseteq \mathcal{Y}$. Neighborhoods are problem-specific, and must respect the structure of the problem, i.e., must maintain solution feasibility. They are typically defined implicitly via a set of allowed *moves* from y. For instance, Table 4 lists example moves for a vehicle routing problem.

Algorithm 1 SA / MH as a layer

80

81

82

```
Inputs: \boldsymbol{\theta} \in \mathbb{R}^d, \boldsymbol{y}^{(0)} \in \mathcal{Y}, (t_k), K \in \mathbb{N}, \mathcal{N}, q
 for k = 0 : K do
          Sample a neighbor in \mathcal{N}(\boldsymbol{y}^{(k)}):
         oldsymbol{y}' \sim q\left(oldsymbol{y}^{(k)},\,\cdot\,
ight)
         \alpha(\boldsymbol{y}^{(k)}, \boldsymbol{y}') \leftarrow 1 \text{ (SA) or } 
\alpha(\boldsymbol{y}^{(k)}, \boldsymbol{y}') \leftarrow \frac{q(\boldsymbol{y}', \boldsymbol{y}^{(k)})}{q(\boldsymbol{y}^{(k)}, \boldsymbol{y}')} \text{ (MH)} 
U \sim \mathcal{U}([0, 1])
Output: \widehat{y}(\theta) \approx y^{(K)} (SA) or \widehat{y}_t(\theta) = \mathbb{E}_{\pi_{\theta,t}}[Y] \approx \frac{1}{K} \sum_{k=1}^{K} y^{(k)} (MH)
```

Algorithm 2 Neighborhood mixture MCMC

```
Inputs: \theta \in \mathbb{R}^d, y^{(0)} \in \mathcal{Y}, t, K \in \mathbb{N}, (\mathcal{N}_s, q_s)_{s=1}^S
                                                                                                                                                                                                  for k = 0 : K \operatorname{do}
                                                                                                                                                                                                            Sample a neighborhood system:
                                                                                                                                                                                                            s \sim \mathcal{U}(Q(\mathbf{y}^{(k)}))
\alpha(\boldsymbol{y}^{(k)},\boldsymbol{y}') \leftarrow \frac{q(\boldsymbol{y}',\boldsymbol{y}^{(k)})}{q(\boldsymbol{y}^{(k)},\boldsymbol{y}')} \text{ (MH)}
U \sim \mathcal{U}([0,1])
\Delta^{(k)} \leftarrow \langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}') - \langle \boldsymbol{\theta}, \boldsymbol{y}^{(k)} \rangle - \varphi(\boldsymbol{y}^{(k)})
p^{(k)} \leftarrow \alpha(\boldsymbol{y}^{(k)}, \boldsymbol{y}') \exp\left(\Delta^{(k)}/t_k\right)
If U \leq p^{(k)}, accept move: \boldsymbol{y}^{(k+1)} \leftarrow \boldsymbol{y}'
If U > p^{(k)}, reject move: \boldsymbol{y}^{(k+1)} \leftarrow \boldsymbol{y}'
If U < p^{(k)}, reject move: \boldsymbol{y}^{(k+1)} \leftarrow \boldsymbol{y}'
If U < p^{(k)}, reject move: \boldsymbol{y}^{(k+1)} \leftarrow \boldsymbol{y}'
                                                                                                                                                                                                           Sample a neighbor in \mathcal{N}_s(\boldsymbol{y}^{(k)}):
                                                                                                                                                                                                          If U \le p^{(k)}, accept move: \boldsymbol{y}^{(k+1)} \leftarrow \boldsymbol{y}'
If U > p^{(k)}, reject move: \boldsymbol{y}^{(k+1)} \leftarrow \boldsymbol{y}^{(k)}
                                                                                                                                                                                                  end for
                                                                                                                                                                                                 Output: \widehat{y}_t(\theta) = \mathbb{E}_{\pi_{\theta,t}}[Y] \approx \frac{1}{K} \sum_{k=1}^{K} y^{(k)}
```

Formally, we denote the neighborhood graph by $G_{\mathcal{N}} := (\mathcal{Y}, E_{\mathcal{N}})$, where edges are defined by \mathcal{N} . We assume the graph is undirected, i.e., $y' \in \mathcal{N}(y)$ if and only if $y \in \mathcal{N}(y')$, and without self-loops – 77 i.e., $y \notin \mathcal{N}(y)$. A stochastic neighbor generating function is also provided, in the form of a proposal 78 distribution $q(\boldsymbol{y},\cdot)$ with support either equal to $\mathcal{N}(\boldsymbol{y})$ or $\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\}$. 79

Link between simulated annealing and Metropolis-Hastings. A well-known example of local search heuristic is simulated annealing (SA) [26]. It is intimately related to Metropolis-Hastings (MH) [21], an instance of a MCMC algorithm. We provide a unified view of both in Algorithm 1.

The difference lies in the acceptance rule, which incorporates a proposal correction ratio for MH, 83 and in the choice of the sequence of temperatures $(t_k)_{k\in\mathbb{N}}$. In the case of SA, it is chosen to verify 84 $t_k \to 0$. In the case of MH, it is such that $t_k \equiv t$. In this case, the iterates $y^{(k)}$ of Algorithm 1 follow 85 a time-homogenous Markov chain on \mathcal{Y} , defined by the following transition kernel:

$$P_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}') = \begin{cases} q(\boldsymbol{y},\boldsymbol{y}') \min\left[1, \frac{q(\boldsymbol{y}',\boldsymbol{y})}{q(\boldsymbol{y},\boldsymbol{y}')} \exp\left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}') - \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle - \varphi(\boldsymbol{y})}{t}\right)\right] & \text{if } \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}), \\ 1 - \sum_{\boldsymbol{y}'' \in \mathcal{N}(\boldsymbol{y})} P_{\boldsymbol{\theta},t_k}(\boldsymbol{y},\boldsymbol{y}'') & \text{if } \boldsymbol{y}' = \boldsymbol{y}, \\ 0 & \text{else.} \end{cases}$$
(3)

In past work, the link between the two algorithms has primarily been used to show that SA converges to the exact MAP solution in the limit of infinite iterations [36, 17]. Under mild conditions – if the 89 neighborhood graph G_N is connected and the chain is aperiodic – the iterates $y^{(k)}$ of Algorithm 1 90 (MH case) converge in distribution to the Gibbs distribution (see Appendix E.1 for a proof): 91

$$\pi_{\boldsymbol{\theta},t}(\boldsymbol{y}) \propto \exp\left(\left[\langle \boldsymbol{\theta}, \, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})\right]/t\right).$$
 (4)

Proposed layer. Algorithm 1 and this result motivate us to define the combinatorial MCMC layer

$$\widehat{\boldsymbol{y}}_t(\boldsymbol{\theta}) \coloneqq \mathbb{E}_{\pi_{\boldsymbol{\theta},t}}[Y], \tag{5}$$

where $\theta \in \mathbb{R}^d$ are logits and t > 0 is a temperature parameter, defaulting to t = 1. Naturally, the 94 estimate of $\hat{y}_t(\theta)$ returned by Algorithm 1 (MH case) is biased, as the Markov chain cannot perfectly 95 mix in a finite number of iterations, except if it is initialized at $\pi_{\theta,t}$. In Section 3, we will show 96 that this does not hinder the convergence of the proposed learning algorithms. The next proposition, 97 proved in Appendix E.2, states some useful properties of the proposed layer.

Proposition 2.1. Let $\theta \in \mathbb{R}^d$. We have the following properties:

$$\widehat{\boldsymbol{y}}_t(\boldsymbol{\theta}) \in \operatorname{relint}(\mathcal{C}), \quad \widehat{\boldsymbol{y}}_t(\boldsymbol{\theta}) \xrightarrow[t \to 0^+]{} \operatorname{argmax}_{\boldsymbol{y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}), \quad and \quad \widehat{\boldsymbol{y}}_t(\boldsymbol{\theta}) \xrightarrow[t \to \infty]{} \frac{1}{|\mathcal{Y}|} \sum_{\boldsymbol{y} \in \mathcal{Y}} \boldsymbol{y} \ .$$

Moreover, \widehat{y}_t is differentiable and its Jacobian matrix is given by $J_{\theta}\widehat{y}_t(\theta) = \frac{1}{t} \cot_{\pi_{\theta,t}}[Y]$.

2.2 Mixing neighborhood systems

100

120

123

134

101 Central to local search algorithms in combinatorial optimization is the use of multiple neighborhood systems to more effectively explore the solution space [37, 10]. We now present a tractable way to incorporate such diversity of neighborhood systems into the proposed layer, while preserving the correct stationary distribution, by mixing the corresponding proposal distributions. A discussion giving intuition on why the proposed method is crucial to get tractable updates is given in Appendix C.1.

Definitions. Let $(\mathcal{N}_s)_{s=1}^S$ be a set of different neighborhood systems. Typically, all neighborhood systems are not defined on all solutions $\boldsymbol{y} \in \mathcal{Y}$, so we note $Q(\boldsymbol{y}) \subseteq [\![1,S]\!]$ the set of neighborhood systems defined on \boldsymbol{y} (i.e., the set of allowed moves on \boldsymbol{y}). Let $(q_s)_{s \in Q(\boldsymbol{y})}$ be the corresponding proposal distributions, such that the support of $q_s(\boldsymbol{y},\cdot)$ is either $\mathcal{N}_s(\boldsymbol{y})$ or $\mathcal{N}_s(\boldsymbol{y}) \cup \{\boldsymbol{y}\}$. Let $\bar{\mathcal{N}}$ be the aggregate neighborhood system defined by $\bar{\mathcal{N}}: \boldsymbol{y} \mapsto \bigcup_{s \in Q(\boldsymbol{y})} \mathcal{N}_s(\boldsymbol{y})$.

We assume that the individual Metropolis correction ratios $\tilde{\alpha}_s(\boldsymbol{y}, \boldsymbol{y}') \coloneqq \frac{q_s(\boldsymbol{y}', \boldsymbol{y})}{q_s(\boldsymbol{y}, \boldsymbol{y}')}$ are tractable. The proposed procedure is summarized in Algorithm 2.

Proposition 2.2. If each neighborhood graph $G_{\mathcal{N}_s}$ is undirected and without self-loops, and the aggregate neighborhood graph $G_{\overline{\mathcal{N}}}$ is connected, the iterations $\boldsymbol{y}^{(k)}$ produced by Algorithm 2 follow a Markov chain with unique stationary distribution $\pi_{\boldsymbol{\theta},t}$.

See Appendix E.3 for the proof. Importantly, only the connectedness of the aggregate neighborhood graph $G_{\widetilde{\mathcal{N}}}$ is required. This allows us to combine neighborhood systems that could not connect \mathcal{Y} if used individually, i.e., an irreducible Markov chain can be obtained by mixing the proposal distributions of reducible ones. Such an example is given with the moves defined in Table 4.

3 Loss functions and theoretical analysis

We now derive and study loss functions for learning models using the proposed layer. The analysis for the case where only one iteration of MCMC is performed (K = 1) is given in Appendix C.2.

3.1 Negative log-likelihood and associated Fenchel-Young loss

We now show that the proposed layer $\widehat{y}_t(\theta)$ can be viewed as the solution of a regularized optimization problem on $\mathcal{C} = \text{conv}(\mathcal{Y})$. Let $A_t(\theta) \coloneqq t \cdot \log \sum_{\boldsymbol{y} \in \mathcal{Y}} \exp\left(\left[\left\langle \boldsymbol{\theta} \,,\, \boldsymbol{y} \right\rangle + \varphi(\boldsymbol{y})\right]/t\right)$ be the cumulant function [51] associated to the exponential family defined by $\pi_{\theta,t}$, scaled by the temperature t. We define the regularization function Ω_t and the corresponding Fenchel-Young loss [8] as:

$$\Omega_t(\boldsymbol{\mu}) \coloneqq A_t^*(\boldsymbol{\mu}) = \sup_{\boldsymbol{\theta} \in \mathbb{R}^d} \langle \boldsymbol{\mu} \,,\, \boldsymbol{\theta} \rangle - A_t(\boldsymbol{\theta}), \quad \text{and} \quad \ell_t(\boldsymbol{\theta} \,; \boldsymbol{y}) \coloneqq (\Omega_t)^*(\boldsymbol{\theta}) + \Omega_t(\boldsymbol{y}) - \langle \boldsymbol{\theta},\, \boldsymbol{y} \rangle.$$

Since $\Omega_t = A_t^*$ is strictly convex on relint(\mathcal{C}) (see Appendix E.4 for a proof) and $\widehat{y}_t(\theta) = \nabla_{\theta} A_t(\theta)$, the proposed layer is the solution of the regularized optimization problem

$$\widehat{\boldsymbol{y}}_{t}(\boldsymbol{\theta}) = \underset{\boldsymbol{\mu} \in \mathcal{C}}{\operatorname{argmax}} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_{t}(\boldsymbol{\mu}) \right\}, \tag{6}$$

the Fenchel-Young loss ℓ_t is differentiable, satisfies $\ell_t(\theta, y) = 0 \Leftrightarrow \widehat{y}_t(\theta) = y$, and has gradient $\nabla_{\theta}\ell_t(\theta; y) = \widehat{y}_t(\theta) - y$ [8]. It is therefore equivalent, up to a constant, to the negative log-likelihood loss, as we have $-\nabla_{\theta}\log \pi_{\theta,t}(y) = (\widehat{y}_t(\theta) - y)/t$. Algorithms 1 and 2 can thus be used to perform maximum likelihood estimation, by returning a (biased) stochastic estimate of the gradient of ℓ_t .

3.2 Empirical risk minimization

In the supervised learning setting, we are given observations $(\boldsymbol{x}_i,\,\boldsymbol{y}_i)_{i=1}^N\in(\mathbb{R}^p\times\mathcal{Y})^N$, and want to fit a model $g_W\colon\mathbb{R}^p\to\mathbb{R}^d$ such that $\widehat{\boldsymbol{y}}_t(g_W(\boldsymbol{x}_i))\approx\boldsymbol{y}_i$. This is motivated by a generative model where, for some weights $W_0\in\mathbb{R}^p$, the data is generated with $\boldsymbol{y}_i\sim\pi_{g_{W_0}(\boldsymbol{x}_i),t}$. We aim at minimizing the empirical risk L_N , defined below along with its exact gradient ∇L_N :

$$L_N(W) \coloneqq \frac{1}{N} \sum_{i=1}^N \ell_t \left(g_W(\boldsymbol{x}_i); \boldsymbol{y}_i \right) \quad \text{and} \quad \nabla_W L_N(W) = \frac{1}{N} \sum_{i=1}^N J_W g_W(\boldsymbol{x}_i) \left(\widehat{\boldsymbol{y}}_t(g_W(\boldsymbol{x}_i)) - \boldsymbol{y}_i \right).$$

Doubly stochastic gradient estimator. In practice, we cannot compute the exact gradient above. Using Algorithm 1 to get a MCMC estimate of $\hat{y}_t(g_W(x_i))$, we propose the following estimator:

$$\nabla_W L_N(W) \approx J_W g_W(\boldsymbol{x}_i) \left(\frac{1}{K} \sum_{k=1}^K \boldsymbol{y}_i^{(k)} - \boldsymbol{y}_i \right),$$

where $y_i^{(k)}$ is the k-th iterate of Algorithm 1 with maximization direction $\theta_i = g_W(x_i)$ and temperature t. This estimator is doubly stochastic, since we sample both data points and iterations of Algorithm 1, and can be seamlessly used with batches. The vector-jacobian product with $J_W g_W(x_i)$ is computed via autodiff. The Markov chain initialization methods for the supervised and unsupervised settings are inspired from the contrastive divergence literature [22] and detailed in Appendix C.5.

3.3 Convergence analysis in the unsupervised setting

144

161

162

163

164

165

In the unsupervised setting, we are only given observations $(y_i)_{i=1}^N \in \mathcal{Y}^N$ and want to fit a model $\pi_{\theta,t}$, motivated by an underlying generative model such that $y_i \sim \pi_{\theta_0,t}$ for an unknown true parameter θ_0 .

We assume here that $\mathcal{C} = \text{conv}(\mathcal{Y})$ is of full dimension in \mathbb{R}^d . We have the following empirical L_N and population L_{θ_0} Fenchel-Young losses:

$$L_N(\boldsymbol{\theta}; \boldsymbol{y}_1, \dots, \boldsymbol{y}_N) \coloneqq \frac{1}{N} \sum_{i=1}^N \ell_t\left(\boldsymbol{\theta}; \, \boldsymbol{y}_i\right) \;, \quad L_{\boldsymbol{\theta}_0}(\boldsymbol{\theta}) \coloneqq \mathbb{E}_{(\boldsymbol{y}_i)_{i=1}^N \sim (\pi_{\boldsymbol{\theta}_0, t})^{\otimes N}} \left[L_N(\boldsymbol{\theta}; \boldsymbol{y}_1, \dots, \boldsymbol{y}_N)\right],$$

which are minimized for $\boldsymbol{\theta}$ such that $\widehat{y}_t(\boldsymbol{\theta}) = \bar{Y}_N \coloneqq \frac{1}{N} \sum_{i=1}^N \boldsymbol{y}_i$, and for $\boldsymbol{\theta}$ such that $\widehat{y}_t(\boldsymbol{\theta}) = \widehat{y}_t(\boldsymbol{\theta}_0)$, respectively. Let $\boldsymbol{\theta}_N^\star$ as the minimizer of the empirical loss L_N . For it to be defined, we assume N is large enough to have $\bar{Y}_N \in \operatorname{relint}(\mathcal{C})$ (which is always possible as $\pi_{\boldsymbol{\theta}_0,t}$ has dense support on \mathcal{Y}). A slight variation on Proposition 4.1 in Berthet et al. [6], proved in Appendix E.5, gives the following asymptotic normality as $N \to \infty$.

Proposition 3.1 (Convergence of the empirical loss minmizer to the true parameter).

$$\sqrt{N}(\boldsymbol{\theta}_{N}^{\star}-\boldsymbol{\theta}_{0})\xrightarrow[N\to\infty]{\mathcal{D}}\mathcal{N}\left(\mathbf{0},\,t^{2}\operatorname{cov}_{\pi_{\boldsymbol{\theta}_{0},t}}\left[Y\right]^{-1}\right).$$

We now consider the sample size as fixed to N samples, and define $\hat{\theta}_n$ as the n-th iterate of the following stochastic gradient algorithm:

$$\hat{\boldsymbol{\theta}}_{n+1} = \hat{\boldsymbol{\theta}}_n + \gamma_{n+1} \left[\bar{Y}_N - \frac{1}{K_{n+1}} \sum_{k=1}^{K_{n+1}} \boldsymbol{y}^{(n+1,k)} \right],$$
 (7)

where $y^{(n+1,k)}$ is the k-th iterate of Algorithm 1 with temperature t, maximization direction $\hat{\theta}_n$, and initialized at $y^{(n+1,1)} = y^{(n,K_n)}$. This initialization corresponds to the persistent contrastive divergences (PCD) algorithm [48], and is further discussed in Appendix C.5.

Proposition 3.2 (Convergence of the stochastic gradient estimate). Suppose the following assumptions on the step sizes $(\gamma_n)_{n>1}$, sample sizes $(K_n)_{n>1}$, and proposal distribution q hold:

•
$$\gamma_n = an^{-b}$$
, with $b \in \left]\frac{1}{2}, 1\right]$ and $a > 0$,

•
$$K_{n+1} > \left| 1 + a' \exp\left(\frac{8R_{\mathcal{C}}}{t} \cdot ||\hat{\boldsymbol{\theta}}_n||\right) \right|$$
 with $a' > 0$ and $R_{\mathcal{C}} = \max_{\boldsymbol{y} \in \mathcal{Y}} ||\boldsymbol{y}||$,

•
$$\frac{1}{\sqrt{K_n}} - \frac{1}{\sqrt{K_{n-1}}} \le a'' n^{-c}$$
, with $a'' > 0$ and $c > 1 - \frac{b}{2}$,

$$\mathbf{\bullet} \ q(\mathbf{y}, \mathbf{y}') = \begin{cases} \frac{1}{2d^*} & \text{if } \mathbf{y}' \in \mathcal{N}(\mathbf{y}), \\ 1 - \frac{d(\mathbf{y})}{2d^*} & \text{if } \mathbf{y}' = \mathbf{y}, \\ 0 & \text{else}, \end{cases}$$

$$\text{where } d(\mathbf{y}) \coloneqq |\mathcal{N}(\mathbf{y})| \text{ is the degree of } \mathbf{y} \text{ in } G_{\mathcal{N}}, \text{ and } d^* \coloneqq \max_{\mathbf{y} \in \mathcal{Y}} d(\mathbf{y}).$$

166 Then, we have the almost sure convergence $\hat{\theta}_n \xrightarrow{a.s.} \theta_N^*$ of the iterates $\hat{\theta}_n$ defined by Eq. (7).

See Appendix E.6 for the proof. The assumptions on q are used for obtaining a closed-form convergence rate bound for the Markov chain, using graph-based geometric bounds [23].

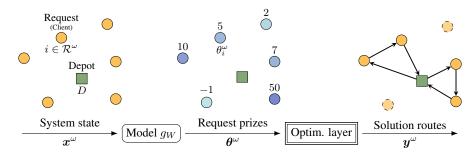


Figure 1: Overview of the vehicle routing pipeline, represented at request wave ω .

4 Experiments on dynamic vehicle routing

We demonstrate our approach on the EURO Meets NeurIPS 2022 Vehicle Routing Competition, a large-scale, ML-enriched dynamic vehicle routing problem with time windows (DVRPTW). A detailed introduction to the challenge with precise notations, together with precise explications on the reduction to supervised learning, the proposed approach, the perturbation-based baseline, full experimental details and additional results, are given in Appendix D.

Approach. We adopt the winning strategy of Baty et al. [4], which reduces the problem to a supervised learning task. This involves decomposing the DVRPTW as multiple prize-collecting problems (PC-VRPTW) for each wave ω , where a model g_W predicts a prize vector $\boldsymbol{\theta}^\omega$ for serving each request. This fits our general problem formulation $\widehat{\boldsymbol{y}}(\boldsymbol{\theta}^\omega) = \operatorname{argmax}_{\boldsymbol{y}}\langle \boldsymbol{\theta}^\omega, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})$, where $\varphi(\boldsymbol{y})$ represents the negative routing cost. The overall pipeline is shown in Fig. 1. To train the model g_W , we use the Fenchel-Young loss associated with our proposed MCMC layer. The proposal distributions for the MCMC sampler are derived from the local search moves used by the state-of-the-art PC-HGS solver $\widetilde{\boldsymbol{y}}$, which are summarized in Table 4. At inference, we use the trained model g_W with the PC-HGS solver, forming the policy $f_W := \widetilde{\boldsymbol{y}} \circ g_W$. We compare our learning algorithms to the perturbation-based baseline from Baty et al. [4], which sacrifices the theoretical guarantees of the general framework from Berthet et al. [6] it instantiates by using an inexact solver $(\widetilde{\boldsymbol{y}})$.

Results. We evaluate performance using the competition's metric: the routing cost relative to an anticipative (oracle) baseline. Fig. 2 shows that initializing the MCMC chain with the ground-truth solution significantly outperforms a random start and that performance improves with more MCMC iterations (K). Table 2 compares our method with the baseline under a fixed time budget for the layer's forward pass. Our approach significantly outperforms the perturbation-based method in low-time-limit regimes (1-100 ms), enabling faster and more efficient training.

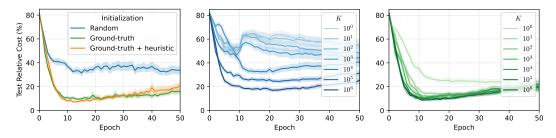


Figure 2: Test relative cost (%). **Left**: varying initialization method. **Center**: varying number of Markov iterations K with random initialization. **Right**: varying K with ground-truth initialization.

Table 2: Best test relative cost (%) for different training methods and time limits.

Tuble 2. Best test relative cost (%) for different training methods and time minus.						
Time limit (ms)	1	5	10	50	100	1000
Perturbed inexact oracle	65.2 ± 5.8	13.1 ± 3.4	8.7 ± 1.9	6.5 ± 1.1	6.3 ± 0.76	5.5 ± 0.4
Proposed $(\boldsymbol{y}^{(0)} = \boldsymbol{y})$		12.0 ± 2.6				
Proposed ($y^{(0)} = y$ +heuristic)	7.8 ± 0.8	7.2 ± 0.6	6.3 ± 0.7	6.2 ± 0.8	5.9 ± 0.7	5.9 ± 0.6

2 References

- [1] Akshay Agrawal, Brandon Amos, Shane Barratt, Stephen Boyd, Steven Diamond, and Zico
 Kolter. Differentiable convex optimization layers, 2019. URL https://arxiv.org/abs/
 1910.12430.
- [2] Kareem Ahmed, Zhe Zeng, Mathias Niepert, and Guy Van den Broeck. SIMPLE: A gradient estimator for \$k\$-subset sampling, 2024. URL http://arxiv.org/abs/2210.01941.
- [3] Brandon Amos and J Zico Kolter. Optnet: Differentiable optimization as a layer in neural
 networks. In *International Conference on Machine Learning*, pages 136–145. PMLR, 2017.
- [4] Léo Baty, Kai Jungel, Patrick S. Klein, Axel Parmentier, and Maximilian Schiffer. Combinatorial
 optimization enriched machine learning to solve the dynamic vehicle routing problem with time
 windows, 2023. URL http://arxiv.org/abs/2304.00789.
- Yoshua Bengio, Andrea Lodi, and Antoine Prouvost. Machine learning for combinatorial optimization: a methodological tour d'horizon, 2020. URL http://arxiv.org/abs/1811.
 06128.
- 206 [6] Quentin Berthet, Mathieu Blondel, Olivier Teboul, Marco Cuturi, Jean-Philippe Vert, and
 207 Francis Bach. Learning with differentiable perturbed optimizers, 2020. URL http://arxiv.
 208 org/abs/2002.08676.
- [7] Mathieu Blondel and Vincent Roulet. The Elements of Differentiable Programming. *arXiv* preprint arXiv:2403.14606, 2024.
- [8] Mathieu Blondel, André F. T. Martins, and Vlad Niculae. Learning with fenchel-young losses, 2020. URL http://arxiv.org/abs/1901.02324.
- [9] Mathieu Blondel, Quentin Berthet, Marco Cuturi, Roy Frostig, Stephan Hoyer, Felipe Llinares-López, Fabian Pedregosa, and Jean-Philippe Vert. Efficient and modular implicit differentiation. Advances in neural information processing systems, 35:5230–5242, 2022.
- 216 [10] Christian Blum and Andrea Roli. Metaheuristics in combinatorial optimization: Overview and conceptual comparison. 35(3):268–308, 2003. ISSN 0360-0300. doi: 10.1145/937503.937505. URL https://doi.org/10.1145/937503.937505.
- [11] Miguel A Carreira-Perpiñán and Geoffrey Hinton. On contrastive divergence learning. In
 International Workshop on Artificial Intelligence and Statistics, pages 33–40. PMLR, 2005.
 URL https://proceedings.mlr.press/r5/carreira-perpinan05a.html.
- 222 [12] Bor-Liang Chen and Ko-Wei Lih. Hamiltonian uniform subset graphs. 42(3):257–263, 1987. ISSN 0095-8956. doi: 10.1016/0095-8956(87)90044-X. URL https://www.sciencedirect.com/science/article/pii/009589568790044X.
- 225 [13] John M. Danskin. The theory of max-min, with applications. 14(4):641–664, 1966. ISSN 226 0036-1399. doi: 10.1137/0114053. URL https://epubs.siam.org/doi/abs/10.1137/227 0114053.
- 228 [14] Priya Donti, Brandon Amos, and J Zico Kolter. Task-based end-to-end model learning in stochastic optimization. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- 231 [15] Yilun Du and Igor Mordatch. Implicit generation and generalization in energy-based models, 232 2020. URL https://arxiv.org/abs/1903.08689.
- 233 [16] Yilun Du, Shuang Li, Joshua Tenenbaum, and Igor Mordatch. Improved contrastive divergence training of energy based models, 2021. URL https://arxiv.org/abs/2012.01316.
- 235 [17] Ulrich Faigle and Rainer Schrader. On the convergence of stationary distributions in sim-236 ulated annealing algorithms. 27(4):189–194, 1988. ISSN 0020-0190. doi: 10.1016/ 237 0020-0190(88)90024-5. URL https://www.sciencedirect.com/science/article/ 238 pii/0020019088900245.

- 239 [18] Ari Freedman. CONVERGENCE THEOREM FOR FINITE MARKOV
 240 CHAINS. 2017. URL https://www.semanticscholar.org/paper/
 241 CONVERGENCE-THEOREM-FOR-FINITE-MARKOV-CHAINS-%E2%8B%82t/
 242 65f7c092bd9c59cbbc88dd69266d39cd79840648.
- ²⁴³ [19] Michel Gendreau, Jean-Yves Potvin, et al. *Handbook of metaheuristics*, volume 2. Springer, 2010.
- [20] Will Grathwohl, Kevin Swersky, Milad Hashemi, David Duvenaud, and Chris J. Maddison.
 Oops i took a gradient: Scalable sampling for discrete distributions, 2021. URL https://arxiv.org/abs/2102.04509.
- 248 [21] W. K. Hastings. Monte carlo sampling methods using markov chains and their applications.
 249 *Biometrika*, 57(1):97–109, 1970. ISSN 00063444, 14643510. URL http://www.jstor.org/
 250 stable/2334940.
- 251 [22] Geoffrey E. Hinton. Training products of experts by minimizing con-252 trastive divergence. 2000. URL https://www.semanticscholar. 253 org/paper/Training-Products-of-Experts-by-Minimizing-Hinton/ 254 9360e5ce9c98166bb179ad479a9d2919ff13d022.
- 255 [23] Salvatore Ingrassia. On the rate of convergence of the metropolis algorithm and gibbs sampler 256 by geometric bounds. 4(2):347–389, 1994. ISSN 1050-5164. URL https://www.jstor. 257 org/stable/2245161.
- 258 [24] Gareth A. Jones. Automorphisms and regular embeddings of merged johnson graphs. 26
 259 (3):417-435, 2005. ISSN 0195-6698. doi: 10.1016/j.ejc.2004.01.012. URL https://www.sciencedirect.com/science/article/pii/S0195669804000630.
- 261 [25] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL http://arxiv.org/abs/1412.6980.
- [26] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*,
 220(4598):671–680, 1983. doi: 10.1126/science.220.4598.671. URL https://www.science.org/doi/abs/10.1126/science.220.4598.671.
- Wouter Kool, Laurens Bliek, Danilo Numeroso, Yingqian Zhang, Tom Catshoek, Kevin Tierney,
 Thibaut Vidal, and Joaquim Gromicho. The EURO meets NeurIPS 2022 vehicle routing
 competition. In *Proceedings of the NeurIPS 2022 Competitions Track*, pages 35–49. PMLR,
 2023. URL https://proceedings.mlr.press/v220/kool23a.html.
- 270 [28] Rahul G. Krishnan, Simon Lacoste-Julien, and David Sontag. Barrier frank-wolfe for marginal inference, 2015. URL https://arxiv.org/abs/1511.02124.
- [29] John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. Conditional random fields:
 Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning*, ICML '01, page 282–289, San Francisco,
 CA, USA, 2001. Morgan Kaufmann Publishers Inc. ISBN 1558607781.
- 276 [30] Yann Lecun, Sumit Chopra, Raia Hadsell, Marc Aurelio Ranzato, and Fu Jie Huang. *A tutorial* 277 on energy-based learning. MIT Press, 2006.
- 278 [31] Zhifei Li and Jason Eisner. First- and second-order expectation semirings with applica-279 tions to minimum-risk training on translation forests. In Philipp Koehn and Rada Mihalcea, 280 editors, *Proceedings of the 2009 Conference on Empirical Methods in Natural Language* 281 *Processing*, pages 40–51. Association for Computational Linguistics, August 2009. URL 282 https://aclanthology.org/D09-1005/.
- 283 [32] Neal Madras and Dana Randall. Markov chain decomposition for convergence rate analysis. 12
 284 (2):581-606, 2002. ISSN 1050-5164, 2168-8737. doi: 10.1214/aoap/1026915617. URL https:
 285 //projecteuclid.org/journals/annals-of-applied-probability/volume-12/
 286 issue-2/Markov-chain-decomposition-for-convergence-rate-analysis/10.
 287 1214/aoap/1026915617.full.

- 288 [33] Jayanta Mandi and Tias Guns. Interior point solving for LP-based prediction+optimisation, 2020. URL http://arxiv.org/abs/2010.13943.
- [34] Jayanta Mandi, James Kotary, Senne Berden, Maxime Mulamba, Victor Bucarey, Tias Guns, and Ferdinando Fioretto. Decision-focused learning: Foundations, state of the art, benchmark and future opportunities. 80:1623–1701, 2024. ISSN 1076-9757. doi: 10.1613/jair.1.15320. URL http://arxiv.org/abs/2307.13565.
- 294 [35] Arthur Mensch and Mathieu Blondel. Differentiable dynamic programming for structured prediction and attention, 2018. URL https://arxiv.org/abs/1802.03676.
- [36] Debasis Mitra, Fabio Romeo, and Alberto Sangiovanni-Vincentelli. Convergence and finite-time
 behavior of simulated annealing. *Advances in Applied Probability*, 18(3):747–771, 1986. ISSN 0001-8678. doi: 10.2307/1427186. URL https://www.jstor.org/stable/1427186.
- 299 [37] Nenad Mladenović and Pierre Hansen. Variable neighborhood search. *Computers & operations* 300 research, 24(11):1097–1100, 1997.
- 301 [38] Volodymyr Mnih, Hugo Larochelle, and Geoffrey E. Hinton. Conditional restricted boltzmann machines for structured output prediction, 2012. URL http://arxiv.org/abs/1202.3748.
- 303 [39] Vlad Niculae, André F. T. Martins, Mathieu Blondel, and Claire Cardie. Sparsemap: Differen-304 tiable sparse structured inference, 2018. URL https://arxiv.org/abs/1802.04223.
- Benjamin Rhodes and Michael Gutmann. Enhanced gradient-based MCMC in discrete spaces, 2022. URL http://arxiv.org/abs/2208.00040.
- 307 [41] Fred J. Rispoli. The graph of the hypersimplex, 2008. URL http://arxiv.org/abs/0811.
 308 2981.
- 309 [42] R. Tyrrell Rockafellar. *Convex Analysis*. Princeton University Press, 1970. ISBN 9780691015866. URL http://www.jstor.org/stable/j.ctt14bs1ff.
- Utsav Sadana, Abhilash Chenreddy, Erick Delage, Alexandre Forel, Emma Frejinger, and
 Thibaut Vidal. A survey of contextual optimization methods for decision making under uncertainty, 2024. URL http://arxiv.org/abs/2306.10374.
- 314 [44] Yang Song and Diederik P. Kingma. How to train your energy-based models, 2021. URL https://arxiv.org/abs/2101.03288.
- Haoran Sun, Hanjun Dai, Bo Dai, Haomin Zhou, and Dale Schuurmans. Discrete langevin sampler via wasserstein gradient flow, 2023. URL http://arxiv.org/abs/2206.14897.
- [46] Haoran Sun, Katayoon Goshvadi, Azade Nova, Dale Schuurmans, and Hanjun Dai. Revisiting sampling for combinatorial optimization. In *Proceedings of the 40th International Conference on Machine Learning*, pages 32859–32874. PMLR, 2023. URL https://proceedings.mlr.press/v202/sun23c.html.
- Ilya Sutskever and Tijmen Tieleman. On the convergence properties of contrastive divergence. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence* and Statistics, pages 789–795. JMLR Workshop and Conference Proceedings, 2010. URL https://proceedings.mlr.press/v9/sutskever10a.html.
- Tijmen Tieleman. Training restricted boltzmann machines using approximations to the likelihood gradient. In *Proceedings of the 25th international conference on Machine learning*, ICML '08, pages 1064–1071. Association for Computing Machinery, 2008. ISBN 9781605582054. doi: 10.1145/1390156.1390290. URL https://doi.org/10.1145/1390156.1390290.
- Thibaut Vidal. Hybrid genetic search for the cvrp: Open-source implementation and swap* neighborhood. *Computers & Operations Research*, 140:105643, April 2022. ISSN 0305-0548. doi: 10.1016/j.cor.2021.105643. URL http://dx.doi.org/10.1016/j.cor.2021.105643.
- [50] Marin Vlastelica, Anselm Paulus, Vít Musil, Georg Martius, and Michal Rolínek. Differentiation
 of blackbox combinatorial solvers, 2020. URL http://arxiv.org/abs/1912.02175.

- Martin J. Wainwright and Michael I. Jordan. Graphical models, exponential families, and variational inference. 1(1):1–305, 2008. ISSN 1935-8237, 1935-8245. doi: 10.1561/2200000001.
 URL https://www.nowpublishers.com/article/Details/MAL-001.
- [52] Laurent Younes. Stochastic gradient estimation strategies for markov ran-339 In Bayesian Inference for Inverse Problems, volume 3459, dom fields. 340 URL https://www. 10.1117/12.323811. pages 315-325. SPIE, 1998. doi: 341 spiedigitallibrary.org/conference-proceedings-of-spie/3459/0000/ 342 Stochastic-gradient-estimation-strategies-for-Markov-random-fields/ 343 10.1117/12.323811.full. 344
- Ruqi Zhang, Xingchao Liu, and Qiang Liu. A langevin-like sampler for discrete distributions, 2022. URL https://arxiv.org/abs/2206.09914.

A Background and Related Work

A.1 Combinatorial optimization as a layer

Because the function in Eq. (1) is piecewise constant and discontinuous, a frequent strategy consists in introducing regularization in the problem so as to obtain a continuous relaxation. In some cases, we may have access to an oracle for directly solving the regularized problem. For instance, when the unregularized problem can be solved by dynamic programming, its entropic regularization can be computed using a change of semi-ring [31] or by algorithmic smoothing [35]. As another example, interior point solvers can be used to compute a logarithmic barrier regularized solution [33].

We focus settings where only a MAP oracle is available for the original, unregularized optimization problem. While many prior works are limited to the linear form in Eq. (1) for the latter, our framework is more general and also handles problems of the form in Eq. (2). Frank-Wolfe-like methods can be used to solve the regularized problem using only MAP oracle calls [39, 28]. Another strategy consists in injecting noise perturbations [6] in the oracle. This approach can be shown to be implicitly using regularization. In both cases, a Fenchel-Young loss can be associated, giving a principled way to learn with the optimization layer. However, formal guarantees only hold if the oracle used is exact, and in practice it is typically called multiple times during the forward pass. Our proposal enjoys guarantees even with inexact solvers and a single call.

Regarding differentiation, several strategies are possible. When the approach only needs to differentiate through a (regularized) max, as is the case of Fenchel-Young losses, we can use Danskin's theorem [13]. When the approach needs to differentiate a (regularized) argmax, we can either use autodiff on the unrolled solver iterations or implicit differentiation [3, 1, 9]. Differently, Vlastelica et al. [50] propose to compute gradients via continuous interpolation of the solver.

A.2 Contrastive divergences

An alternative approach to learning in combinatorial spaces is to use energy-based models (EBMs) [30], which define a distribution over outputs via a parameterized energy function E_{θ} :

$$p_{\theta}(y) \propto \exp(E_{\theta}(y)), \text{ with } \nabla_{\theta} \log p_{\theta}(y) = \nabla_{\theta} E_{\theta}(y) - \mathbb{E}_{Y \sim p_{\theta}} \left[\nabla_{\theta} E_{\theta}(Y) \right].$$

Therefore, we can perform maximum likelihood estimation (MLE) if we can sample from p_{θ} , but this is hard both in continuous and combinatorial settings, due to its intractable normalization constant. Contrastive divergences [22, 11, 44] address this by using MCMC to obtain (biased) stochastic gradients. Originally developed for restricted Boltzmann machines with $\mathcal{Y} = \{0, 1\}^d$ and a Gibbs sampler, they have also been applied in continuous domains via Langevin dynamics [15, 16].

MCMC in discrete spaces. Contrastive divergences rely on MCMC to sample from the current model distribution. Unfortunately, designing an MCMC sampler is usually case by case, and MCMC on discrete domains has received comparatively less attention than continuous domains. Recent efforts adapt continuous techniques, such as Langevin dynamics [53, 45] or gradient-informed proposals [20, 40], to discrete settings. However, these works typically assume simple state space structure, like the hypercube or categorical codebooks, and do not handle complex constraints that are ubiquitous in operations research problems. Sun et al. [46] allow structured state spaces via relaxed constraints in the energy function, yet ignore these structures in the proposal supports. Notably, we emphasize that all these works focus on sampling, not on designing differentiable MCMC layers.

B Experiments on empirical convergence of gradients and parameters

In this section, we evaluate the proposed approach on two discrete output spaces: sets and κ -subsets. These output spaces are for instance useful for multilabel classification. We focus on these output spaces because the exact MAP and marginal inference oracles are available, allowing us to compare our gradient estimators to exact gradients.

389 B.1 Polytopes and corresponding oracles

The vertex set of the first polytope is the set of binary vectors in \mathbb{R}^d , which we denote $\mathcal{Y}^d := \{0,1\}^d$, and $conv(\mathcal{Y}^d) = [0,1]^d$ is the "hypercube". The vertex set of the second is the set of binary vectors

with exactly κ ones and $d - \kappa$ zeros (with $0 < \kappa < d$),

$$\mathcal{Y}_{\kappa}^{d} := \{ \boldsymbol{y} \in \{0,1\}^{d} : \langle \boldsymbol{y}, \boldsymbol{1} \rangle = \kappa \},$$

and $conv(\mathcal{Y}_{\kappa}^d)$ is referred to as "top- κ polytope" or "hypersimplex". Although these polytopes would not provide relevant use cases of the proposed approach in practice, since exact marginal inference oracles are available (see below), they allow us to compare the Fenchel-Young loss value and gradient estimated by our algorithm to their true value.

394 Marginal inference. For the hypercube, we have:

$$\begin{split} \mathbb{E}_{\pi_{\boldsymbol{\theta},t}}\left[Y_{i}\right] &= \sum_{\boldsymbol{y} \in \mathcal{Y}^{d}} \frac{\exp\left(\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle / t\right)}{\sum_{\boldsymbol{y}' \in \mathcal{Y}^{d}} \exp\left(\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle / t\right)} y_{i} \\ &= \sum_{\boldsymbol{y} \in \{0,1\}^{d}} \frac{\exp\left(\sum_{j=1}^{d} \theta_{j} y_{j} / t\right)}{\sum_{\boldsymbol{y}' \in \{0,1\}^{d}} \exp\left(\sum_{j=1}^{d} \theta_{j} y_{j} / t\right)} y_{i} \\ &= \sum_{\boldsymbol{y}_{i} \in \{0,1\}} \sum_{\boldsymbol{y}_{-i} \in \{0,1\}^{d-1}} \frac{\exp\left(\theta_{i} y_{i} / t + \sum_{j \neq i} \theta_{j} y_{j} / t\right)}{\sum_{\boldsymbol{y}'_{i} \in \{0,1\}} \sum_{\boldsymbol{y}'_{-i} \in \{0,1\}^{d-1}} \exp\left(\theta_{i} y_{i} / t + \sum_{j \neq i} \theta_{j} y_{j} / t\right)} y_{i} \\ &= \sum_{\boldsymbol{y}_{i} \in \{0,1\}} \frac{\exp\left(\theta_{i} y_{i} / t\right)}{\sum_{\boldsymbol{y}'_{i} \in \{0,1\}} \exp\left(\theta_{i} y_{i} / t\right)} y_{i} \sum_{\boldsymbol{y}_{-i} \in \{0,1\}^{d-1}} \frac{\exp\left(\sum_{j \neq i} \theta_{j} y_{j} / t\right)}{\sum_{\boldsymbol{y}'_{i} \in \{0,1\}} \exp\left(\theta_{i} y_{i} / t\right)} y_{i} \\ &= \sum_{\boldsymbol{y}_{i} \in \{0,1\}} \frac{\exp\left(\theta_{i} y_{i} / t\right)}{\sum_{\boldsymbol{y}'_{i} \in \{0,1\}} \exp\left(\theta_{i} y_{i} / t\right)} y_{i} \\ &= \frac{0 \cdot \exp\left(0\right) + 1 \cdot \exp\left(\theta_{i} / t\right)}{\exp\left(0\right) + \exp\left(\theta_{i} / t\right)} \\ &= \frac{\exp\left(\theta_{i} / t\right)}{1 + \exp\left(\theta_{i} / t\right)} \\ &= \sigma\left(\frac{\theta_{i}}{t}\right), \end{split}$$

which gives $\mathbb{E}_{\pi_{\theta,t}}[Y] = \sigma\left(\frac{\theta}{t}\right)$, where the logistic sigmoid function σ is applied component-wise. The cumulant function is also tractable, as we have

$$\log \sum_{\boldsymbol{y} \in \mathcal{Y}^d} \exp\left(\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle / t\right) = \log \sum_{\boldsymbol{y} \in \{0,1\}^d} \exp\left(\sum_{i=1}^d \theta_i y_i / t\right)$$

$$= \log \sum_{\boldsymbol{y}_1 = 0}^1 \sum_{\boldsymbol{y}_2 = 0}^1 \cdots \sum_{\boldsymbol{y}_d = 0}^1 \exp\left(\sum_{i=1}^d \theta_i y_i / t\right)$$

$$= \log \prod_{i=1}^d \sum_{\boldsymbol{y}_i = 0}^1 \exp\left(\theta_i y_i / t\right)$$

$$= \log \prod_{i=1}^d \left(\exp(0) + \exp\left(\theta_i / t\right)\right)$$

$$= \log \prod_{i=1}^d \left(1 + \exp\left(\theta_i / t\right)\right)$$

$$= \sum_{i=1}^d \log\left(1 + \exp\left(\theta_i / t\right)\right).$$

Another way to derive this is via the Fenchel conjugate.

398

For the top- κ polytope, such closed-form formulas do not exist for the cumulant and its gradient. However, we implement them with dynamic programming, by viewing the top- κ MAP problem as a 0/1-knapsack problem with constant item weights, and by changing the $(\max, +)$ semiring into a (LSE, +) semiring. This returns the cumulant function, and we leverage PyTorch's automatic differentiation framework to compute its gradient. This simple implementation allows us to compute true Fenchel-Young losses values and their gradients in $\mathcal{O}(d\kappa)$ time and space complexity.

Sampling. For the hypercube, sampling from the Gibbs distribution on \mathcal{Y}^d has closed form. Indeed, the latter is fully factorized, and we can sample $\boldsymbol{y} \sim \pi_{\boldsymbol{\theta},t}$ by sampling independently each component with $y_i \sim \text{Bern}\left(\sigma(\theta_i/t)\right)$. Sampling from $\pi_{\boldsymbol{\theta},t}$ is also possible on \mathcal{Y}^d_{κ} , by sampling coordinates iteratively using the dynamic programming table used to compute the cumulant function (see, e.g., Algorithm 2 in Ahmed et al. [2] for a detailed explanation).

410 B.2 Neighborhood graphs

Hypercube. On \mathcal{Y}^d , we use a family of neighborhood systems \mathcal{N}^r_{\leq} parameterized by a Hamming distance radius $r \in [d-1]$. The graph is defined by:

$$\forall \boldsymbol{y}, \boldsymbol{y}' \in \mathcal{Y}^d: \ \boldsymbol{y}' \in \mathcal{N}_{<}^r(\boldsymbol{y}) \Leftrightarrow 1 \leq d_H(\boldsymbol{y}, \ \boldsymbol{y}') \leq r.$$

That is, two vertices are neighbors if their Hamming distance is at most r. This graph is regular, with degree $|\mathcal{N}_{\leq}^r(\boldsymbol{y})| = \sum_{i=1}^r {d \choose i}$. This graph is naturally connected, as any binary vector \boldsymbol{y}' can be reached from any other binary vector \boldsymbol{y} in $||\boldsymbol{y}'-\boldsymbol{y}||_1$ moves, by flipping each bit where $y_i' \neq y_i$, iteratively. Indeed, this trajectory consists in moves between vertices with Hamming distance equal to 1, and are therefore along edges of the neighborhood graph, regardless of the value of r.

We also use a slight variation on this family of neighborhood systems, the graphs $\mathcal{N}_{=}^{r}$, defined by:

$$\forall \boldsymbol{y}, \boldsymbol{y}' \in \mathcal{Y}^d: \ \boldsymbol{y}' \in \mathcal{N}^r_{=}(\boldsymbol{y}) \Leftrightarrow d_H(\boldsymbol{y}, \ \boldsymbol{y}') = r.$$

These graphs, on the contrary, are not always connected: indeed, if r is even, they contain two connected components (binary vectors with an even sum, and binary vectors with an odd sum).
We only use such graphs when experimenting with neighborhood mixtures (see Algorithm 2), by aggregating them into a connected graph.

Top- κ polytope. On \mathcal{Y}_{κ}^d , we use a family of neighborhoods systems \mathcal{N}^s parameterized by a number of "swaps" $s \in [\![1, \min(\kappa, d - \kappa)]\!]$. The graph is defined by

$$\forall \boldsymbol{y}, \boldsymbol{y}' \in \mathcal{Y}_{\kappa}^{d}: \ \boldsymbol{y}' \in \mathcal{N}^{s}(\boldsymbol{y}) \Leftrightarrow d_{H}(\boldsymbol{y}, \ \boldsymbol{y}') = 2s.$$

That is, two vertices are neighbors if one can be reached from the other by performing s "swaps", each swap corresponding to flipping a 1 to a 0 and vice-versa. This ensures that the resulting vector is still in \mathcal{Y}_{κ}^d . All s swaps must be performed on distinct components. The resulting graph is known as the *Generalized Johnson Graph* $J(d, \kappa, \kappa - s)$, or *Uniform Subset Graph* [12]. It is a regular graph, with degree $|\mathcal{N}^s(y)| = {\kappa \choose s} {d-\kappa \choose s}$. It is proved to be connected in Jones [24], except if $d=2\kappa$ and $s=\kappa$ (in this case, it consists in $\frac{1}{2} {d \choose \kappa}$ disjoint edges).

When s=1, the neighborhood graph is the Johnson Graph $J(d,\kappa)$, which coincides with the graph associated to the polytope $\text{conv}(\mathcal{Y}_{\kappa}^d) = \Delta_{d,\kappa}$ [41].

B.3 Convergence to exact gradients

In this section, we conduct experiments on the convergence of the MCMC estimators to the exact corresponding expectation (that is, convergence of the stochastic gradient estimator to the true gradient). The estimators are defined as

$$\widehat{\boldsymbol{y}}_t(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta},t}}\left[Y\right] \approx \frac{1}{K - K_0} \sum_{k=K_0+1}^{K} \boldsymbol{y}^{(k)},$$

where $\boldsymbol{y}^{(k)}$ is the k-th iterate of Algorithm 1 with maximization direction $\boldsymbol{\theta}$, final temperature t, and K_0 is a number of burn-in (or warm-up) iterations. The obtained estimator is compared to the exact expectation by performing marginal inference as described in Appendix B.1 (with a closed-form formula in the case of \mathcal{Y}^d , and by dynamic programming in the case of \mathcal{Y}^d).

438

443

444

445

Setup. For $T > K_0$, let $\tilde{\mathbb{E}}(\boldsymbol{\theta}, T) \coloneqq \frac{1}{T - K_0} \sum_{k=K_0+1}^T \boldsymbol{y}^{(k)}$ be the stochastic estimate of the expectation at step T. We proceed by first randomly generating $\boldsymbol{\Theta} \in \mathbb{R}^{M \times d}$, with M being the number of instances, by sampling $\boldsymbol{\Theta}_{i,j} \sim \mathcal{N}(0,1)$ independently. Then, we evaluate the impact of the following hyperparameters on the estimation of $\mathbb{E}_{\pi_{\boldsymbol{\Theta}_i},t}[Y]$, for $i \in [M]$:

- 1. K_0 , the number of burn-in iterations,
- 2. t, the temperature parameter,
 - 3. C, the number of parallel Markov chains.

Metric. The metric used is the squared Euclidean distance to the exact expectation, averaged on the M instances

$$\frac{1}{M} \sum_{i=1}^{M} ||\mathbb{E}_{\pi_{\boldsymbol{\Theta}_{i}}, t}[Y] - \tilde{\mathbb{E}}(\boldsymbol{\Theta}_{i}, T)||_{2}^{2},$$

which we measure for $T \in [K_0 + 1, K]$.

Polytopes. For the hypercube \mathcal{Y}^d and its neighborhood system \mathcal{N}^r_{κ} , we use d=10 and r=1, which gives $|\mathcal{Y}^d|=2^{10}$ and $|\mathcal{N}^r_{\leq}(\boldsymbol{y})|=10$. For the top- κ polytope \mathcal{Y}^d_{κ} and its neighborhood system \mathcal{N}^s , we use d=10, $\kappa=3$ and s=1, which gives $|\mathcal{Y}^d_{\kappa}|=120$ and $|\mathcal{N}^s(\boldsymbol{y})|=30$. We also use a larger scale for both polytopes in order to highlight the varying impact of the temperature t depending on the combinatorial size of the problem, in the second experiment. For the large scale, we use d=1000 and r=10 for the hypercube, which give $|\mathcal{Y}^d|=2^{1000}\approx 10^{301}$ and $|\mathcal{N}^r_{\leq}(\boldsymbol{y})|\approx 2.7\times 10^{23}$, and we use d=1000, $\kappa=50$ and s=10 for the top- κ polytope, which give $|\mathcal{Y}^d|\approx 9.5\times 10^{84}$ and $|\mathcal{N}^s(\boldsymbol{y})|\approx 1.6\times 10^{33}$.

Hyperparameters. For each experiment, we use K=3000. We average over M=1000 problem instances for statistical significance. We use $K_0=0$, except for the first experiment, where it varies. We use a final temperature t=1, except for the second experiment, where it varies. We use an initial temperature t=1 (leading to a constant temperature schedule), except for the first experiment, where it depends on K_0 . We use only one Markov chain and thus have C=1, except for the third experiment, where it varies.

- (1) Impact of burn-in. First, we evaluate the impact of K_0 , the number of burn-in iterations. We use a truncated geometric cooling schedule $t_k = \max(\gamma^k \cdot t_0, t)$ with $\gamma = 0.995$. The initial temperature t_0 is set to $1/(\gamma^{K_0})$, so that $\forall k \geq K_0 + 1, \ t_k = t = 1$. The results are gathered in Fig. 3.
- 464 **(2) Impact of temperature.** We then evaluate the impact of the final temperature t on the difficulty 465 of the estimation task (different temperatures lead to different target expectations). The results for the 466 small scale are gathered in Fig. 4, and the results for the large scale are gathered in Fig. 5.
- (3) Impact of the number of parallel Markov chains. Finally, we evaluate the impact of the number of parallel Markov chains C on the estimation. The results are gathered in Fig. 6.

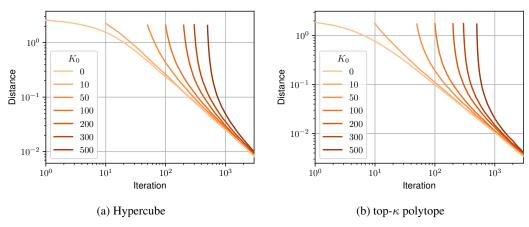


Figure 3: Convergence to exact expectation on the hypercube and the top- κ polytope, for varying number of burn-in iterations K_0 . We conclude that burn-in is not beneficial to the estimation, and taking $K_0 = 0$ is the best option.

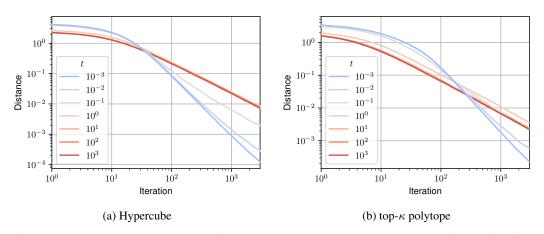


Figure 4: Convergence to exact expectation on the hypercube and the top- κ polytope, for varying final temperature t (small scale experiment). We conclude that lower temperatures facilitate the estimation.

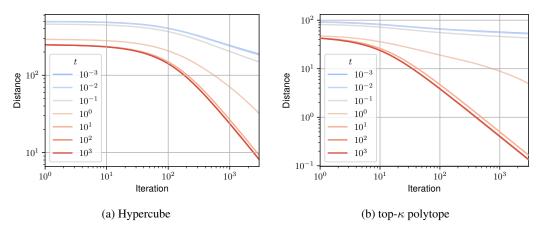


Figure 5: Convergence to exact expectation on the hypercube and the top- κ polytope, for varying final temperature t (large scale experiment). Contrary to the small scale case, larger temperatures are beneficial to the estimation when the solution set is combinatorially large.

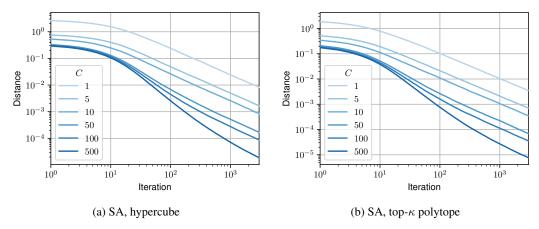


Figure 6: Convergence to exact expectation on the hypercube and the top- κ polytope, for varying number of parallel Markov chains C. Running 10 times more chains in parallel provides roughly the same benefit as extending each chain by 10 times more iterations, highlighting the advantage of massively parallelized estimation.

469 B.4 Convergence to exact parameters

In this section, we conduct experiments in the unsupervised setting described in Section 3.3. As a reminder, the empirical L_N and population L_{θ_0} Fenchel-Young losses are given by:

$$L_{N}(\boldsymbol{\theta}; \boldsymbol{y}_{1}, \dots, \boldsymbol{y}_{N}) \coloneqq \frac{1}{N} \sum_{i=1}^{N} \ell_{t} (\boldsymbol{\theta}; \boldsymbol{y}_{i})$$

$$= A_{t}(\boldsymbol{\theta}) + \frac{1}{N} \sum_{i=1}^{N} \Omega_{t}(\boldsymbol{y}_{i}) - \langle \boldsymbol{\theta}, \bar{Y}_{N} \rangle$$

$$= \ell_{t}(\boldsymbol{\theta}; \bar{Y}_{N}) + C_{1}(Y), \tag{8}$$

472 and

$$L_{\boldsymbol{\theta}_{0}}(\boldsymbol{\theta}) := \mathbb{E}_{(\boldsymbol{y}_{i})_{i=1}^{N} \sim (\pi_{\boldsymbol{\theta}_{0},t})^{\otimes N}} [L_{N}(\boldsymbol{\theta}; \boldsymbol{y}_{1}, \dots, \boldsymbol{y}_{N})]$$

$$= A_{t}(\boldsymbol{\theta}) + \mathbb{E}_{\pi_{\boldsymbol{\theta}_{0},t}} [\Omega_{t}(Y)] - \langle \boldsymbol{\theta}, \, \widehat{\boldsymbol{y}}_{t}(\boldsymbol{\theta}_{0}) \rangle$$

$$= \ell_{t}(\boldsymbol{\theta}; \, \widehat{\boldsymbol{y}}_{t}(\boldsymbol{\theta}_{0})) + C_{2}(\boldsymbol{\theta}_{0}), \tag{9}$$

where the constants $C_1(Y) = \frac{1}{N} \sum_{i=1}^{N} \Omega_t(\boldsymbol{y}_i) - \Omega_t(\bar{Y}_N)$ and $C_2(\boldsymbol{\theta}_0) = \mathbb{E}_{\pi_{\boldsymbol{\theta}_0,t}} \left[\Omega_t(Y)\right] - \Omega_t\left(\widehat{\boldsymbol{y}}_t(\boldsymbol{\theta}_0)\right)$ do not depend on $\boldsymbol{\theta}$. As Jensen gaps, they are non-negative by convexity of Ω_t .

2D visualization. As an introductory example, we display stochastic gradient trajectories in Fig. 7. The parameter $\theta \in \mathbb{R}^d$ is updated following Eq. (7) to minimize the population loss L_{θ_0} defined in Eq. (9), with $\theta_0 = (1/2, 1/2)$. The polytope used is the 2-dimensional hypercube \mathcal{Y}^2 , with neighborhood graph \mathcal{N}_1 (neighbors are adjacent vertices of the square). We present trajectories obtained using MCMC-sampled gradients, comparing results from both 1 and 100 Markov chain iterations with Algorithm 1. For comparison, we include trajectories obtained using Monte Carlosampled (i.e., unbiased) gradients, using 1 and 100 samples.

General setup. We proceed by first randomly generating true parameters $\Theta_0 \in \mathbb{R}^{M \times d}$, with M being a number of problem instances we average on (in order to reduce noise in our observations), by sampling $\Theta_{i,j} \sim \mathcal{N}(0,1)$ independently. The goal is to learn each parameter vector $(\Theta_0)_i \in \mathbb{R}^d$, $i \in [M]$, as M independent problems. The model is randomly initialized at $\hat{\Theta}_0$, and updated with Adam [25] to minimize the loss. In order to better separate noise due to the optimization process and noise due to the sampling process, we use the population loss $L_{(\Theta_0)_i}$ for general experiments, and use the empirical loss L_N only when focusing on the impact of the dataset size N. In this case,

Simulated Annealing (Blue) vs. Monte-Carlo (Red) SA-1 SA-100 MC-1 MC-100 -1 -2 -3

Figure 7: Comparison of stochastic gradient trajectories for a SA / M-H oracle on \mathcal{Y}^2 and unbiased stochastic gradients obtained via Monte Carlo sampling. Increasing the number of Markov chain iterations yields smoother trajectories, similar to the effect of using more Monte Carlo samples in the case of perturbation-based methods [6].

0

-1

we create a dataset $Y \in \mathbb{R}^{M \times N \times d}$, with N being the number of samples, by sampling independently $Y_{i,j} \sim \pi_{(\Theta_0)_i}, \ \forall i \in [M], \ \forall j \in [N].$

We study the impact of the following hyperparameters on learning:

-2

1. K, the number of Markov chain iterations,

492

493

495

505

-3

- 2. C, the number of parallel Markov chains,
- 3. the initialization method used for the chains (either random, persistent, or data-based),
- 4. N, the number of samples in the dataset.

Metrics. The first metric used is the objective function actually minimized, i.e., the population loss, averaged on the M instances:

$$\frac{1}{M} \sum_{i=1}^{M} L_{(\boldsymbol{\Theta}_0)_i}((\hat{\boldsymbol{\Theta}}_n)_i),$$

where $(\hat{\Theta}_n)_i$ is the *n*-th iterate of the optimization process for the problem instance $i \in [M]$. We measure this loss for $n \in [n_{\text{max}}]$, with n_{max} the total number of gradient iterations. For the fourth experiment, where we evaluate the impact of the number of samples N, we measure instead the empirical Fenchel-Young loss:

$$\frac{1}{M}\sum_{i=1}^{M}L_{N}((\hat{\boldsymbol{\Theta}}_{n})_{i};Y_{i,1},\ldots Y_{i,N})$$

In both cases, the best loss value that can be reached is positive but cannot be computed: it corresponds to the constants C_1 and C_2 in Eq. (8) and Eq. (9). Thus, we also provide "stretched" figures, where we plot the loss minus the best loss found during the optimization process.

The second metric used is the squared euclidean distance of the estimate to the true parameter, also averaged on the M instances:

$$\frac{1}{M} \sum_{i=1}^{M} ||(\mathbf{\Theta}_0)_i - (\hat{\mathbf{\Theta}}_n)_i||_2^2.$$

As the top- κ polytope is of dimension d-1, the model is only specified up to vectors orthogonal to the direction of the smallest affine subspace it spans. Thus, in this case, we measure instead:

$$\frac{1}{M} \sum_{i=1}^{M} || \mathsf{P}_{D}^{\perp} ((\Theta_{0})_{i}) - \mathsf{P}_{D}^{\perp} ((\hat{\Theta}_{n})_{i}) ||_{2}^{2},$$

- where $\mathsf{P}^{\perp}{}_D$ is the orthogonal projector on the hyperplane $D=\{\boldsymbol{x}\in\mathbb{R}^d:\langle\mathbf{1},\,\boldsymbol{x}\rangle=0\}$, which is the corresponding direction.
- Polytopes. For the hypercube \mathcal{Y}^d and its neighborhood system \mathcal{N}^r_{\leq} , we use d=10 and r=1, except in the fifth experiment, where we use a mixture of $\mathcal{N}^r_{=}$ neighborhoods (detailed below). For the top- κ polytope \mathcal{Y}^d_{κ} and its neighborhood system \mathcal{N}^s , we use d=10, $\kappa=3$ and s=1.
- **Hyperparameters.** For each experiment, we perform 1000 gradient steps. We use $K_0 = 0$, final 515 temperature t = 1 and initial temperature $t_0 = t = 1$ (leading to a constant temperature schedule). 516 We use K = 1000 Markov chain iterations, except in the first experiment, where it varies. We use 517 only one Markov chain and thus have C=1, except for the second experiment, where it varies. We 518 use a persistent initialization method for the Markov chains, except in the third experiment, where we 519 compare the three different methods. For statistical significance, we average over M=100 problem 520 instances for each experiment, except in the third experiment, where we use M=1000. We work in 521 the limit case $N \to \infty$, except in the fourth experiment, where N varies. 522
- (1) Impact of the length of Markov chains. First, we evaluate the impact of K, the number of inner iterations, i.e., the length of each Markov chain. The results are gathered in Fig. 8.
- 525 (2) Impact of the number of parallel Markov chains. We now evaluate the impact of the number of Markov chains C run in parallel to perform each gradient estimation on the learning process. The results are gathered in Fig. 9.
- (3) Impact of the initialization method. Then, we evaluate the impact of the method to initialize each Markov chain used for gradient estimation. The persistent method consists in setting $\boldsymbol{y}^{(n+1,0)} = \boldsymbol{y}^{(n,K)}$, the data-based method consists in setting $\boldsymbol{y}^{(n+1,0)} = \boldsymbol{y}_i$ with $i \sim \mathcal{U}([N])$, and the random method consists in setting $\boldsymbol{y}^{(n+1,0)} \sim \mathcal{U}(\mathcal{Y})$ (see Appendix C.5 and Table 3 for a detailed explanation). The results are gathered in Fig. 10.
- 533 **(4) Impact of the dataset size.** We now evaluate the impact of the number of samples N from 534 π_{θ_0} (i.e., the size of the dataset $(y_i)_{i=1}^N$) on the estimation of the true parameter θ_0 . The results are 535 gathered in Fig. 11.
- 536 **(5) Impact of neigborhood mixtures.** Finally, we evaluate the impact of the use of neighborhood mixtures. To do so, we use mixtures $\{\mathcal{N}_{=}^{r_s}\}_{s=1}^{S}$, once with $\{r_s\}_{s=1}^{S}=\{5\}$ opposed to $\{r_s\}_{s=1}^{S}=\{5\}$ opposed to $\{r_s\}_{s=1}^{S}=\{5\}$ once with $\{r_s\}_{s=1}^{S}=\{5\}$ opposed to $\{r_s\}_{s=1}^{S}=\{5\}$

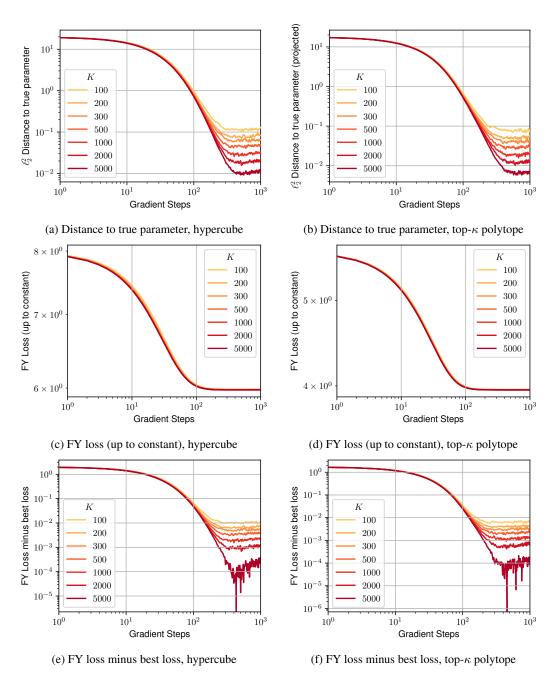


Figure 8: Convergence to the true parameter on the hypercube (left) and the top- κ polytope (right), for varying number of Markov chain iterations K. Longer chains improve learning.

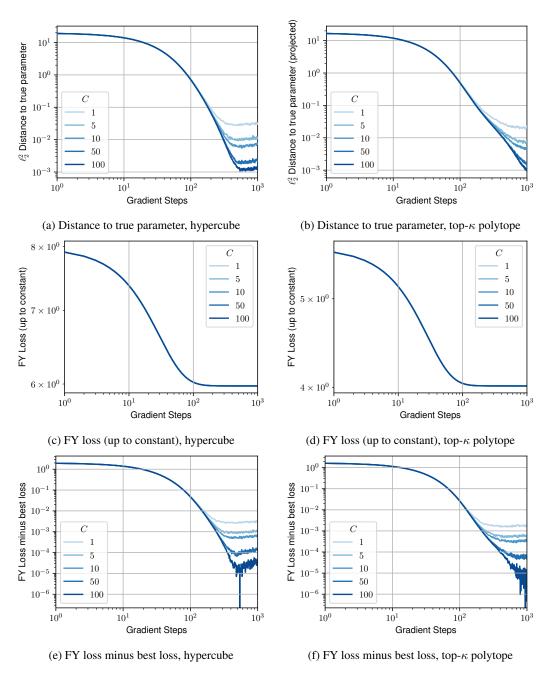


Figure 9: Convergence to the true parameter on the hypercube (left) and the top- κ polytope (right), for varying number of parallel Markov chains C. Adding Markov chains improves estimation.

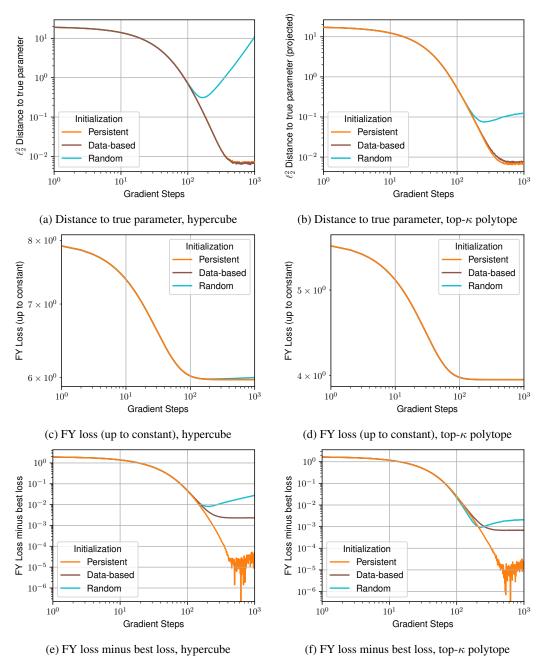


Figure 10: Convergence to the true parameter on the hypercube (left) and the top- κ polytope (right), for varying Markov chain initialization method. The persistent and data-based initialization methods significantly outperform the random initialization method.

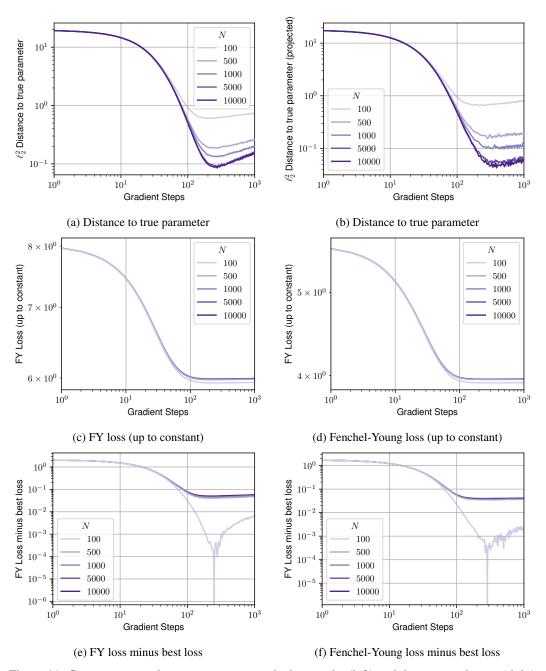


Figure 11: Convergence to the true parameter on the hypercube (left) and the top- κ polytope (right), for varying number of samples N in the dataset. As the dataset is different for each task, the empirical Fenchel-Young loss L_N , which is the minimized objective function (contrary to other experiments, where we minimize L_{θ_0}), also varies. Although empirical Fenchel-Young losses associated to smaller datasets appear easier to minimize, increasing the dataset size reduces the bias and thus the distance to θ_0 , as expected.

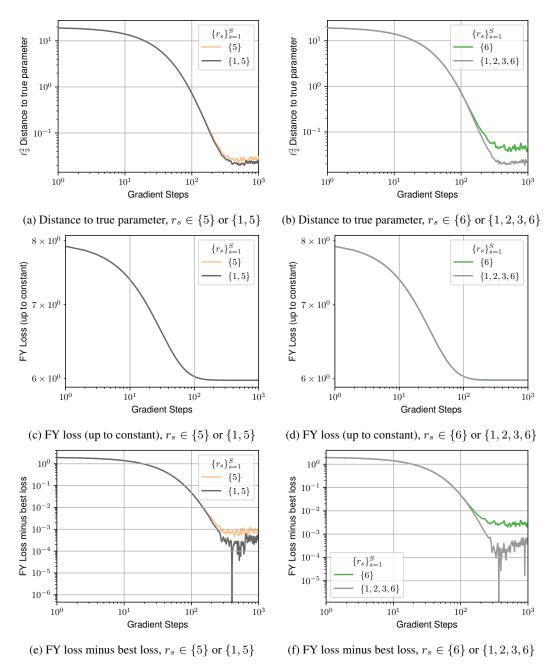


Figure 12: Convergence to the true parameter on the hypercube, with different mixtures of neighborhood systems $\{\mathcal{N}_{=}^{r_s}\}_{s=1}^S$: comparing $r_s \in \{5\}$ to $r_s \in \{1,5\}$ (left), and comparing $r_s \in \{6\}$ to $r_s \in \{1,2,3,6\}$ (right). Using more neighborhoods in the mixture improves learning.

Additional material

553

554

555

556

557

Mixing neighborhood systems: a discussion

In this section, we give intuition on why the update proposed in Section 2.2 and Algorithm 2 is crucial 543 as a tractable way to mix different neighborhood systems.

A naive way to combine these neighborhoods $(\mathcal{N}_s)_{s=1}^S$ and proposals $(q_s)_{s=1}^S$ would be to use Algorithm 1 by defining a unique aggregated proposal $q(\boldsymbol{y},\,\cdot)$ with support $\mathcal{N}(\boldsymbol{y})$ or $\bar{\mathcal{N}}(\boldsymbol{y})\cup\{\boldsymbol{y}\}$ as,

$$q(\boldsymbol{y}, \boldsymbol{y}') \coloneqq \frac{1}{|Q(\boldsymbol{y})|} \sum_{s \in Q(\boldsymbol{y})} q_s(\boldsymbol{y}, \boldsymbol{y}').$$

However, this would lead to non-tractable updates because of the computation of the Metropolis-Hastings correction ratio. Indeed, the latter would be equal to:

$$\alpha(\boldsymbol{y}, \boldsymbol{y}') = \frac{|Q(\boldsymbol{y})|}{|Q(\boldsymbol{y}')|} \cdot \frac{\sum_{s \in Q(\boldsymbol{y}')} q_s(\boldsymbol{y}', \boldsymbol{y})}{\sum_{s \in Q(\boldsymbol{y})} q_s(\boldsymbol{y}, \boldsymbol{y}')}.$$

This calculation is prohibitively expensive because it involves summing the forward proposal probabilities for all move types in Q(y) and the reverse probabilities for all move types in Q(y'). The main difficulty is that multiple, distinct proposal types can generate the same solution y' from y. For 547 example, in the vehicle routing application presented in Section 4, relocating a pair of clients (using 548 the relocate pair move from Table 4) before the first one in a route of 3 gives the same solution y'549 as relocating the first client (with the relocate move) at the last position. Identifying and calculating 550 all these potential forward and reverse pathways for every step is a significant computational hurdle. 551 In contrast, the update we propose in Algorithm 2 only requires computing the single individual ratio 552 $\frac{q_s({m y}',{m y})}{q_s({m y},{m y}')}$ for the unique move type s that was actually sampled.

C.2 Associated Fenchel-Young loss with a single MCMC iteration

To obtain an unbiased gradient estimator for the Fenchel-Young loss ℓ_t associated with \hat{y}_t , the MCMC sampler must be run until it reaches its stationary distribution $\pi_{\theta,t}$. This requirement makes any practical estimator with a finite number of steps K inherently biased.

Although our convergence analysis in Section 3.3 shows that this bias does not hinder the convergence 558 of the proposed learning algorithms, we now demonstrate that when a single MCMC iteration is used (K=1), there exists another target-dependent Fenchel-Young loss such that the stochastic gradient estimator is *unbiased* with respect to that loss. See Appendix E.7 for the construction of Ω_u and the 561 proof. 562

Proposition C.1 (Existence of a Fenchel-Young loss when K=1). Let $p_{\theta,y}^{(1)}$ denote the distribution of the first iterate of the Markov chain defined by the Markov transition kernel given in Eq. (3), with proposal distribution q and initialized at ground-truth $m{y} \in \mathcal{Y}$. There exists a target-dependent regularization function $\Omega_{\boldsymbol{y}}$ with the following properties: $\Omega_{\boldsymbol{y}}$ is $t/\mathbb{E}_{q(\boldsymbol{y},\cdot)}||Y-\boldsymbol{y}||_2^2$ -strongly convex; it is such that

$$\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}\left[Y\right] = \operatorname*{argmax}_{\boldsymbol{\mu} \in \mathsf{conv}(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\})} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_{\boldsymbol{y}}(\boldsymbol{\mu}) \right\}$$

and the Fenchel-Young loss ℓ_{Ω_y} generated by Ω_y is $\mathbb{E}_{q(y,\cdot)}||Y-y||_2^2/t$ -smooth in its first argument, 563 and such that $\nabla_{\boldsymbol{\theta}} \ell_{\Omega_{\boldsymbol{y}}}(\boldsymbol{\theta}; \boldsymbol{y}) = \mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta}, \boldsymbol{y}}^{(1)}}[Y] - \boldsymbol{y}.$ 564

A similar result in the unsupervised setting with data-based initialization is given in Proposition C.2. 565 Interestingly, theses results contrast with prior work on the expected CD-1 update. Indeed, when 566 applied with Gibbs sampling to train restricted Boltzmann machines, the latter was shown in Sutskever 567 and Tieleman [47] not to be the gradient of any function – let alone a convex one. 568

Note that, similarly to the regularization Ω_t , the target-dependent Ω_y extends the influence of φ from 569 the set $\mathcal{N}(y) \cup \{y\}$ to its convex hull in a principled way. As a verification, we give properties of the regularized maximizer $\mathbb{E}_{p_{\theta,y}^{(1)}}[Y]$ in Proposition C.3.

Fig. C.3 Fenchel-Young loss for K=1 in the unsupervised setting

This proposition is analogous to Proposition C.1, but in the unsupervised setting, when using a data-based initialization method – i.e., the original CD initialization scheme, without persistent Markov chains. See Appendix C.5 for a detailed discussion about this.

Proposition C.2. Let $p_{\theta, \bar{Y}_N}^{(1)}$ denote the distribution of the first iterate of the Markov chain defined by the Markov transition kernel given in Eq. (3), with proposal distribution q and initialized by $\mathbf{y}^{(0)} = \mathbf{y}_i$, with $i \sim \mathcal{U}(\llbracket 1, N \rrbracket)$. There exists a dataset-dependent regularization $\Omega_{\bar{Y}_N}$ with the following properties: $\Omega_{\bar{Y}_N}$ is $tN/\sum_{i=1}^N \mathbb{E}_{q(\mathbf{y}_i,\cdot)}||Y-\mathbf{y}_i||_2^2$ -strongly convex; it is such that:

$$\mathbb{E}_{p_{\theta,\bar{Y}_{N}}^{(1)}}\left[Y\right] = \operatorname*{argmax}_{\boldsymbol{\mu} \in \operatorname{conv}\left(\bigcup_{i=1}^{N} \left\{\mathcal{N}(y_{i}) \cup \left\{y_{i}\right\}\right\}\right)} \left\{\left\langle\boldsymbol{\theta},\boldsymbol{\mu}\right\rangle - \Omega_{\bar{Y}_{N}}(\boldsymbol{\mu})\right\};$$

and the Fenchel-Young loss $L_{\Omega_{\widetilde{Y}_N}}$ generated by $\Omega_{\widetilde{Y}_N}$ is $\frac{1}{N}\sum_{i=1}^N \mathbb{E}_{q(\boldsymbol{y}_i,\cdot)}||Y-\boldsymbol{y}_i||_2^2/t$ -smooth in its first argument, and such that $\nabla_{\boldsymbol{\theta}} L_{\Omega_{\widetilde{Y}_N}}(\boldsymbol{\theta}\,;\boldsymbol{y}) = \mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\widetilde{Y}_N}^{(1)}}[Y]-\boldsymbol{y}.$

The proof is given in Appendix E.7.

779 C.4 Properties of the expected first iterate

Proposition C.3. Let $\theta \in \mathbb{R}^d$, $y \in \mathcal{Y}$. Let

$$\mathcal{N}_{better}(\boldsymbol{y}) \coloneqq \{ \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}) \mid \langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}') > \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}) \}$$

denote the set of improving neighbors of y for the unregularized objective function. We have the following properties:

$$\begin{split} & \mathbb{E}_{p_{\theta, \boldsymbol{y}}^{(1)}}[Y] \in \mathsf{conv}\left(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\}\right), \\ & \mathbb{E}_{p_{\theta, \boldsymbol{y}}^{(1)}}[Y] \xrightarrow[t \to 0^+]{} \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}_{better}(\boldsymbol{y})} q(\boldsymbol{y}, \boldsymbol{y}') \cdot (\boldsymbol{y}' - \boldsymbol{y}), \\ & \text{and} \quad \mathbb{E}_{p_{\theta, \boldsymbol{y}}^{(1)}}[Y] \xrightarrow[t \to +\infty]{} \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} \min\left[q(\boldsymbol{y}, \boldsymbol{y}'), q(\boldsymbol{y}', \boldsymbol{y})\right] \cdot (\boldsymbol{y}' - \boldsymbol{y}). \end{split}$$

The proof is given in Appendix E.8. Thus, as the set $\mathcal{N}_{\text{better}}$ is defined according the value of the original, unregularized objective function $\boldsymbol{y} \mapsto \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})$, the low temperature behavior of the regularized maximizer $\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y]$ effectively reflects the fact that the regularization function $\Omega_{\boldsymbol{y}}$ extends the influence of φ from the vertices $\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\}$ to their convex hull.

586 C.5 Markov chain initialization

In contrastive divergence (CD) learning, the intractable expectation in the log-likelihood gradient is approximated by short-run MCMC, initialized at the data distribution [22] (using a Gibbs sampler in the setting of Restricted Boltzmann Machines).

Here, we note, at the n-th iteration of gradient descent:

$$\nabla_W L_N(\hat{W}_n) \approx \frac{1}{|B_n|} \sum_{i \in B_n} J_W g_{\hat{W}_n}(\boldsymbol{x}_i) \left(\frac{1}{K} \sum_{k=1}^K \boldsymbol{y}_i^{(n+1,k)} - \boldsymbol{y}_i \right),$$

for the supervised setting, with B_n being the mini-batch (or full batch) used at iteration n, y_i the ground-truth structure associated to x_i in the dataset, and $y_i^{(n+1,k)}$ the k-th iterate of Algorithm 1, with maximization direction $g_{\hat{W}_n}(x_i)$, and initialization point $y_i^{(n+1,0)}$. We also note:

$$\nabla_{\boldsymbol{\theta}} L_N(\hat{\boldsymbol{\theta}}_n) \approx \frac{1}{K} \sum_{k=1}^K \boldsymbol{y}^{(n+1,\,k)} - \bar{Y}_N$$

for the unsupervised setting, with $y^{(n+1,k)}$ being the k-th iterate of Algorithm 1, with maximization direction $\hat{\theta}_n$, and initialization point $y^{(n+1,0)}$.

In CD learning of unconditional EBMs (i.e., in our unsupervised setting), the Markov Chain is initialized at the empirical data distribution [22, 11], as explained earlier. Persistent Contrastive Divergence (PCD) learning [48] modifies CD by maintaining a persistent Markov chain. Thus, instead of initializing the chain from the data distribution in each iteration, the chain continues from its last state in the previous iteration, by setting $\boldsymbol{y}^{(n+1,0)} = \boldsymbol{y}^{(n,K)}$. This approach aims to provide a better approximation of the model distribution and to reduce the bias introduced by the initialization of the Markov chain in CD. These are two types of informative initialization methods, which aim at reducing the mixing times of the Markov Chains.

However, neither of these can be applied to the supervised (or conditional) setting, as observed in [38] in the context of conditional Restricted Boltzmann Machines (which are a type of EBMs). Indeed, on the one hand, PCD takes advantage of the fact that the parameter $\hat{\theta}$ does not change too much from one iteration to the next, so that a Markov Chain that has reached equilibrium on $\hat{\theta}_n$ is not far from equilibrium on $\hat{\theta}_{n+1}$. This does not hold in the supervised setting, as each x_i leads to a different $\hat{\theta}_i = g_{\hat{W}}(x_i)$. On the other hand, the data-based initialization method in CD would amount to initialize the chains at the empirical marginal data distribution on \mathcal{Y} , and would be irrelevant in a supervised setting, since the distribution we want each Markov Chain to approximate is conditioned on the input x_i .

An option is to use persistent chains if training for multiple epochs, and to initialize the Markov Chain associated to (x_i, y_i) for epoch j at the final state of the one associated to the same data point (x_i, y_i) at epoch j-1. However, this method is relevant than PCD in the unsupervised setting, as \hat{w} changes a lot more in a full epoch than $\hat{\theta}$ in just one gradient step in the unsupervised setting. It might be relevant, however, if each epoch consists in a single, full-batch gradient step. Nevertheless, it would require to store a significant number of states $y_i^{(n,K)}$ (one for each point in the dataset). The solution we propose, for both full-batch and mini-batch settings, is to initialize the chains at the empirical data distribution conditioned on the input x_i , which amounts to initialize them at the ground-truth y_i .

This discussion is summed up in Table 3.

Table 3: Possible Markov Chain Initialization Methods under each Learning Setting

Setting Init. Method	Unsupervised	Supervised, Batch	Supervised, Mini-Batch
Persistent	${m y}^{(n+1,0)} = {m y}^{(n,K)}$	$oldsymbol{y}_i^{(n+1,0)} = oldsymbol{y}_i^{(n,K)}$	/
Data-Based	$oldsymbol{y}^{(n+1,0)} = oldsymbol{y}_j, ext{ with } \ j \sim \mathcal{U}(\llbracket 1, N rbracket)$	$oldsymbol{y}_i^{(n+1,0)} = oldsymbol{y}_i$	$oldsymbol{y}_i^{(n+1,0)} = oldsymbol{y}_i$
Random	$oldsymbol{y}^{(n+1,0)} \sim \mathcal{U}(\mathcal{Y})$	$oldsymbol{y}_i^{(n+1,0)} \sim \mathcal{U}(\mathcal{Y})$	$oldsymbol{y}_i^{(n+1,0)} \sim \mathcal{U}(\mathcal{Y})$

Remark C.4. The use of uniform distributions on \mathcal{Y} for the random initialization method can naturally be replaced by any other different prior distribution.

C.6 Proposal distribution design for the DVRPTW

Original deterministic moves. The selected moves, designed for Local Search algorithms on vehicle routing problems (specifically for the PC-VRPTW for serve request and remove request), are given in Table 4.

Name	Description
relocate	removes request i from its route and re-inserts it before or after request j
relocate pair	removes pair of requests $(i, next(i))$ from their route and re-inserts them before or after request j
swap	exchanges the position of requests i and j in the solution
swap pair	exchanges the positions of the pairs $(i, next(i))$ and $(j, next(j))$ in the solution
2-opt	reverses the route segment between i and j
serve request	inserts currently undispatched request i before or after request j
remove request	removes currently dispatched request i from the solution

Table 4: PC-VRPTW Local search moves

Move	$V_s^1(oldsymbol{y})$	$V_s^2(oldsymbol{y})[i]$
relocate	$\mathcal{D}(oldsymbol{y}) \setminus \mathcal{D}_1(oldsymbol{y})$	$\mathcal{D}(oldsymbol{y})$
relocate pair	$\mathcal{D}(oldsymbol{y}) \setminus \left\{ \mathcal{D}_2(oldsymbol{y}) \cup \mathcal{D}^{ ext{last}}(oldsymbol{y}) ight\}$	$\mathcal{D}(oldsymbol{y})\setminus\{next(i)\}$
swap	$\mathcal{D}(oldsymbol{y})$	$\mathcal{D}(oldsymbol{y})$
swap pair	$\mathcal{D}(oldsymbol{y}) \setminus \mathcal{D}^{ ext{last}}(oldsymbol{y})$	$\mathcal{D}(oldsymbol{y}) \setminus ig\{\mathcal{D}^{ ext{last}}(oldsymbol{y}) \cup \{prev(i), next(i)\}ig\}$
2-opt	$\mathcal{D}(oldsymbol{y}) \setminus \mathcal{D}_2(oldsymbol{y})$	$\mathcal{D}(oldsymbol{y}) \setminus \hat{\mathcal{D}}_2(oldsymbol{y})$
serve request	$\overline{\mathcal{D}}(oldsymbol{y})$	$\mathcal{D}(oldsymbol{y}) \cup \mathcal{I}_D(oldsymbol{y})$
remove request	$ig\{\mathcal{D}(oldsymbol{y})\setminus\mathcal{D}_1(oldsymbol{y})ig\}\cup\mathcal{I}_1(oldsymbol{y})$	

Table 5: Sets of valid clients for each move. $\mathcal{D}(\boldsymbol{y})$ contains all dispatched clients in solution \boldsymbol{y} . $\mathcal{D}_1(\boldsymbol{y})$ contains all dispatched clients that are the only client in their route. $\mathcal{D}_2(\boldsymbol{y})$ contains all dispatched clients that are in a route with 2 clients or less. $\mathcal{D}^{\text{last}}(\boldsymbol{y})$ contains all dispatched clients that are the last of their route. $\overline{\mathcal{D}}(\boldsymbol{y})$ contains all non-dispatched clients. $\mathcal{I}_D(\boldsymbol{y})$ contains the depot of the first empty route, if it exists (all routes may be non-empty), or else is the empty set. $\mathcal{I}_1(\boldsymbol{y})$ contains the only client in the last non-empty route if it contains exactly one client, or else is the empty set.

All of these moves (except for remove request) involve selecting two clients i and j from the request set \mathcal{R}^{ω} (for example, the relocate move relocates client i after client j in the solution).

In the Local Search part of the PC-HGS algorithm from Vidal [49], they are implemented as deterministic functions used within a quadratic loop over clients, and are performed only if they improve the solution's objective value. The search is narrowed down to client pairs (i,j) such that d(i,j) is among the N_{prox} lowest values in $\left\{d(i,k) \mid k \in \mathcal{R}^{\omega} \setminus \{D,i\}\right\}$, where d is a problem-specific heuristic distance measure between clients, based on spatial features and time windows, and N_{prox} is a hyperparameter. These distances are independent from the chosen solution routes (they are computed once at the start of the algorithm, from the problem features), non-negative, and symmetric: d(i,j) = d(j,i).

Randomization. In order to transform these deterministic moves into proposals, we first adapt the choice of clients i and j, by sampling i uniformly from $V^1_s(\boldsymbol{y})$, which contains the set of valid choices of client i for move s from solution \boldsymbol{y} . Then, we sample j from $V^2_s(\boldsymbol{y})[i] \setminus \{i\}$ using the following softmax distribution: $P_s(j \mid i) = \frac{\exp[-d(i,j)/\beta]}{\sum_{k \in V^2_s(\boldsymbol{y})[i] \setminus \{i\}} \exp[-d(i,k)/\beta]}$, where $\beta > 0$ is a neighborhood sampling temperature. The set $V^2_s(\boldsymbol{y})[i]$ contains all valid choices of client j for move s from solution \boldsymbol{y} , and is precised along with $V^1_s(\boldsymbol{y})$ in Table 5. We normalize the distance measures inside the softmax, by dividing them by the maximum distance: $d(i,\cdot) \leftarrow d(i,\cdot)/\max_{k \in V^2(\boldsymbol{y})[i] \setminus \{i\}} d(i,k)$.

Neighborhood graph symmetrization. Then, we ensure that each individual neighborhood graph \mathcal{N}_s is undirected. This is already the case for the moves swap, swap pair and 2-opt, as they are actually involutions (applying the same move on the same couple (i,j) from y' will result in y). However, this is obviously not the case for serve request and remove request. Indeed, if solution y' is obtained from y by removing a dispatched client (respectively serving an non-dispatched one), y cannot be obtained by removing another one (respectively, serving another one). To fix this, we merge these two moves into a single one. First, it evaluates which of the two moves are allowed (i.e., if they are such that $V_s^1(y) \neq \emptyset$). Then, it samples one (the probability of selecting "remove" is chosen to be equal to the number of removable clients divided by the number of removable clients plus the number of servable clients) in the case where both are possible, or else simply performs the only move allowed. Thus, the corresponding neighborhood graph is undirected as it is always

possible to perform the reverse operation (as when removing a client, it becomes unserved, thus allowing the serve request move from y', and vice-versa). We also allow the serve request move to insert a client after the depot of the first empty route, to allow the creation of new routes. In consequence, we allow the remove request move to remove the only client in the last non-empty route if it contains exactly one client (to maintain symmetry of the neighborhood graph).

For the relocate and relocate pair moves, the non-reversibility comes from the fact that they only relocate client i (or clients i and next(i) in the pair case) after client j, so that if client i was the first in its route, relocating it back would be impossible (the depot, which is the start of the route, cannot be selected as j). Thus, we allow insertions before clients too, and add a random choice with probability $(\frac{1}{2}, \frac{1}{2})$ to determine if the relocated client(s) will be inserted before or after j. We also add this feature to the serve request move.

Correction ratio computation. Next, we implement the computation of the individual correction ratio $\tilde{\alpha}_s(\boldsymbol{y}, \boldsymbol{y}') = \frac{q_s(\boldsymbol{y}', \boldsymbol{y})}{q_s(\boldsymbol{y}, \boldsymbol{y}')}$ for each proposal q_s .

• In the case of swap and 2-opt, we have $\tilde{\alpha}_s(\boldsymbol{y}, \boldsymbol{y}') = 1$. Indeed, let \boldsymbol{y}' be the result of applying one of these moves s on \boldsymbol{y} when sampling $i \in V_s^1(\boldsymbol{y})$ and $j \in V_s^2(\boldsymbol{y})[i] \setminus \{i\}$. We then have:

$$q_s(\boldsymbol{y}, \boldsymbol{y}') = \frac{1}{|V_s^1(\boldsymbol{y})|} \cdot \frac{\exp\left[-d(i, j)/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y})[i] \setminus \{i\}} \exp\left[-d(i, k)/\beta\right]} + \frac{1}{|V_s^1(\boldsymbol{y})|} \cdot \frac{\exp\left[-d(j, i)/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y})[j] \setminus \{j\}} \exp\left[-d(j, k)/\beta\right]},$$

where the first term accounts for the probability of selecting i then j and the second term accounts for that of selecting j then i (one can easily check that these two cases are the only way of sampling \mathbf{y}' from \mathbf{y}). Then, noticing that we have $|V_s^1(\mathbf{y}')| = |V_s^1(\mathbf{y})|$, that these moves are involutions (selecting (i,j) or (j,i) from \mathbf{y}' is also the only way to sample \mathbf{y}), and that we have the equalities $V_s^2(\mathbf{y})[i] = V_s^2(\mathbf{y}')[i]$ and $V_s^2(\mathbf{y})[j] = V_s^2(\mathbf{y}')[j]$, we actually have $q_s(\mathbf{y}',\mathbf{y}) = q_s(\mathbf{y},\mathbf{y}')$.

• For swap pair, the same arguments hold (leading to the same form for q_s), except for the equalities $V_s^2(\boldsymbol{y})[i] = V_s^2(\boldsymbol{y}')[i]$ and $V_s^2(\boldsymbol{y})[j] = V_s^2(\boldsymbol{y}')[j]$. Thus, we have the following form for the correction ratio:

$$\frac{q_s(\boldsymbol{y}', \boldsymbol{y})}{q_s(\boldsymbol{y}, \boldsymbol{y}')} = \frac{\sum_{k \in V_s^2(\boldsymbol{y})[i] \setminus \{i\}} \exp\left[-d(i, k) / \beta\right] + \sum_{k \in V_s^2(\boldsymbol{y})[j] \setminus \{j\}} \exp\left[-d(j, k) / \beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y}')[i] \setminus \{i\}} \exp\left[-d(i, k) / \beta\right] + \sum_{k \in V_s^2(\boldsymbol{y}')[j] \setminus \{j\}} \exp\left[-d(j, k) / \beta\right]}$$

• In the case of relocate, let j' denote $\operatorname{next}(j)$ if the selected insertion type was "after", and $\operatorname{prev}(j)$ if it was "before" – where $\operatorname{next}(j) \in \mathcal{R}^\omega$ denotes the request following j in solution \boldsymbol{y} , i.e., the only index k such that $\boldsymbol{y}_{j,k} = 1$, and $\operatorname{prev}(j)$ is the one preceding it, i.e., the only k such that $\boldsymbol{y}_{k,j} = 1$. We have:

$$q_s(\boldsymbol{y}, \boldsymbol{y}') = \frac{1}{2} \cdot \frac{1}{|V_s^1(\boldsymbol{y})|} \cdot \frac{\exp\left[-d(i, j)/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y})[i] \setminus \{i\}} \exp\left[-d(i, k)/\beta\right]} + \frac{1}{2} \cdot \frac{1}{|V_s^1(\boldsymbol{y})|} \cdot \frac{\exp\left[-d(i, j')/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y})[i] \setminus \{i\}} \exp\left[-d(i, k)/\beta\right]}$$

Indeed, if i was relocated $after\ j$, the same solution y' could have been obtained by relocating i before $j' = \mathsf{next}(j)$. Similarly, if i was relocated before j, the same solution y' could have been obtained by relocating i after $j' = \mathsf{prev}(j)$. For the reverse move probability, the way of obtaining y from y' is either to select $(i, \mathsf{prev}(i))$ in the after-type insertion case, or $(i, \mathsf{next}(i))$ in the before-type insertion case (where prev and next are taken w.r.t. y, i.e., before applying the move). Thus, we have:

$$\begin{split} q_s(\boldsymbol{y}', \boldsymbol{y}) &= \frac{1}{2} \cdot \frac{1}{|V_s^1(\boldsymbol{y}')|} \cdot \frac{\exp\left[-d(i, \mathsf{prev}(i)/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y}')[i] \backslash \{i\}} \exp\left[-d(i, k)/\beta\right]} \\ &+ \frac{1}{2} \cdot \frac{1}{|V_s^1(\boldsymbol{y}')|} \cdot \frac{\exp\left[-d(i, \mathsf{next}(i))/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y}')[i] \backslash \{i\}} \exp\left[-d(i, k)/\beta\right]}. \end{split}$$

 For the relocate pair move, the exact same reasoning and proposal probability form hold for the forward move, but we have for the reverse direction:

$$\begin{split} q_s(\boldsymbol{y}', \boldsymbol{y}) &= \frac{1}{2} \cdot \frac{1}{|V_s^1(\boldsymbol{y}')|} \cdot \frac{\exp\left[-d(i, \mathsf{prev}(i)/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y}')[i] \backslash \{i\}} \exp\left[-d(i, k)/\beta\right]} \\ &+ \frac{1}{2} \cdot \frac{1}{|V_s^1(\boldsymbol{y}')|} \cdot \frac{\exp\left[-d(i, \mathsf{next}(\mathsf{next}(i)))/\beta\right]}{\sum_{k \in V_s^2(\boldsymbol{y}')[i] \backslash \{i\}} \exp\left[-d(i, k)/\beta\right]}, \end{split}$$

as client next(i) is also relocated.

• For the serve request / remove request move, we have the forward probability:

$$q_{s}(\boldsymbol{y}, \boldsymbol{y}') = \frac{|\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_{1}(\boldsymbol{y})\} \cup \mathcal{I}_{1}(\boldsymbol{y})|}{|\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_{1}(\boldsymbol{y})\} \cup \mathcal{I}_{1}(\boldsymbol{y})| + |\overline{\mathcal{D}}(\boldsymbol{y})|} \times \frac{1}{|\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_{1}(\boldsymbol{y})\} \cup \mathcal{I}_{1}(\boldsymbol{y})|}$$
$$= \frac{1}{|\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_{1}(\boldsymbol{y})\} \cup \mathcal{I}_{1}(\boldsymbol{y})| + |\overline{\mathcal{D}}(\boldsymbol{y})|}$$

if the chosen move is remove request. The expression corresponds to the composition of move choice sampling and uniform sampling over removable clients.

Still in the same case (remove request is chosen) and if the removed request i was in $\mathcal{I}_1(y)$ (i.e., was the only client in the last non-empty route if the latter contained exactly one client), we have the reverse move probability:

$$q_{s}(\mathbf{y}', \mathbf{y}) = \frac{1}{|\{\mathcal{D}(\mathbf{y}') \setminus \mathcal{D}_{1}(\mathbf{y}')\} \cup \mathcal{I}_{1}(\mathbf{y}')| + |\overline{\mathcal{D}}(\mathbf{y}')|} \times \frac{\exp\left[-\overline{d}(i)/\beta\right]}{\exp\left[-\overline{d}(i)/\beta\right] + \sum_{k \in \mathcal{D}(\mathbf{y}')} \exp\left[-d(i, k)/\beta\right]} = \frac{1}{|\{\mathcal{D}(\mathbf{y}) \setminus \mathcal{D}_{1}(\mathbf{y})\} \cup \mathcal{I}_{1}(\mathbf{y})| + |\overline{\mathcal{D}}(\mathbf{y})|} \times \frac{\exp\left[-\overline{d}(i)/\beta\right]}{\exp\left[-\overline{d}(i)/\beta\right] + \sum_{\substack{k \in \mathcal{D}(\mathbf{y}) \\ k \neq i}} \exp\left[-d(i, k)/\beta\right]}.$$

The expression corresponds to the composition of move choice sampling and softmax sampling of the depot of the first empty route (which was the route of the removed client i, so that $\mathcal{I}_D(\boldsymbol{y}') \neq \emptyset$ in this case). We use the average distance to dispatched clients $\bar{d}(i) \coloneqq \frac{1}{|\mathcal{D}(\boldsymbol{y}')|} \sum_{k \in \mathcal{D}(\boldsymbol{y}')} d(i,k)$ as distance to the depot.

In the case where the removed request i was not in $\mathcal{I}_1(y)$, we have instead:

$$\begin{split} q_s(\boldsymbol{y}',\boldsymbol{y}) &= \frac{1}{|\left\{\mathcal{D}(\boldsymbol{y}') \setminus \mathcal{D}_1(\boldsymbol{y}')\right\} \cup \mathcal{I}_1(\boldsymbol{y}')| + |\overline{\mathcal{D}}(\boldsymbol{y}')|} \\ &\times \frac{\frac{1}{2} \cdot \exp\left[-d(i,\operatorname{prev}(i))\right] + \frac{1}{2} \cdot \exp\left[-d(i,\operatorname{next}(i))\right]}{\mathbf{1}_{\left\{\mathcal{I}_D(\boldsymbol{y}') \neq \emptyset\right\}} \cdot \exp\left[-\overline{d}(i)/\beta\right] + \sum_{k \in \mathcal{D}(\boldsymbol{y}')} \exp\left[-d(i,k)/\beta\right]} \\ &= \frac{1}{|\left\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_1(\boldsymbol{y})\right\} \cup \mathcal{I}_1(\boldsymbol{y})| + |\overline{\mathcal{D}}(\boldsymbol{y})|} \\ &\times \frac{\frac{1}{2} \cdot \exp\left[-d(i,\operatorname{prev}(i))\right] + \frac{1}{2} \cdot \exp\left[-d(i,\operatorname{next}(i))\right]}{\mathbf{1}_{\left\{\mathcal{I}_D(\boldsymbol{y}') \neq \emptyset\right\}} \cdot \exp\left[-\overline{d}(i)/\beta\right] + \sum_{k \in \mathcal{D}(\boldsymbol{y})} \exp\left[-d(i,k)/\beta\right]}. \end{split}$$

The right term corresponds to softmax sampling of the previous node with "after" insertion type (which has probability 1/2) and of the next node with "before" insertion type. The non-emptiness of $\mathcal{I}_D(\boldsymbol{y}')$ is not guaranteed anymore, as all routes might be non-empty (indeed, we did not create an empty one by removing i, as $i \in \mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_1(\boldsymbol{y})$ in this case).

Similarly, if the chosen move is serve request, we have the forward probability:

$$q_s(\boldsymbol{y}, \boldsymbol{y}') = \frac{|\overline{\mathcal{D}}(\boldsymbol{y})|}{|\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_1(\boldsymbol{y})\} \cup \mathcal{I}_1(\boldsymbol{y})| + |\overline{\mathcal{D}}(\boldsymbol{y})|} \times \frac{\frac{1}{2} \cdot \exp\left[-d(i, j)\right] + \frac{1}{2} \cdot \exp\left[-d(i, j')\right]}{\mathbf{1}_{\{\mathcal{I}_D(\boldsymbol{y}) \neq \emptyset\}} \cdot \exp\left[-\overline{d}(i)/\beta\right] + \sum_{k \in \mathcal{D}(\boldsymbol{y})} \exp\left[-d(i, k)/\beta\right]}$$

if the selected insertion node j is not in $\mathcal{I}_D(\boldsymbol{y})$ (i.e., is not the depot of the first empty route in \boldsymbol{y}), where j' = prev(j) if the insertion type selected was "before" (which has probability 1/2), and j' = next(j) if it was "after".

We have instead the forward probability:

714

718

719

720

730

$$q_s(\boldsymbol{y}, \boldsymbol{y}') = \frac{1}{|\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_1(\boldsymbol{y})\} \cup \mathcal{I}_1(\boldsymbol{y})| + |\overline{\mathcal{D}}(\boldsymbol{y})|} \times \frac{\exp\left[-\overline{d}(i)/\beta\right]}{\exp\left[-\overline{d}(i)/\beta\right] + \sum_{k \in \mathcal{D}(\boldsymbol{y})} \exp\left[-d(i, k)/\beta\right]}$$

if the selected insertion node j is in $\mathcal{I}_D(y)$ (i.e., is the depot of the first empty route in y). In every case, we have the reverse move probability:

$$q_s(\boldsymbol{y}',\boldsymbol{y}) = \frac{1}{|\{\mathcal{D}(\boldsymbol{y}) \setminus \mathcal{D}_1(\boldsymbol{y})\} \cup \mathcal{I}_1(\boldsymbol{y})| + |\overline{\mathcal{D}}(\boldsymbol{y})|}.$$

In each case, we set $d(i,D)=+\infty$ to account for the fact that the depot can never be sampled during the process (except in the serve request / remove request move, where we allow the depot of the first empty route / last non-empty route to be selected, for which we use the average distance to other requests as explained earlier) – in fact, the distance measure from a client to the depot is not even defined in the original HGS implementation.

The second correction factor needed is $\frac{|Q(\boldsymbol{y})|}{|Q(\boldsymbol{y}')|}$ (see Algorithm 2). We compute it by checking if each move is allowed, i.e., if there exists at least one $i \in V_s^1(\boldsymbol{y})$ such that $V_s^2(\boldsymbol{y})[i] \setminus \{i\} \neq \emptyset$. This can be determined in $\mathcal{O}(\mathcal{R}^\omega)$ for each move.

729 **D Details on the DVRPTW**

D.1 Overview of the challenge.

We evaluate the proposed approach on a large-scale, ML-enriched combinatorial optimization problem: the EURO Meets NeurIPS 2022 Vehicle Routing Competition [27]. In this dynamic vehicle routing problem with time windows (DVRPTW), requests arrive continuously throughout a planning horizon, which is partitioned into a series of delivery waves $\mathcal{W} = \{ [\tau_0, \tau_1], [\tau_1, \tau_2], \dots, [\tau_{|\mathcal{W}|-1}, \tau_{|\mathcal{W}|}] \}$. At the start of each wave ω , a dispatching and vehicle routing problem must be solved for the set of requests \mathcal{R}^{ω} specific to that wave (in which we include the depot D), encoded into the system state

 x^{ω} . We note $\mathcal{Y}(x^{\omega})$ the set of feasible decisions associated to state x^{ω} . 737 A feasible solution $y^{\omega} \in \mathcal{Y}(x^{\omega})$ must contain all requests that must be dispatched before τ_{ω} (the rest 738 are postponable), allow each of its routes to visit the requests they dispatch within their respective 739 time windows, and be such that the cumulative customer demand on each of its routes does not exceed 740 a given vehicle capacity. It is encoded thanks to a vector $\left(y_{i,j}^{\omega}\right)_{i,j\in\mathcal{R}^{\omega}}$, where $y_{i,j}^{\omega}=1$ if the solution 741 contains the directed route segment from i to j, and $y_{i,j}^{\omega}=0$ otherwise. The set of requests $\mathcal{R}^{\omega+1}$ is 742 obtained by removing all requests dispatched by the chosen solution y^ω from \mathcal{R}^ω and adding all new 743 requests which arrived between τ_{ω} and $\tau_{\omega+1}$. 744

The aim of the challenge is to find an optimal policy $f: \mathcal{X} \to \mathcal{Y}$ assigning decisions $\mathbf{y}^{\omega} \in \mathcal{Y}(\mathbf{x}^{\omega})$ to system states $\mathbf{x}^{\omega} \in \mathcal{X}$. This can be cast as a reinforcement learning problem:

$$\min_{f} \mathbb{E}\left[c_{\mathcal{W}}(f)\right], \quad \text{with} \quad c_{\mathcal{W}}(f) \coloneqq \sum_{\omega \in \mathcal{W}} c(f(\boldsymbol{x}^{\omega})),$$

where $c: \boldsymbol{y}^{\omega} \mapsto \sum_{i,j \in \mathcal{R}^{\omega}} c_{i,j} y_{i,j}^{\omega}$ gives the routing cost of $\boldsymbol{y}^{\omega} \in \mathcal{Y}^{\omega}$ and where $c_{i,j} \geq 0$ is the routing cost from i to j. The expectation is taken over full problem instances.

747 D.2 Reduction to supervised learning.

We follow the method of [4], which was the winning approach for the challenge. Central to this approach is the concept of prize-collecting dynamic vehicle routing problem with time windows (PC-VRPTW). In this setting, each request $i \in \mathcal{R}^{\omega}$ is assigned an artificial $prize \ \theta_i^{\omega} \in \mathbb{R}$, that reflects the benefit of serving it. The prize of the depot D is set to $\theta_D^{\omega} = 0$. The objective is then to identify a set of routes that maximizes the total prize collected while minimizing the associated travel costs. The model g_W predicts the prize vector $\theta^{\omega} = g_W(\boldsymbol{x}^{\omega})$. Denoting $\varphi(\boldsymbol{y}) := -\langle \boldsymbol{c}, \boldsymbol{y} \rangle$, the corresponding optimization problem can be written as

$$\widehat{\boldsymbol{y}}(\boldsymbol{\theta}^{\omega}) = \underset{\boldsymbol{y} \in \mathcal{Y}(\boldsymbol{x}^{\omega})}{\operatorname{argmax}} \sum_{i,j \in \mathcal{R}^{\omega}} \theta_{j}^{\omega} y_{i,j} - \sum_{i,j \in \mathcal{R}^{\omega}} c_{i,j} y_{i,j} = \langle \boldsymbol{\theta}^{\omega}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}). \tag{10}$$

The overall pipeline is summarized in Fig. 1. Following [4], we approximately solve the problem 755 in Eq. (10) using the prize-collecting HGS heuristic (PC-HGS), a variant of hybrid genetic search 756 (HGS) [49]. We denote this approximate solver $\widetilde{y} \approx \widehat{y}$, so that their proposed policy decomposes 757 as $f_W := \widetilde{y} \circ g_W$. The ground-truth routes are created by using an anticipative strategy, i.e., by 758 solving multiple instances where all future information is revealed from the start, and the requests' 759 arrival times information is translated into time windows (thus removing the dynamic aspect of the 760 problem). This anticipative policy, which we note f^* (which cannot be attained as it needs unavailable 761 information) is thus the target policy imitated by the model – see Appendix D.7 for more details. 762

D.3 Perturbation-based baseline.

In [4], a perturbation-based method [6] was used. This method is based on injecting noise in the PC-HGS solver \widetilde{y} . Similarly to our approach, the parameters W can then be learned using a Fenchel-Young loss, since the loss is minimized when the perturbed solver correctly predicts the ground truth. However, since \widetilde{y} is not an exact solver, all theoretical learning guarantees associated with this method (e.g., correctness of the gradients) no longer hold.

769 D.4 Proposed approach.

763

Our proposed approach instead uses the Fenchel-Young loss associated with the proposed layer, which is minimized when the proposed layer correctly predicts the ground-truth. At inference time, however, we use $f_W \coloneqq \widetilde{\boldsymbol{y}} \circ g_W$. We use a mixture of proposals, as defined in Algorithm 2. To design each proposal q_s , we build randomized versions of moves specifically designed for the prize-collecting dynamic vehicle routing problem with time windows. More precisely, we base our proposals on moves used in the local search part of the PC-HGS algorithm, which are summarized in Table 4. The details of turning these moves into proposal distributions with tractable individual correction ratios are given in Appendix C.6.

We evaluate three different initialization methods: (i) initialize $y^{(0)}$ by constructing routes dispatching random requests, (ii) initialize $y^{(0)}$ to the ground-truth solution, (iii) initialize $y^{(0)}$ by starting from the dataset ground-truth and applying a heuristic initialization algorithm to improve it. This heuristic initialization, similar to a short local search, is also used by the PC-HGS algorithm \tilde{y} , and is set to take up to half the time allocated to the layer (a limit it does not reach in practice).

D.5 Performance metric.

As the Fenchel-Young loss ℓ_t actually minimized is intractable to compute exactly, we only use the challenge metric. More precisely, we measure the cost relative to that of the anticipative baseline, $\frac{c_{\mathcal{W}}(f_W)-c_{\mathcal{W}}(f^*)}{c_{\mathcal{W}}(f^*)}$, which we average over a test dataset of unseen instances.

787 D.6 Results.

783

In Fig. 2, we observe that the initialization method plays an important role, and the ground-truth-based ones greatly outperform the random one.

We observe that the number of Markov iterations K is an important performance factor. Interestingly, the ground-truth initialization significantly improves the learning process for small K.

In Table 2, we compare training methods with fixed compute time budget for the layer (perturbed solver or proposed MCMC approach), which is by far the main computational bottleneck. This parameter limits the time allowed for a single forward pass through the combinatorial optimization layer (be it the perturbed inexact oracle or the proposed method). In both cases, the backward pass through the layer is immediate, as a property of the expression of the gradient of Fenchel-Young losses. The models are selected using a validation set and evaluated on the test set. We observe that the proposed approach significantly outperforms the perturbation-based method [6] using \tilde{y} in low time limit regimes, thus allowing for faster and more efficient training.

Full experimental details and additional results on the impact of temperature are given in Appendix D.7.

D.7 Additional experimental details and results for Section 4

Model, features, dataset, hyperparameters, compute. Following Baty et al. [4], the differentiable ML model g_W is implemented as a sparse graph neural network. We also use the same feature set, which represents the system state x^ω as a vector comprising request-level features, such as coordinates, time windows, demands, travel time to the depot, and quantiles from the distribution of the travel time to all other requests (named complete feature set, and described in the Table 4 of their paper). We use the same training, validation, and testing datasets, which are created from 30, 15 and 25 problem instances respectively. The training set uses a sample size of 50 requests per wave, while the rest use 100. The solutions in the training dataset, i.e., the examples from the anticipative strategy f^* imitated by the model, are obtained by solving the corresponding offline VRPTWs using HGS [49] with a time limit of 3600 seconds. During evaluation, the PC-HGS solver \tilde{y} is used with a constant time limit of 60 seconds for all models. We use Adam [25] together with the proposed stochastic gradient estimators, with a learning rate of $5 \cdot 10^{-3}$. Each training is performed using only a single CPU worker. For Fig. 2, we use a temperature $t = 10^2$. For Table 2, we use 1 Monte-Carlo sample for the perturbation-based method and 1 Markov chain for the proposed approach (in order to have a fair comparison: an equal number of oracle calls / equal compute).

Statistical significance. Each training is performed times with the same parameters and different random seeds. Then, the learning curves are averaged, and plotted with a 95% confidence interval. For the results in Table 2, we report the performance of the best model iteration (selected with respect to the validation set) on the test set. This procedure is also averaged over 50 trainings, and reported with 95% confidence intervals.

Additional results. In Fig. 13, we report model performance for varying temperature t. Interestingly, lower temperatures perform better when using random initialization. In the ground-truth initialization setting, a sweet spot is found at $t=10^2$, but lower temperatures do not particularly decrease performance.

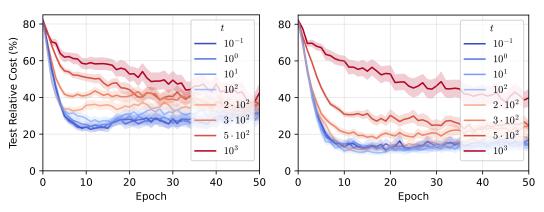


Figure 13: Test relative cost (%). **Left**: varying temperature t with random initialization. **Right**: varying temperature t with ground-truth initialization.

827 E Proofs

828 E.1 Proof of Eq. (4)

Proof. At fixed temperature $t_k = t$, the iterates of Algorithm 1 (MH case) follow a time-homogenous Markov chain, defined by the following transition kernel $P_{\theta,t}$:

$$P_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}') = \begin{cases} q(\boldsymbol{y},\boldsymbol{y}') \min\left[1, \frac{q(\boldsymbol{y}',\boldsymbol{y})}{q(\boldsymbol{y},\boldsymbol{y}')} \exp\left(\frac{\langle \boldsymbol{\theta},\boldsymbol{y}'\rangle + \varphi(\boldsymbol{y}') - \langle \boldsymbol{\theta},\boldsymbol{y}\rangle - \varphi(\boldsymbol{y})}{t}\right)\right] & \text{if } \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}), \\ 1 - \sum_{\boldsymbol{y}'' \in \mathcal{N}(\boldsymbol{y})} P_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}'') & \text{if } \boldsymbol{y}' = \boldsymbol{y}, \\ 0 & \text{else.} \end{cases}$$

Irreducibility. As we assumed the neighborhood graph $G_{\mathcal{N}}$ to be connected and undirected, the Markov Chain is irreducible as we have $\forall y \in \mathcal{Y}, \forall y' \in \mathcal{N}(y), P_{\theta,t}(y,y') > 0$.

Aperiodicity. For simplicity, we directly assumed aperiodicity in the main text. Here, we show that this is a mild condition, which is verified for instance if there is a solution $y \in \mathcal{Y}$ such that q(y, y) > 0. Indeed, we then have:

$$P_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}) = 1 - \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} P_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}')$$

$$= 1 - \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} q(\boldsymbol{y},\boldsymbol{y}') \min \left[1, \frac{q(\boldsymbol{y}',\boldsymbol{y})}{q(\boldsymbol{y},\boldsymbol{y}')} \exp \left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}') - \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle - \varphi(\boldsymbol{y})}{t} \right) \right]$$

$$\geq 1 - \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} q(\boldsymbol{y},\boldsymbol{y}')$$

$$\geq q(\boldsymbol{y},\boldsymbol{y}')$$

$$\geq 0.$$

Thus, we have $P_{\theta,t}(y,y) > 0$, which implies that the chain is aperiodic. As an irreducible and aperiodic Markov Chain on a finite state space, it converges to its stationary distribution and the latter is unique [18]. Finally, one can easily check that the detailed balance equation is satisfied for $\pi_{\theta,t}$, i.e.:

$$\forall \boldsymbol{y}, \boldsymbol{y}' \in \mathcal{Y}, \, \pi_{\boldsymbol{\theta},t}(\boldsymbol{y}) P_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}') = \pi_{\boldsymbol{\theta},t}(\boldsymbol{y}') P_{\boldsymbol{\theta},t}(\boldsymbol{y}',\boldsymbol{y}),$$

giving that $\pi_{\theta,t}$ is indeed the stationary distribution of the chain, which concludes the proof.

837 E.2 Proof of Proposition 2.1

Proof. Let $\theta \in \mathbb{R}^d$ and t > 0. The fact that $\hat{y}_t(\theta) \in \text{relint}(\mathcal{C}) = \text{relint}(\text{conv}(\mathcal{Y}))$ follows directly from the fact that $\hat{y}_t(\theta)$ is a convex combination of the elements of \mathcal{Y} with positive coefficients, as

840 $\forall \boldsymbol{y} \in \mathcal{Y}, \, \pi_{\boldsymbol{\theta},t}(\boldsymbol{y}) > 0.$

Low temperature limit. Let $y^* := \operatorname{argmax}_{y \in \mathcal{Y}} \langle \theta, y \rangle + \varphi(y)$. The argmax is assumed to be single-valued. Let $y \in \mathcal{Y} \setminus \{y^*\}$. We have:

$$\pi_{\boldsymbol{\theta},t}(\boldsymbol{y}) = \frac{\exp\left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t}\right)}{\sum_{\boldsymbol{y}' \in \mathcal{Y}} \exp\left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}')}{t}\right)}$$

$$\leq \frac{\exp\left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t}\right)}{\exp\left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}^{\star})}{t}\right)}$$

$$\leq \exp\left(\frac{(\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})) - (\langle \boldsymbol{\theta}, \boldsymbol{y}^{\star} \rangle + \varphi(\boldsymbol{y}^{\star}))}{t}\right)$$

$$\xrightarrow{t \to 0^{+}} 0,$$

as $\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}) < \langle \boldsymbol{\theta}, \boldsymbol{y}^* \rangle + \varphi(\boldsymbol{y}^*)$ by definition of \boldsymbol{y}^* . Thus, we have:

Thus, the expectation of $\pi_{\theta,t}$ converges to y^* . Naturally, if the argmax is not unique, the distribution converges to a uniform distribution on the maximizing structures.

High temperature limit. For all $y \in \mathcal{Y}$, we have:

845

854

$$\pi_{\boldsymbol{\theta},t}(\boldsymbol{y}) = \frac{\exp\left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t}\right)}{\sum_{\boldsymbol{y}' \in \mathcal{Y}} \exp\left(\frac{\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}')}{t}\right)}$$

$$\xrightarrow[t \to +\infty]{} \frac{1}{|\mathcal{Y}|},$$

as $\exp(x/t) \xrightarrow[t \to +\infty]{} 1$ for all $x \in \mathbb{R}$. Thus, $\pi_{\theta,t}$ converges to the uniform distribution on \mathcal{Y} , and its expectation converges to the average of all structures.

Expression of the Jacobian. Let $A_t: \theta \mapsto t \cdot \log \sum_{y \in \mathcal{Y}} \exp\left(\langle \theta, y \rangle + \varphi(y)\right)$ be the cumulant function of the exponential family defined by $\pi_{\theta,t}$, scaled by t. One can easily check that we have $\nabla_{\theta} A_t(\theta) = \widehat{y}_t(\theta)$. Thus, we have $J_{\theta} \widehat{y}_t(\theta) = \nabla^2_{\theta} A_t(\theta)$. However, we also have that the hessian matrix of the cumulant function $\theta \mapsto \frac{1}{t} A_t(\theta)$ is equal to the covariance matrix of the random vector under $\pi_{\theta,t}$ [51]. Thus, we have:

$$\begin{split} J_{\theta} \widehat{\boldsymbol{y}}_t(\boldsymbol{\theta}) &= \nabla_{\boldsymbol{\theta}}^2 A_t(\boldsymbol{\theta}) \\ &= t \cdot \nabla_{\boldsymbol{\theta}}^2 \left(\frac{1}{t} A_t(\boldsymbol{\theta}) \right) \\ &= t \cdot \mathsf{cov}_{\pi_{\boldsymbol{\theta},t}} \left[\frac{Y}{t} \right] \\ &= \frac{1}{t} \, \mathsf{cov}_{\pi_{\boldsymbol{\theta},t}} \left[Y \right]. \end{split}$$

855 E.3 Proof of Proposition 2.2

856 *Proof.* Let $K_{\theta,t}$ be the Markov transition kernel associated to Algorithm 2, which can be written as:

$$K_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}') = \begin{cases} \sum_{\substack{s \in Q(\boldsymbol{y}) \\ \text{s.t. } q_s(\boldsymbol{y},\boldsymbol{y}') > 0}} \frac{1}{|Q(\boldsymbol{y})|} q_s(\boldsymbol{y},\boldsymbol{y}') \min\left(1, \frac{|Q(\boldsymbol{y})|}{|Q(\boldsymbol{y}')|} \cdot \frac{q_s(\boldsymbol{y}',\boldsymbol{y})\pi_{\boldsymbol{\theta},t}(\boldsymbol{y}')}{q_s(\boldsymbol{y},\boldsymbol{y}')\pi_{\boldsymbol{\theta},t}(\boldsymbol{y})}\right) & \text{if } \boldsymbol{y}' \in \bar{\mathcal{N}}(\boldsymbol{y}), \\ 1 - \sum_{\boldsymbol{y}'' \in \bar{\mathcal{N}}(\boldsymbol{y})} K_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}'') & & \text{if } \boldsymbol{y}' = \boldsymbol{y}, \\ 0 & & \text{else.} \end{cases}$$

As $\forall y \in \mathcal{Y}, \forall y' \in \bar{\mathcal{N}}(y), K_{\theta,t}(y,y') > 0$, the irreducibility of the chain on \mathcal{Y} is directly implied by the connectedness of $G_{\bar{\mathcal{N}}}$.

Thus, we only have to check that the detailed balance equation $\pi_{\theta,t}(y)K_{\theta,t}(y,y') = \pi_{\theta,t}(y')K_{\theta,t}(y',y)$ is satisfied for all $y' \in \bar{\mathcal{N}}(y)$. We have:

$$\pi_{\theta,t}(\boldsymbol{y})K_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}') = \sum_{\substack{s \in Q(\boldsymbol{y})\\ \text{s.t. } q_s(\boldsymbol{y},\boldsymbol{y}') > 0}} \left[\frac{q_s(\boldsymbol{y},\boldsymbol{y}')\pi_{\theta,t}(\boldsymbol{y})}{|Q(\boldsymbol{y})|} \min\left(1, \frac{|Q(\boldsymbol{y})|}{|Q(\boldsymbol{y}')|} \cdot \frac{q_s(\boldsymbol{y}',\boldsymbol{y})\pi_{\theta,t}(\boldsymbol{y}')}{q_s(\boldsymbol{y},\boldsymbol{y}')\pi_{\theta,t}(\boldsymbol{y})}\right) \right].$$

The main point consists in noticing that the undirectedness assumption for each neighborhood graph $G_{\mathcal{N}_s}$ implies:

$${s \in Q(y) : q_s(y, y') > 0} = {s \in Q(y') : q_s(y', y) > 0}.$$

Thus, a simple case analysis on how $|Q(y)|q_s(y',y)\pi_{\theta,t}(y')$ and $|Q(y')|q_s(y,y')\pi_{\theta,t}(y)$ compare allows us to observe that the expression of $\pi_{\theta,t}(y)K_{\theta,t}(y,y')$ is symmetric in y and y', which concludes the proof.

862 E.4 Proof of strict convexity

Proof. As A_t is a differentiable convex function on \mathbb{R}^d (as the log-sum-exp of such functions), it is an essentially smooth closed proper convex function. Thus, it is such that

relint
$$(dom((A_t)^*)) \subseteq \nabla A_t(\mathbb{R}^d) \subseteq dom((A_t)^*),$$

and we have that the restriction of $(A_t)^*$ to $\nabla A_t(\mathbb{R}^d)$ is strictly convex on every convex subset of $\nabla A_t(\mathbb{R}^d)$ (corollary 26.4.1 in Rockafellar [42]). As the range of the gradient of the cumulant function $\theta \mapsto A_t(\theta)/t$ is exactly the relative interior of the marginal polytope conv $(\{y/t, y \in \mathcal{Y}\})$ (see appendix B.1 in Wainwright and Jordan [51]), and $(A_t)^* =: \Omega_t$, we actually have that

$$\operatorname{relint}(\operatorname{dom}(\Omega_t)) \subseteq \operatorname{relint}(\mathcal{C}) \subseteq \operatorname{dom}(\Omega_t),$$

and that Ω_t is strictly convex on every convex subset of relint(\mathcal{C}), i.e., strictly convex on relint(\mathcal{C}) (as relint(\mathcal{C}) is itself convex).

As A_t is closed proper convex, it is equal to its biconjugate by the Fenchel-Moreau theorem. Thus, we have:

$$A_t(\boldsymbol{\theta}) = \sup_{\boldsymbol{\mu} \in \mathbb{R}^d} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - (A_t)^*(\boldsymbol{\mu}) \right\} = \sup_{\boldsymbol{\mu} \in \mathbb{R}^d} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_t(\boldsymbol{\mu}) \right\}.$$

Moreover, as $\nabla A_t(\mathbb{R}^d) = \operatorname{relint}(\mathcal{C})$, we have $||\nabla A_t(\boldsymbol{\theta})|| \leq R_{\mathcal{C}} := \max_{\boldsymbol{\mu} \in \mathcal{C}} ||\boldsymbol{\mu}||$, which gives $\operatorname{dom}(\Omega_t) \subset B(\mathbf{0}, R_{\mathcal{C}})$. Thus we can actually write:

$$A_t(\boldsymbol{\theta}) = \sup_{\boldsymbol{\mu} \in B(\mathbf{0}, R_{\mathcal{C}})} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_t(\boldsymbol{\mu}) \right\},\,$$

and now apply Danksin's theorem as $B(\mathbf{0}, R_{\mathcal{C}})$ is compact, which further gives:

$$\partial A_t(\boldsymbol{\theta}) = \mathop{\mathrm{argmax}}_{\boldsymbol{\mu} \in B(\mathbf{0}, R_{\mathcal{C}})} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_t(\boldsymbol{\mu}) \right\},$$

and the fact that A_t is differentiable gives that both sides are single-valued. Moreover, as $\nabla A_t(\mathbb{R}^d)$ = relint(\mathcal{C}), we know that the right hand side is maximized in \mathcal{C} , and we can actually write:

$$\nabla A_t(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\mu} \in \mathcal{C}} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_t(\boldsymbol{\mu}) \right\}.$$

We end this proof by noting that a simple calculation yields $\nabla A_t(\boldsymbol{\theta}) = \mathbb{E}_{\pi_{\boldsymbol{\theta},t}}[Y] = \widehat{\boldsymbol{y}}_t(\boldsymbol{\theta})$. The expression of $\nabla_{\boldsymbol{\theta}} \ell_t(\boldsymbol{\theta}; \boldsymbol{y})$ follows.

Remark E.1. The proposed Fenchel-Young loss can also be obtained via distribution-space regularization. Let $s_{\theta} \coloneqq (\langle \theta, y \rangle + \varphi(y))_{y \in \mathcal{Y}} \in \mathbb{R}^{|\mathcal{Y}|}$ be a vector containing the score of all structures, and

869 $L_{-tH}: \mathbb{R}^{|\mathcal{Y}|} \times \Delta^{|\mathcal{Y}|} \to \mathbb{R}$ be the Fenchel-Young loss generated by -tH, where H is the Shannon

entropy. We have $\nabla_{s_{\theta}}(-tH)^*(s_{\theta}) = \pi_{\theta,t}$. The chain rule further gives $\nabla_{\theta}(-tH)^*(s_{\theta}) = \mathbb{E}_{\pi_{\theta,t}}[Y]$.

Thus, we have $\nabla_{\theta} L_{-tH}(s_{\theta}; p_y) = \nabla_{\theta} \ell_t(\theta; y)$, where p_y is the dirac distribution on y. In the case

where $\varphi \equiv 0$ and t=1, we have $\Omega_t(\mu) = -\left(\max_{p \in \Delta^{|\mathcal{Y}|}} H^s(p) \text{ s.t. } \mathbb{E}_p[Y] = \mu\right)$, with H^s the

Shannon entropy [8], and ℓ_t is known as the CRF loss [29].

874 E.5 Proof of Proposition 3.1

Proof. The proof is exactly the proof of Proposition 4.1 in Berthet et al. [6], in which the setting is similar, and all the same arguments hold (we also have that π_{θ_0} is dense on \mathcal{Y} , giving $\bar{Y}_N \in \operatorname{relint}(\mathcal{C})$ for N large enough). The only difference is the choice of regularization function, and we have to prove that it is also convex and smooth in our case. While the convexity of Ω_t is directly implied by its definition as a Fenchel conjugate, the fact that is is smooth is due to Theorem 26.3 in Rockafellar [42] and the essential strict convexity of A_t (which is itself closed proper convex). The latter relies on the fact that \mathcal{C} is assumed to be of full-dimension (otherwise A_t would be linear when restricted to any affine subspace of direction equal to the subspace orthogonal to the direction of the smallest affine subspace spanned by \mathcal{C}), which in turn implies that A_t is strictly convex on \mathbb{R}^d . Thus, Proposition 4.1 in Berthet et al. [6] gives the asymptotic normality:

$$\sqrt{N}(\boldsymbol{\theta}_{N}^{\star}-\boldsymbol{\theta}_{0})\xrightarrow[N\to\infty]{\mathcal{D}}\mathcal{N}\left(\mathbf{0},\,\left(\nabla_{\boldsymbol{\theta}}^{2}A_{t}(\boldsymbol{\theta}_{0})\right)^{-1}\operatorname{cov}_{\pi_{\boldsymbol{\theta}_{0},t}}\left[Y\right]\left(\nabla_{\boldsymbol{\theta}}^{2}A_{t}(\boldsymbol{\theta}_{0})\right)^{-1}\right).$$

Moreover, we already derived $\nabla^2_{\theta} A_t(\theta_0) = \frac{1}{t} \operatorname{cov}_{\pi_{\theta_0,t}}[Y]$ in Appendix E.2, leading to the simplified asymptotic normality given in the proposition.

877

878 E.6 Proof of Proposition 3.2

Proof. The proof consists in bounding the convergence rate of the Markov chain $(y^{(k)})_{k \in \mathbb{N}}$ (which has transition kernel $P_{\theta,t}$) for all θ , in order to apply Theorem 4.1 in Younes [52]. It is defined as the smallest constant λ_{θ} such that:

$$\exists A > 0 : \forall y \in \mathcal{Y}, |\mathbb{P}(y^{(k)} = y) - \pi_{\theta,t}(y)| \leq A\lambda_{\theta}^{k}$$

More precisely, we must find a constant D such that $\exists B>0: \lambda_{\theta} \leq 1-Be^{-D||\theta||}$, in order to impose $K_{n+1}>\left\lfloor 1+a'\exp\left(2D||\hat{\theta}_n||\right)\right\rfloor$.

A known result gives $\lambda_{\boldsymbol{\theta}} \leq \rho(\boldsymbol{\theta})$ with $\rho(\boldsymbol{\theta}) = \max_{\lambda \in S_{\boldsymbol{\theta}} \setminus \{1\}} |\lambda|$ [32], where $S_{\boldsymbol{\theta}}$ is the spectrum of the transition kernel $P_{\boldsymbol{\theta},t}$ (here, $1 - \rho(\boldsymbol{\theta})$ is known as the *spectral gap* of the Markov chain). To bound $\rho(\boldsymbol{\theta})$, we use the results of Ingrassia [23], which study the Markov chain with transition kernel $P'_{\boldsymbol{\theta},t}$, such that $P_{\boldsymbol{\theta},t} = \frac{1}{2} \left(I + P'_{\boldsymbol{\theta},t} \right)$. It corresponds to the same algorithm, but with a proposal distribution q' defined as:

$$q'(\boldsymbol{y}, \boldsymbol{y}') = \begin{cases} \frac{1}{d^*} & \text{if } \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}), \\ 1 - \frac{d(\boldsymbol{y})}{d^*} & \text{if } \boldsymbol{y}' = \boldsymbol{y}, \\ 0 & \text{else.} \end{cases}$$

As $P'_{\theta,t}$ is a row-stochastic matrix, Gershgorin's circle theorem gives that its spectrum is included in the complex unit disc. Moreover, one can easily check that the associated Markov chain is also reversible with respect to $\pi_{\theta,t}$, and the corresponding detailed balance equation gives:

$$\forall \boldsymbol{y}, \boldsymbol{y}' \in \mathcal{Y}, \ \pi_{\boldsymbol{\theta},t}(\boldsymbol{y}) P'_{\boldsymbol{\theta},t}(\boldsymbol{y}, \boldsymbol{y}') = \pi_{\boldsymbol{\theta},t}(\boldsymbol{y}') P'_{\boldsymbol{\theta},t}(\boldsymbol{y}', \boldsymbol{y}),$$

893 which is equivalent to:

$$\forall \boldsymbol{y}, \boldsymbol{y}' \in \mathcal{Y}, \ \sqrt{\frac{\pi_{\boldsymbol{\theta},t}(\boldsymbol{y})}{\pi_{\boldsymbol{\theta},t}(\boldsymbol{y}')}} P_{\boldsymbol{\theta},t}'(\boldsymbol{y},\boldsymbol{y}') = \sqrt{\frac{\pi_{\boldsymbol{\theta},t}(\boldsymbol{y}')}{\pi_{\boldsymbol{\theta},t}(\boldsymbol{y})}} P_{\boldsymbol{\theta},t}'(\boldsymbol{y}',\boldsymbol{y})$$

as $\pi_{\theta,t}$ has full support on \mathcal{Y} , which can be further written in matrix form as:

$$\Pi_{\boldsymbol{\theta}}^{1/2} P_{\boldsymbol{\theta},t}' \Pi_{\boldsymbol{\theta}}^{-1/2} = \Pi_{\boldsymbol{\theta}}^{-1/2} P_{\boldsymbol{\theta},t}'^{\top} \Pi_{\boldsymbol{\theta}}^{1/2},$$

where $\Pi_{\theta}=\mathrm{diag}(\pi_{\theta;t})$. Thus, the matrix $\Pi_{\theta}^{1/2}P'_{\theta,t}\Pi_{\theta}^{-1/2}$ is symmetric, and the spectral theorem ensures its eigenvalues are real. As it is similar to the transition kernel $P'_{\theta,t}$ (with change of basis matrix $\Pi_{\theta}^{-1/2}$), they share the same spectrum S'_{θ} , and we have $S'_{\theta}\subset [-1,1]$. Let us order S'_{θ} as $-1\leq \lambda'_{\min}\leq \cdots \leq \lambda'_2\leq \lambda'_1=1$. As $P_{\theta,t}=\frac{1}{2}\left(I+P'_{\theta,t}\right)$, we clearly have $\rho(\theta)=\frac{1+\lambda'_2}{2}$. Thus, we can use Theorem 4.1 of Ingrassia [23], which gives $\lambda'_2\leq 1-G\cdot Z(\theta)\exp(-m(\theta))$ (we keep their notations for Z and m, and add the dependency in θ for clarity), where G is a constant depending only on the graph $G_{\mathcal{N}}$, and with:

$$\begin{split} Z(\boldsymbol{\theta}) &= \sum_{\boldsymbol{y} \in \mathcal{Y}} \exp\left(\frac{\langle \boldsymbol{\theta}, \, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t} - \max_{\boldsymbol{y}' \in \mathcal{Y}} \left[\frac{\langle \boldsymbol{\theta}, \, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}')}{t}\right]\right) \\ &\geq |\mathcal{Y}| \exp\left(\frac{1}{t} \left[\min_{\boldsymbol{y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \, \boldsymbol{y} \rangle + \min_{\boldsymbol{y} \in \mathcal{Y}} \varphi(\boldsymbol{y}) - \max_{\boldsymbol{y}' \in \mathcal{Y}} \langle \boldsymbol{\theta}, \, \boldsymbol{y}' \rangle - \max_{\boldsymbol{y}' \in \mathcal{Y}} \varphi(\boldsymbol{y}')\right]\right) \\ &\geq |\mathcal{Y}| \exp\left(-\frac{2R_{\mathcal{C}}}{t} ||\boldsymbol{\theta}|| - \frac{2R_{\varphi}}{t}\right), \end{split}$$

902 and:

$$\begin{split} m(\boldsymbol{\theta}) &\leq \max_{\boldsymbol{y} \in \mathcal{Y}} \left\{ \max_{\boldsymbol{y}' \in \mathcal{Y}} \left[\frac{\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}')}{t} \right] - \frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t} \right\} - 2 \min_{\boldsymbol{y} \in \mathcal{Y}} \left\{ \max_{\boldsymbol{y}' \in \mathcal{Y}} \left[\frac{\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}')}{t} \right] - \frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t} \right\} \\ &= \max_{\boldsymbol{y}' \in \mathcal{Y}} \left[\frac{\langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}')}{t} \right] - \min_{\boldsymbol{y} \in \mathcal{Y}} \left[\frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t} \right] \\ &\leq \frac{1}{t} \left(\max_{\boldsymbol{y}' \in \mathcal{Y}} \langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \max_{\boldsymbol{y}' \in \mathcal{Y}} \varphi(\boldsymbol{y}') - \min_{\boldsymbol{y} \in \mathcal{Y}} \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle - \min_{\boldsymbol{y} \in \mathcal{Y}} \varphi(\boldsymbol{y}) \right) \\ &\leq \frac{2R_{\mathcal{C}}}{t} ||\boldsymbol{\theta}|| + \frac{2R_{\varphi}}{t}, \end{split}$$

where $R_{\mathcal{C}} = \max_{y \in \mathcal{Y}} ||y||$ and $R_{\varphi} = \max_{y \in \mathcal{Y}} |\varphi(y)|$. Thus, we have:

$$\lambda_2' \le 1 - G|\mathcal{Y}| \exp\left(-\frac{4R_{\varphi}}{t}\right) \exp\left(-\frac{4R_{\mathcal{C}}}{t}||\boldsymbol{\theta}||\right),$$

903 and finally:

$$\lambda_{\theta} \leq 1 - \frac{G|\mathcal{Y}|\exp\left(-\frac{4R_{\varphi}}{t}\right)}{2}\exp\left(-\frac{4R_{\mathcal{C}}}{t}||\boldsymbol{\theta}||\right),$$

so taking $D = 4R_{\mathcal{C}}/t$ concludes the proof.

905

Remark E.2. The stationary distribution in Ingrassia [23] is defined as proportional to $\exp(-H(y))$, with the assumption that the function H is such that $\min_{y \in \mathcal{Y}} H(y) = 0$. Thus, we apply their results with

$$H(\boldsymbol{y}) \coloneqq \max_{\boldsymbol{y}' \in \mathcal{Y}} \left[\frac{\langle \boldsymbol{\theta}, \, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}')}{t} \right] - \frac{\langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y})}{t}$$

(which gives correct distribution $\pi_{\theta,t}$ and respects this assumption), hence the obtained forms for $Z(\theta)$ and the upper bound on $m(\theta)$.

908 E.7 Proofs of Proposition C.1 and Proposition C.2

Proposition C.1. The distribution of the first iterate of the Markov chain with transition kernel defined in Eq. (3) and initialized at the ground-truth structure y is given by:

$$\begin{split} (\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)})(\boldsymbol{y}') &= P_{\boldsymbol{\theta},t}(\boldsymbol{y},\boldsymbol{y}') \\ &= \begin{cases} q(\boldsymbol{y},\boldsymbol{y}') \min\left[1,\frac{q(\boldsymbol{y}',\boldsymbol{y})}{q(\boldsymbol{y},\boldsymbol{y}')} \exp\left(\left[\langle \boldsymbol{\theta},\boldsymbol{y}'-\boldsymbol{y}\rangle + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y})\right]/t\right)\right] & \text{if } \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}), \\ 1 - \sum_{\boldsymbol{y}'' \in \mathcal{N}(\boldsymbol{y})} (\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)})(\boldsymbol{y}'') & \text{if } \boldsymbol{y}' = \boldsymbol{y}, \\ 0 & \text{else.} \end{split}$$

1911 Let $\alpha_{\boldsymbol{y}}(\boldsymbol{\theta}, \boldsymbol{y}') \coloneqq \frac{q(\boldsymbol{y}', \boldsymbol{y})}{q(\boldsymbol{y}, \boldsymbol{y}')} \exp\left(\left[\left\langle \boldsymbol{\theta}, \boldsymbol{y}' - \boldsymbol{y} \right\rangle + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y})\right]/t\right)$. Define also the following sets:

$$\mathcal{N}_{\boldsymbol{y}}^{-}(\boldsymbol{\theta}) = \left\{ \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}) \mid \alpha_{\boldsymbol{y}}(\boldsymbol{\theta}, \boldsymbol{y}') \leq 1 \right\}, \quad \mathcal{N}_{\boldsymbol{y}}^{+}(\boldsymbol{\theta}) = \left\{ \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}) \mid \alpha_{\boldsymbol{y}}(\boldsymbol{\theta}, \boldsymbol{y}') > 1 \right\}.$$

The expectation of the first iterate is then given by:

$$\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] = \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} (\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)})(\boldsymbol{y}') \cdot \boldsymbol{y}' + \left(1 - \sum_{\boldsymbol{y}'' \in \mathcal{N}(\boldsymbol{y})} (\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)})(\boldsymbol{y}'')\right) \cdot \boldsymbol{y}$$

$$= \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} (\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)})(\boldsymbol{y}') \cdot (\boldsymbol{y}' - \boldsymbol{y})$$

$$= \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}_{\boldsymbol{y}}^{(1)}(\boldsymbol{\theta})} q(\boldsymbol{y}',\boldsymbol{y}) \exp\left(\left[\langle \boldsymbol{\theta}, \boldsymbol{y}' - \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y})\right]/t\right) \cdot (\boldsymbol{y}' - \boldsymbol{y}) + \sum_{\boldsymbol{y}' \in \mathcal{N}_{\boldsymbol{y}}^{(1)}(\boldsymbol{\theta})} q(\boldsymbol{y},\boldsymbol{y}') \cdot (\boldsymbol{y}' - \boldsymbol{y}).$$

913 Let now $f_{\boldsymbol{y}}: \mathbb{R}^d \times \mathcal{N}(\boldsymbol{y}) \to \mathbb{R}$ be defined as:

914

915

916

917

$$f_{\boldsymbol{y}}: (\boldsymbol{\theta}; \boldsymbol{y}') \mapsto \begin{cases} t \cdot q(\boldsymbol{y}', \boldsymbol{y}) \exp\left(\left[\langle \boldsymbol{\theta}, \boldsymbol{y}' - \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y})\right] / t & \text{if } \alpha_{\boldsymbol{y}}(\boldsymbol{\theta}, \boldsymbol{y}') \leq 1, \\ t \cdot q(\boldsymbol{y}, \boldsymbol{y}') \left(\left[\langle \boldsymbol{\theta}, \boldsymbol{y}' - \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y})\right] / t + 1 - \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} \right) & \text{if } \alpha_{\boldsymbol{y}}(\boldsymbol{\theta}, \boldsymbol{y}') > 1. \end{cases}$$

Let $F_y: \theta \mapsto \langle \theta, y \rangle + \sum_{y' \in \mathcal{N}(y)} f_y(\theta; y')$. We define the target-dependent regularization function Ω_y and the corresponding Fenchel-Young loss as:

$$\Omega_{\boldsymbol{y}}: \boldsymbol{\mu} \mapsto (F_{\boldsymbol{y}})^*(\boldsymbol{\mu}), \qquad L_{\Omega_{\boldsymbol{y}}}(\boldsymbol{\theta}; \boldsymbol{y}) \coloneqq (\Omega_{\boldsymbol{y}})^*(\boldsymbol{\theta}) + \Omega_{\boldsymbol{y}}(\boldsymbol{y}) - \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle.$$

•
$$\Omega_{\boldsymbol{y}}$$
 is $t/\mathbb{E}_{q(\boldsymbol{y},\cdot)}||Y-\boldsymbol{y}||_2^2$ -strongly convex:

One can easily check that $f_y(\cdot; y')$ is continuous for all $y' \in \mathcal{N}(y)$, as it is defined piecewise as continuous functions that match on the junction affine hyperplane defined by:

$$\left\{\boldsymbol{\theta} \in \mathbb{R}^d \mid \alpha_{\boldsymbol{y}}(\boldsymbol{\theta}; \boldsymbol{y}') = 1\right\} = \left\{\boldsymbol{\theta} \in \mathbb{R}^d \mid \langle \boldsymbol{\theta}, \boldsymbol{y}' - \boldsymbol{y} \rangle = t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')\right\}.$$

Moreover, we have that $f_{\boldsymbol{y}}(\cdot; \boldsymbol{y}')$ is actually differentiable everywhere as its gradient can be continuously extended to the junction affine hyperplane with constant value equal to $q(\boldsymbol{y}, \boldsymbol{y}')(\boldsymbol{y}' - \boldsymbol{y})$. We now show that $f_{\boldsymbol{y}}(\cdot; \boldsymbol{y}')$ is $\frac{1}{t}q(\boldsymbol{y}, \boldsymbol{y}')\cdot||\boldsymbol{y}'-\boldsymbol{y}||^2$ -smooth. Indeed, it is defined as the composition of the linear form $\boldsymbol{\theta}\mapsto\langle\boldsymbol{\theta},\boldsymbol{y}'-\boldsymbol{y}\rangle$ and the function $g:\mathbb{R}\to\mathbb{R}$ given by:

$$g: x \mapsto \begin{cases} t \cdot q(\mathbf{y}', \mathbf{y}) \exp\left(\left[x + \varphi(\mathbf{y}') - \varphi(\mathbf{y})\right]/t\right) & \text{if } x \le t \log \frac{q(\mathbf{y}, \mathbf{y}')}{q(\mathbf{y}', \mathbf{y})} + \varphi(\mathbf{y}) - \varphi(\mathbf{y}'), \\ t \cdot q(\mathbf{y}, \mathbf{y}') \left(\left[x + \varphi(\mathbf{y}') - \varphi(\mathbf{y})\right]/t + 1 - \log \frac{q(\mathbf{y}, \mathbf{y}')}{q(\mathbf{y}', \mathbf{y})}\right) & \text{if } x > t \log \frac{q(\mathbf{y}, \mathbf{y}')}{q(\mathbf{y}', \mathbf{y})} + \varphi(\mathbf{y}) - \varphi(\mathbf{y}'). \end{cases}$$

We begin by showing that g is $\frac{1}{t}q(\boldsymbol{y},\boldsymbol{y}')$ -smooth. We have:

$$g': x \mapsto \begin{cases} q(\mathbf{y}', \mathbf{y}) \exp\left(\left[x + \varphi(\mathbf{y}') - \varphi(\mathbf{y})\right]/t\right) & \text{if } x \le t \log \frac{q(\mathbf{y}, \mathbf{y}')}{q(\mathbf{y}', \mathbf{y})} + \varphi(\mathbf{y}) - \varphi(\mathbf{y}'), \\ q(\mathbf{y}, \mathbf{y}') & \text{if } x > t \log \frac{q(\mathbf{y}, \mathbf{y}')}{q(\mathbf{y}', \mathbf{y})} + \varphi(\mathbf{y}) - \varphi(\mathbf{y}'). \end{cases}$$

Thus, g' is continuous, and differentiable everywhere except in $x_0 := t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')$. Its derivative is given by:

$$g'': x \mapsto \begin{cases} \frac{1}{t}q(\mathbf{y}', \mathbf{y}) \exp\left(\left[x + \varphi(\mathbf{y}') - \varphi(\mathbf{y})\right]/t\right) & \text{if } x \le t \log \frac{q(\mathbf{y}, \mathbf{y}')}{q(\mathbf{y}', \mathbf{y})} + \varphi(\mathbf{y}) - \varphi(\mathbf{y}'), \\ 0 & \text{if } x > t \log \frac{q(\mathbf{y}, \mathbf{y}')}{q(\mathbf{y}', \mathbf{y})} + \varphi(\mathbf{y}) - \varphi(\mathbf{y}'). \end{cases}$$

• For $x_1, x_2 \leq x_0$, we have:

$$|g'(x_1) - g'(x_2)| \le |x_1 - x_2| \sup_{\substack{x \in]-\infty, x_0[\\ x < x_0}} |g''(x)|$$

$$= |x_1 - x_2| \lim_{\substack{x \to x_0\\ x < x_0}} |g''(x)|$$

$$= \frac{1}{t} q(\mathbf{y}, \mathbf{y}') \cdot |x_1 - x_2|.$$

• For $x_1, x_2 \ge x_0$, we trivially have $|g'(x_1) - g'(x_2)| = 0$.

• For $x_1 \leq x_0 \leq x_2$, we have:

$$|g'(x_1) - g'(x_2)| = |(g'(x_1) - g'(x_0)) - (g'(x_2) - g'(x_0))|$$

$$\leq |g'(x_1) - g'(x_0)| + |g'(x_2) - g'(x_0)|$$

$$\leq \frac{1}{t}q(\mathbf{y}, \mathbf{y}') \cdot |x_1 - x_0|$$

$$\leq \frac{1}{t}q(\mathbf{y}, \mathbf{y}') \cdot |x_1 - x_2|.$$

Thus, we have:

$$\forall x_1, x_2 \in \mathbb{R}, |g'(x_1) - g'(x_2)| \le \frac{1}{t} q(\boldsymbol{y}, \boldsymbol{y}') \cdot |x_1 - x_2|,$$

and g is $\frac{1}{t}q(\boldsymbol{y},\boldsymbol{y}')$ -smooth. Nevertheless, we have $f_{\boldsymbol{y}}(\,\cdot\,,\boldsymbol{y}')=g(\langle\,\cdot\,,\boldsymbol{y}'-\boldsymbol{y}\rangle)$. Thus, we have, for

$$\begin{split} ||\nabla_{\boldsymbol{\theta}} f_{\boldsymbol{y}}(\boldsymbol{\theta}_{1}, \boldsymbol{y}') - \nabla_{\boldsymbol{\theta}} f_{\boldsymbol{y}}(\boldsymbol{\theta}_{2}, \boldsymbol{y}')|| &= ||g'(\langle \boldsymbol{\theta}_{1}, \boldsymbol{y}' - \boldsymbol{y} \rangle)(\boldsymbol{y}' - \boldsymbol{y}) - g'(\langle \boldsymbol{\theta}_{2}, \boldsymbol{y}' - \boldsymbol{y} \rangle)(\boldsymbol{y}' - \boldsymbol{y})|| \\ &= |g'(\langle \boldsymbol{\theta}_{1}, \boldsymbol{y}' - \boldsymbol{y} \rangle) - g'(\langle \boldsymbol{\theta}_{2}, \boldsymbol{y}' - \boldsymbol{y} \rangle)| \cdot ||\boldsymbol{y}' - \boldsymbol{y}|| \\ &\leq \frac{1}{t} q(\boldsymbol{y}, \boldsymbol{y}') \cdot |\langle \boldsymbol{\theta}_{1}, \boldsymbol{y}' - \boldsymbol{y} \rangle - \langle \boldsymbol{\theta}_{2}, \boldsymbol{y}' - \boldsymbol{y} \rangle| \cdot ||\boldsymbol{y}' - \boldsymbol{y}|| \\ &\leq \frac{1}{t} q(\boldsymbol{y}, \boldsymbol{y}') \cdot ||\boldsymbol{y}' - \boldsymbol{y}||^{2} \cdot ||\boldsymbol{\theta}_{1} - \boldsymbol{\theta}_{2}||, \end{split}$$

and $f_y(\cdot, y')$ is $\frac{1}{t}q(y, y') \cdot ||y' - y||^2$ -smooth. Thus, recalling that F_y is defined as

$$F_{m{y}}: m{ heta} \mapsto \langle m{ heta}, m{y}
angle + \sum_{m{y}' \in \mathcal{N}(m{y})} f_{m{y}}(m{ heta}; m{y}'),$$

we have that $F_{\boldsymbol{y}}$ is $\sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} \frac{1}{t} q(\boldsymbol{y}, \boldsymbol{y}') \cdot ||\boldsymbol{y}' - \boldsymbol{y}||^2 = \mathbb{E}_{q(\boldsymbol{y}, \cdot)} ||Y - \boldsymbol{y}||_2^2 / t$ -smooth. Finally, as $\Omega_{\boldsymbol{y}} := (F_{\boldsymbol{y}})^*$, Fenchel duality theory gives that $\Omega_{\boldsymbol{y}}$ is $t/\mathbb{E}_{q(\boldsymbol{y}, \cdot)} ||Y - \boldsymbol{y}||_2^2$ -strongly convex.

$$\bullet \ \mathbb{E}_{p_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] = \mathrm{argmax}_{\boldsymbol{\mu} \in \mathsf{conv}(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\})} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_{\boldsymbol{y}}(\boldsymbol{\mu}) \right\} :$$

Noticing that g is continuous on \mathbb{R} , convex on $\left]-\infty, t\log\frac{q(y,y')}{q(y',y)} + \varphi(y) - \varphi(y')\right[$ and on $t \log \frac{q(y,y')}{q(y',y)} + \varphi(y) - \varphi(y'), +\infty$, and with matching derivatives on the junction:

$$g'(t) \xrightarrow[t < t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}') \xrightarrow[t < t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}) - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y}')} \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y})} + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y}')] \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y}')} + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y}')] \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y}')} + \varphi(\boldsymbol{y}')] \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y}')} + \varphi(\boldsymbol{y}')] \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y}')} + \varphi(\boldsymbol{y}')] \xrightarrow[t > t \log \frac{q(\boldsymbol{y}, \boldsymbol{y}')}{q(\boldsymbol{y}', \boldsymbol{y}')} + \varphi(\boldsymbol{y}', \boldsymbol{y$$

gives that g is convex on \mathbb{R} . Thus, $f_y(\cdot; y')$ is convex on \mathbb{R}^d by composition. Thus,

$$F_{\boldsymbol{y}}: \boldsymbol{\theta} \mapsto \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} f_{\boldsymbol{y}}(\boldsymbol{\theta}; \boldsymbol{y}')$$

is closed proper convex as the sum of such functions. The Fenchel-Moreau theorem then gives that it is equal to its biconjugate. Thus, we have:

$$F_{\boldsymbol{y}}(\boldsymbol{\theta}) = \sup_{\boldsymbol{\mu} \in \mathbb{R}^d} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - (F_{\boldsymbol{y}})^*(\boldsymbol{\mu}) \right\} = \sup_{\boldsymbol{\mu} \in \mathbb{R}^d} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_{\boldsymbol{y}}(\boldsymbol{\mu}) \right\}.$$

Nonetheless, the gradient of $F_{\boldsymbol{y}}$ is given by:

$$\nabla_{\boldsymbol{\theta}} F_{\boldsymbol{y}}(\boldsymbol{\theta}) = \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}_{\boldsymbol{y}}^{-}(\boldsymbol{\theta})} q(\boldsymbol{y}', \boldsymbol{y}) \exp\left(\left[\langle \boldsymbol{\theta}, \boldsymbol{y}' - \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y})\right] / t\right) \cdot (\boldsymbol{y}' - \boldsymbol{y}) + \sum_{\boldsymbol{y}' \in \mathcal{N}_{\boldsymbol{y}}^{+}(\boldsymbol{\theta})} q(\boldsymbol{y}, \boldsymbol{y}') \cdot (\boldsymbol{y}' - \boldsymbol{y})$$

$$= \mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta}, \boldsymbol{y}}^{(1)}} [Y].$$

Thus, we have $\nabla F_{\boldsymbol{y}}(\mathbb{R}^d) \subset \operatorname{conv}(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\})$, which gives:

$$\forall \boldsymbol{\theta} \in \mathbb{R}^d, \, ||\nabla F_{\boldsymbol{y}}(\boldsymbol{\theta})|| \leq R_{\mathcal{N}(\boldsymbol{y})} \coloneqq \max_{\boldsymbol{\mu} \in \mathsf{conv}(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\})} ||\boldsymbol{\mu}||,$$

so that we have $\mathsf{dom}(\Omega_{\pmb{y}}) \subset B(\pmb{0}, R_{\mathcal{N}(\pmb{y})}).$ Thus we can actually write:

$$F_{\boldsymbol{y}}(\boldsymbol{\theta}) = \sup_{\boldsymbol{\mu} \in B(0, R_{\boldsymbol{\lambda}(\boldsymbol{\mu})})} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_{\boldsymbol{y}}(\boldsymbol{\mu}) \right\},$$

and now apply Danksin's theorem as $B(\mathbf{0}, R_{\mathcal{N}(y)})$ is compact, which further gives:

$$\partial F_{\boldsymbol{y}}(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\mu} \in B(\boldsymbol{0}, R_{\mathcal{N}(\boldsymbol{y})})} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_{\boldsymbol{y}}(\boldsymbol{\mu}) \right\},$$

and the fact that $F_{\boldsymbol{y}}$ is differentiable gives that both sides are single-valued. Moreover, as $\nabla F_{\boldsymbol{y}}(\mathbb{R}^d) \subset \operatorname{conv}(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\})$, we know that the right hand side is maximized in $\operatorname{conv}(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\})$, and we can actually write:

$$\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] = \nabla F_{\boldsymbol{y}}(\boldsymbol{\theta}) = \operatorname*{argmax}_{\boldsymbol{\mu} \in \mathsf{conv}(\mathcal{N}(\boldsymbol{y}) \cup \{\boldsymbol{y}\})} \left\{ \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle - \Omega_{\boldsymbol{y}}(\boldsymbol{\mu}) \right\}.$$

• Smoothness of $L_{\Omega_{\boldsymbol{y}}}(\cdot;\boldsymbol{y})$ and expression of its gradient:

Based on the above, we have:

924

$$L_{\Omega_{\boldsymbol{y}}}(\boldsymbol{\theta}; \boldsymbol{y}) = F_{\boldsymbol{y}}(\boldsymbol{\theta}) + \Omega_{\boldsymbol{y}}(\boldsymbol{y}) - \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle.$$

Thus, the $\mathbb{E}_{q(\boldsymbol{y},\cdot)}||Y-\boldsymbol{y}||_2^2/t$ -smoothness of $L_{\Omega_{\boldsymbol{y}}}(\cdot;\boldsymbol{y})$ follows directly from the previously established $\mathbb{E}_{q(\boldsymbol{y},\cdot)}||Y-\boldsymbol{y}||_2^2/t$ -smoothness of $F_{\boldsymbol{y}}$. Similarly, the expression of $\nabla_{\boldsymbol{\theta}}L_{\Omega_{\boldsymbol{y}}}(\boldsymbol{\theta};\boldsymbol{y})$ follows from the previously established expression of $\nabla_{\boldsymbol{\theta}}F_{\boldsymbol{y}}(\boldsymbol{\theta})$, and we have:

$$\nabla_{\boldsymbol{\theta}} L_{\Omega_{\boldsymbol{y}}}(\boldsymbol{\theta}\,;\boldsymbol{y}) = \nabla_{\boldsymbol{\theta}} F_{\boldsymbol{y}}(\boldsymbol{\theta}) - \boldsymbol{y} = \mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] - \boldsymbol{y}.$$

925

Proposition C.2. In the unsupervised setting, given a dataset $(y_i)_{i=1}^N$, the distribution of the first iterate of the Markov chain with transition kernel defined in Eq. (3) and initialized by $y^{(0)} = y_i$, with $i \sim \mathcal{U}(\lceil 1, N \rceil)$, is given by:

$$(\boldsymbol{p}_{\boldsymbol{\theta},\bar{Y}_{N}}^{(1)})(\boldsymbol{y}) = \sum_{\boldsymbol{y}' \in \mathcal{Y}} \left(\sum_{i=1}^{N} \mathbf{1}_{\{\boldsymbol{y}_{i} = \boldsymbol{y}'\}} \cdot \frac{1}{N} \right) P_{\boldsymbol{\theta},t}(\boldsymbol{y}',\boldsymbol{y})$$

$$= \sum_{\boldsymbol{y}' \in \mathcal{Y}} \left(\sum_{i=1}^{N} \mathbf{1}_{\{\boldsymbol{y}_{i} = \boldsymbol{y}'\}} \cdot \frac{1}{N} \right) \boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}'}^{(1)}(\boldsymbol{y})$$

$$= \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}_{i}}^{(1)}(\boldsymbol{y}).$$

Thus, keeping the same notations as in the previous proof, previous calculations give:

$$\begin{split} \mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\bar{Y}_{N}}^{(1)}}[Y] &= \frac{1}{N} \sum_{i=1}^{N} \mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}_{i}}^{(1)}}[Y] \\ &= \frac{1}{N} \sum_{i=1}^{N} \nabla_{\boldsymbol{\theta}} F_{\boldsymbol{y}_{i}}(\boldsymbol{\theta}) \\ &= \nabla_{\boldsymbol{\theta}} \left(\frac{1}{N} \sum_{i=1}^{N} F_{\boldsymbol{y}_{i}} \right) (\boldsymbol{\theta}). \end{split}$$

Let $F_{\bar{Y}_N} \coloneqq \frac{1}{N} \sum_{i=1}^N F_{m{y}_i}$ Then, the exact same arguments as in the supervised case hold, and the results of Proposition C.2 are obtained by replacing $F_{m{y}}$ by $F_{\bar{Y}_N}$ in the proof of Proposition C.1, and noticing that the previously shown $\mathbb{E}_{q(m{y}_i,\cdot)}||Y-m{y}_i||_2^2/t$ -smoothness of $F_{m{y}_i}$ gives that $F_{\bar{Y}_N}$ is $\frac{1}{N}\sum_{i=1}^N \mathbb{E}_{q(m{y}_i,\cdot)}||Y-m{y}_i||_2^2/t$ -smooth. Similar arguments also hold for the regularized optimization formulation, by noting that this time we have $\nabla F_{\bar{Y}_N}(\mathbb{R}^d) \subset \text{conv}\left(\bigcup_{i=1}^N \{\mathcal{N}(m{y}_i) \cup \{m{y}_i\}\}\right)$.

E.8 Proof of Proposition C.3

935

Proof. The first point is directly given by the fact that $\mathbb{E}_{p_{\theta,y}^{(1)}}[Y]$ is the expectation of a distribution

over $\mathcal{N}(y) \cup \{y\}$. For the second and third points, as derived in Appendix E.7, we have:

$$\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] = \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}_{\boldsymbol{y}}^{-}(\boldsymbol{\theta})} q(\boldsymbol{y}',\boldsymbol{y}) \exp\left(\left[\langle \boldsymbol{\theta}, \boldsymbol{y}' - \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}') - \varphi(\boldsymbol{y})\right]/t\right) \cdot (\boldsymbol{y}' - \boldsymbol{y}) \\ + \sum_{\boldsymbol{y}' \in \mathcal{N}_{\boldsymbol{y}}^{+}(\boldsymbol{\theta})} q(\boldsymbol{y}, \boldsymbol{y}') \cdot (\boldsymbol{y}' - \boldsymbol{y}) \ .$$

938 Define then:

$$\mathcal{N}_{\text{better}}(\boldsymbol{y}) \coloneqq \left\{ \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}) \mid \langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}') > \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}) \right\},$$

$$\mathcal{N}_{\text{worse}}(\boldsymbol{y}) \coloneqq \left\{ \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}) \mid \langle \boldsymbol{\theta}, \boldsymbol{y}' \rangle + \varphi(\boldsymbol{y}') < \langle \boldsymbol{\theta}, \boldsymbol{y} \rangle + \varphi(\boldsymbol{y}) \right\}$$

as the sets of improving and worsening neighbors of y respectively (assuming no neighbor of y has exactly equal objective value for simplicity, which is true almost everywhere w.r.t. $\theta \in \mathbb{R}^d$).

941 **Low temperature limit.** We have:

$$\mathcal{N}^+_{m{y}}(m{ heta}) \xrightarrow[t o 0^+]{} \mathcal{N}_{ ext{better}}(m{y}), \quad ext{and} \quad \mathcal{N}^-_{m{y}}(m{ heta}) \xrightarrow[t o 0^+]{} \mathcal{N}_{ ext{worse}}(m{y}).$$

Then, as $x < 0 \implies \exp(x/t) \xrightarrow[t \to 0^+]{} 0$, we have effectively

$$\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] \xrightarrow[t \to 0^+]{} \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}_{\text{better}}(\boldsymbol{y})} q(\boldsymbol{y},\boldsymbol{y}') \cdot (\boldsymbol{y}' - \boldsymbol{y}).$$

942 **High temperature limit.** As $\forall x \in \mathbb{R}, \ \exp(x/t) \xrightarrow[t \to +\infty]{} 1$, we have:

$$\mathcal{N}^+_{\boldsymbol{y}}(\boldsymbol{\theta}) \xrightarrow[t \to +\infty]{} \left\{ \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}) \mid q(\boldsymbol{y}', \boldsymbol{y}) > q(\boldsymbol{y}, \boldsymbol{y}') \right\}, \quad \text{and} \quad \mathcal{N}^-_{\boldsymbol{y}}(\boldsymbol{\theta}) \xrightarrow[t \to +\infty]{} \left\{ \boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y}) \mid q(\boldsymbol{y}', \boldsymbol{y}) \leq (\boldsymbol{y}, \boldsymbol{y}') \right\}.$$

943 Thus, we have:

945

$$\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] \xrightarrow[t \to +\infty]{} \boldsymbol{y} + \sum_{\boldsymbol{y}' \mid q(\boldsymbol{y}',\boldsymbol{y}) \leq (\boldsymbol{y},\boldsymbol{y}')} q(\boldsymbol{y}',\boldsymbol{y}) \cdot (\boldsymbol{y}'-\boldsymbol{y}) + \sum_{\boldsymbol{y}' \mid q(\boldsymbol{y}',\boldsymbol{y}) > (\boldsymbol{y},\boldsymbol{y}')} q(\boldsymbol{y},\boldsymbol{y}') \cdot (\boldsymbol{y}'-\boldsymbol{y}),$$

944 which gives effectively:

$$\mathbb{E}_{\boldsymbol{p}_{\boldsymbol{\theta},\boldsymbol{y}}^{(1)}}[Y] \xrightarrow[t \to +\infty]{} \boldsymbol{y} + \sum_{\boldsymbol{y}' \in \mathcal{N}(\boldsymbol{y})} \min\left[q(\boldsymbol{y},\boldsymbol{y}'),q(\boldsymbol{y}',\boldsymbol{y})\right] \cdot (\boldsymbol{y}'-\boldsymbol{y}).$$

46 NeurIPS Paper Checklist

1. Claims

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962 963

964

965

966

967

968

969

970

971

972

973

974

975

976

977 978

980

981

982

983

984

985

986

987

988

990

992

993

994

995

996

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: In the abstract and introduction, we clearly outline the specific settings addressed in this paper and the corresponding contributions.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the
 contributions made in the paper and important assumptions and limitations. A No or
 NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: We include a discussion on the limitations of this work in the conclusion.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was
 only tested on a few datasets or with a few runs. In general, empirical results often
 depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

997 Answer: [Yes]

Justification: All proofs are included in the appendix, and provide the full set of needed assumptions.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if
 they appear in the supplemental material, the authors are encouraged to provide a short
 proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented
 by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: The main contributions of the paper are high-level training algorithms, which are described in detail. Additionally, the numerical experiments are carefully documented to ensure they are as reproducible as possible.

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The data of the *EURO Meets NeurIPS 2022 Vehicle Routing Competition* is accessible online. Apart from this, we only use synthetic data, for which the generation process is detailed. We will release the code upon acceptance of the paper.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
 including code, unless this is central to the contribution (e.g., for a new open-source
 benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new
 proposed method and baselines. If only a subset of experiments are reproducible, they
 should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyperparameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We present the experimental setting with sufficient level of detail to fully appreciate the results in the core of the paper, and give full details in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We average experiments over multiple runs with varying seeds and report error bars and statistical significance statements in consequence.

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how
 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

1103

1104

1105

1106

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

Justification: We include sufficient relevant information on the compute resources.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: We have carefully reviewed the NeurIPS Code of Ethics and ensured that all aspects of the research comply with its guidelines.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: As this is a foundational research paper, we do not foresee direct negative societal impacts in its current form. Potential negative impacts would depend on downstream applications, which are beyond the scope of this work.

- The answer NA means that there is no societal impact of the work performed.
 - If the authors answer NA or No, they should explain why their work has no societal
 impact or why the paper does not address societal impact.
 - Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
 - The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: this work does not pose such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We properly credit original papers for the code we use (Vidal [49] and Baty et al. [4]), and proper credit will be also given in the released code.

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

1206

1207

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230 1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

1242

1243

1244

1245

1246

1247

1248

1250

1251

1252

1253

1254

1255

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No assets are released.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: This paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects. Depending on the country in which research is conducted, IRB approval (or equivalent)
 may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 We recognize that the procedures for this may vary significantly between institutions

1260

1261

1262

1263

- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.