

# On counterfactual inference with unobserved confounding

Abhin Shah<sup>1</sup> Raaz Dwivedi<sup>2</sup> Devavrat Shah<sup>1</sup> Gregory W. Wornell<sup>1</sup>  
abhin@mit.edu, dwivedi@cornell.edu, {devavrat, gww}@mit.edu

<sup>1</sup>MIT and <sup>2</sup>Cornell Tech  
September 18, 2023

## Abstract

Given an observational study with  $n$  independent but heterogeneous units, our goal is to learn the counterfactual distribution for each unit using only one  $p$ -dimensional sample per unit containing covariates, interventions, and outcomes. Specifically, we allow for unobserved confounding that introduces statistical biases between interventions and outcomes as well as exacerbates the heterogeneity across units. Modeling the conditional distribution of the outcomes as an exponential family, we reduce learning the unit-level counterfactual distributions to learning  $n$  exponential family distributions with heterogeneous parameters and only one sample per distribution. We introduce a convex objective that pools all  $n$  samples to jointly learn all  $n$  parameter vectors, and provide a unit-wise mean squared error bound that scales linearly with the metric entropy of the parameter space. For example, when the parameters are  $s$ -sparse linear combination of  $k$  known vectors, the error is  $O(s \log k/p)$ . En route, we derive sufficient conditions for compactly supported distributions to satisfy the logarithmic Sobolev inequality. As an application of the framework, our results enable consistent imputation of sparsely missing covariates.

## 1 Introduction

We are interested in the problem of unit-level counterfactual inference owing to the increasing importance of personalized decision-making in many domains. As a motivating example, consider an observational dataset corresponding to an interaction between a recommender system and a user over time. At each time, the user was exposed to a product based on observed demographic factors as well as factors that are not observed in the dataset, e.g., user’s energy level (i.e., whether they’re feeling energetic or tired). Additionally, at each time, the user’s engagement level, which could have sequentially depended on the prior interaction in addition to the ongoing interaction, was also recorded. Also, the system could have sequentially adapted its recommendation. Given such data of many heterogeneous users (e.g., a movie recommender system for a streaming media platform), we want to infer each user’s average engagement level if it were exposed to a different sequence of products while the observed and the unobserved factors remain unchanged. This task is challenging since: (a) the unobserved factors could give rise to spurious associations, (b) the users could be heterogeneous in that they may have different responses to same sequence of products, and (c) each user provides a single interaction trajectory.

More generally, to address problems of this kind, we consider an observational setting where a unit undergoes multiple interventions (or treatments) denoted by  $\mathbf{a}$ . We denote the outcomes of interest by  $\mathbf{y}$ , and allow the interventions  $\mathbf{a}$  and the outcomes  $\mathbf{y}$  to be confounded by observed covariates  $\mathbf{v}$  as well as unobserved covariates  $\mathbf{z}$ . The graphical structure shown in Figure. 1(a) captures these

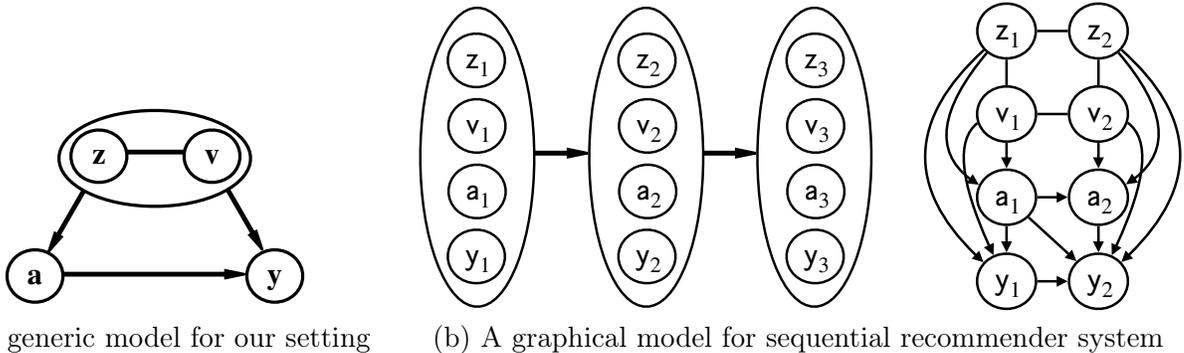


Figure 1: Graphical models covered by our methodology. Directed arrows denote causation and undirected arrows denote association. Thin arrows denote low-level causal links and thick arrows denote high-level causal links, i.e., aggregated thin arrows. Our methodology does not assume knowledge of low-level causal links and is applicable to any graphical model with high-level causal links between variables as in panel (a). Panel (b) presents an example of a sequential recommender system (consistent with the model in panel (a)) interacting with a user at 3 time points where  $z_t$ ,  $v_t$ ,  $a_t$ , and  $y_t$  denote the user’s unobserved energy levels, observed demographic factors, the product exposed to the user, and the user’s engagement level, respectively, at time  $t$ . The left subplot illustrates the high-level dependency between the variables while the right subplot expands on it for time 1 and 2.

interactions and is at the heart of our problem. In the recommender system example above, a unit corresponds to a user,  $\mathbf{a}$  corresponds to the products recommended,  $\mathbf{y}$  corresponds to the engagement levels,  $\mathbf{v}$  corresponds to the observed demographic factors, and  $\mathbf{z}$  corresponds to the unobserved energy levels (see Figure. 1(b)). We consider  $n$  heterogeneous and independent units indexed by  $i \in [n] \triangleq \{1, \dots, n\}$ , and assume access to one observation per unit with  $(\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)})$  denoting the realizations of  $(\mathbf{v}, \mathbf{a}, \mathbf{y})$  for unit  $i$ .

We operate within the Neyman-Rubin potential outcomes framework (Neyman, 1923; Rubin, 1974) and denote the potential outcome of unit  $i \in [n]$  under interventions  $\mathbf{a}$  by  $\mathbf{y}^{(i)}(\mathbf{a})$ . Given the realizations  $\{(\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$ , our goal is to answer counterfactual questions for these  $n$  units. For example, what would the potential outcomes  $\mathbf{y}^{(i)}(\tilde{\mathbf{a}}^{(i)})$  for interventions  $\tilde{\mathbf{a}}^{(i)} \neq \mathbf{a}^{(i)}$  be, while the observed and unobserved covariates remain unchanged? Under the graphical model in Figure. 1(a) and the stable unit treatment value assumption (SUTVA), i.e., the potential outcomes of unit  $i$  are not affected by the interventions at other units, learning unit-level counterfactual distributions is equivalent to learning unit-level conditional distributions

$$\left\{ f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}(\mathbf{y} = \cdot | \mathbf{a} = \cdot, \mathbf{z}^{(i)}, \mathbf{v}^{(i)}) \right\}_{i=1}^n. \quad (1)$$

Here, the  $i$ -th distribution represents the conditional distribution for the outcomes  $\mathbf{y}$  as a function of the interventions  $\mathbf{a}$ , while keeping the observed covariates  $\mathbf{v}$  and the unobserved covariates  $\mathbf{z}$  fixed at the corresponding realizations for unit  $i$ , i.e.,  $\mathbf{v}^{(i)}$  and  $\mathbf{z}^{(i)}$ , respectively.

Such questions cannot be answered without structural assumptions due to two key challenges: (a) unobserved confounding and (b) single observation per unit. First, the unobserved covariates  $\mathbf{z}$  introduce spurious statistical dependence between interventions and outcomes, termed unobserved confounding, which results in biased estimates. Second, we only observe one realization, namely the

outcomes  $\mathbf{y}^{(i)}(\mathbf{a}^{(i)})$  under the interventions  $\mathbf{a}^{(i)}$ , that is consistent with the unit-level conditional distribution  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}(\mathbf{y}|\mathbf{a},\mathbf{z}^{(i)},\mathbf{v}^{(i)})$ . As a result, we need to learn  $n$  heterogeneous conditional distributions while having access to only one sample from each of them.

In this work, we model the conditional distribution of the outcomes of interest conditioned on the unobserved covariates, the observed covariates, the intervention as an exponential family distribution motivated by the principle of maximum entropy.<sup>1</sup> With this model structure, we show that both the aforementioned challenges can be tackled. In particular, we show that the  $n$  unit-level conditional distributions in (1) lead to  $n$  distributions from the same exponential family, albeit with parameters that vary across units. The parameter corresponding to the  $i^{\text{th}}$  unit, for brevity in terminology denoted by  $\gamma^{(i)}$  (defined later), captures the effect of  $\mathbf{z}^{(i)}$  and helps tackle the challenge of unobserved confounding. However, the challenge still remains to learn  $n$  heterogeneous exponential family distributions with one sample per distribution. This challenge has been addressed in two specific scenarios in the literature: (a) if the unobserved confounding is identical across units, i.e., the parameters  $\{\gamma^{(i)}\}_{i=1}^n$  were all equal, then the challenge boils down to learning parameters of a single exponential family distribution from  $n$  samples, which has been well-studied (cf. [Shah et al. \(2021b\)](#) for an overview); (b) if  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  take binary values and have pairwise interactions, then the challenge boils down to learning parameters of an Ising model (a special sub-class of exponential family defined later) with one sample. This specific challenge has been studied under restricted settings: (i) where the dependencies between the variables are known (e.g., [Kandiros et al. \(2021\)](#); [Mukherjee et al. \(2021\)](#)) and (ii) where a specific subset of the parameters are known ([Dagan et al., 2021](#)). In this work, we consider a generalized setting where  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  can be either discrete, continuous, or both, and do not assume that the underlying dependencies or a specific subset of parameters are known.

**Summary of contributions** This work introduces a method to learn unit-level counterfactual distributions from observational studies, in the presence of unobserved confounding, with one sample per unit, using exponential family modeling. For every unit  $i \in [n]$ , we reduce learning its counterfactual distribution to learning the unit-specific parameter  $\gamma^{(i)}$  with access to one sample  $(\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)})$  from unit  $i$ . Here,  $\{\gamma^{(1)}, \dots, \gamma^{(n)}\}$  are parameters of  $n$  different distributions from the same exponential family. The specific technical contributions are as follows:

1. We introduce a convex (and strictly proper) loss function (Definition. 1) that pools the data  $\{(\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^n$  across all  $n$  samples to jointly learn all  $n$  parameters  $\{\gamma^{(i)}\}_{i=1}^n$ .
2. For every unit  $i$ , we prove that the mean squared errors of our estimates of (a)  $\gamma^{(i)}$  (Theorem. 1) and (b) the expected potential outcomes under alternate interventions (Theorem. 2) scale linearly with the metric entropy of the underlying parameter space. For instance, when  $\gamma^{(i)}$  is  $s$ -sparse linear combination of  $k$  known vectors (Corollary. 1), the error—just with one sample—decays as  $O(s \log k/p)$ , where  $p$  is the dimension of the tuple  $(\mathbf{v}, \mathbf{a}, \mathbf{y})$ .
3. We apply our method to impute missing covariates when they are sparse. Formally, we consider a setup (with no systematically unobserved covariates) where the observed covariates are entirely missing for some fixed fraction of the units. Specifically, for unit  $i$  with missing covariates, only  $(\mathbf{a}^{(i)}, \mathbf{y}^{(i)})$  is observed. For every such unit, we show that our method can recover the

---

<sup>1</sup>Exponential family distributions are the maximum entropy distributions given linear constraints on distributions such as specifying the moments (see [Jaynes \(1957\)](#)).

missing covariates with the mean squared error decaying as  $O(p_v^2/p)$ , where  $p_v$  and  $p$  are the dimensions of  $\mathbf{v}$  and  $(\mathbf{v}, \mathbf{a}, \mathbf{y})$ , respectively (Proposition. 2).

4. Methodologically, our work advances three threads: (a) learning Ising models (and their extensions to discrete, continuous, or mixed variables) from a single sample, where we learn the dependencies between variables, generalizing prior work [Kandiros et al. \(2021\)](#); [Dagan et al. \(2021\)](#), (b) learning Markov random fields (a sub-class of exponential family) from multiple independent but non-identical samples, generalizing prior work [Vuffray et al. \(2016, 2022\)](#); [Shah et al. \(2021a\)](#), and (c) learning counterfactual outcomes with an exponential family model, allowing each unit to have different unobserved covariates and providing unit-level guarantees instead of average-level, generalizing [Arkhangelsky and Imbens \(2018\)](#).
5. In our analysis, we (a) derive sufficient conditions for a continuous random vector supported on a compact set to satisfy the logarithmic Sobolev inequality (Proposition. F.1) and (b) provide new concentration bounds for arbitrary functions of a continuous random vector that satisfies the logarithmic Sobolev inequality (Proposition. F.2). These results may be of independent interest.

**Outline** Section. 2 discusses background and related work. We discuss our formulation and algorithm in Section. 3 and present their analysis in Section. 4. We develop an application of our methodology to impute missing covariates in Section. 6. We sketch the proof of our main result in Section. 7 with detailed proofs deferred to the appendices. We conclude with a discussion in Section. 8.

**Notation** For any positive integer  $n$ , let  $[n] := \{1, \dots, n\}$ . For a deterministic sequence  $u_1, \dots, u_n$ , we let  $\mathbf{u} := (u_1, \dots, u_n)$ . For a random sequence  $u_1, \dots, u_n$ , we let  $\mathbf{u} := (u_1, \dots, u_n)$ . For a vector  $\mathbf{u} \in \mathbb{R}^p$ , we use  $u_t$  to denote its  $t^{\text{th}}$  coordinate and  $u_{-t} \in \mathbb{R}^{p-1}$  to denote the vector after deleting the  $t^{\text{th}}$  coordinate. We denote the  $\ell_0$ ,  $\ell_p$  ( $p \geq 1$ ), and  $\ell_\infty$  norms of a vector  $\mathbf{v}$  by  $\|\mathbf{v}\|_0$ ,  $\|\mathbf{v}\|_p$ , and  $\|\mathbf{v}\|_\infty$ , respectively. For a matrix  $\mathbf{M} \in \mathbb{R}^{p \times p}$ , we denote the element in  $t^{\text{th}}$  row and  $u^{\text{th}}$  column by  $\mathbf{M}_{tu}$ , the  $t^{\text{th}}$  row by  $\mathbf{M}_t$ , and the vector obtained after deleting  $\mathbf{M}_{tt}$  from  $\mathbf{M}_t$  by  $\mathbf{M}_{t,-t}$ . Further, we denote the matrix maximum norm by  $\|\mathbf{M}\|_{\max}$ , the Frobenius norm by  $\|\mathbf{M}\|_F$ , the spectral norm (operator 2-norm) by  $\|\mathbf{M}\|_{\text{op}}$ , the induced 1-norm (operator 1-norm) by  $\|\mathbf{M}\|_1$ , the induced  $\infty$ -norm (operator  $\infty$ -norm) by  $\|\mathbf{M}\|_\infty$ , and the  $(2, \infty)$ -norm by  $\|\mathbf{M}\|_{2,\infty}$ . Finally, for vectors  $\hat{\mathbf{u}} \in \mathbb{R}^p$  and  $\tilde{\mathbf{u}} \in \mathbb{R}^p$ , the mean squared error between  $\hat{\mathbf{u}}$  and  $\tilde{\mathbf{u}}$  is defined as  $\text{MSE}(\hat{\mathbf{u}}, \tilde{\mathbf{u}}) \triangleq p^{-1} \sum_{t \in [p]} (\hat{u}_t - \tilde{u}_t)^2$ .

## 2 Background and related work

This work builds on two vast bodies of literature: exponential family learning and unit-level counterfactual inference with unobserved confounding. For a detailed literature overview of the former, we refer the readers to [Bresler \(2015\)](#); [Klivans and Meka \(2017\)](#); [Vuffray et al. \(2022\)](#); [Shah et al. \(2021a\)](#) (for a special sub-class, Markov random fields (MRFs)<sup>2</sup>) and [Shah et al. \(2021b\)](#) for general exponential families. For an introduction to counterfactual inference, see the books [Imbens and Rubin \(2015\)](#); [Hernán and Robins \(2020\)](#) for settings with no unobserved confounding and [Pearl \(2009\)](#); [Pearl et al. \(2016\)](#) for settings with known causal mechanism (in the form of a causal graph).

<sup>2</sup>MRFs can be naturally represented as exponential family distributions with certain sparsity constraints on the parameters via the principle of maximum entropy ([Wainwright et al., 2008](#)).

**Exponential family learning** There is a series of works for learning Ising models, a special MRF with binary variables and an instance of a pair-wise exponential family, from a single sample. Such a model has two distinct sets of parameters capturing the contribution of nodes and edges in the underlying undirected graph, referred to as the external field and the interaction matrix.<sup>3</sup> Many strategies exist for learning such a model when the interaction matrix is known up to a constant and under varying assumptions on the external field; see, e.g., Chatterjee (2007); Bhattacharya and Mukherjee (2018); Daskalakis et al. (2019); Ghosal and Mukherjee (2020); Kandiros et al. (2021); Mukherjee et al. (2021). More recently, Dagan et al. (2021) provide guarantees for learning the interaction matrix from a single sample when the external field is known. Kandiros et al. (2021) and Mukherjee et al. (2021) extend the tools in Dagan et al. (2021) to learn the external field for an Ising model with a known interaction matrix (up to a scalar multiple). Notably, all of these works are based on the pseudo-likelihood estimation (Besag, 1975). Our work extends the techniques and results from Dagan et al. (2021) to learn the external field from one sample of continuous variables with an estimated interaction matrix.

Vuffray et al. (2016) introduced a novel M-estimation-based loss function for learning Ising models from many independent and identically distributed samples. Vuffray et al. (2022) and Shah et al. (2021a) generalize it to learn general MRFs with multi-ary discrete and continuous variables, respectively. Ren et al. (2021) showed that this loss function has superior numerical performance compared to the ones based on pseudo-likelihood. We contribute to this line of work by generalizing that loss function further to learn MRFs with discrete, continuous, and mixed variables with independent but not identically distributed samples.

For settings closer to our work, namely, exponential families with unobserved variables, the two common modeling approaches include restricted Boltzmann machines (Bresler et al., 2019; Goel, 2020; Bresler and Buhai, 2020) and latent variable Gaussian graphical models; see, e.g., Chandrasekaran et al. (2012); Ma et al. (2013); Vinyes and Obozinski (2018); Wang et al. (2023). While the former assumes a bipartite structure with edges only across observed and unobserved variables, the latter imposes a Gaussian generative model. In this thread, most related to our set-up is the work by Taeb et al. (2020) as they model the conditional distribution of the observed variables conditioned on the unobserved variables as an exponential family similar to us. They provide empirically promising results for recovering the underlying graph and the number of unobserved variables (assumed to be small), albeit with limited theoretical guarantees. In contrast, here we provide parameter estimation error in the presence of unobserved variables (notably, we cover all the models they considered).

**Unit-level counterfactual inference** Recent years have seen an active interest in developing different strategies for unit-level inference with unobserved confounding.

For the settings with univariate outcomes for each unit, a common approach to deal with unobserved confounding is the instrumental variable (IV) method (Imbens and Angrist, 1994) when one has access to a variable—the IV—that induces changes in intervention assignment but has no independent effect on outcomes allowing causal effect estimation. Recent works for IV methods with unit-level inference include Hartford et al. (2017); Athey et al. (2019); Syrgkanis et al. (2019); Singh et al. (2019); Xu et al. (2020); Semenova and Chernozhukov (2021); Wang et al. (2022). Another approach for univariate outcomes, called causal sensitivity analysis (Rosenbaum and Rubin, 1983), estimates the worst-case effect on the causal estimand as a function of the extent of unobserved

---

<sup>3</sup>E.g., in our model (defined later in (2)),  $\phi$  and  $\Phi$  correspond to the external field and the interaction matrix, respectively.

confounding in a given dataset under varying assumptions on the generative model. For such analysis with unit-level guarantees, see, e.g., [Yadlowsky et al. \(2022\)](#); [Kallus et al. \(2019\)](#); [Yin et al. \(2022\)](#); [Jin et al. \(2023\)](#); [Jesson et al. \(2021\)](#).

Closer to our work are those on panel or longitudinal data settings, where one observes multiple outcomes for each unit. For linear panel data settings, a common approach is factor modeling, where potential outcomes and interventions (binary or multi-ary) are assumed to be independent conditional on some latent factors. See, e.g., difference-in-difference methods ([Bertrand et al., 2004](#); [Angrist and Pischke, 2009](#)), synthetic control ([Abadie and Gardeazabal, 2003](#); [Abadie et al., 2010](#)), its variants [Arkhangelsky et al. \(2021\)](#); [Dwivedi et al. \(2022b\)](#), and extensions to multi-ary interventions in synthetic interventions ([Agarwal et al., 2020](#)) and sequential experiments ([Dwivedi et al., 2022a](#)). For non-linear panel data settings, the most commonly used models include probit, logit, Poisson, negative binomial, proportional hazard, and tobit models (see [Fernández-Val and Weidner \(2018\)](#) for an overview) where some parametric model characterises the distribution of the outcomes conditional on the unobserved covariates, the observed covariates, and the interventions. Notably, these works on linear and non-linear panel data directly estimate effects (averaged over all observed and unobserved covariates or unit-level for given observed and unobserved covariates) for finitely many interventions when the intervention assignment has special structure, while we focus on learning the counterfactual distributions while allowing for multi-ary discrete and continuous interventions without any special structure. In this thread, our work is most related to [Arkhangelsky and Imbens \(2018\)](#), who also use an exponential family to model the unit-wise distribution of the observed covariates and interventions conditioned on the unobserved covariates. They connect this model to the commonly used fixed effects model for the outcomes in latent factor modeling ([Angrist and Pischke, 2009](#)), and provide estimates for the average treatment effect given multiple units with the same set of unobserved covariates. Our work generalizes their set-up by allowing each unit to have a different set of unobserved covariates and provides the first unit-level counterfactual inference guarantee with an exponential family model.

### 3 Problem formulation and algorithm

This section formalizes the problem, specifies our model, and defines the inference tasks of interest.

#### 3.1 Underlying causal mechanism and counterfactual distributions

We consider a counterfactual inference task where units go through  $p_a \geq 1$  interventions. For every unit, we observe  $p_y \geq 1$  outcomes of interest. The interventions and the outcomes could be confounded by  $p_v \geq 0$  observed covariates as well as  $p_z \geq 0$  unobserved covariates. Additionally, the observed covariates and the unobserved covariates could be arbitrarily associated. We denote the random vector associated with the interventions, the outcomes, the observed covariates, and the unobserved covariates by  $\mathbf{a} \triangleq (a_1, \dots, a_{p_a}) \in \mathcal{A}^{p_a}$ ,  $\mathbf{y} = (y_1, \dots, y_{p_y}) \in \mathcal{Y}^{p_y}$ ,  $\mathbf{v} \triangleq (v_1, \dots, v_{p_v}) \in \mathcal{V}^{p_v}$ , and  $\mathbf{z} \triangleq (z_1, \dots, z_{p_z}) \in \mathcal{Z}^{p_z}$ , respectively, where  $\mathcal{A}, \mathcal{Y}, \mathcal{V}$ , and  $\mathcal{Z}$  denote the support of interventions, outcomes, observed covariates, and unobserved covariates, respectively. We allow these sets to contain discrete, continuous, or mixed values.

**Causal mechanism** We summarize the causal relationship between the random vectors  $\mathbf{z}$ ,  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  in Figure. 1(a) where we denote the arbitrary association between  $\mathbf{z}$  and  $\mathbf{v}$  by a undirected arrow, and the causal association between (i)  $(\mathbf{z}, \mathbf{v})$  and  $\mathbf{a}$ , (ii)  $(\mathbf{z}, \mathbf{v})$  and  $\mathbf{y}$ , and (iii)  $\mathbf{a}$  and  $\mathbf{y}$  by directed arrows. More generally, we are interested in any setup consistent with the graphical model in

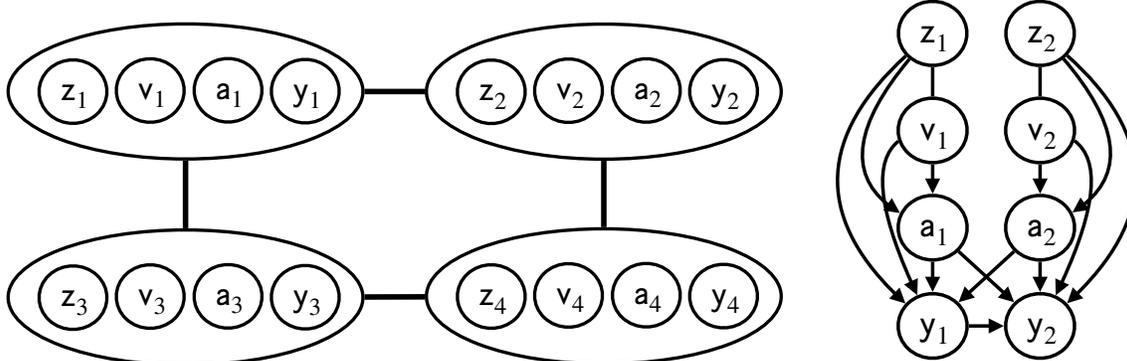


Figure 2: A graphical model for a single unit in the network setting with 4 users; arrows have same meaning as in Figure. 1. Here  $v_t$ ,  $z_t$ ,  $a_t$ , and  $y_t$  denote user  $t$ 's observed factors, unobserved factors, exposed product, and engagement level, respectively. The left plot illustrates the high-level dependency between the variables of different users in the network, and the right plot expands on it for (user 1, user 2) pair. Analogous dependencies exist for (user 1, user 3), (user 2, user 4), and (user 3, user 4) pairs.

Figure. 1(a). We assume access to  $n$  independent realizations indexed by  $i \in [n]$ :  $\mathbf{v}^{(i)}$ ,  $\mathbf{a}^{(i)}$ , and  $\mathbf{y}^{(i)}$  denote the realizations of  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  for unit  $i$ , respectively. For every realized tuple  $(\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)})$ , there is a corresponding realization  $\mathbf{z}^{(i)}$  of the unobserved covariates  $\mathbf{z}$  that is unobserved. Next, we discuss some examples covered by our framework.

**Examples: sequential and network settings** While Figure. 1(a) exhibits the high-level causal links between  $\mathbf{z}$ ,  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$ , there could be complex low-level causal links between elements of these vectors. We do not assume any knowledge of such low-level causal links. In Figure. 1(b), we provide an instance of a sequential setting covered by our work where every unit's (i)  $a_{t+1}$  depends on  $a_t$  in addition to  $v_{t+1}$  and  $\mathbf{z}$ , and (ii)  $y_{t+1}$  depends on  $a_t$  and  $y_t$  in addition to  $a_{t+1}$ ,  $v_{t+1}$  and  $\mathbf{z}$ . Another classical example covered by our framework includes the network setting where a unit represents a social network where users are linked to each other by interpersonal relationships as shown in Figure. 2. Similar to the sequential recommender system, every user was exposed to a product based on observed demographic factors as well as certain unobserved factors, and the user's engagement level was recorded. The engagement level of user  $t$ , i.e.,  $y_t$ , depended its observed demographic factors  $v_t$ , its unobserved factors  $z_t$ , its exposed product  $a_t$  as well as on the product exposed to its neighbor  $u$ , i.e.,  $a_u$ . Further,  $y_t$  could have been associated with  $y_u$ .

**Unit-level counterfactual distributions** We denote the Neyman-Rubin potential outcomes of unit  $i \in [n]$  under interventions  $\mathbf{a} \in \mathcal{A}^{p_a}$  by  $\mathbf{y}^{(i)}(\mathbf{a})$ . We make the stable unit treatment value assumption (SUTVA) (Rubin, 1980) for the observed outcome, i.e.,  $\mathbf{y}^{(i)} = \mathbf{y}^{(i)}(\mathbf{a}^{(i)})$  for all  $i \in [n]$ . For independent units with the causal mechanism and SUTVA assumed here, the unit-level counterfactual distributions are equivalent to certain unit-level conditional distributions as we now argue. Consider unit  $i \in [n]$  and fix the observed covariates and the unobserved covariates at  $\mathbf{v}^{(i)}$  and  $\mathbf{z}^{(i)}$ , respectively. Then, let  $\tilde{\mathbf{y}}^{(i)}$  be a realization of  $\mathbf{y}$  when  $\mathbf{a} = \tilde{\mathbf{a}}^{(i)}$ . We are interested in the distribution of the potential outcomes of unit  $i$  for interventions  $\tilde{\mathbf{a}}^{(i)}$ , i.e., the distribution of  $\mathbf{y}^{(i)}(\tilde{\mathbf{a}}^{(i)})$  given  $\mathbf{z} = \mathbf{z}^{(i)}$ ,  $\mathbf{v} = \mathbf{v}^{(i)}$ . Under the causal framework considered here (see Figure. 1(a)), it is equivalent to the distribution

of  $\mathbf{y}^{(i)}(\tilde{\mathbf{a}}^{(i)})$  given  $\mathbf{a} = \tilde{\mathbf{a}}^{(i)}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}$  since  $(\mathbf{z}, \mathbf{v})$  satisfy ignorability (Pearl, 2009; Imbens and Rubin, 2015), i.e., the potential outcomes are independent of the interventions given  $(\mathbf{z}, \mathbf{v})$ . Further, under SUTVA, it is equivalent to the distribution of  $\tilde{\mathbf{y}}^{(i)}$  given  $\mathbf{a} = \tilde{\mathbf{a}}^{(i)}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}$ , i.e.,  $f_{\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}}(\mathbf{y} = \cdot | \mathbf{a} = \tilde{\mathbf{a}}^{(i)}, \mathbf{z}^{(i)}, \mathbf{v}^{(i)})$ . Therefore, our goal is to learn the  $n$  unit-level conditional distributions in (1). Now, we proceed to the modeling details.

### 3.2 Exponential family modeling and its consequences

Let  $\mathbf{w} \triangleq (\mathbf{z}, \mathbf{v}, \mathbf{a}, \mathbf{y})$  be the  $\tilde{p}$ -dimensional random vector obtained by concatenating  $\mathbf{z}, \mathbf{v}, \mathbf{a}$  and  $\mathbf{y}$  where  $\tilde{p} \triangleq p_z + p_v + p_a + p_y$ . For notational convenience, we start by modeling the joint probability distribution  $f_{\mathbf{w}}$  as an exponential family and relax this model to the conditional distribution of the outcomes in Section 5.1. In particular, we parameterize  $f_{\mathbf{w}}$  with natural parameters  $\phi \in \mathbb{R}^{\tilde{p} \times 1}$  and  $\Phi \in \mathbb{R}^{\tilde{p} \times \tilde{p}}$ , and natural statistics  $\mathbf{w}$  and  $\mathbf{w}\mathbf{w}^\top$  so that

$$f_{\mathbf{w}}(\mathbf{w}; \phi, \Phi) \propto \exp\left(\phi^\top \mathbf{w} + \mathbf{w}^\top \Phi \mathbf{w}\right), \quad \text{where } \mathbf{w} \triangleq (\mathbf{z}, \mathbf{v}, \mathbf{a}, \mathbf{y}), \quad (2)$$

and  $\mathbf{z} \triangleq (z_1, \dots, z_{p_z}), \mathbf{v} \triangleq (v_1, \dots, v_{p_v}), \mathbf{a} \triangleq (a_1, \dots, a_{p_a})$ , and  $\mathbf{y} \triangleq (y_1, \dots, y_{p_y})$  denote realizations of  $\mathbf{z}, \mathbf{v}, \mathbf{a}$ , and  $\mathbf{y}$ , respectively. Without loss of generality, we can assume  $\Phi$  to be a symmetric matrix. Next, we show that with this modeling assumption, learning unit-level counterfactual distribution can be reduced to learning a suitable exponential family model.

Under the exponential family in (2), the unit-level conditional distribution of  $\mathbf{y}$  conditioned on  $\mathbf{a} = \mathbf{a}, \mathbf{z} = \mathbf{z}$ , and  $\mathbf{v} = \mathbf{v}$  is an exponential family model with natural statistics  $\mathbf{y}$  and  $\mathbf{y}\mathbf{y}^\top$  and

$$f_{\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}}(\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}) \propto \exp\left([\phi^{(y)}]^\top + 2\mathbf{z}^\top \Phi^{(z,y)} + 2\mathbf{v}^\top \Phi^{(v,y)} + 2\mathbf{a}^\top \Phi^{(a,y)}\right] \mathbf{y} + \mathbf{y}^\top \Phi^{(y,y)} \mathbf{y}\right), \quad (3)$$

where  $\phi^{(y)} \in \mathbb{R}^{p \times 1}$  is the component of  $\phi$  corresponding to  $\mathbf{y}$  and  $\Phi^{(u,y)} \in \mathbb{R}^{p_u \times p_y}$  is the component of  $\Phi$  corresponding to  $\mathbf{u}$  and  $\mathbf{y}$  for all  $\mathbf{u} \in \{\mathbf{z}, \mathbf{v}, \mathbf{a}, \mathbf{y}\}$ .<sup>4</sup> We make two key observations: (a) the term  $\Phi^{(z,y)\top} \mathbf{z}$  captures the effect of unobserved covariates  $\mathbf{z}$  on  $f_{\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}}(\mathbf{y} = \cdot | \mathbf{a} = \cdot, \mathbf{z}, \mathbf{v})$  and (b) the task of learning  $f_{\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}}(\mathbf{y} = \cdot | \mathbf{a} = \cdot, \mathbf{z}, \mathbf{v})$  in (3) as a function of  $\mathbf{a}$  reduces to learning

$$(i) \phi^{(y)} + 2\Phi^{(z,y)\top} \mathbf{z} + 2\Phi^{(v,y)\top} \mathbf{v}, \quad (ii) \Phi^{(a,y)}, \quad \text{and} \quad (iii) \Phi^{(y,y)}. \quad (4)$$

That is, learning the unit-level conditional distribution for unit  $i$  is equivalent to learning

$$\gamma^{(i)} = \{\phi^{(y)} + 2\Phi^{(z,y)\top} \mathbf{z}^{(i)} + 2\Phi^{(v,y)\top} \mathbf{v}^{(i)}, \Phi^{(a,y)}, \Phi^{(y,y)}\}, \quad (5)$$

where the notation  $\gamma^{(i)}$  is the same as in Section 1. We note that, given  $\mathbf{a} = \mathbf{a}, \mathbf{z} = \mathbf{z}$ , and  $\mathbf{v} = \mathbf{v}, \mathbf{y} = \mathbf{a} + \mathbf{z} + \mathbf{v} + \boldsymbol{\eta}$  is one plausible data generating process (DGP) consistent with (3) when the noise variable  $\boldsymbol{\eta}$  has an exponential family distribution. More specifically, this DGP, with  $\boldsymbol{\eta}$  such that  $f(\boldsymbol{\eta}) \propto \exp(\phi^{(y)\top} \boldsymbol{\eta} + \boldsymbol{\eta}^\top \Phi^{(y,y)} \boldsymbol{\eta})$ , results in the conditional distribution in (3) with  $\Phi^{(z,y)} = \Phi^{(v,y)} = \Phi^{(a,y)} = \Phi^{(y,y)}$ .

Next, we argue that learning the three quantities in (4) is subsumed in learning the parameters of the (unit-level) conditional distribution  $f_{\mathbf{x}|\mathbf{z}}$  of the random vector  $\mathbf{x} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y})$  conditioned on  $\mathbf{z} = \mathbf{z}$ . Note that  $f_{\mathbf{x}|\mathbf{z}}$  belongs to an exponential family with natural statistics  $\mathbf{x}$  and  $\mathbf{x}\mathbf{x}^\top$ . For all

<sup>4</sup>The exponential family in (3) is same as the one considered in Taeb et al. (2020, Equation 1.3).

$\mathbf{u} \in \{\mathbf{v}, \mathbf{a}, \mathbf{y}\}$ , let  $\phi^{(u)} \in \mathbb{R}^{p_u \times 1}$  be the component of  $\phi$  corresponding to  $\mathbf{u}$ , and  $\Phi^{(z,u)} \in \mathbb{R}^{p_z \times p_u}$  be the component of  $\Phi$  corresponding to  $\mathbf{z}$  and  $\mathbf{u}$ . Then  $f_{\mathbf{x}|\mathbf{z}}$  can be parameterized as follows:

$$f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta(\mathbf{z}), \Theta) \propto \exp\left([\theta(\mathbf{z})]^\top \mathbf{x} + \mathbf{x}^\top \Theta \mathbf{x}\right), \text{ where } \theta(\mathbf{z}) \triangleq \begin{bmatrix} \phi^{(v)} + 2\Phi^{(z,v)\top} \mathbf{z} \\ \phi^{(a)} + 2\Phi^{(z,a)\top} \mathbf{z} \\ \phi^{(y)} + 2\Phi^{(z,y)\top} \mathbf{z} \end{bmatrix} \in \mathbb{R}^{p \times 1}, \quad (6)$$

$\mathbf{x} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y})$ ,  $p \triangleq p_v + p_a + p_y$  and  $\Theta \in \mathbb{R}^{p \times p}$  denotes the component of  $\Phi$  corresponding to  $\mathbf{x}$ . Given some estimates for  $\theta(\mathbf{z})$  and  $\Theta$ , using their appropriate components also yields an estimate of the three quantities in (4) for any  $\mathbf{v} = \mathbf{v}$ . To summarize, the spurious associations or unobserved confounding between  $\mathbf{a}$  and  $\mathbf{y}$  introduced due to unobserved  $\mathbf{z}$  are fully captured by  $\Phi^{(z,y)\top} \mathbf{z}$  or equivalently by  $\theta(\mathbf{z})$ ; thereby, learning unit-level counterfactual distributions require us to learn these unit-level parameters.

### 3.2.1 Reduced inference task and modeling constraints

Let  $f_{\mathbf{w}}(\cdot; \phi^*, \Phi^*)$  denote the true data generating distribution of  $\mathbf{w}$  in (2), and let  $f_{\mathbf{x}|\mathbf{z}}(\cdot | \mathbf{z}; \theta^*(\mathbf{z}), \Theta^*)$  denote the true distribution of  $\mathbf{x}$  conditioned on  $\mathbf{z} = \mathbf{z}$  in (6). Then, for all  $i \in [n]$ , we note that the realization  $\mathbf{x}^{(i)} \triangleq (\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)})$  is consistent with the conditional distribution  $f_{\mathbf{x}|\mathbf{z}}(\cdot | \mathbf{z}^{(i)}; \theta^*(\mathbf{z}^{(i)}), \Theta^*)$  where we do not observe  $\mathbf{z}^{(i)}$ . Our primary goal is to learn the  $n$  unit-level counterfactual distributions, which as noted above simplifies to estimating the following parameters:

$$(i) \text{ Unit-level } \theta^{*(i)} \triangleq \theta^*(\mathbf{z}^{(i)}) \text{ for } i \in [n], \quad \text{and (ii) Population-level } \Theta^*. \quad (7)$$

Our secondary goal is to estimate the expected potential outcomes for any given unit  $i$  (with  $\mathbf{z} = \mathbf{z}^{(i)}$ ,  $\mathbf{v} = \mathbf{v}^{(i)}$ ) and an alternate intervention  $\tilde{\mathbf{a}}^{(i)}$ :

$$\mu^{(i)}(\tilde{\mathbf{a}}^{(i)}) \triangleq \mathbb{E}[\mathbf{y}^{(i)}(\tilde{\mathbf{a}}^{(i)}) | \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}], \quad (8)$$

where  $\mathbf{y}^{(i)}(\tilde{\mathbf{a}}^{(i)})$  denotes the potential outcomes for unit  $i \in [n]$  under interventions  $\tilde{\mathbf{a}}^{(i)} \in \mathcal{A}^{p_a}$ .

For ease of exposition, we consider bounded continuous sets  $\mathcal{V}$ ,  $\mathcal{A}$ , and  $\mathcal{Y}$  with  $\mathcal{V} = \mathcal{A} = \mathcal{Y} \triangleq \mathcal{X} = [-x_{\max}, x_{\max}]$  for a given  $x_{\max}$ . In Section. 5.3, we consider compact discrete and mixed sets. Throughout this paper, it is convenient to further constrain the model as follows:

**Assumption 1** (Bounded and sparse parameters). *The true model parameters (7) satisfy*

$$\theta^{*(i)} \in \Lambda_\theta \triangleq \{\theta \in \mathbb{R}^{p \times 1} : \|\theta\|_\infty \leq \alpha\} \text{ for all } i \in [n], \quad (9)$$

and

$$\Theta^* \in \Lambda_\Theta \triangleq \{\Theta \in \mathbb{R}^{p \times p} : \Theta = \Theta^\top, \|\Theta\|_{\max} \leq \alpha, \|\Theta\|_\infty \leq \beta\}. \quad (10)$$

While (9) bounds the unit-level parameters (a necessary condition for model identifiability (Santhanam and Wainwright, 2012)), (10) bounds the  $\ell_1$  norm of the interaction of each  $\mathbf{x}_t \in \mathbf{x}$  with every  $\mathbf{x}_u \in \mathbf{x}$  in (6). As a result, Assumption. 1 implies that the exponential family in (6) corresponds to MRFs (see Section. 2), also known as undirected graphical models (defined in Appendix. G). We note that Assumption. 1 is standard in the literature on learning MRFs (Bresler, 2015; Vuffray et al., 2016; Klivans and Meka, 2017; Vuffray et al., 2022; Shah et al., 2021a). We are now ready to state our algorithm.

### 3.3 An efficient algorithm via a convex objective

We first describe our strategy to estimate the parameters in (7). Then, we use the estimated parameters to estimate the expected potential outcomes in (8). We remark that for exponential families considered here, maximum likelihood for parameter estimation is not computationally tractable (Wainwright et al., 2008; Shah et al., 2021b). As a result, we resort to an alternative objective function inspired by the convex loss functions used in Vuffray et al. (2016, 2022); Shah et al. (2021a) as they do not depend on the partition function of the distribution. These loss functions are designed in a specific way (see below for details): (i) the sufficient statistics of the conditional distribution of a variable given all other variables are *centered* by adding appropriate constants, (ii) the loss function is an empirical average of the sum of the inverses of all of these conditional distributions (without the partition function) with *centered* sufficient statistics.

#### 3.3.1 Parameter estimation

Our convex objective function jointly learns all the parameters of interest by pooling the observations across all  $n$  units and exploiting the exponential family structure of  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  conditioned on  $\mathbf{z} = \mathbf{z}$  in (6), i.e., the objective explicitly utilizes the fact that the population-level parameter  $\Theta^*$  is shared across units. In particular, we use the following two steps.

**Centering sufficient statistics of the conditional distribution of a variable** Consider the conditional distribution  $f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}$  of the random variable  $x_t$  conditioned on  $\mathbf{x}_{-t} = \mathbf{x}_{-t}$  and  $\mathbf{z} = \mathbf{z}$  for any  $t \in [p]$ :

$$f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}(x_t|\mathbf{x}_{-t}, \mathbf{z}; \theta_t(\mathbf{z}), \Theta_t) \propto \exp\left([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t + \Theta_{tt}x_t^2\right), \quad (11)$$

where  $\theta_t(\mathbf{z})$  is the  $t^{\text{th}}$  element of  $\theta(\mathbf{z})$ ,  $\Theta_t$  is the  $t^{\text{th}}$  row of  $\Theta$ ,  $\Theta_{tt}$  is the  $t^{\text{th}}$  element of  $\Theta_t$ , and  $\Theta_{t,-t} \triangleq \Theta_t \setminus \Theta_{tt} \in \mathbb{R}^{p-1}$  is the vector obtained after deleting  $\Theta_{tt}$  from  $\Theta_t$ . Then, the sufficient statistics in (11), namely  $x_t$  and  $x_t^2$ , are centered by subtracting their expected value with respect to the uniform distribution on  $\mathcal{X}$  resulting in

$$f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}(x_t|\mathbf{x}_{-t}, \mathbf{z}; \theta_t(\mathbf{z}), \Theta_t) \propto \exp\left([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t + \Theta_{tt}\left(x_t^2 - \frac{x_{\max}^2}{3}\right)\right), \quad (12)$$

as the integral of  $x_t$  and  $x_t^2$  with respect to the uniform distribution on  $\mathcal{X}$  is 0 and  $x_{\max}^2/3$ , respectively. As we see later (in Proposition. 1), this centering ensures that our loss function is a proper loss function as well as leads to connections with the surrogate likelihood (Shah et al., 2021a, Proposition. 4.1). We emphasize that the term  $x_{\max}^2/3$  inside the exponent in (12) is vacuous (as it is a constant) and the distribution in (12) is equivalent to the one in (11).

**Constructing the loss function** Next, the loss function (defined below) is designed to be an empirical average of the sum over  $t \in [p]$  of the inverse of the term in the right hand side of (12).

**Definition 1 (Loss function).** Given the samples  $\{\mathbf{x}^{(i)}\}_{i \in [n]}$ , the loss  $\mathcal{L} : \mathbb{R}^{p \times (n+p)} \rightarrow \mathbb{R}$  is given by

$$\mathcal{L}(\underline{\Theta}) = \frac{1}{n} \sum_{t \in [p]} \sum_{i \in [n]} \exp\left(-[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \Theta_{tt}\left([x_t^{(i)}]^2 - \frac{x_{\max}^2}{3}\right)\right) \quad \text{where} \quad \underline{\Theta} \triangleq \begin{bmatrix} \Theta_1^\top \\ \vdots \\ \Theta_p^\top \end{bmatrix}, \quad (13)$$

and  $\underline{\Theta}_t \triangleq \{\theta_t^{(1)}, \dots, \theta_t^{(n)}, \Theta_t\}$  for  $t \in [p]$ .

Our estimate of  $\underline{\Theta}^*$  (defined analogous to  $\underline{\Theta}$ ) is given by

$$\widehat{\underline{\Theta}} \in \arg \min_{\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta} \mathcal{L}(\underline{\Theta}). \quad (14)$$

We note (14) is a convex optimization problem, and a projected gradient descent algorithm (see Appendix. A.2) returns an  $\epsilon$ -optimal estimate with  $\tau = O(p/\epsilon)$  iterations<sup>5</sup> where  $\widehat{\underline{\Theta}}_\epsilon$  is said to be an  $\epsilon$ -optimal estimate if  $\mathcal{L}(\widehat{\underline{\Theta}}_\epsilon) \leq \mathcal{L}(\widehat{\underline{\Theta}}) + \epsilon$  for any  $\epsilon > 0$ . The loss function  $\mathcal{L}$  admits a notable property (see Appendix. A.1 for the proof).

**Proposition 1 (Proper loss function).** *The loss function  $\mathcal{L}$  is strictly proper, i.e.,  $\underline{\Theta}^* = \arg \min_{\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta} \mathbb{E}_{\mathbf{x}|\mathbf{z}}[\mathcal{L}(\underline{\Theta})]$ .*

Proposition. 1 shows that the solution of the idealized convex program  $\min_{\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta} \mathbb{E}_{\mathbf{x}|\mathbf{z}}[\mathcal{L}(\underline{\Theta})]$  is unique and equal to  $\underline{\Theta}^*$ . In this idealized convex program, conditioned on the realized values of the unobserved covariates of the  $n$  units  $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(n)}$ , the loss function is averaged over all the randomness in the observed covariates, the interventions, and the outcomes. In other words, for every  $i \in [n]$ , the idealized convex program has infinite samples from  $f_{\mathbf{x}|\mathbf{z}}$  with unobserved covariates  $\mathbf{z}$  conditioned to be  $\mathbf{z}^{(i)}$ . Thus, the convex program in (14) can be seen as a single sample version of this idealized program, thereby providing an intuitive justification of our loss function (instead of a maximum likelihood objective, which is not tractable here). As we show later in our proofs (see Section. 7 for an overview), different partial averages on the RHS of (13) also admit useful properties and are critical to our analyses.

We note that loss function in (13) is a generalization of the loss functions used in Vuffray et al. (2016, 2022); Shah et al. (2021a). In particular, if the unobserved confounding is identical across units, i.e.,  $\theta^{*(1)} = \dots = \theta^{*(n)}$ , then  $\mathcal{L}(\underline{\Theta})$  in (13) can be decomposed into  $p$  independent loss functions, one for every  $t \in [p]$ . These decomposed loss functions are identical to the ones used in these prior works.

### 3.3.2 Causal estimate

Given the estimate  $\widehat{\underline{\Theta}}$ , our estimate of the expected potential outcome  $\mu^{(i)}(\tilde{\mathbf{a}}^{(i)})$  under an alternate intervention  $\tilde{\mathbf{a}}^{(i)} \in \mathcal{A}^{p_a}$  (8) is derived as follows: First, we identify  $\widehat{\Phi}^{(u,y)} \in \mathbb{R}^{p_u \times p_y}$  to be the component of  $\widehat{\underline{\Theta}}$  corresponding to  $\mathbf{u}$  and  $\mathbf{y}$  for all  $\mathbf{u} \in \{\mathbf{v}, \mathbf{a}, \mathbf{y}\}$  and  $\widehat{\theta}^{(i,y)} \in \mathbb{R}^{p_y}$  to be the component of  $\widehat{\theta}^{(i)}$  corresponding to  $\mathbf{y}$ . Next, we estimate the conditional distribution of  $\mathbf{y}$  for unit  $i$  as a function of the interventions  $\mathbf{a}$ , while keeping  $\mathbf{v} = \mathbf{v}^{(i)}$  and  $\mathbf{z} = \mathbf{z}^{(i)}$  fixed as

$$\widehat{f}_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a}) \propto \exp\left([\widehat{\theta}^{(i,y)} + 2\mathbf{v}^{(i)\top} \widehat{\Phi}^{(v,y)} + 2\mathbf{a}^\top \widehat{\Phi}^{(a,y)}] \mathbf{y} + \mathbf{y}^\top \widehat{\Phi}^{(y,y)} \mathbf{y}\right). \quad (15)$$

Finally, we estimate  $\mu^{(i)}(\tilde{\mathbf{a}}^{(i)})$  as the mean under the above conditional distribution, given by

$$\widehat{\mu}^{(i)}(\tilde{\mathbf{a}}^{(i)}) \triangleq \mathbb{E}_{\widehat{f}_{\mathbf{y}|\mathbf{a}}^{(i)}}[\mathbf{y}|\mathbf{a} = \tilde{\mathbf{a}}^{(i)}], \quad (16)$$

which can be computed by standard algorithms for estimating marginals of graphical models, e.g., via the junction tree algorithm (Wainwright et al., 2008) or message-passing algorithms.<sup>6</sup>

<sup>5</sup>This follows from (Bubeck et al., 2015, Theorem. 3.7) by noting that  $\mathcal{L}(\underline{\Theta})$  is  $O(p)$  smooth function of  $\underline{\Theta}$ .

<sup>6</sup>In general, estimating the marginals exactly is computationally hard for undirected graphical models. While the junction tree algorithm works well for graphical models with small treewidth (Wainwright et al., 2008, Section. 2.5),

## 4 Main results

In this section, we analyze our estimates. First, we provide our guarantee on estimating the unit-level and the population-level parameters in Section. 4.1. Next, we provide our guarantee on estimating the causal estimand of interest in Section. 4.2. Before stating our main results, we define a standard notion of complexity of the set  $\Lambda_\theta$ , namely metric entropy (defined below) that our guarantees rely on.

**Definition 2** ( $\varepsilon$ -covering number and metric entropy). *Given a set  $\mathcal{V} \subset \mathbb{R}^p$  and a scalar  $\varepsilon > 0$ , we use  $\mathcal{C}(\mathcal{V}, \varepsilon)$  to denote the  $\varepsilon$ -covering number of  $\mathcal{V}$  with respect to  $\|\cdot\|_1$ , i.e.,  $\mathcal{C}(\mathcal{V}, \varepsilon)$  denotes the minimum cardinality over all possible subsets  $\mathcal{U} \subset \mathcal{V}$  that satisfy  $\mathcal{V} \subset \cup_{u \in \mathcal{U}} \mathcal{B}(u; \varepsilon)$ , where  $\mathcal{B}(u; \varepsilon) \triangleq \{v \in \mathbb{R}^p : \|u - v\|_1 \leq \varepsilon\}$ . We let  $\mathcal{M}_\theta(\varepsilon) \triangleq \log \mathcal{C}(\Lambda_\theta, \varepsilon)$  denote the metric entropy of  $\Lambda_\theta$ , and  $\mathcal{M}_{\theta, n}(\varepsilon) \triangleq n\mathcal{M}_\theta(n\varepsilon)$  denote a scaled version of it.*

Next, we state two settings with upper bounds on the metric entropy, and we use them as running examples to unpack our general results throughout this paper.

**Example 1** (Linear combination). *Consider a set  $\Lambda_\theta$  containing vectors with bounded entries that are also a linear combination of  $k$  known vectors in  $\mathbb{R}^p$  collected as  $\mathbf{B} \in \mathbb{R}^{p \times k}$ , i.e.,  $\Lambda_\theta = \{\mathbf{B}\mathbf{a} : \mathbf{a} \in \mathbb{R}^k, \|\mathbf{B}\mathbf{a}\|_\infty \leq \alpha\}$ . Then, [Dagan et al. \(2021, Lemma. 11\)](#) implies that  $\mathcal{M}_\theta(\eta) = O(k \log(1 + \frac{\alpha}{\eta}))$ . Further,  $\mathcal{M}_{\theta, n}(\eta) = O(\frac{\alpha k}{\eta})$ .*

**Example 2** (Sparse linear combination). *Consider a set  $\Lambda_\theta$  containing vectors with bounded entries that are also a  $s$ -sparse linear combination of  $k$  known vectors in  $\mathbb{R}^p$  collected as  $\mathbf{B} \in \mathbb{R}^{p \times k}$ , i.e.,  $\Lambda_\theta = \{\mathbf{B}\mathbf{a} : \mathbf{a} \in \mathbb{R}^k, \|a\|_0 \leq s, \|\mathbf{B}\mathbf{a}\|_\infty \leq \alpha\}$ . Then [Dagan et al. \(2021, Corollary. 4\)](#) implies that  $\mathcal{M}_\theta(\eta) = O(s \log k \log(1 + \frac{\alpha}{\eta}))$ . Further,  $\mathcal{M}_{\theta, n}(\eta) = O(\frac{\alpha s \log k}{\eta})$ .*

### 4.1 Guarantee on quality of parameter estimate

Our non-asymptotic guarantees use an assumption of a lower bound on the smallest eigenvalue of a suitable set of autocorrelation matrices.

**Assumption 2.** *For any  $\mathbf{z} \in \mathcal{Z}^{p_z}$  and  $t \in [p]$ , let  $\lambda_{\min}(\mathbf{z}, t)$  denote the smallest eigenvalue of the matrix  $\mathbb{E}_{\mathbf{x}|\mathbf{z}}[\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top | \mathbf{z} = \mathbf{z}]$  where  $\tilde{\mathbf{x}} \triangleq (\mathbf{x}_t, 2\mathbf{x}_{-t}\mathbf{x}_t, \mathbf{x}_t^2 - x_{\max}^2/3) \in \mathbb{R}^{p+1}$ . We assume  $\lambda_{\min} \triangleq \min_{\mathbf{z} \in \mathcal{Z}^{p_z}, t \in [p]} \lambda_{\min}(\mathbf{z}, t)$  is strictly positive.*

We note that all eigenvalues of any autocorrelation matrix are non-negative implying  $\lambda_{\min}(\mathbf{z}, t) \geq 0$  for all  $\mathbf{z} \in \mathcal{Z}^{p_z}, t \in [p]$ . Assumption. 2 requires  $\lambda_{\min}(\mathbf{z}, t) > 0$  for all  $\mathbf{z} \in \mathcal{Z}^{p_z}, t \in [p]$  and serves as a sufficient condition to rule out certain singular distributions ([Shah et al., 2021b](#), Section. 5).<sup>7</sup> In Appendix. B.2, we show that  $\lambda_{\min} = \Omega(e^{-c\beta})$  when  $\Theta_{tt}^* = 0$  for all  $t \in [p]$  as in Ising model where  $x_t^2 = 1$  for all  $t \in [p]$ .

We are now ready to state our main result that characterizes a high probability bound on the estimation error for the estimate  $\hat{\Theta}$  computed via (14). To simplify the presentation, we use  $c$  and  $c'$  to denote universal constants or constants that depend on the parameters  $\alpha, x_{\max}$ , and  $\lambda_{\min}$  and can take a different value in each appearance.

e.g., for trees or chains as in hidden Markov models or state-space models, message-passing algorithms are the default choice for computing approximate marginals for complex graphs, especially with cycles. However, message-passing algorithms may induce additional approximations, which we do not discuss here.

<sup>7</sup>Essentially, we use this assumption to lower bound the variance of a non-constant random variable (Appendix. B.1).

**Theorem 1** (Guarantee on quality of parameter estimate). *Suppose Assumptions. 1 and 2 hold. Fix an  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , and define*

$$R(\varepsilon, \delta) \triangleq \max\{ce^{c'\beta}\sqrt{\log(\log p/\delta) + \mathcal{M}_\theta(ce^{-c'\beta})}, \varepsilon\gamma\} \text{ with } \gamma \triangleq \max_{\theta, \bar{\theta} \in \Lambda_\theta} \frac{\|\theta - \bar{\theta}\|_1}{\|\theta - \bar{\theta}\|_2} \quad (17)$$

and

$$\widetilde{\mathcal{M}}_{\theta, n}(\varepsilon, \delta) \triangleq \mathcal{M}_{\theta, n}\left(\frac{\varepsilon^2}{p}\right) + p\mathcal{M}_\theta(R^2(\varepsilon, \delta)). \quad (18)$$

Then, with probability at least  $1 - \delta$ , the estimates  $\widehat{\Theta}, \widehat{\theta}^{(1)}, \dots, \widehat{\theta}^{(n)}$  defined in (14) satisfy

$$\|\widehat{\Theta} - \Theta^*\|_{2, \infty} \leq \varepsilon \quad \text{when } n \geq \frac{ce^{c'\beta}p^2\left(p \log \frac{p}{\delta\varepsilon^2} + \mathcal{M}_{\theta, n}(\varepsilon^2)\right)}{\varepsilon^4} \quad (19)$$

and

$$\max_{i \in [n]} \|\widehat{\theta}^{(i)} - \theta^{*(i)}\|_2 \leq R\left(\varepsilon, \frac{\delta}{n}\right) \quad \text{when } n \geq \frac{ce^{c'\beta}p^4\left(p \log \frac{np^2}{\delta\varepsilon^2} + \widetilde{\mathcal{M}}_{\theta, n}\left(\varepsilon, \frac{\delta}{n}\right)\right)}{\varepsilon^4}. \quad (20)$$

We split the proof into two parts: First, we establish the bound (19) in Appendix. B, which we then use to establish the bound (20) in Appendix. C.

Our guarantee in (19) provides a non-asymptotic error bound of order  $\frac{p^2(p \log p + \mathcal{M}_{\theta, n}(n^{-1/2}))}{n^{1/4}}$  (where we treat  $\beta$  as a constant) for estimating  $\Theta^*$  although the  $n$  samples have different unit-level parameters  $\{\theta^{*(i)}\}_{i=1}^n$ . On the other hand, after squaring both sides and dividing by  $p$ , the guarantee (20) for the unit-level parameters can be simplified as follows:<sup>8</sup> whenever  $n \geq c'\varepsilon^{-4}p^4\left(p \log \frac{p^2}{\delta\varepsilon^2} + \mathcal{M}_{\theta, n}(\varepsilon^2/p) + p\mathcal{M}_\theta(c)\right)$ , we have

$$\text{MSE}(\widehat{\theta}^{(i)}, \theta^{*(i)}) \leq \max\left\{\varepsilon^2, \frac{\mathcal{M}_\theta(c) + \log(\log \frac{p}{\delta})}{p}\right\}, \quad (21)$$

where we use  $\gamma \leq \sqrt{p}$  in (17) and treat  $\beta$  as a constant. For large  $n$  so that  $\varepsilon$  is small, this error scales linearly with the metric entropy  $\mathcal{M}_\theta$ —the error becomes worse as the unit-level parameter set  $\Lambda_\theta$  becomes more complex.

The next corollary (stated without proof) provides a formal version of the population-level guarantee in (19) and the unit-level guarantee in (21) for the two examples discussed earlier. We treat  $\beta$  as a constant and note that the dependence is exponential as in Theorem. 1.

**Corollary 1** (Consequences for examples). *Suppose Assumptions. 1 and 2 hold. Then, for any fixed  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , the following results hold with probability at least  $1 - \delta$ .*

(a) Linear combination: *If  $\Lambda_\theta$  is as in Example. 1, then for all  $i \in [n]$ ,*

$$\begin{aligned} \|\widehat{\Theta} - \Theta^*\|_{2, \infty} \leq \varepsilon & \quad \text{for } n \geq \frac{cp^2\left(p \log \frac{p}{\delta\varepsilon^2} + \frac{k}{\varepsilon^2}\right)}{\varepsilon^4} \\ \text{MSE}(\widehat{\theta}^{(i)}, \theta^{*(i)}) \leq \max\left\{\varepsilon^2, \frac{c\left(k + \log(\log \frac{p}{\delta})\right)}{p}\right\} & \quad \text{for } n \geq \frac{cp^5\left(\log \frac{p^2}{\delta\varepsilon^2} + k + \frac{k}{\varepsilon^2}\right)}{\varepsilon^4}. \end{aligned}$$

<sup>8</sup>We replace  $\delta/n$  in (20) by  $\delta$  as we do not require a union bound over  $i \in [n]$  for unit-wise guarantees.

(b) Sparse linear combination: If  $\Lambda_\theta$  is as in Example. 2, then for all  $i \in [n]$ ,

$$\begin{aligned} \|\widehat{\Theta} - \Theta^*\|_{2,\infty} \leq \varepsilon & \quad \text{for } n \geq \frac{cp^2(p \log \frac{p}{\delta \varepsilon^2} + \frac{s \log k}{\varepsilon^2})}{\varepsilon^4} \\ \text{MSE}(\widehat{\theta}^{(i)}, \theta^{*(i)}) \leq \max \left\{ \varepsilon^2, \frac{c(s \log k + \log(\log \frac{p}{\delta}))}{p} \right\} & \quad \text{for } n \geq \frac{cp^5(\log \frac{p^2}{\delta \varepsilon^2} + s \log k + \frac{s \log k}{\varepsilon^2})}{\varepsilon^4}. \end{aligned}$$

Corollary. 1 states that, as long as  $n$  is polynomially large in  $p$ , our strategy learns the unit-level parameters (on average in terms of mean square error across coordinates) for each user if  $p$  is large compared to either the number of vectors  $k$  (Example. 1) or the sparsity parameter  $s$  (Example. 2).

**Sharpness of guarantees and generalization of prior results** The exponential dependence on  $\beta$  in Theorem. 1 is unavoidable given the lower bounds for learning exponential families even with i.i.d. samples (Santhanam and Wainwright, 2012). Regarding the dependence on error tolerance  $\varepsilon$ , prior works with suitable analogs of our loss function provide two different error scaling: (i)  $1/\varepsilon^4$  in Vuffray et al. (2022); Shah et al. (2021a,b) and (ii)  $1/\varepsilon^2$  in Vuffray et al. (2016) and Shah et al. (2023). The works in category (ii) use techniques from Negahban et al. (2012), and it remains an interesting future direction to see whether similar ideas could be used to sharpen the error scaling of  $1/\varepsilon^4$  to the parametric rate of  $1/\varepsilon^2$  in Theorem. 1. We note that improving the dependence on  $\varepsilon$  in (19) improves the dependence on  $\varepsilon$  as well as  $p$  in (20). In the special case of equal unit-level parameters ( $\theta^{*(1)} = \dots = \theta^{*(n)}$ ), the analysis in Appendix. B to establish the bound (19) can be modified to recover (up to constants) prior guarantee (Shah et al., 2021a, Lemma. 9.1) on learning exponential family from  $n$  i.i.d. samples. Further, the guarantee (20) recovers the prior guarantee (Kandiros et al., 2021, Theorem. 6) as a special case where the authors consider learning an Ising model from one sample when the population-level parameter is known up to a scaling factor.

## 4.2 Guarantee on quality of outcome estimate

Our non-asymptotic guarantee on outcome estimate assumes that the following matrices are suitably stable under small perturbation in the parameters: (i) the covariance matrix of  $\mathbf{y}$  conditioned on  $\mathbf{a}$ ,  $\mathbf{z}$ , and  $\mathbf{v}$  and (ii) the cross-covariance matrix of  $\mathbf{y}$  and  $y_t \mathbf{y}$  conditioned on  $\mathbf{a}$ ,  $\mathbf{z}$ , and  $\mathbf{v}$  for all  $t \in [p_y]$ .

**Assumption 3.** For any set  $\mathbb{B}$  containing  $\theta, \Theta$ , there exists a constant  $C(\mathbb{B})$  such that

$$\sup_{\theta, \Theta \in \mathbb{B}} \max \left\{ \|\text{Cov}_{\theta, \Theta}(\mathbf{y}, \mathbf{y} | \mathbf{a}, \mathbf{z}, \mathbf{v})\|_{\text{op}}, \max_{t \in [p_y]} \|\text{Cov}_{\theta, \Theta}(\mathbf{y}, y_t \mathbf{y} | \mathbf{a}, \mathbf{z}, \mathbf{v})\|_{\text{op}} \right\} \leq C(\mathbb{B}), \quad (22)$$

almost surely. The expectation in (22) is with respect to the distribution of  $\mathbf{y}$  conditioned on  $\mathbf{a} = \mathbf{a}$ ,  $\mathbf{z} = \mathbf{z}$ , and  $\mathbf{v} = \mathbf{v}$  which is fully parameterized by  $\theta$  and  $\Theta$ , and can be obtained from (6) after replacing  $\theta(\mathbf{z})$  by  $\theta$ .

In Appendix. D.1, we show that  $C(\mathbb{B})$  is a constant for a class of distributions. We note that this assumption is common in the literature on learning Gaussian graphical models to rule out singular distributions (Won and Kim, 2006; Zhou et al., 2011; Ma and Michailidis, 2016).

We are now ready to state our guarantee for the estimate  $\widehat{\mu}^{(i)}(\widetilde{\mathbf{a}}^{(i)})$  (see (16)) of the expected potential outcomes for any unit  $i \in [n]$  under an alternate intervention  $\widetilde{\mathbf{a}}^{(i)} \in \mathcal{A}^{p_a}$ . We assume  $p_v = p_a = p_y$  for brevity. See the proof in Appendix. D where we also state a more general result.

**Theorem 2** (Guarantee on quality of outcome estimate). *Suppose Assumptions. 1 to 3 hold. Then for any fixed  $\varepsilon > 0$  and  $\delta \in (0, 1)$ , the estimates  $\{\hat{\mu}^{(i)}(\tilde{\mathbf{a}}^{(i)})\}_{i=1}^n$  defined in (16) for any  $\{\tilde{\mathbf{a}}^{(i)} \in \mathcal{A}^{p_a}\}_{i=1}^n$  satisfy*

$$\max_{i \in [n]} \frac{\|\mu^{(i)}(\tilde{\mathbf{a}}^{(i)}) - \hat{\mu}^{(i)}(\tilde{\mathbf{a}}^{(i)})\|_2}{C(\mathbb{B}_i)} \leq R\left(\varepsilon, \frac{\delta}{n}\right) + p\varepsilon \quad \text{for } n \geq \frac{ce^{c'\beta} p^4 (p \log \frac{np^2}{\delta\varepsilon^2} + \widetilde{\mathcal{M}}_{\theta,n}(\varepsilon, \frac{\delta}{n}))}{\varepsilon^4}, \quad (23)$$

with probability at least  $1 - \delta$ , where  $R(\varepsilon, \delta)$  was defined in (17),  $\widetilde{\mathcal{M}}_{\theta,n}(\varepsilon, \delta)$  was defined in (18),  $C(\mathbb{B})$  was defined in (22), and

$$\mathbb{B}_i \triangleq \left\{ \theta \in \Lambda_\theta : \|\theta - \theta^{*(i)}\|_2 \leq R\left(\varepsilon, \frac{\delta}{n}\right) \right\} \times \left\{ \Theta \in \Lambda_\Theta : \max_{t \in [p]} \|\Theta_t - \Theta_t^*\|_2 \leq \varepsilon \right\}.$$

Repeating the algebra as in (21) and treating  $C(\mathbb{B}_i)$  as a constant, the bound (23) yields the following simplified bound for the MSE of our mean outcome estimate  $\mu^{(i)}(\tilde{\mathbf{a}}^{(i)})$  for unit  $i \in [n]$  under treatment  $\tilde{\mathbf{a}}^{(i)} \in \mathcal{A}^{p_a}$ : whenever  $n \geq c'\varepsilon^{-4} p^4 (p \log \frac{p^2}{\delta\varepsilon^2} + \mathcal{M}_{\theta,n}(\varepsilon^2/p) + p\mathcal{M}_\theta(c))$ , we have

$$\text{MSE}(\mu^{(i)}(\tilde{\mathbf{a}}^{(i)}), \hat{\mu}^{(i)}(\tilde{\mathbf{a}}^{(i)})) \leq \varepsilon^2 + \frac{\mathcal{M}_\theta(c) + \log(\log \frac{p}{\delta})}{p}.$$

This bound is of the same order as in (21) and can be formalized for the two examples (Examples. 1 and 2) by deriving a suitable analog of Corollary. 1. In a nutshell, in both settings, the unit-level expected potential outcomes can be estimated well when the total number of units  $n$  is large and the observations for each unit are high dimensional compared to the number of vectors  $k$  in Example. 1 or the sparsity parameter  $s$  in Example. 2. We omit a formal statement for brevity.

Finally, we also note that as in Theorem. 1, the exponential dependence on  $\beta$  is expected to be unavoidable due to the principle of conjugate duality (Wainwright et al., 2008), i.e., the existence of a unique mapping from the parameters to the means and vice versa for the exponential family. Moreover, as in the discussion after Corollary. 1, the sharpness of the rate of  $1/\varepsilon^4$  is left for future work. Improving the dependency on  $\varepsilon$  in (23) would also improve the dependency on  $p$ .

## 5 Possible extensions

We now discuss how to extend our theoretical results with various relaxations of the exponential family modeling.

### 5.1 Modeling only the conditional distribution as exponential family

Our framework and analysis can be extended to the setting where, instead of the joint distribution  $f_{\mathbf{w}}$  of  $\mathbf{w} = (\mathbf{z}, \mathbf{v}, \mathbf{a}, \mathbf{y})$ , we model only the conditional distribution  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}$  of  $\mathbf{y}$  conditioned on  $\mathbf{a}$ ,  $\mathbf{z}$ , and  $\mathbf{v}$  as an exponential family. Note that when the joint distribution  $f_{\mathbf{w}}$  is an exponential family, the conditional distribution  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}$  is also an exponential family, however a vice versa implication does not hold so that the setting considered here is a strict generalization of our previous setting. In fact, the conditional distribution  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}$  being an exponential family puts no restrictions on the marginal distribution  $f_{\mathbf{z},\mathbf{v},\mathbf{a}}$  of the unobserved covariates, the observed covariates, and the interventions as is the case with non-linear panel data models (Section. 2).

To estimate the expected potential outcomes  $\mu^{(i)}(\tilde{\mathbf{a}}^{(i)})$  in (8) for any given unit  $i$  and any alternate intervention  $\tilde{\mathbf{a}}^{(i)}$ , it suffices to estimate the conditional distribution of  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}(\cdot|\mathbf{v}=\mathbf{a},\mathbf{v}=\mathbf{v}^{(i)}\mathbf{z}=\mathbf{z}^{(i)})$   $\mathbf{y}$  for unit  $i$  as a function of the intervention  $\mathbf{a}$  (as in (15)). This task is equivalent to estimating  $\gamma^{(i)}$  in (5) under the exponential family models in (2) or (3).

In Section. 3.2, under the exponential family in (2), we argued (for analytical convenience) that learning  $\gamma^{(i)}$  is subsumed in learning the parameters corresponding to the conditional distribution  $f_{\mathbf{x}|\mathbf{z}}$  of  $\mathbf{x} = (\mathbf{v}, \mathbf{a}, \mathbf{y})$  conditioned on  $\mathbf{z}$  (which also belongs to an exponential family with linear and quadratic interactions) as in (6). Then, we set the goal of estimating the parameters in (7) and designed a loss function to do so. The loss function depended on the conditional distribution  $f_{x_t|\mathbf{x}_{-t},\mathbf{z}}$  (11) of the random variable  $x_t$  conditioned on  $\mathbf{x}_{-t} = \mathbf{x}_{-t}$  and  $\mathbf{z} = \mathbf{z}$  for every  $t \in [p]$ .

Under the exponential family in (3), we focus on directly learning the components of (7) relevant to learning  $\gamma^{(i)}$ , i.e.,

$$\theta_t^*(\mathbf{z}^{(i)}) = \phi^{*(y)} + 2\Phi^{*(z,y)\top} \mathbf{z}^{(i)} \in \mathbb{R}^{p_y \times 1}, \quad \text{for all } t \in \{p_v + p_a + 1, \dots, p_v + p_a + p_y\} \quad (24)$$

$$\Theta_t^* = (\Phi^{*(v,y)}, \Phi^{*(a,y)}, \Phi^{*(y,y)}) \in \mathbb{R}^{p \times 1} \quad \text{for all } t \in \{p_v + p_a + 1, \dots, p_v + p_a + p_y\}. \quad (25)$$

We note that the conditional distribution  $f_{y_t|\mathbf{y}_{-t},\mathbf{v},\mathbf{a},\mathbf{z}}$  of the random variable  $y_t$  conditioned on  $\mathbf{y}_{-t} = \mathbf{y}_{-t}$ ,  $\mathbf{v} = \mathbf{v}$ ,  $\mathbf{a} = \mathbf{a}$ , and  $\mathbf{z} = \mathbf{z}$  for every  $t \in [p_y]$  is consistent with the conditional distribution  $f_{x_{t'}|\mathbf{x}_{-t'},\mathbf{z}}$  in (11) for every  $t' \in \{p_v + p_a + 1, \dots, p_v + p_a + p_y\}$ . As a result, we can adapt the loss function in (13) to learn the parameters in (24) and (25) by summing over  $t \in \{p_v + p_a + 1, \dots, p_v + p_a + p_y\}$  instead of  $t \in [p]$ . Consequently, the guarantees in Section. 4 continue to hold with  $p$  replaced by  $p_y$ .

## 5.2 Higher order terms in the conditional exponential family

In Section. 5.1, we described how our framework and results apply when only the conditional distribution  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}$  is modeled as the exponential family distribution in (3) where the term inside the exponent is linear in  $(\mathbf{z}, \mathbf{v}, \mathbf{a})$  and quadratic in  $\mathbf{y}$ . We now describe how our framework and results are applicable when the conditional distribution  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}$  is modeled as the following exponential family distribution

$$f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}(\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}) \propto \exp(q_\Phi(\mathbf{v}, \mathbf{a}, \mathbf{y})) \exp(2\mathbf{z}^\top \Phi^{(z,y)} \mathbf{y}), \quad (26)$$

where  $q_\Phi(\mathbf{v}, \mathbf{a}, \mathbf{y})$  is some bounded degree polynomial in  $(\mathbf{v}, \mathbf{a}, \mathbf{y})$  parameterized by  $\Phi$ , i.e., the term inside the exponent is linear in  $\mathbf{z}$  and arbitrary bounded degree polynomial in  $(\mathbf{v}, \mathbf{a}, \mathbf{y})$ . We note that every term in  $q_\Phi(\mathbf{v}, \mathbf{a}, \mathbf{y})$  needs to depend on  $\mathbf{y}$  for it to contribute to  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}$  in (26). For convenience, hereon, we ignore any dependence on  $\mathbf{v}$ , and abuse notation to let  $q_\Phi(\mathbf{a}, \mathbf{y}) = q_\Phi(\mathbf{v}, \mathbf{a}, \mathbf{y})$ . Then, in (3),  $q_\Phi(\mathbf{a}, \mathbf{y})$  was a polynomial of degree 2, i.e.,

$$q_\Phi(\mathbf{a}, \mathbf{y}) = q_\Phi^{(2)}(\mathbf{a}, \mathbf{y}) \triangleq \text{Sum}\left(\phi^{(y)} \odot \mathbf{y} + 2\Phi^{(a,y)} \odot (\mathbf{a} \otimes \mathbf{y}) + \Phi^{(y,y)} \odot (\mathbf{y} \otimes \mathbf{y})\right),$$

where  $\odot$  denotes the Hadamard product,  $\otimes$  denotes the Kronecker product,  $\Phi = (\phi^{(y)}, \Phi^{(a,y)}, \Phi^{(y,y)})$  with  $\Phi^{(y,y)}$  being symmetric, and  $\text{Sum}(s_1 + \dots + s_h) \in \mathbb{R}$  sums, over all  $i \in [h]$ , all the entries of  $s_i$  which could be a real number/vector/matrix/tensor. To explain how the loss function in (13) needs to be modified for general  $q_\Phi(\mathbf{a}, \mathbf{y})$ , we consider a polynomial of degree 3:

$$q_\Phi(\mathbf{a}, \mathbf{y}) = q_\Phi^{(2)}(\mathbf{a}, \mathbf{y}) + \text{Sum}\left(\sum_{(u_1, u_2) \in \{(a,a), (a,y), (y,y)\}} c_{u_1, u_2} \cdot \Phi^{(u_1, u_2, y)} \odot (\mathbf{u}_1 \otimes \mathbf{u}_2 \otimes \mathbf{y})\right),$$

where  $c_{a,a} = c_{a,y} = 3$ ,  $c_{y,y} = 1$  are constants chosen for consistency, and  $\Phi^{(u_1, u_2, y)} \in \mathbb{R}^{p_{u_1} \times p_{u_2} \times p_y}$  is symmetric with respect to indices that are repeated for every  $(u_1, u_2) \in \{(a, a), (a, y), (y, y)\}$ . We illustrate the two steps from Section. 3.3.1 below.

**Centering sufficient statistics of the conditional distribution of a variable** The conditional distribution  $f_{y_t|\mathbf{y}_{-t}, \mathbf{a}, \mathbf{z}}$  of the random variable  $y_t$  conditioned on  $\mathbf{y}_{-t} = \mathbf{y}_{-t}$ ,  $\mathbf{a} = \mathbf{a}$ , and  $\mathbf{z} = \mathbf{z}$  for every  $t \in [p_y]$  is given by

$$f_{y_t|\mathbf{y}_{-t}, \mathbf{a}, \mathbf{z}}(y_t|\mathbf{y}_{-t}, \mathbf{a}, \mathbf{z}) \propto \exp\left(\text{Sum}\left(\left[\phi_t(\mathbf{z}) + \sum_{u \in \{\mathbf{y}_{-t}, \mathbf{a}\}} 2\Phi^{(u, y_t)} \odot \mathbf{u} + \sum_{(u_1, u_2) \in \{(a, a), (a, y-t), (y-t, y-t)\}} c_{u_1, u_2} \Phi^{(u_1, u_2, y_t)} \odot (\mathbf{u}_1 \otimes \mathbf{u}_2)\right] y_t + \left[\Phi^{(y_t, y_t)} + \sum_{u \in \{\mathbf{y}_{-t}, \mathbf{a}\}} 3\Phi^{(u, y_t, y_t)} \odot \mathbf{u}\right] \left(y_t^2 - \frac{x_{\max}^2}{3}\right) + \Phi^{(y_t, y_t, y_t)} y_t^3\right)\right),$$

where  $\phi_t(\mathbf{z}) \triangleq \phi^{(y_t)} + 2\Phi^{(z, y_t)} \odot \mathbf{z}$ ,  $c_{y-t, y-t} = 3$ , and  $c_{a, y-t} = 6$ . Let  $\Phi_t$  denote the concatenation of all the remaining parameters. As in (12), the term  $x_{\max}^2/3$  inside the exponent is vacuous and centers the sufficient statistics  $y_t^2$ . The other sufficient statistics, i.e.,  $\mathbf{x}_t$  and  $\mathbf{x}_t^3$ , are naturally centered as their integrals with respect to the uniform distribution on  $\mathcal{X}$  are both zeros.

**Constructing the loss function** Now, it is easy to see that the corresponding loss  $\mathcal{L}$  is given by

$$\mathcal{L} = \frac{1}{n} \sum_{t \in [p_y]} \sum_{i \in [n]} \exp\left(-\text{Sum}\left(\left[\phi_t^{(i)} + \sum_{u \in \{\mathbf{y}_{-t}, \mathbf{a}\}} 2\Phi^{(u, y_t)} \odot \mathbf{u}^{(i)} + \sum_{(u_1, u_2) \in \{(a, a), (a, y-t), (y-t, y-t)\}} c_{u_1, u_2} \Phi^{(u_1, u_2, y_t)} \odot (\mathbf{u}_1^{(i)} \otimes \mathbf{u}_2^{(i)})\right] y_t^{(i)} + \left[\Phi^{(y_t, y_t)} + \sum_{u \in \{\mathbf{y}_{-t}, \mathbf{a}\}} 3\Phi^{(u, y_t, y_t)} \odot \mathbf{u}^{(i)}\right] \left([y_t^{(i)}]^2 - \frac{x_{\max}^2}{3}\right) + \Phi^{(y_t, y_t, y_t)} [y_t^{(i)}]^3\right)\right),$$

and minimizing this convex loss results in the estimates of  $\{\phi_t^{(i)}\}_{i \in [n]}$  and  $\{\Phi_t\}_{t \in [p_y]}$ . Consequently, the guarantees in Section. 4 continue to hold with  $p$  replaced by  $p_y$  as long as Assumptions. 1 to 3 are appropriately generalized.

**Tilting the base distribution** We note that the exponential family in (3) can be rewritten as

$$f_{\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}}(\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}) \propto \exp(2\mathbf{z}^\top \Phi^{(z, y)} \mathbf{y}) \exp(2\mathbf{v}^\top \Phi^{(v, y)} \mathbf{y}) \exp(2\mathbf{a}^\top \Phi^{(a, y)} \mathbf{y}) \exp(\phi^{(y)}^\top \mathbf{y} + \mathbf{y}^\top \Phi^{(y, y)} \mathbf{y}),$$

where  $\exp(\phi^{(y)}^\top \mathbf{y} + \mathbf{y}^\top \Phi^{(y, y)} \mathbf{y})$  stands for a base distribution on  $\mathbf{y}$  which is exponentially tilted by  $\mathbf{z}$ ,  $\mathbf{v}$ , and  $\mathbf{a}$ , i.e., by  $\exp(2\mathbf{z}^\top \Phi^{(z, y)} \mathbf{y})$ ,  $\exp(2\mathbf{v}^\top \Phi^{(v, y)} \mathbf{y})$ , and  $\exp(2\mathbf{a}^\top \Phi^{(a, y)} \mathbf{y})$ , respectively. Then, generalizing the exponential family in (3) to the one in (26) is equivalent to saying that our approach and results continue to apply when (a) the base distribution on  $\mathbf{y}$  is an exponential family distribution where the term inside the exponent is arbitrary bounded degree polynomial (instead of quadratic) and (b) the exponent of the exponential tilting of this base distribution by  $(\mathbf{v}, \mathbf{a})$  is arbitrary bounded degree polynomial (instead of linear).

### 5.3 Discrete and mixed variables

In Section. 3.2, we described how our framework and results are applicable when the support of  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  are bounded continuous sets, i.e.,  $\mathcal{V} = \mathcal{A} = \mathcal{Y} = [-x_{\max}, x_{\max}]$ . In Section. 5.1, we showed that

it suffices to only model the conditional distribution  $f_{\mathbf{y}|\mathbf{a},\mathbf{z},\mathbf{v}}$  as an exponential family distribution implying that we do not need any restrictions on the support of  $\mathbf{v}$  and  $\mathbf{a}$ . Now, we describe how to adapt our loss function when  $\mathbf{y} = (y_1, \dots, y_{p_y}) \in \mathcal{Y}_1 \times \dots \times \mathcal{Y}_{p_y}$  where  $\mathcal{Y}_t$  is either a discrete compact set or a continuous compact set for  $t \in [p_y]$ .

We note that the conditional distribution  $f_{y_t|\mathbf{y}_{-t},\mathbf{v},\mathbf{a},\mathbf{z}}$  of the random variable  $y_t$  conditioned on  $\mathbf{y}_{-t} = \mathbf{y}_{-t}$ ,  $\mathbf{v} = \mathbf{v}$ ,  $\mathbf{a} = \mathbf{a}$ , and  $\mathbf{z} = \mathbf{z}$  for every  $t \in [p_y]$  is still consistent with the conditional distribution  $f_{x_{t'}|\mathbf{x}_{-t'},\mathbf{z}}$  in (11) for every  $t' \in \{p_v + p_a + 1, \dots, p_v + p_a + p_y\}$ . However, the constants used to center the sufficient statistics in (12) may change. More precisely, for any  $t \in [p]$ , the sufficient statistics  $x_t$  and  $x_t^2$  are centered by subtracting  $\mathbb{E}_{\mathcal{U}_t}[x_t]$  and  $\mathbb{E}_{\mathcal{U}_t}[x_t^2]$ , respectively where  $\mathcal{U}_t$  denotes the uniform distribution supported over  $\mathcal{Y}_t$ . Consequently, the loss function in (13) as well as Assumption. 2 can be adapted, and the guarantees in Section. 4 continue to hold.

## 6 Application: Imputing missing covariates

Consider a setting with no systematically unobserved covariates  $\mathbf{z}$ ; instead, elements of  $(\mathbf{v}, \mathbf{a}, \mathbf{y})$  are missing or have measurement error for some fraction of the units. Our goal is to impute these missing values or denoise the measurement error in the observed values.

**Problem setup** For the ease of exposition, we assume the observed covariates  $\mathbf{v}$  can have measurement error<sup>9</sup> but the interventions and the outcomes do not have any measurement error. More concretely, for every unit  $i \in [n]$ , along with the interventions  $\mathbf{a}^{(i)}$  and the outcomes  $\mathbf{y}^{(i)}$ , we observe  $\bar{\mathbf{v}}^{(i)} = \mathbf{v}^{(i)} + \Delta\mathbf{v}^{(i)}$  instead of true covariates  $\mathbf{v}^{(i)}$  where  $\Delta\mathbf{v}^{(i)}$  denotes (unobserved) bounded measurement error. We assume that a certain number of units (known to us) have no measurement error: say,  $\Delta\mathbf{v}^{(i)} = 0$  for all  $i \in \{n/2 + 1, \dots, n\}$ .

**Questions of interest** Besides counterfactual estimates, our goal is to estimate  $\Delta\mathbf{v}^{(i)}$  for units with measurement error.

### 6.1 A theoretical guarantee

Our methodology can be applied to estimate these measurement errors when the joint distribution of the true covariates  $\mathbf{v} \in \mathcal{X}^{p_v}$ , the interventions  $\mathbf{a} \in \mathcal{X}^{p_a}$ , and the observed outcomes  $\mathbf{y} \in \mathcal{X}^{p_y}$  can be modeled as an exponential family, parameterized by a vector  $\phi \in \mathbb{R}^p$  and a symmetric matrix  $\Phi \in \mathbb{R}^{p \times p}$  where  $p \triangleq p_v + p_a + p_y$ , i.e., with  $\mathbf{w} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y})$

$$f_{\mathbf{w}}(\mathbf{w}; \phi, \Phi) \propto \exp\left(\phi^\top \mathbf{w} + \mathbf{w}^\top \Phi \mathbf{w}\right), \quad \text{where } \mathbf{w} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y}), \quad (27)$$

and  $\mathbf{v} \triangleq (v_1, \dots, v_{p_v})$ ,  $\mathbf{a} \triangleq (a_1, \dots, a_{p_a})$ , and  $\mathbf{y} \triangleq (y_1, \dots, y_{p_y})$  denote realizations of  $\mathbf{v}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$ , respectively. To estimate the counterfactual distribution, we decompose  $\mathbf{v}$  into  $\bar{\mathbf{v}}$  and  $\Delta\mathbf{v}$ , and obtain the distribution of the observed quantities  $\mathbf{x} \triangleq (\bar{\mathbf{v}}, \mathbf{a}, \mathbf{y})$  conditioned on  $\Delta\mathbf{v} = \Delta\mathbf{v}$  as follows (see Appendix. E for details)

$$f_{\mathbf{x}|\Delta\mathbf{v}}(\mathbf{x}|\Delta\mathbf{v}; \theta(\Delta\mathbf{v}), \Theta) \propto \exp\left([\theta(\Delta\mathbf{v})]^\top \mathbf{x} + \mathbf{x}^\top \Theta \mathbf{x}\right) \quad \text{where } \theta(\Delta\mathbf{v}) \triangleq \begin{bmatrix} \phi^{(v)} - 2\Phi^{(v,v)\top} \Delta\mathbf{v} \\ \phi^{(a)} - 2\Phi^{(v,a)\top} \Delta\mathbf{v} \\ \phi^{(y)} - 2\Phi^{(v,y)\top} \Delta\mathbf{v} \end{bmatrix}, \quad (28)$$

<sup>9</sup>Our analysis remains the same when observed covariates  $\mathbf{v}$  are missing instead of having measurement error.

$\mathbf{x} \triangleq (\bar{\mathbf{v}}, \mathbf{a}, \mathbf{y})$ ,  $\Theta \triangleq \Phi$ , and  $\bar{\mathbf{v}}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  denote realizations of  $\bar{\mathbf{v}}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$ , respectively. As in Section. 3.2, to estimate the counterfactual distribution, it suffices to learn  $\theta(\Delta\mathbf{v}) \in \mathbb{R}^{p \times 1}$  and  $\Theta \in \mathbb{R}^{p \times p}$ .

Let  $f_{\mathbf{w}}(\cdot; \phi^*, \Phi^*)$  denote the true data generating distribution of  $\mathbf{w}$  in (27) and let  $f_{\mathbf{x}|\Delta\mathbf{v}}(\cdot | \Delta\mathbf{v}; \theta^*(\Delta\mathbf{v}), \Theta^*)$  denote the true distribution of  $\mathbf{x}$  conditioned on  $\Delta\mathbf{v} = \Delta\mathbf{v}$ . We assume (a)  $\max\{\|\Delta\mathbf{v}\|_\infty, \|\phi^*\|_\infty, \|\Phi^*\|_{\max}\} \leq \alpha$  and (b)  $\|\Phi^*\|_\infty \leq \beta$  analogous to Assumption. 1 where the row-wise  $\ell_1$  sparsity in (b) is assumed to be induced by row-wise  $\ell_0$  sparsity, i.e.,  $\|\Phi_t^*\|_0 \leq \beta/\alpha$  for all  $t \in [p]$ . Then, given realizations  $\{\mathbf{x}^{(i)}\}_{i=1}^n$  consistent with  $f_{\mathbf{x}|\Delta\mathbf{v}}(\cdot | \Delta\mathbf{v}^{(i)}; \theta^*(\Delta\mathbf{v}^{(i)}), \Theta^*)$ , first, we estimate the parameters  $\phi^*$  and  $\Phi^* = \Theta^*$  using the realizations for units  $\{n/2 + 1, \dots, n\}$ . Next, we exploit the structure in the problem to show that  $\theta^{*(i)} \triangleq \theta^*(\Delta\mathbf{v}^{(i)})$  can be written as a linear combination of known vectors with some error, for every unit  $i \in \{1, \dots, n/2\}$ . Then, we use (14) to estimate  $\{\theta^{*(i)}\}_{i=1}^n$  and obtain estimates of  $\{\Delta\mathbf{v}^{(i)}\}_{i=1}^n$  as by-products. In particular, the estimate of the coefficients associated with the aforementioned linear combination for  $\theta^{*(i)}$  turn out to be our estimate of the measurement error  $\Delta\mathbf{v}^{(i)}$  for every  $i \in \{1, \dots, n/2\}$ . For  $i \in \{n/2 + 1, \dots, n\}$ , estimating  $\theta^{*(i)}$  and  $\Delta\mathbf{v}^{(i)}$  is straightforward since  $\theta^{*(i)} = \phi^*$  and  $\Delta\mathbf{v}^{(i)} = 0$ . We provide our guarantee on estimating  $\Theta^*$ ,  $\theta^{*(i)}$  for  $i \in [n]$ , and  $\Delta\mathbf{v}^{(i)}$  for  $i \in [n]$  below with a proof in Appendix. E.

**Proposition 2** (Impute missing covariates). *Suppose the eigenvalues of  $\mathbf{B}^\top \mathbf{B}$  are lower bounded by  $\kappa p$  for some  $\kappa > 0$  where  $\mathbf{B} \triangleq [\phi^*, -2\Phi_1^*, \dots, -2\Phi_{p_v}^*] \in \mathbb{R}^{p \times (p_v+1)}$ . Then, for any fixed  $\varepsilon_1 > 0$  and  $\delta \in (0, 1)$ , there exists estimates  $\widehat{\Theta}$  and  $\{\widehat{\theta}^{(i)}\}_{i=1}^n$  such that, with probability at least  $1 - \delta$ ,*

$$\|\widehat{\Theta} - \Theta^*\|_{2,\infty} \leq \varepsilon_1 \quad \text{when } n \geq \frac{ce^{c'\beta} \log \frac{p}{\sqrt{\delta}}}{\varepsilon_1^2},$$

and

$$\max_{i \in [n]} \text{MSE}(\widehat{\theta}^{(i)}, \theta^{*(i)}) \leq \max \left\{ \varepsilon_1^2, \frac{ce^{c'\beta} (p_v + \log(\log \frac{np}{\delta}))}{p} \right\} \quad \text{when } n \geq \frac{ce^{c'\beta} (\log \frac{\sqrt{np}}{\sqrt{\delta}} + p_v)}{\varepsilon_1^2}.$$

Further, for any fixed  $\varepsilon_2 > 0$ , if  $\varepsilon_2 \leq \frac{1}{8} \sqrt{\frac{p}{p_v+1}}$ , there exist estimates  $\{\widehat{\Delta\mathbf{v}}^{(i)}\}_{i=1}^n$  such that,

$$\max_{i \in [n]} \|\widehat{\Delta\mathbf{v}}^{(i)} - \Delta\mathbf{v}^{(i)}\|_2^2 \leq \max \left\{ \frac{c_1 \varepsilon_2^2 \kappa}{p_v + 1}, \frac{ce^{c'\beta} (p_v + \log(\log \frac{np}{\delta}))}{p\kappa} \right\} + \varepsilon_2^2 \kappa,$$

with probability at least  $1 - \delta$ , whenever  $n \geq ce^{c'\beta} \kappa^{-2} \varepsilon_2^{-2} (p_v + 1) (\log \frac{\sqrt{np}}{\sqrt{\delta}} + p_v)$ .

The above guarantees can be simplified as follows by treating  $\beta$  and  $\kappa$  as constants as well as ignoring the constants, and the logarithmic factors in  $n$  and  $\delta$  (denoted by  $\lesssim$  and  $\gtrsim$ ): for any  $\varepsilon_1 > 0$  and  $\frac{1}{8} \sqrt{\frac{p}{p_v+1}} \geq \varepsilon_2 > 0$

$$\|\widehat{\Theta} - \Theta^*\|_{2,\infty} \leq \varepsilon_1 \quad \text{when } n \gtrsim \frac{\log p}{\varepsilon_1^2}, \quad (29)$$

$$\max_{i \in [n]} \text{MSE}(\widehat{\theta}^{(i)}, \theta^{*(i)}) \lesssim \max \left\{ \varepsilon_1^2, \frac{p_v}{p} \right\} \quad \text{when } n \gtrsim \frac{\log p + p_v}{\varepsilon_1^2}, \quad (30)$$

and

$$\max_{i \in [n]} \|\widehat{\Delta \mathbf{v}}^{(i)} - \Delta \mathbf{v}^{(i)}\|_2^2 \lesssim \max \left\{ \frac{\varepsilon_2^2}{p_v}, \frac{p_v}{p} \right\} + \varepsilon_2^2 \quad \text{when } n \gtrsim \frac{p_v(\log p + p_v)}{\varepsilon_2^2}. \quad (31)$$

For large  $n$ , whenever,  $\max \left\{ \varepsilon_1^2, \frac{p_v}{p} \right\} = \frac{p_v}{p}$  and  $\max \left\{ \frac{\varepsilon_2^2}{p_v}, \frac{p_v}{p} \right\} = \frac{p_v}{p}$ , the guarantees in (30) and (31) can be written as

$$\max_{i \in [n]} \text{MSE}(\widehat{\theta}^{(i)}, \theta^{*(i)}) \lesssim \frac{p_v}{p} \quad \text{when } n \gtrsim \frac{p \log p}{p_v}, \quad (32)$$

and

$$\max_{i \in [n]} \|\widehat{\Delta \mathbf{v}}^{(i)} - \Delta \mathbf{v}^{(i)}\|_2^2 \lesssim \frac{p_v^2}{p} \quad \text{when } n \gtrsim \frac{p \log p}{p_v}. \quad (33)$$

**Remark** The measurement errors can be recovered well as long as enough units with no measurement error are observed (i.e.,  $n/2$  is large) and the observation per unit is high dimensional (i.e.,  $p$  is large compared to  $p_v^2$ ). We note that the quadratic dependence (on  $p_v$ ) in (33) arises because of the error in expressing  $\theta^{*(i)}$  as a linear combination of known vectors. In contrast, we get a linear dependence (on  $k$ ) in Corollary 1(a) where there is no error in expressing  $\theta^{*(i)}$  as a linear combination of known vectors (via Example 1).

## 6.2 Simulations

We now present some simulation results to empirically evaluate the error scaling of our parameter estimates with three key aspects of the application above: number of units  $n$ , dimension  $p$ , and dimension  $p_v$  of covariates with measurement error.

**Data generation** We choose  $\mathcal{X} = [-1, 1]$  and  $p_a = p_y = (p - p_v)/2$ . The true joint distribution (27) of  $\mathbf{w} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y})$  is set as a truncated Gaussian distribution with the parameters  $\phi^* = \mathbf{1} \in \mathbb{R}^p$  and a positive definite  $\Phi^* \in \mathbb{R}^{p \times p}$  generated using *sklearn* package (Pedregosa et al., 2011) such that  $\alpha = 6$ ,  $\beta = 4$ , and  $\kappa = 0.15$ . We draw  $n$  i.i.d. samples  $\{\mathbf{w}^{(i)}\}_{i=1}^n$  from this true distribution using *tmvtnorm* package (Wilhelm and Manjunath, 2010). Next, we generate  $\Delta \mathbf{v}^{(i)}$  uniformly from  $[0.9, 1]^{p_v}$  for units  $i \in \{1, \dots, n/2\}$  while setting  $\Delta \mathbf{v}^{(i)} = \mathbf{0}$  for other units. Combining  $\{\mathbf{w}^{(i)}\}_{i=1}^n$  and  $\{\Delta \mathbf{v}^{(i)}\}_{i=1}^n$  yields  $\{\mathbf{x}^{(i)}\}_{i=1}^n$  (see (28)).

**Plot details** In Figure 3, we plot the scaling of errors in our estimates for  $\Theta^*$  in the top row,  $\{\theta^{*(i)}\}_{i=1}^n$  in the middle row, and  $\{\Delta \mathbf{v}^{(i)}\}_{i=1}^n$  in the bottom row. In particular, we present how the error scales as the dimension  $n$  grows for various  $p$  and  $p_v$ . We plot the averaged error across 50 independent trials along with  $\pm 1$  standard error (the standard error is too small to be visible in our results).

To help see the error scaling, we provide the least squares fit on the log-log scale (log error vs log x-axis). We display the best linear fit and mention an empirical decay rate in the legend based on the slope of that fit, e.g., for a slope of  $-0.56$  for estimating  $\Theta^*$  when  $p = 16$  and  $p_v = 4$ , we report an empirical rate of  $n^{-0.56}$  for the averaged error. In the middle row and the bottom row of Figure 3, the rates vary from  $n^{0.00}$  to  $n^{-0.17}$ , and we omit these weak dependencies in the legend to reduce clutter.

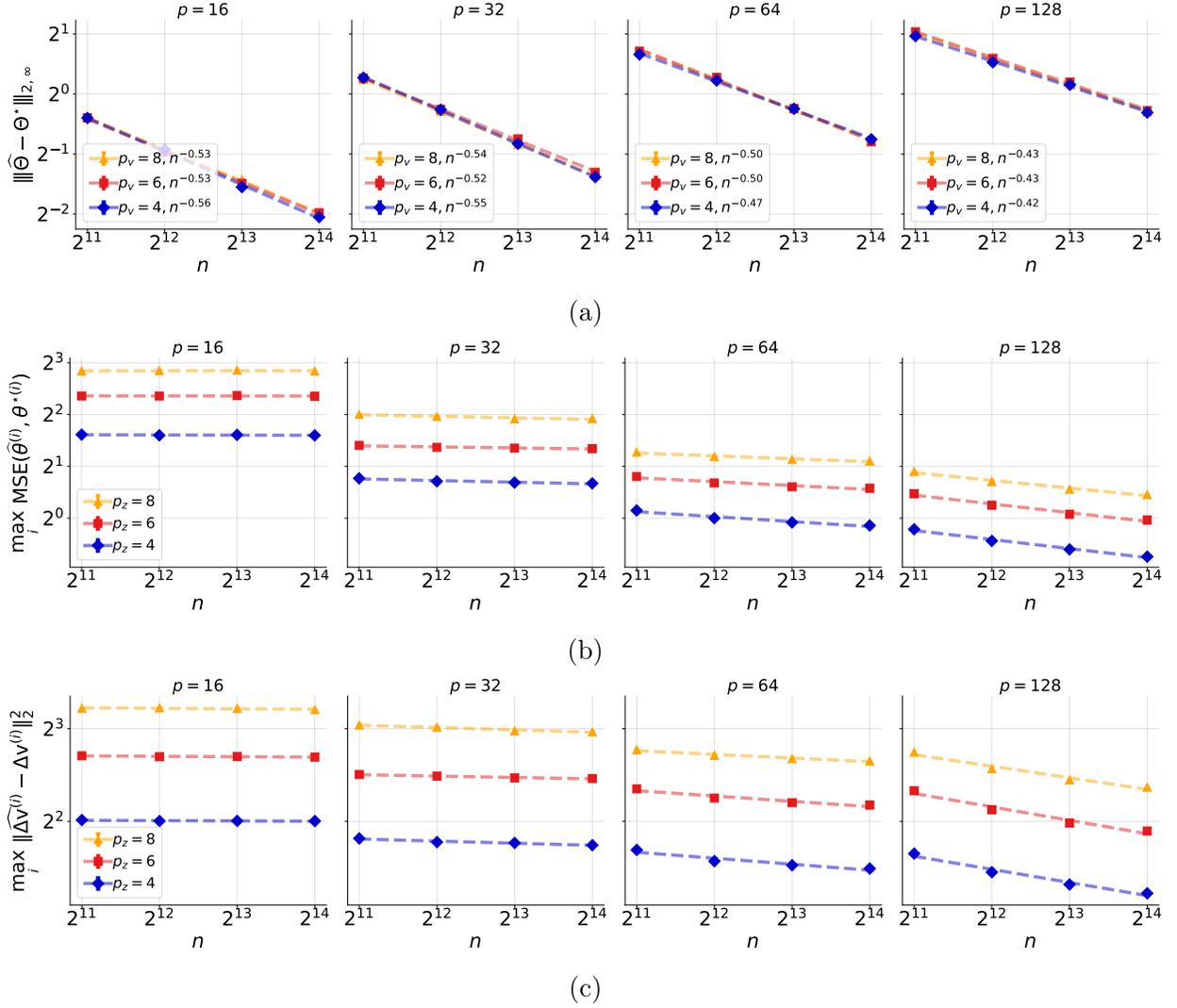


Figure 3: Error scaling with number of units  $n$ , for various  $p$  and  $p_v$ , for our estimates of  $\Theta^*$  (top row),  $\{\theta^{*(i)}\}_{i=1}^n$  (middle row), and  $\{\Delta v^{(i)}\}_{i=1}^n$  (bottom row).

**Error scaling for  $\hat{\Theta}$**  From the first row of Figure 3, we observe that the error  $\|\hat{\Theta} - \Theta^*\|_{2,\infty}$  admits a scaling of between  $n^{-0.56}$  and  $n^{-0.42}$  for various  $p$  and  $p_v$ . These empirical rates indicate a parametric error rate of  $n^{-0.5}$  for  $\|\hat{\Theta} - \Theta^*\|_{2,\infty}$ , consistent with the scaling of  $\varepsilon^{-2}$  in (29). Further, as expected, the error  $\|\hat{\Theta} - \Theta^*\|_{2,\infty}$  does not depend on  $p_v$  but increases with an increase in  $p$ .

**Error scaling for  $\hat{\theta}^{(i)}$**  In the middle row of Figure 3, we see the error  $\max_{i \in [n]} \text{MSE}(\hat{\theta}^{(i)}, \theta^{*(i)})$  has a weak dependence on  $n$  for a fixed  $p$  and  $p_v$ , decreases with an increase in  $p$  for any fixed  $n$  and  $p_v$ , and increases with an increase in  $p_v$  for any fixed  $n$  and  $p$ . This is consistent with (30) when  $\max\{\varepsilon_1^2, \frac{p_v}{p}\} = \frac{p_v}{p}$  (see (32)). Further, we note that the decay of the error with  $p$  is slower for smaller  $n$  (cf.  $n = 2^{11}$  vs  $n = 2^{14}$ ). This is expected from (30) where the  $n$  required to ensure  $\max\{\varepsilon_1^2, \frac{p_v}{p}\} = \frac{p_v}{p}$  increases with an increase in  $p$ . As a result, for larger  $p$ ,  $\varepsilon_1^2$  comes into the picture

explaining the increased dependence of the error on  $n$  (cf.  $p = 16$  vs  $p = 128$ ).

**Error scaling for  $\widehat{\Delta \mathbf{v}}^{(i)}$**  The trends in the error  $\max_{i \in [n]} \|\widehat{\Delta \mathbf{v}}^{(i)} - \Delta \mathbf{v}^{(i)}\|_2^2$  are similar to the error  $\max_{i \in [n]} \text{MSE}(\widehat{\theta}^{(i)}, \theta^{*(i)})$ . In the bottom row of Figure. 3, we see  $\max_{i \in [n]} \|\widehat{\Delta \mathbf{v}}^{(i)} - \Delta \mathbf{v}^{(i)}\|_2^2$  has a weak dependence on  $n$  for a fixed  $p$  and  $p_v$ , decreases with an increase in  $p$  for any fixed  $n$  and  $p_v$ , and increases with an increase in  $p_v$  for any fixed  $n$  and  $p$ . This is consistent with (31) when  $\max\{\frac{\varepsilon^2}{p_v}, \frac{p_v}{p}\} = \frac{p_v}{p}$  (see (33)). For the same reason mentioned in the previous paragraph, we see a slower decay in the error with  $p$  for smaller  $n$  (cf.  $n = 2^{11}$  vs  $n = 2^{14}$ ), and a higher dependence of the error on  $n$  for larger  $p$  (cf.  $p = 16$  vs  $p = 128$ ).

## 7 Proof Sketch for Theorem. 1: Guarantee on quality of parameter estimate

Our proof of Theorem. 1 proceeds in two stages (see Figure. 4 for an overview). First, we establish (19) for estimating  $\Theta^*$ . Next, we use this guarantee to establish the unit-level guarantee (20) for each of  $\{\theta^{*(1)}, \dots, \theta^{*(n)}\}$  by substituting  $\Theta = \widehat{\Theta}$  in (14), i.e., analyzing the following convex optimization problem:

$$\{\widehat{\theta}^{(1)}, \dots, \widehat{\theta}^{(n)}\} \in \arg \min_{\{\theta^{(1)}, \dots, \theta^{(n)}\} \in \Lambda_\theta^n} \mathcal{L}(\widehat{\Theta}, \theta^{(1)}, \dots, \theta^{(n)}). \quad (34)$$

### 7.1 Estimating the population-level parameter

In the first part, we show that all points  $\underline{\Theta} \in \Lambda_\Theta \times \Lambda_\theta^n$ , such that  $\|\Theta_t - \Theta_t^*\|_2 \geq \varepsilon$  for at least one  $t \in [p]$ , uniformly satisfy

$$\mathcal{L}(\underline{\Theta}) \geq \mathcal{L}(\underline{\Theta}^*) + \Omega(\varepsilon^2) \text{ for } n \geq \frac{ce^{c'\beta} p^2}{\varepsilon^4} \cdot \left( p \log \frac{p}{\delta \varepsilon^2} + \mathcal{M}_{\theta, n}(\varepsilon^2) \right), \quad (35)$$

with probability at least  $1 - \delta$ . Then, we conclude the proof using contraposition.

To prove (35), we first decompose the convex (and positive) objective  $\mathcal{L}(\underline{\Theta})$  in (13) as a sum of  $p$  convex (and positive) auxiliary objectives  $\mathcal{L}_t$ , namely,  $\mathcal{L}(\underline{\Theta}) = \sum_{t \in [p]} \mathcal{L}_t(\underline{\Theta}_t)$  where

$$\mathcal{L}_t(\underline{\Theta}_t) \triangleq \frac{1}{n} \sum_{i \in [n]} \exp \left( - [\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top x_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \left[ [x_t^{(i)}]^2 - \frac{x_{\max}^2}{3} \right] \right). \quad (36)$$

Next, for any fixed  $t \in [p]$ ,  $\varepsilon > 0$ , and  $\underline{\Theta} \in \Lambda_\Theta^n \times \Lambda_\Theta$  with  $\|\Theta_t - \Theta_t^*\|_2 \geq \varepsilon$ , we show (see Lemma. B.1)

$$\mathcal{L}_t(\underline{\Theta}_t) \geq \mathcal{L}_t(\underline{\Theta}_t^*) + \Omega(\varepsilon^2) - \varepsilon_1 \quad \text{whenever } n \geq \frac{ce^{c'\beta} \log \frac{p}{\delta}}{\varepsilon_1^2}, \quad (37)$$

and then establish the same bound uniformly for all  $t \in [p]$  with probability  $1 - \delta$ . Taking a sum over  $t$  on both sides of (37), we conclude that for any fixed  $\underline{\Theta}$  with  $\|\Theta_t - \Theta_t^*\|_2 \geq \varepsilon$  for some  $t \in [p]$ ,

$$\mathcal{L}(\underline{\Theta}) \geq \mathcal{L}(\underline{\Theta}^*) + \Omega(\varepsilon^2) \quad \text{whenever } n \geq \frac{ce^{c'\beta} p^2 \log \frac{p}{\delta}}{\varepsilon^4}, \quad (38)$$

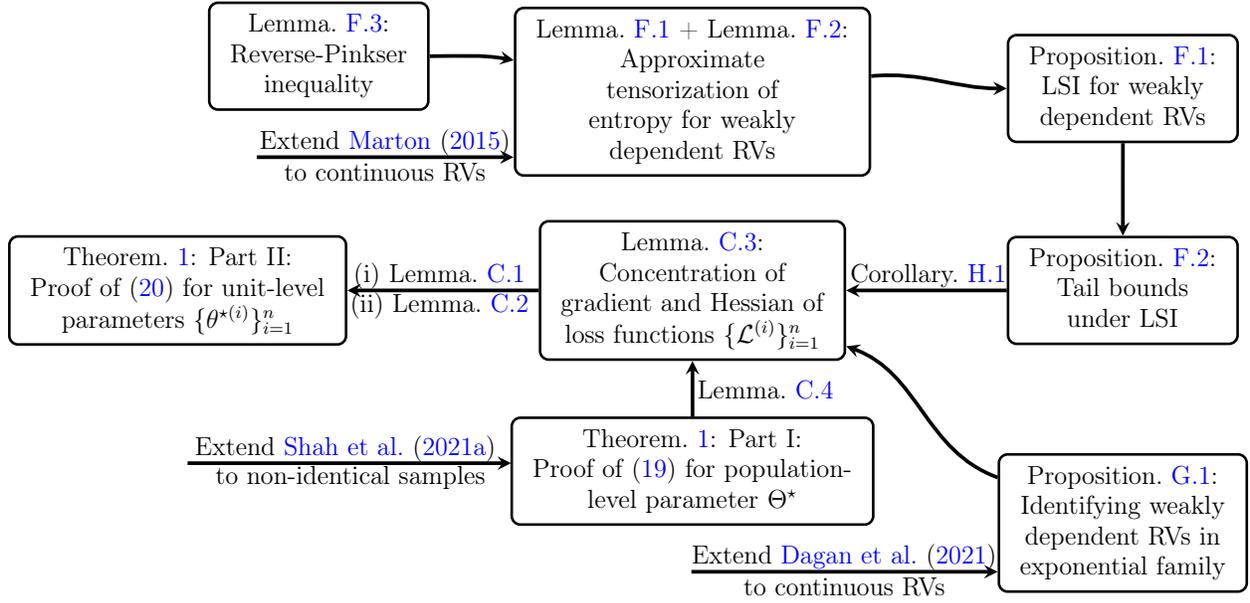


Figure 4: **Sketch diagram of the results and the proof techniques for Theorem 1.** First, we establish (19) for estimating  $\Theta^*$  by extending Shah et al. (2021a, Proposition I.1, Proposition I.2) for i.i.d. data to non-identical samples. Next, we use (19) to establish (20) for the unit-level parameters  $\{\theta^{*(i)}\}_{i=1}^n$  via suitable concentration results for derivatives of the auxiliary loss functions in k(39). En route, we establish three results of independent interest: (i) Proposition. F.1 that shows that weakly dependent and bounded random variables satisfy logarithmic Sobolev inequality (LSI) by both extending Marton (2015, Theorem. 1, Theorem. 2) and establishing a reverse-Pinkser inequality to continuous random vectors; (ii) Proposition. F.2 that extends the tail bounds Dagan et al. (2021, Theorem. 6) to continuous distributions satisfying LSI; and (iii) Proposition. G.1 that extends the conditioning trick Dagan et al. (2021, Lemma. 2) for identifying a weakly dependent subset to continuous random vectors.

with probability at least  $1-\delta$  where we substituted  $\varepsilon_1 = c\varepsilon^2/p$ . Finally, we conclude (35) by using (38), the Lipschitzness of  $\mathcal{L}$  (see Lemma. B.2), and a covering number argument (see Appendix. B).

We establish (37) (Lemma. B.1) via Lemma. B.3, which provides suitable concentration and anti-concentration results for the first-order and second-order derivatives, respectively, for the auxiliary objective  $\mathcal{L}_t$  in (36). We prove Lemma. B.3 by extending the results from Shah et al. (2021a) to the setting with non-identical but independent samples  $\{\mathbf{x}^{(i)} \sim f_{\mathbf{x}|\mathbf{z}}(\cdot | \mathbf{z}^{(i)}; \theta^*(\mathbf{z}^{(i)}), \Theta^*)\}_{i=1}^n$ .

## 7.2 Estimating the unit-level parameters

In the second part, we decompose the convex optimization problem in (34) into  $n$  convex optimization problems:

$$\mathcal{L}^{(i)}(\theta^{(i)}) \triangleq \sum_{t \in [p]} \exp \left( - [\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} ([x_t^{(i)}]^2 - \frac{x_{\max}^2}{3}) \right) \text{ for } i \in [n]. \quad (39)$$

Noting that the set  $\Lambda_\theta^n$  places independent constraints on the  $n$  unit-level parameters, namely  $\theta^{(i)} \in \Lambda_\theta$ , independently for all  $i \in [n]$  and combining (13) and (34), we find that

$$\min_{\{\theta^{(1)}, \dots, \theta^{(n)}\} \in \Lambda_\theta^n} \mathcal{L}(\widehat{\Theta}, \theta^{(1)}, \dots, \theta^{(n)}) \stackrel{(39)}{=} \frac{1}{n} \sum_{i \in [n]} \min_{\theta^{(i)} \in \Lambda_\theta} \mathcal{L}^{(i)}(\theta^{(i)}) \implies \widehat{\theta}^{(i)} \in \arg \min_{\theta^{(i)} \in \Lambda_\theta} \mathcal{L}^{(i)}(\theta^{(i)}),$$

for each  $i \in [n]$ . Next, we establish that with probability at least  $1 - \delta$ ,

$$\mathcal{L}^{(i)}(\theta^{(i)}) \geq \mathcal{L}^{(i)}(\theta^{\star(i)}) + R^2(\varepsilon, \delta) \text{ when } n \geq \frac{ce^{c'\beta}p^4}{\varepsilon^4} \left( p \log \frac{p^2}{\delta\varepsilon^2} + \widetilde{\mathcal{M}}_{\theta, n}(\varepsilon, \delta) \right), \quad (40)$$

uniformly for all points  $\theta^{(i)} \in \Lambda_\theta$  with  $\|\theta^{(i)} - \theta^{\star(i)}\|_2 \geq R(\varepsilon, \delta)$  (see (17)). We conclude the proof by contraposition with the basic inequality  $\mathcal{L}^{(i)}(\widehat{\theta}^{(i)}) \leq \mathcal{L}^{(i)}(\theta^{\star(i)})$  and a standard union bound over all  $i \in [n]$ .

The proof of (40) mimics the same road map as that for (35). Lemma. C.1 shows that for any fixed  $\theta^{(i)} \in \Lambda_\theta$ , if  $\theta^{(i)}$  is far from  $\theta^{\star(i)}$ , then with high probability  $\mathcal{L}^{(i)}(\theta^{(i)})$  is significantly larger than  $\mathcal{L}^{(i)}(\theta^{\star(i)})$ . We prove Lemma. C.1 via concentration of derivatives of  $\mathcal{L}^{(i)}$  (39) in Lemma. C.3, this objective's Lipschitznes in Lemma. C.2, and a covering number argument (see Appendix. C).

The proof of Lemma. C.3 involves several novel arguments: First, for a  $\tau$ -Sparse Graphical Model (Definition. G.1), i.e., a generalization of the random vector  $\mathbf{w}$  in (2), Proposition. G.1 identifies a subset that satisfies Dobrushin's uniqueness condition (Definition. F.2) after conditioning on the complementary subset. Second, Proposition. F.1 shows that a bounded and weakly dependent continuous random vector (defined using Dobrushin's uniqueness condition) satisfies the logarithmic Sobolev inequality (LSI). Third, Proposition. F.2 establishes tail bounds for arbitrary functions of a continuous random vector that satisfies LSI. Putting together these results and a robustness result (Lemma. C.4) while invoking concentration results to account for the estimation error for  $\Theta^*$ , yields Lemma. C.3.

## 8 Discussion

We introduce an exponential family approach to learn unit-level counterfactual distributions from a single sample per unit even when there is unobserved confounding. By conditioning on the latent confounders and using a novel convex loss function, we estimate the parameters of unit-level counterfactual distributions given the information about what actually happened. The resulting estimates of unit-level counterfactual distributions enable us to estimate any functional of each unit's potential outcomes under alternate interventions. We analyze each unit's expected potential outcomes under alternate interventions, thereby providing a guarantee on unit-level counterfactual effects, i.e., individual treatment effects. We note that our approach makes only macro-level assumptions about the underlying causal graph and does not assume the knowledge of the micro-level causal graph.

A side product of our results is a strategy for answering interventional questions, e.g., to estimate average treatment effects. These questions are equivalent to estimating distributions of the form  $f_{\mathbf{y}|\text{do}(\mathbf{a})}(\mathbf{y}|\text{do}(\mathbf{a} = \mathbf{a}))$  where the do-operator (Pearl, 2009) forces  $\mathbf{a}$  to be  $\mathbf{a}$ . Under the causal framework considered (Figure. 1(b)), we have  $f_{\mathbf{y}|\text{do}(\mathbf{a})}(\mathbf{y}|\text{do}(\mathbf{a} = \mathbf{a})) = \mathbb{E}_{\mathbf{v}, \mathbf{z}}[f_{\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v}}(\mathbf{y}|\mathbf{a}, \mathbf{z}, \mathbf{v})]$ . Consequently, the mixture distribution  $n^{-1} \sum_{i \in [n]} \widehat{f}_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a})$  with  $\widehat{f}_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a})$  defined in (15), serves as a natural estimate via our strategy. Investigating the efficacy of this estimator is an interesting future direction.

In this work, the conditional exponential family distribution of  $\mathbf{y}$  in Section. 3.2 or in Section. 5.2 was such that the effect of unobserved covariates  $\mathbf{z}$ —after conditioning on them—was captured by a

first-order interaction term varying with the realized value of  $\mathbf{z}$  for each unit, e.g.,  $\{\theta(\mathbf{z}^{(i)})\}_{i=1}^n$  for the conditional distribution in Section. 3.2. Focusing on Section. 3.2, when one considers higher-order interaction terms in the joint distribution, the conditional distributions would also have higher-order interaction terms (the highest order in the conditional distribution is one less than the highest order in the joint distribution) that vary with  $\mathbf{z}$ . Focusing on Section. 5.2, the exponent of the exponential tilting of the base distribution of the outcomes by the unobserved covariates could have higher-order terms. For such cases, while our analysis for population-level parameters (Theorem. 1 Part I’s proof in Appendix. B) is likely to extend easily, new arguments for analyzing quadratic (or higher-order) interaction terms that vary for each unit seem necessary. Developing these results, e.g., suitable analogs of Dobrushin’s condition for higher-order exponential family, present an exciting future venue for research.

Our methodology can be useful for a class of multi-task learning problems (Caruana, 1997), e.g., when we have multiple logistic regression tasks with some commonalities. For a logistic regression task, the exponential family model (6) has been used by Dagan et al. (2021) to allow dependencies between the labels via the parameter  $\Theta$  (instead of assuming independence between the labels), e.g., for spatio-temporal data. They consider a single regression task and assume that the dependency matrix  $\Theta$  is known up to a constant and learn a task-specific parameter  $\theta(\mathbf{z})$  (where  $\mathbf{z}$  denotes a task). Our model and methodology apply to the case of fully unknown  $\Theta$  given multiple datasets that share the same dependency parameter  $\Theta$  but have varying task-specific parameters  $\theta(\mathbf{z})$ ; and provide a tractable way to estimate all these parameters together. In fact, our framework and results also apply beyond the quadratic dependencies captured by  $\Theta$  as described in Section. 5.2. Analyzing whether our methodology can be extended beyond logistic regression models for multi-task learning is a question worthy of further investigation.

## Acknowledgments

The authors thank Thomas Courtade, Yuzhou Gu, Anuran Makur, Wenlong Mou, Felix Otto, and Yury Polyanskiy for helpful pointers regarding Logarithmic Sobolev inequalities. The authors thank Yuval Dagan and Anthimos Vardis Kandiros for helpful discussion about Dagan et al. (2021) and Kandiros et al. (2021). The authors also thank Alberto Abadie, Avi Feller, and Martin Wainwright for helpful comments. Lastly, the authors thank the anonymous reviewers of Workshop on Causality for Real-world Impact (NeurIPS 2022) for their comments and suggestions.

## Funding

This work was supported, in part, by NSF under Grant No. DMS-2023528 as part of the Foundations of Data Science Institute (FODSI), the MIT-IBM Watson AI Lab under Agreement No. W1771646, MIT-IBM projects on Time Series and Causal Inference as well as project with DSO National Laboratory.

# Appendix

<b>A</b>	<b>Proper loss function and projected gradient descent</b>	<b>26</b>
A.1	Proof of Proposition 1 . . . . .	26
A.2	Algorithm . . . . .	27

<b>B</b>	<b>Proof of Theorem 1 Part I: Recovering population-level parameter</b>	<b>28</b>
B.1	Proof of Lemma. B.1: Gap between the loss function for a fixed parameter . . . . .	30
B.2	Example for Assumption. 2 . . . . .	37
B.3	Proof of Lemma. B.2: Lipschitzness of the loss function . . . . .	39
<b>C</b>	<b>Proof of Theorem 1 Part II: Recovering unit-level parameters</b>	<b>40</b>
C.1	Proof of Lemma. C.1: Gap between the loss function for a fixed parameter . . . . .	42
C.2	Proof of Lemma. C.2: Lipschitzness of the loss function . . . . .	51
<b>D</b>	<b>Proof of Theorem 2: Guarantee on quality of outcome estimate</b>	<b>52</b>
D.1	Bounded operator norms for perturbations in the parameters . . . . .	55
<b>E</b>	<b>Proof of Proposition 2: Impute missing covariates</b>	<b>55</b>
<b>F</b>	<b>Logarithmic Sobolev inequality and tail bounds</b>	<b>59</b>
F.1	Proof of Proposition. F.1: Logarithmic Sobolev inequality . . . . .	60
F.2	Proof of Proposition. F.2: Tail bounds for arbitrary functions under LSI . . . . .	68
<b>G</b>	<b>Identifying weakly dependent random variables</b>	<b>72</b>
<b>H</b>	<b>Supporting concentration results</b>	<b>75</b>
H.1	Proof of Lemma. H.1: Logarithmic Sobolev inequality for $x_t \mathbf{x}_{-t}, \mathbf{z}$ . . . . .	75
H.2	Proof of Corollary. H.1: Supporting concentration bounds . . . . .	76

## A Proper loss function and projected gradient descent

In this section, we prove Proposition. 1 showing that the loss function in (13) is a proper loss function. We also provide an algorithm to obtain an  $\epsilon$ -optimal estimate of  $\hat{\Theta}$ .

### A.1 Proof of Proposition 1

Fix any  $\mathbf{z} \in \mathcal{Z}^{p_z}$ . For every  $t \in [p]$ , define the following parametric distribution

$$u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t(\mathbf{z}), \Theta_t) \propto \frac{f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*)}{f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}(x_t|\mathbf{x}_{-t}, \mathbf{z}; \theta_t(\mathbf{z}), \Theta_t)}, \quad (41)$$

where  $f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*)$  is as defined in (6) and  $f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}(x_t|\mathbf{x}_{-t}, \mathbf{z}; \theta_t(\mathbf{z}), \Theta_t)$  is as defined in (12). Letting  $\bar{x}_t \triangleq x_t^2 - x_{\max}^2/3$  and using (12), we can write  $u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t(\mathbf{z}), \Theta_t)$  in (41) as

$$u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t(\mathbf{z}), \Theta_t) \propto f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*) \exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t).$$

Then, we have

$$\begin{aligned} u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t(\mathbf{z}), \Theta_t) &= \frac{f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*) \exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t)}{\int_{\mathbf{x} \in \mathcal{X}^p} f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*) \exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t) d\mathbf{x}} \\ &= \frac{f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*) \exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t)}{\mathbb{E}_{\mathbf{x}|\mathbf{z}}[\exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t)]}. \end{aligned} \quad (42)$$

Further, for  $\theta_t(\mathbf{z}) = \theta_t^*(\mathbf{z})$ , and  $\Theta_t = \Theta_t^*$ , we can write an expression for  $u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*)$  which does not depend on  $x_t$  functionally. From (12), we have

$$u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) \propto f_{\mathbf{x}_{-t}|\mathbf{z}}(\mathbf{x}_{-t}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*). \quad (43)$$

Now, consider the difference between  $\text{KL}(u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) \| u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t(\mathbf{z}), \Theta_t))$  and  $\text{KL}(u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) \| f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*))$ . We have

$$\begin{aligned} & \text{KL}(u_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) \| u_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}; \theta_t(\mathbf{z}), \Theta_t)) \\ & \quad - \text{KL}(u_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) \| f_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*)) \\ \stackrel{(a)}{=} & \int_{\mathbf{x} \in \mathcal{X}^p} u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) \log \frac{f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*)}{u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t(\mathbf{z}), \Theta_t)} d\mathbf{x} \\ \stackrel{(42)}{=} & \int_{\mathbf{x} \in \mathcal{X}^p} u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) \log \frac{\mathbb{E}_{\mathbf{x}|\mathbf{z}} \left[ \exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t) \right]}{\exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t)} d\mathbf{x} \\ = & \log \mathbb{E}_{\mathbf{x}|\mathbf{z}} \left[ \exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t) \right] \\ & \quad - \int_{\mathbf{x} \in \mathcal{X}^p} u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*) ([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t + \Theta_{tt}\bar{x}_t) d\mathbf{x} \\ \stackrel{(b)}{=} & \log \mathbb{E}_{\mathbf{x}|\mathbf{z}} \left[ \exp(-[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t) \right], \end{aligned} \quad (44)$$

where (a) follows from the definition of KL-divergence and (b) follows because integral is zero since  $u_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta_t^*(\mathbf{z}), \Theta_t^*)$  does not functionally depend on  $x_t$  as in (43), and  $\int_{x_t \in \mathcal{X}} x_t dx_t = 0$  and  $\int_{x_t \in \mathcal{X}} \bar{x}_t dx_t = 0$ . Now, we can write

$$\begin{aligned} \mathbb{E}_{\mathbf{x}|\mathbf{z}}[\mathcal{L}(\underline{\Theta})] &= \frac{1}{n} \sum_{t \in [p]} \sum_{i \in [n]} \mathbb{E}_{\mathbf{x}|\mathbf{z}} \left[ \exp(-[\theta_t(\mathbf{z}^{(i)}) + \Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \Theta_{tt}\bar{x}_t^{(i)}) \right] \\ \stackrel{(44)}{=} & \frac{1}{n} \sum_{t \in [p]} \sum_{i \in [n]} \exp \left( \text{KL} \left( u_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}^{(i)}; \theta_t^*(\mathbf{z}^{(i)}), \Theta_t^*) \| u_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}^{(i)}; \theta_t(\mathbf{z}^{(i)}), \Theta_t) \right) \right. \\ & \quad \left. - \text{KL} \left( u_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}^{(i)}; \theta_t^*(\mathbf{z}^{(i)}), \Theta_t^*) \| f_{\mathbf{x}|\mathbf{z}}(\cdot|\mathbf{z}^{(i)}; \theta^*(\mathbf{z}^{(i)}), \Theta^*) \right) \right). \end{aligned} \quad (45)$$

We note that the parameters only show up in the first KL-divergence term in the right-hand-side of (45). Therefore, it is easy to see that  $\mathbb{E}_{\mathbf{x}|\mathbf{z}}[\mathcal{L}(\underline{\Theta})]$  is minimized uniquely when  $\theta_t(\mathbf{z}^{(i)}) = \theta_t^*(\mathbf{z}^{(i)})$  and  $\Theta_t = \Theta_t^*$  for all  $t \in [p]$  and all  $i \in [n]$ , i.e., when  $\underline{\Theta} = \underline{\Theta}^*$ .

## A.2 Algorithm

In this section, we provide a projected gradient descent algorithm to return an  $\epsilon$ -optimal estimate of the convex optimization in (14). We note that alternative algorithms (including Frank-Wolfe) can also be used.

We note that, in general, projecting onto the space  $\Lambda_\theta^n \times \Lambda_\Theta$  may not be easy depending on the specific form of  $\Lambda_\theta$ . For Examples 1 and 2, projecting on  $\Lambda_\theta$  is equivalent to projecting onto the  $k$ -dimensional vector  $\mathbf{a}$ . For Example 2, the  $\ell_0$ -sparsity is relaxed to  $\ell_1$  sparsity. We also do not focus on any issues that may arise due to the choice of the step size  $\eta$ .

---

**Algorithm 1:** Projected Gradient Descent

---

**Input:** number of iterations  $\tau$ , step size  $\eta$ ,  $\epsilon$ , parameter sets  $\Lambda_\theta$  and  $\Lambda_\Theta$   
**Output:**  $\epsilon$ -optimal estimate  $\widehat{\Theta}_\epsilon$   
**Initialization:**  $\underline{\Theta}^{(0)} = \mathbf{0}$   
**1** for  $j = 0, \dots, \tau$  do  
**2**     $\underline{\Theta}^{(j+1)} \leftarrow \arg \min_{\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta} \|\underline{\Theta}^{(j)} - \eta \nabla \mathcal{L}(\underline{\Theta}^{(j)}) - \underline{\Theta}\|_2$   
**3**     $\widehat{\Theta}_\epsilon \leftarrow \underline{\Theta}^{(\tau+1)}$

---

## B Proof of Theorem 1 Part I: Recovering population-level parameter

To prove this part, it is sufficient to show that all points  $\underline{\Theta} \in \Lambda_\Theta \times \Lambda_\theta^n$ , such that  $\|\Theta_t - \Theta_t^*\|_2 \geq \epsilon$  for at least one  $t \in [p]$ , uniformly satisfy

$$\mathcal{L}(\underline{\Theta}) \geq \mathcal{L}(\underline{\Theta}^*) + \Omega(\epsilon^2) \text{ for } n \geq \frac{ce^{c'\beta}p^2}{\epsilon^4} \cdot \left( p \log \frac{p}{\delta} + \mathcal{M}_{\theta,n}(\epsilon^2) \right), \quad (46)$$

with probability at least  $1 - \delta$ . Then, the guarantee in Theorem. 1 follows from (14) by contraposition.

To that end, we decompose  $\mathcal{L}(\underline{\Theta})$  in (13) as a sum of  $p$  convex (and positive) auxiliary objectives  $\mathcal{L}_t(\underline{\Theta}_t)$ , i.e.,  $\mathcal{L}(\underline{\Theta}) = \sum_{t \in [p]} \mathcal{L}_t(\underline{\Theta}_t)$  where

$$\mathcal{L}_t(\underline{\Theta}_t) \triangleq \frac{1}{n} \sum_{i \in [n]} \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top]x_t^{(i)} - \Theta_{tt}\bar{x}_t^{(i)} \right), \quad (47)$$

with  $\bar{x}_t^{(i)} = [x_t^{(i)}]^2 - x_{\max}^2/3$  and  $\underline{\Theta}_t = \{\theta_t^{(1)}, \dots, \theta_t^{(n)}, \Theta_t\}$  as defined in (13). The lemma below, proven in Appendix. B.1, shows that for any fixed and feasible  $\underline{\Theta}_t$ , if  $\Theta_t$  is far from  $\Theta_t^*$ , then with high probability  $\mathcal{L}_t(\underline{\Theta}_t)$  is significantly larger than  $\mathcal{L}_t(\underline{\Theta}_t^*)$ . The lemma uses the following constants that depend on model parameters  $\tau \triangleq (\alpha, \beta, x_{\max}, \Theta)$ :

$$C_{1,\tau} \triangleq \alpha + 4\beta x_{\max} \quad \text{and} \quad C_{2,\tau} \triangleq \exp(x_{\max}(\alpha + 2\beta x_{\max})). \quad (48)$$

**Lemma B.1** (Gap between the loss function for a fixed parameter). *Consider any  $\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta$ . Fix any  $\delta \in (0, 1)$ . Then, we have uniformly for all  $t \in [p]$*

$$\mathcal{L}_t(\underline{\Theta}_t) \geq \mathcal{L}_t(\underline{\Theta}_t^*) + \frac{\lambda_{\min} \|\Theta_t - \Theta_t^*\|_2^2}{2C_{2,\tau}} - \epsilon \quad \text{for } n \geq \frac{ce^{c'\beta} \log(p/\delta)}{\epsilon^2},$$

with probability at least  $1 - \delta$ , where  $C_{2,\tau}$  was defined in (48).

Next, we show that the loss function  $\mathcal{L}$  is Lipschitz (see Appendix. B.3 for the proof).

**Lemma B.2** (Lipschitzness of the loss function). *Consider any  $\underline{\Theta}, \tilde{\Theta} \in \Lambda_\Theta$ . Then, the loss function  $\mathcal{L}$  is  $2x_{\max}^2 C_{2,\tau}$ -Lipschitz in a suitably-adjusted  $\ell_1$  norm:*

$$|\mathcal{L}(\tilde{\Theta}) - \mathcal{L}(\underline{\Theta})| \leq 2x_{\max}^2 C_{2,\tau} \left( \sum_{t \in [p]} \|\tilde{\Theta}_t - \Theta_t\|_1 + \frac{1}{n} \sum_{i \in [n]} \|\tilde{\theta}^{(i)} - \theta^{(i)}\|_1 \right), \quad (49)$$

where the constant  $C_{2,\tau}$  was defined in (48).

Given these lemmas, we now proceed with the proof.

**Proof strategy** We want to show that all points  $\underline{\Theta} \in \Lambda_{\Theta} \times \Lambda_{\theta}^n$ , such that  $\|\Theta_t - \Theta_t^*\|_2 \geq \varepsilon$  for at least one  $t \in [p]$ , uniformly satisfy (46) with probability at least  $1 - \delta$ . To do so, we consider the set of feasible  $\underline{\Theta}$  such that the distance of  $\Theta_t$  from  $\Theta_t^*$  is at least  $\varepsilon > 0$  in  $\ell_2$  norm for some  $t \in [p]$ , and denote the set by  $\Lambda_{\Theta}^{\varepsilon} \times \Lambda_{\theta}^n$  (see (50) and (9)). Then, using an appropriate covering set of  $\Lambda_{\Theta}^{\varepsilon} \times \Lambda_{\theta}^n$  and the Lipschitzness of  $\mathcal{L}$ , we show that the value of  $\mathcal{L}$  at all points in  $\Lambda_{\Theta}^{\varepsilon} \times \Lambda_{\theta}^n$  is uniformly  $\Omega(\varepsilon^2)$  larger than the value of  $\mathcal{L}$  at  $\underline{\Theta}^*$  with high probability.

**Arguments for points in the covering set** Define the set

$$\Lambda_{\Theta}^{\varepsilon} \triangleq \left\{ \Theta \in \mathbb{R}^{p \times p} : \Theta = \Theta^{\top}, \|\Theta\|_{\max} \leq \alpha, \|\Theta\|_{\infty} \leq \beta, \max_{t \in [p]} \|\Theta_t^* - \Theta_t\|_2 \geq \varepsilon \right\}. \quad (50)$$

Let  $\mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon')$  be the  $\varepsilon'$ -cover of smallest size for the set  $\Lambda_{\Theta}^{\varepsilon}$  with respect to  $\|\cdot\|_1$  (see Definition. 2) and let  $\mathcal{C}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') = |\mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon')|$  be the  $\varepsilon'$ -covering number. Similarly, let  $\mathcal{U}(\Lambda_{\theta}^n, \varepsilon'')$  be the  $\varepsilon''$ -cover of the smallest size for the set  $\Lambda_{\theta}^n$  with respect to  $\|\cdot\|_1$  and let  $\mathcal{C}(\Lambda_{\theta}^n, \varepsilon'') = |\mathcal{U}(\Lambda_{\theta}^n, \varepsilon'')|$  be the  $\varepsilon''$ -covering number. We choose

$$\varepsilon' \triangleq \frac{\lambda_{\min} \varepsilon^2}{32x_{\max}^2 C_{2,\tau}^2} \quad \text{and} \quad \varepsilon'' \triangleq \frac{\lambda_{\min} \varepsilon^2 n}{32x_{\max}^2 C_{2,\tau}^2}. \quad (51)$$

Now, we argue by a union bound that the value of  $\mathcal{L}$  at all points in  $\mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \times \mathcal{U}(\Lambda_{\theta}^n, \varepsilon'')$  is uniformly  $\Omega(\varepsilon^2)$  larger than  $\mathcal{L}(\underline{\Theta}^*)$  with high probability. For any  $\underline{\Theta} \in \mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \times \mathcal{U}(\Lambda_{\theta}^n, \varepsilon'')$ , we have

$$\sum_{t \in [p]} \|\Theta_t^* - \Theta_t\|_2^2 \stackrel{(a)}{\geq} \varepsilon^2, \quad (52)$$

where (a) follows because  $\mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \subseteq \Lambda_{\Theta}^{\varepsilon}$ . Now, applying Lemma. B.1 with  $\varepsilon \leftarrow \lambda_{\min} \varepsilon^2 / 4C_{2,\tau} p$  and  $\delta \leftarrow \delta / (\mathcal{C}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') + \mathcal{C}(\Lambda_{\theta}^n, \varepsilon''))$  and summing over  $t \in [p]$ , we find that

$$\begin{aligned} \sum_{t \in [p]} \mathcal{L}_t(\underline{\Theta}_t) &\geq \sum_{t \in [p]} \left( \mathcal{L}_t(\underline{\Theta}_t^*) + \frac{\lambda_{\min} \|\Theta_t - \Theta_t^*\|_2^2}{2C_{2,\tau}} - \frac{\lambda_{\min} \varepsilon^2}{4C_{2,\tau} p} \right) \\ \implies \mathcal{L}(\underline{\Theta}) &\geq \mathcal{L}(\underline{\Theta}^*) + \frac{\lambda_{\min}}{2C_{2,\tau}} \sum_{t \in [p]} \|\Theta_t^* - \Theta_t\|_2^2 - \frac{\lambda_{\min} \varepsilon^2}{4C_{2,\tau}} \\ &\stackrel{(52)}{\geq} \mathcal{L}(\underline{\Theta}^*) + \frac{\lambda_{\min} \varepsilon^2}{4C_{2,\tau}}, \end{aligned}$$

with probability at least  $1 - \delta / (\mathcal{C}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') + \mathcal{C}(\Lambda_{\theta}^n, \varepsilon''))$  whenever

$$n \geq \frac{ce^{c'\beta} p^2 \log((\mathcal{C}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \times \mathcal{C}(\Lambda_{\theta}^n, \varepsilon'')) \cdot p / \delta)}{\lambda_{\min}^2 \varepsilon^4}. \quad (53)$$

By applying the union bound over  $\mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \times \mathcal{U}(\Lambda_{\theta}^n, \varepsilon'')$ , as long as  $n$  satisfies (53), we have

$$\mathcal{L}(\underline{\Theta}) \geq \mathcal{L}(\underline{\Theta}^*) + \frac{\lambda_{\min} \varepsilon^2}{4C_{2,\tau}} \quad \text{uniformly for every } \underline{\Theta} \in \mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \times \mathcal{U}(\Lambda_{\theta}^n, \varepsilon''), \quad (54)$$

with probability at least  $1 - \delta$ .

**Arguments for points outside the covering set** Next, we establish the claim (46) for an arbitrary  $\tilde{\Theta} \in \Lambda_{\Theta}^{\varepsilon} \times \Lambda_{\theta}^n$  conditional on the event that (54) holds. Given a fixed  $\tilde{\Theta} \in \Lambda_{\Theta}^{\varepsilon} \times \Lambda_{\theta}^n$ , let  $\underline{\Theta}$  be (one of) the point(s) in the cover  $\mathcal{U}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \times \mathcal{U}(\Lambda_{\theta}^n, \varepsilon'')$  that satisfies  $\sum_{t \in [p]} \|\tilde{\Theta}_t - \Theta_t\|_1 \leq \varepsilon'$  and  $\sum_{i \in [n]} \|\tilde{\theta}^{(i)} - \theta^{(i)}\|_1 \leq \varepsilon''$  (there exists such a point by Definition. 2). Then, the choices (51) and Lemma. B.2 put together imply that

$$\begin{aligned} \mathcal{L}(\tilde{\Theta}) &\geq \mathcal{L}(\underline{\Theta}) - 2x_{\max}^2 C_{2,\tau} \left( \sum_{t \in [p]} \|\tilde{\Theta}_t - \Theta_t\|_1 + \frac{1}{n} \sum_{i \in [n]} \|\tilde{\theta}^{(i)} - \theta^{(i)}\|_1 \right) \\ &\geq \mathcal{L}(\underline{\Theta}) - 2x_{\max}^2 C_{2,\tau} \left( \varepsilon' + \frac{\varepsilon''}{n} \right) \stackrel{(51)}{\geq} \mathcal{L}(\underline{\Theta}) - \frac{\lambda_{\min} \varepsilon^2}{8C_{2,\tau}} \stackrel{(54)}{\geq} \mathcal{L}(\Theta^*) + \frac{\lambda_{\min} \varepsilon^2}{8C_{2,\tau}}. \end{aligned}$$

**Bounding  $n$**  Using  $\Lambda_{\Theta}^{\varepsilon} \subseteq \Lambda_{\Theta}$  and the outer product definition of  $\theta^n$ , we find that

$$\mathcal{C}(\Lambda_{\Theta}^{\varepsilon}, \varepsilon') \leq \mathcal{C}(\Lambda_{\Theta}, \varepsilon') \quad \text{and} \quad \mathcal{C}(\Lambda_{\theta}^n, \varepsilon'') = (\mathcal{C}(\Lambda_{\theta}, \varepsilon''))^n. \quad (55)$$

Putting together (51) and (55), the lower bound (53) can be replaced by

$$n \geq \frac{ce^{c'\beta} p^2}{\lambda_{\min}^2 \varepsilon^4} \cdot \left( \log \frac{p}{\delta} + \log \mathcal{C} \left( \Lambda_{\Theta}, \frac{\lambda_{\min} \varepsilon^2}{ce^{c'\beta}} \right) + n \log \mathcal{C} \left( \Lambda_{\theta}, \frac{\lambda_{\min} n \varepsilon^2}{ce^{c'\beta}} \right) \right),$$

which yields the claim immediately after noting that

$$\log \mathcal{C} \left( \Lambda_{\Theta}, \frac{\lambda_{\min} \varepsilon^2}{ce^{c'\beta}} \right) = O \left( \beta^2 p \log \left( \frac{1}{\lambda_{\min} \varepsilon^2} \right) \right) \quad \text{and} \quad \log \mathcal{C} \left( \Lambda_{\theta}, \frac{\lambda_{\min} n \varepsilon^2}{ce^{c'\beta}} \right) = \mathcal{M}_{\theta} \left( \frac{\lambda_{\min} n \varepsilon^2}{ce^{c'\beta}} \right).$$

## B.1 Proof of Lemma. B.1: Gap between the loss function for a fixed parameter

Fix any  $\varepsilon > 0$ , any  $\delta \in (0, 1)$ , and  $t \in [p]$ . Consider any direction  $\underline{\Omega}_t \triangleq \{\omega_t^{(1)}, \dots, \omega_t^{(n)}, \Omega_t\} \in \mathbb{R}^{n+p}$  along the parameter  $\underline{\Theta}_t$ , i.e.,

$$\underline{\Omega}_t = \underline{\Theta}_t - \underline{\Theta}_t^*, \quad \text{and} \quad \Omega_t = \Theta_t - \Theta_t^*. \quad (56)$$

We denote the first-order and the second-order directional derivatives of the loss function  $\mathcal{L}_t$  in (47) along the direction  $\underline{\Omega}_t$  evaluated at  $\underline{\Theta}_t$  by  $\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t)$  and  $\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t)$ , respectively. Below, we state a lemma (with proof divided across Appendix. B.1.1 and Appendix. B.1.2) that provides us a control on  $\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t)$  and  $\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t)$ . The assumptions of Lemma. B.1 remain in force.

**Lemma B.3** (Control on first and second directional derivatives). *For any fixed  $\varepsilon_1, \varepsilon_2 > 0$ ,  $\delta_1, \delta_2 \in (0, 1)$ ,  $t \in [p]$ ,  $\underline{\Theta} \in \Lambda_{\theta}^n \times \Lambda_{\Theta}$  defined in (13) and  $\Omega_t$  defined in (56), we have the following:*

(a) Concentration of first directional derivative: *with probability at least  $1 - \delta_1$ ,*

$$|\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*)| \leq \varepsilon_1 \quad \text{for } n \geq \frac{8C_{1,\tau}^2 C_{2,\tau}^2 x_{\max}^2 \log \frac{2p}{\delta_1}}{\varepsilon_1^2} \quad \text{and uniformly for all } t \in [p].$$

(b) Anti-concentration of second directional derivative: *with probability at least  $1 - \delta_2$ ,*

$$\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t) \geq \frac{\lambda_{\min} \|\Omega_t\|_2^2}{C_{2,\tau}} - \varepsilon_2 \quad \text{for } n \geq \frac{32C_{1,\tau}^4 x_{\max}^4 \log \frac{2p}{\delta_2}}{\varepsilon_2^2 C_{2,\tau}^2} \quad \text{and uniformly for all } t \in [p].$$

Given this lemma, we now proceed with the proof. Define a function  $g : [0, 1] \rightarrow \mathbb{R}^{n+p}$

$$g(a) \triangleq \underline{\Theta}_t^* + a(\underline{\Theta}_t - \underline{\Theta}_t^*).$$

Notice that  $g(0) = \underline{\Theta}_t^*$  and  $g(1) = \underline{\Theta}_t$  as well as

$$\frac{d\mathcal{L}_t(g(a))}{da} = \partial_{\underline{\Omega}_t} \mathcal{L}_t(\tilde{\Theta}_t) \Big|_{\tilde{\Theta}_t=g(a)} \quad \text{and} \quad \frac{d^2\mathcal{L}_t(g(a))}{da^2} = \partial_{\underline{\Omega}_t^2}^2 \mathcal{L}_t(\tilde{\Theta}_t) \Big|_{\tilde{\Theta}_t=g(a)}. \quad (57)$$

By the fundamental theorem of calculus, we have

$$\frac{d\mathcal{L}_t(g(a))}{da} \geq \frac{d\mathcal{L}_t(g(a))}{da} \Big|_{a=0} + a \min_{a \in (0,1)} \frac{d^2\mathcal{L}_t(g(a))}{da^2}. \quad (58)$$

Integrating both sides of (58) with respect to  $a$ , we obtain

$$\begin{aligned} \mathcal{L}_t(g(a)) - \mathcal{L}_t(g(0)) &\geq a \frac{d\mathcal{L}_t(g(a))}{da} \Big|_{a=0} + \frac{a^2}{2} \min_{a \in (0,1)} \frac{d^2\mathcal{L}_t(g(a))}{da^2} \\ &\stackrel{(57)}{=} a \partial_{\underline{\Omega}_t} \mathcal{L}_t(\tilde{\Theta}_t) \Big|_{\tilde{\Theta}_t=g(0)} + \frac{a^2}{2} \min_{a \in (0,1)} \partial_{\underline{\Omega}_t^2}^2 \mathcal{L}_t(\tilde{\Theta}_t) \Big|_{\tilde{\Theta}_t=g(a)} \\ &\stackrel{(a)}{=} a \partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*) + \frac{a^2}{2} \min_{a \in (0,1)} \partial_{\underline{\Omega}_t^2}^2 \mathcal{L}_t(\tilde{\Theta}_t) \Big|_{\tilde{\Theta}_t=g(a)} \\ &\stackrel{(b)}{\geq} -a |\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*)| + \frac{a^2}{2} \min_{a \in (0,1)} \partial_{\underline{\Omega}_t^2}^2 \mathcal{L}_t(\tilde{\Theta}_t) \Big|_{\tilde{\Theta}_t=g(a)}, \end{aligned} \quad (59)$$

where (a) follows because  $g(0) = \underline{\Theta}_t^*$  and (b) follows by the triangle inequality. Plugging in  $a = 1$  in (59) as well as using  $g(0) = \underline{\Theta}_t^*$  and  $g(1) = \underline{\Theta}_t$ , we find that

$$\mathcal{L}_t(\underline{\Theta}_t) - \mathcal{L}_t(\underline{\Theta}_t^*) \geq -|\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*)| + \frac{1}{2} \min_{a \in (0,1)} \partial_{\underline{\Omega}_t^2}^2 \mathcal{L}_t(\tilde{\Theta}_t) \Big|_{\tilde{\Theta}_t=g(a)}.$$

Now, we use Lemma. B.3 with

$$\varepsilon_1 \leftarrow \frac{\varepsilon}{2}, \quad \delta_1 \leftarrow \frac{\delta}{2}, \quad \varepsilon_2 \leftarrow \varepsilon, \quad \text{and} \quad \delta_2 \leftarrow \frac{\delta}{2}.$$

Thus for  $n \geq \frac{ce^{c'\beta} \log(p/\delta)}{\varepsilon^2}$ , we have

$$\mathcal{L}_t(\underline{\Theta}_t) - \mathcal{L}_t(\underline{\Theta}_t^*) \geq -\frac{\varepsilon}{2} + \frac{1}{2} \left( \frac{\lambda_{\min} \|\Omega_t\|_2^2}{C_{2,\tau}} - \varepsilon \right) = \frac{\lambda_{\min} \|\Omega_t\|_2^2}{2C_{2,\tau}} - \varepsilon,$$

uniformly for all  $t \in [p]$ , with probability at least  $1 - \delta$ .

### B.1.1 Proof of Lemma. B.3(a): Concentration of first directional derivative

For every  $t \in [p]$  with  $\underline{\Omega}_t$  defined in (56), we claim that the first-order directional derivative of the loss function defined in (47) is given by

$$\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t) = -\frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right) \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right), \quad (60)$$

where  $\Delta_t^{(i)} \triangleq \begin{bmatrix} \omega_t^{(i)} \\ \Omega_{t,-t}^\top \\ \Omega_{tt} \end{bmatrix} \in \mathbb{R}^{p+1}$  and  $\tilde{\mathbf{x}}^{(i)} \triangleq \begin{bmatrix} x_t^{(i)} \\ 2\mathbf{x}_{-t}^{(i)}x_t^{(i)} \\ \bar{x}_t^{(i)} \end{bmatrix} \in \mathbb{R}^{p+1}$  for all  $i \in [n]$  with  $\bar{x}_t^{(i)} = [x_t^{(i)}]^2 - x_{\max}^2/3$ . We provide a proof at the end.

Next, we claim that the mean of the first-order directional derivative evaluated at the true parameter is zero. We provide a proof at the end.

**Lemma B.4** (Zero-meanness of first directional derivative). *For every  $t \in [p]$  with  $\underline{\Omega}_t$  defined in (56), we have  $\mathbb{E}[\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*)] = 0$ .*

Given these, we proceed to show the concentration of the first-order directional derivative evaluated at the true parameter. Fix any  $t \in [p]$ . From (60), we have

$$\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*) \stackrel{(60)}{=} -\frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right) \exp \left( -[\theta_t^{*(i)} + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right).$$

Each term in the above summation is an independent random variable and is bounded as follows

$$\begin{aligned} & \left| \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right) \times \exp \left( -[\theta_t^{*(i)} + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) \right| \\ & \stackrel{(a)}{=} \left| \left( \omega_t^{(i)} x_t^{(i)} + 2\Omega_{t,-t}^\top \mathbf{x}_{-t}^{(i)} x_t^{(i)} + \Omega_{tt} \bar{x}_t^{(i)} \right) \times \exp \left( -[\theta_t^{*(i)} + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) \right| \\ & \stackrel{(b)}{\leq} \left| \omega_t^{(i)} \right| + 2\|\Omega_{t,-t}\|_1 \|\mathbf{x}^{(i)}\|_\infty \times x_{\max} \times \exp \left( (|\theta_t^{*(i)}| + 2\|\Theta_{t,-t}^*\|_1 \|\mathbf{x}^{(i)}\|_\infty) x_{\max} \right) \\ & \stackrel{(c)}{\leq} (2\alpha + 8\beta x_{\max}) \times x_{\max} \times \exp \left( (\alpha + 2\beta x_{\max}) x_{\max} \right) \stackrel{(48)}{=} 2C_{1,\tau} C_{2,\tau} x_{\max}, \end{aligned}$$

where (a) follows by plugging in  $\Delta_t^{(i)}$  and  $\tilde{\mathbf{x}}^{(i)}$ , (b) follows from triangle inequality, Cauchy–Schwarz inequality, and because  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , and (c) follows because  $\theta^{*(i)} \in \Lambda_\theta$  for all  $i \in [n]$ ,  $\Theta^* \in \Lambda_\Theta$ ,  $\omega^{(i)} \in 2\Lambda_\theta$  for all  $i \in [n]$ ,  $\Omega \in 2\Lambda_\Theta$ , and  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ .

Further, from Lemma B.4, we have  $\mathbb{E}[\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*)] = 0$ . Therefore, using the Hoeffding’s inequality results in

$$\mathbb{P} \left( \left| \partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t^*) \right| > \varepsilon_1 \right) < 2 \exp \left( -\frac{n\varepsilon_1^2}{8C_{1,\tau}^2 C_{2,\tau}^2 x_{\max}^2} \right).$$

The proof follows by using the union bound over all  $t \in [p]$ .

**Proof of (60): Expression for first directional derivative** Fix any  $t \in [p]$ . The first-order partial derivatives of  $\mathcal{L}_t$  with respect to entries of  $\underline{\Theta}_t$  defined in (47) are given by

$$\begin{aligned} \frac{\partial \mathcal{L}_t(\underline{\Theta}_t)}{\partial \theta_t^{(i)}} &= \frac{-1}{n} x_t^{(i)} \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \text{ for all } i \in [n], \quad \text{and} \\ \frac{\partial \mathcal{L}_t(\underline{\Theta}_t)}{\partial \Theta_{tu}} &= \begin{cases} \frac{-2}{n} \sum_{i \in [n]} x_t^{(i)} x_u^{(i)} \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \text{ for all } u \in [p] \setminus \{t\}. \\ \frac{-1}{n} \sum_{i \in [n]} \bar{x}_t^{(i)} \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \text{ for } u = t. \end{cases} \end{aligned}$$

Now, we can write the first-order directional derivative of  $\mathcal{L}_t$  as

$$\begin{aligned}
\partial_{\Omega_t} \mathcal{L}_t(\Theta_t) &\triangleq \lim_{h \rightarrow 0} \frac{\mathcal{L}_t(\Theta_t + h\Omega_t) - \mathcal{L}_t(\Theta_t)}{h} = \sum_{i \in [n]} \omega_t^{(i)} \frac{\partial \mathcal{L}_t(\Theta_t)}{\partial \theta_t^{(i)}} + \sum_{u \in [p]} \Omega_{tu} \frac{\partial \mathcal{L}_t(\Theta_t)}{\partial \Theta_{tu}} \\
&= \frac{-1}{n} \sum_{i \in [n]} \left( \omega_t^{(i)} x_t^{(i)} + 2 \sum_{u \in [p] \setminus \{t\}} \Omega_{tu} x_t^{(i)} x_u^{(i)} + \Omega_{tt} \bar{x}_t^{(i)} \right) \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \\
&= \frac{-1}{n} \sum_{i \in [n]} \left( \omega_t^{(i)} x_t^{(i)} + 2\Omega_{t,-t}^\top \mathbf{x}_{-t}^{(i)} x_t^{(i)} + \Omega_{tt} \bar{x}_t^{(i)} \right) \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \\
&\stackrel{(a)}{=} \frac{-1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right) \exp \left( -[\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right),
\end{aligned}$$

where (a) follows from the definitions of  $\Delta_t^{(i)}$  and  $\tilde{\mathbf{x}}^{(i)}$ .

**Proof of Lemma. B.4: Zero-meanness of first directional derivative** Fix any  $t \in [p]$ . From (60), we have

$$\begin{aligned}
&\mathbb{E}[\partial_{\Omega_t} \mathcal{L}_t(\Theta_t^*)] \\
&\stackrel{(60)}{=} -\frac{1}{n} \sum_{i \in [n]} \mathbb{E}_{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right) \exp \left( -[\theta_t^{*(i)} + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) \right] \\
&\stackrel{(a)}{=} -\frac{1}{n} \sum_{i \in [n]} \sum_{u \in [p+1]} \mathbb{E}_{\mathbf{z}^{(i)}} \left[ \Delta_{tu}^{(i)} \mathbb{E}_{\mathbf{x}^{(i)} | \mathbf{z}^{(i)}} \left[ \tilde{\mathbf{x}}_u^{(i)} \exp \left( -[\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) \right] \right],
\end{aligned}$$

where (a) follows by linearity of expectation and by plugging in  $\theta_t^{*(i)} = \theta_t^*(\mathbf{z}^{(i)})$ . Now to complete the proof, we show that for any  $i \in [n]$ ,  $u \in [p+1]$  and  $\mathbf{z}^{(i)} \in \mathcal{Z}^{pz}$ , we have

$$\mathbb{E}_{\mathbf{x}^{(i)} | \mathbf{z}^{(i)}} \left[ \tilde{\mathbf{x}}_u^{(i)} \exp \left( -[\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) \right] = 0.$$

Fix any  $i \in [n]$ ,  $u \in [p+1]$  and  $\mathbf{z}^{(i)} \in \mathcal{Z}^{pz}$ . We have

$$\begin{aligned}
&\mathbb{E}_{\mathbf{x}^{(i)} | \mathbf{z}^{(i)}} \left[ \tilde{\mathbf{x}}_u^{(i)} \exp \left( -[\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) \right] \\
&= \int_{\mathcal{X}^p} \tilde{\mathbf{x}}_u^{(i)} \exp \left( -[\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) f_{\mathbf{x} | \mathbf{z}}(\mathbf{x}^{(i)} | \mathbf{z}^{(i)}) d\mathbf{x}^{(i)} \\
&= \int_{\mathcal{X}^p} \tilde{\mathbf{x}}_u^{(i)} \exp \left( -[\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt}^* \bar{x}_t^{(i)} \right) f_{\mathbf{x}_{-t} | \mathbf{z}}(\mathbf{x}_{-t}^{(i)} | \mathbf{z}^{(i)}) \times \\
&\quad f_{x_t | \mathbf{x}_{-t}, \mathbf{z}}(x_t^{(i)} | \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}; \theta_t^*(\mathbf{z}^{(i)}), \Theta_t^*) d\mathbf{x}^{(i)} \\
&\stackrel{(a)}{=} \int_{\mathcal{X}^p} \frac{\tilde{\mathbf{x}}_u^{(i)} f_{\mathbf{x}_{-t} | \mathbf{z}}(\mathbf{x}_{-t}^{(i)} | \mathbf{z}^{(i)}) d\mathbf{x}^{(i)}}{\int_{\mathcal{X}} \exp \left( [\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)} \right) dx_t^{(i)}} \\
&= \int_{\mathcal{X}^{p-1}} \left[ \int_{\mathcal{X}} \tilde{\mathbf{x}}_u^{(i)} dx_t^{(i)} \right] \frac{f_{\mathbf{x}_{-t} | \mathbf{z}}(\mathbf{x}_{-t}^{(i)} | \mathbf{z}^{(i)}) d\mathbf{x}_{-t}^{(i)}}{\int_{\mathcal{X}} \exp \left( [\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)} \right) dx_t^{(i)}}
\end{aligned}$$

$$\stackrel{(b)}{=} 0,$$

where (a) follows by plugging in  $f_{x_t|x_{-t},z}(x_t^{(i)}|\mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}; \theta_t^*(\mathbf{z}^{(i)}), \Theta_t^*)$  from (12) and (b) follows because  $\int_{\mathcal{X}} x_t^{(i)} dx_t^{(i)} = 0$  and  $\int_{\mathcal{X}} \bar{x}_t^{(i)} dx_t^{(i)} = 0$ .

### B.1.2 Proof of Lemma. B.3(b): Anti-concentration of second directional derivative

We start by claiming that the second-order directional derivative can be lower bounded by a quadratic form. We provide a proof in Appendix. B.1.2.

**Lemma B.5** (Lower bound on the second directional derivative). *For every  $t \in [p]$  with  $\underline{\Omega}_t$  defined in (56), we have*

$$\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t) \geq \frac{1}{nC_{2,\tau}} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2,$$

where  $\Delta_t^{(i)} \triangleq \begin{bmatrix} \omega_t^{(i)} \\ \Omega_{t,-t}^\top \\ \Omega_{tt} \end{bmatrix} \in \mathbb{R}^{p+1}$  and  $\tilde{\mathbf{x}}^{(i)} \triangleq \begin{bmatrix} x_t^{(i)} \\ 2\mathbf{x}_{-t}^{(i)} x_t^{(i)} \\ \bar{x}_t^{(i)} \end{bmatrix} \in \mathbb{R}^{p+1}$  for all  $i \in [n]$  with  $\bar{x}_t^{(i)} = [x_t^{(i)}]^2 - x_{\max}^2/3$  and the constant  $C_{2,\tau}$  was defined in (48).

Given this, we proceed to show the anti-concentration of the second-order directional derivative. Fix any  $t \in [p]$  and any  $\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta$ . From Lemma. B.5, we have

$$\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t) \geq \frac{1}{nC_{2,\tau}} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2. \quad (61)$$

First, using the Hoeffding's inequality, let us show concentration of  $\frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2$  around its mean. We observe that each term in the summation is an independent random variable and is bounded as follows

$$\begin{aligned} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 &\stackrel{(a)}{=} \left( \omega_t^{(i)} x_t^{(i)} + 2\Omega_{t,-t}^\top \mathbf{x}_{-t}^{(i)} x_t^{(i)} + \Omega_{tt} \bar{x}_t^{(i)} \right)^2 \\ &\stackrel{(b)}{\leq} \left( |\omega_t^{(i)}| + 2\|\Omega_{t,-t}\|_1 \|\mathbf{x}^{(i)}\|_\infty \right)^2 x_{\max}^2 \stackrel{(c)}{\leq} (2\alpha + 8\beta x_{\max})^2 x_{\max}^2 \stackrel{(48)}{=} 4C_{1,\tau}^2 x_{\max}^2, \end{aligned}$$

where (a) follows by plugging in  $\Delta_t^{(i)}$  and  $\tilde{\mathbf{x}}^{(i)}$ , (b) follows from triangle inequality, Cauchy–Schwarz inequality and because  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , and (c) follows because  $\Omega \in 2\Lambda_\Theta$ ,  $\omega^{(i)} \in 2\Lambda_\theta$ , and  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ . Then, from the Hoeffding's inequality, for any  $\varepsilon > 0$  we have

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 - \frac{1}{n} \sum_{i \in [n]} \mathbb{E} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right] \right| > \varepsilon \right) < 2 \exp \left( - \frac{n\varepsilon^2}{32C_{1,\tau}^4 x_{\max}^4} \right).$$

Applying the union bound over all  $t \in [p]$ , for any  $\delta \in (0, 1)$  and uniformly for all  $t \in [p]$ , we have

$$\frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \geq \frac{1}{n} \sum_{i \in [n]} \mathbb{E} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right] - \varepsilon, \quad (62)$$

with probability at least  $1 - \delta$  as long as

$$n \geq \frac{32C_{1,\tau}^4 x_{\max}^4}{\varepsilon^2} \log \left( \frac{2p}{\delta} \right).$$

Now, we lower bound  $\mathbb{E} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right]$  for every  $t \in [p]$  and every  $i \in [n]$ . Fix any  $t \in [p]$  and  $i \in [n]$ . We have

$$\begin{aligned} \mathbb{E}_{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right] &= \mathbb{E}_{\mathbf{z}^{(i)}} \left[ [\Delta_t^{(i)}]^\top \mathbb{E}_{\mathbf{x}^{(i)} | \mathbf{z}^{(i)}} \left[ \tilde{\mathbf{x}}^{(i)} \tilde{\mathbf{x}}^{(i)\top} | \mathbf{z}^{(i)} \right] \Delta_t^{(i)} \right] \\ &\stackrel{(a)}{\geq} \lambda_{\min} \mathbb{E}_{\mathbf{z}^{(i)}} \left[ \|\Delta_t^{(i)}\|_2^2 \right] \stackrel{(b)}{\geq} \lambda_{\min} \|\Omega_t\|_2^2, \end{aligned} \quad (63)$$

where (a) follows from Assumption. 2 and (b) follows from the definition of  $\Delta_t^{(i)}$ . Combining (61) to (63), for any  $\delta \in (0, 1)$  and uniformly for all  $t \in [p]$ , we have

$$\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t) \geq \frac{1}{C_{2,\tau}} \left( \lambda_{\min} \|\Omega_t\|_2^2 - \varepsilon \right),$$

with probability at least  $1 - \delta$  as long as

$$n \geq \frac{32C_{1,\tau}^4 x_{\max}^4}{\varepsilon^2} \log \left( \frac{2p}{\delta} \right).$$

Choosing  $\varepsilon = \varepsilon_2 C_{2,\tau}$  and  $\delta = \delta_2$  yields the claim.

**Proof of Lemma. B.5: Lower bound on the second directional derivative** For every  $t \in [p]$  with  $\underline{\Omega}_t$  defined in (56), we claim that the second-order directional derivative of the loss function defined in (47) is given by

$$\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t) = \frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right), \quad (64)$$

where  $\Delta_t^{(i)} \triangleq \begin{bmatrix} \omega_t^{(i)} \\ \Omega_{t,-t}^\top \\ \Omega_{tt} \end{bmatrix} \in \mathbb{R}^{p+1}$  and  $\tilde{\mathbf{x}}^{(i)} \triangleq \begin{bmatrix} x_t^{(i)} \\ 2\mathbf{x}_{-t}^{(i)} x_t^{(i)} \\ \bar{x}_t^{(i)} \end{bmatrix} \in \mathbb{R}^{p+1}$  for all  $i \in [n]$  with  $\bar{x}_t^{(i)} = [x_t^{(i)}]^2 - x_{\max}^2/3$ . We provide a proof at the end.

Given this claim, we proceed to prove the lower bound on the second directional derivative. Fix any  $t \in [p]$ . From (64), we have

$$\begin{aligned} \partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t) &= \frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \times \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \\ &\stackrel{(a)}{\geq} \frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \times \exp \left( - (|\theta_t^{(i)}| + 2\|\Theta_{t,-t}\|_1 \|\mathbf{x}^{(i)}\|_\infty) x_{\max} \right) \end{aligned}$$

$$\begin{aligned}
&\stackrel{(b)}{\geq} \frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \times \exp \left( - (\alpha + 2\beta x_{\max}) x_{\max} \right) \\
&\stackrel{(48)}{=} \frac{1}{C_{2,\tau} n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2,
\end{aligned}$$

where (a) follows from triangle inequality, Cauchy–Schwarz inequality and because  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , and (b) follows because  $\theta^{(i)} \in \Lambda_\theta$  for all  $i \in [n]$ ,  $\Theta \in \Lambda_\Theta$ , and  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ .

**Proof of (64): Expression for second directional derivative** Fix any  $t \in [p]$ . The second-order partial derivatives of  $\mathcal{L}_t$  with respect to entries of  $\underline{\Theta}_t$  defined in (13) are given by

$$\begin{aligned}
\frac{\partial^2 \mathcal{L}_t(\underline{\Theta}_t)}{\partial [\theta_t^{(i)}]^2} &= \frac{1}{n} [x_t^{(i)}]^2 \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \quad \text{for all } i \in [n], \\
\frac{\partial^2 \mathcal{L}_t(\underline{\Theta}_t)}{\partial \Theta_{tu} \Theta_{tv}} &= \begin{cases} \frac{4}{n} \sum_{i \in [n]} [x_t^{(i)}]^2 x_u^{(i)} x_v^{(i)} \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) & \text{for all } u, v \in [p] \setminus \{t\}. \\ \frac{2}{n} \sum_{i \in [n]} \bar{x}_t^{(i)} x_t^{(i)} x_u^{(i)} \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) & \text{for all } u \in [p] \setminus \{t\} \text{ and } v=t. \\ \frac{2}{n} \sum_{i \in [n]} \bar{x}_t^{(i)} x_t^{(i)} x_v^{(i)} \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) & \text{for all } v \in [p] \setminus \{t\} \text{ and } u=t. \\ \frac{1}{n} \sum_{i \in [n]} [\bar{x}_t^{(i)}]^2 \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) & \text{for } v=t \text{ and } u=t. \end{cases} \\
\frac{\partial^2 \mathcal{L}_t(\underline{\Theta}_t)}{\partial \Theta_{tu} \theta_t^{(i)}} &= \frac{\partial^2 \mathcal{L}_t(\underline{\Theta}_t)}{\partial \theta_t^{(i)} \Theta_{tu}} = \begin{cases} \frac{2}{n} [x_t^{(i)}]^2 x_u^{(i)} \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) & \text{for all } i \in [n], u \in [p] \setminus \{t\}. \\ \frac{1}{n} x_t^{(i)} \bar{x}_t^{(i)} \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) & \text{for all } i \in [n], u=t. \end{cases}
\end{aligned}$$

Now, we can write the second-order directional derivative of  $\mathcal{L}_t$  as

$$\begin{aligned}
\partial_{\underline{\Omega}_t}^2 \mathcal{L}_t(\underline{\Theta}_t) &\triangleq \lim_{h \rightarrow 0} \frac{\partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t + h \underline{\Omega}_t) - \partial_{\underline{\Omega}_t} \mathcal{L}_t(\underline{\Theta}_t)}{h} \\
&= \sum_{i \in [n]} [\omega_t^{(i)}]^2 \frac{\partial^2 \mathcal{L}_t(\underline{\Theta}_t)}{\partial [\theta_t^{(i)}]^2} + \sum_{u \in [p]} \sum_{v \in [p]} \Omega_{tu} \Omega_{tv} \frac{\partial^2 \mathcal{L}_t(\underline{\Theta}_t)}{\partial \Theta_{tu} \Theta_{tv}} + 2 \sum_{i \in [n]} \sum_{u \in [p]} \omega_t^{(i)} \Omega_{tu} \frac{\partial^2 \mathcal{L}_t(\underline{\Theta}_t)}{\partial \Theta_{tu} \theta_t^{(i)}} \\
&= \frac{1}{n} \sum_{i \in [n]} \left( [\omega_t^{(i)} x_t^{(i)}]^2 + 4 \sum_{u \in [p]} \Omega_{tu} x_t^{(i)} x_u^{(i)} \sum_{v \in [p]} \Omega_{tv} x_t^{(i)} x_v^{(i)} + 4 \Omega_{tt} \bar{x}_t^{(i)} \sum_{u \in [p]} \Omega_{tu} x_t^{(i)} x_u^{(i)} + [\Omega_{tt} \bar{x}_t^{(i)}]^2 \right. \\
&\quad \left. + 4 \omega_t^{(i)} x_t^{(i)} \sum_{u \in [p]} \Omega_{tu} x_t^{(i)} x_u^{(i)} + 2 \omega_t^{(i)} x_t^{(i)} [\Omega_{tt} \bar{x}_t^{(i)}] \right) \times \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right) \\
&= \frac{1}{n} \sum_{i \in [n]} \left( \omega_t^{(i)} x_t^{(i)} + 2 \Omega_{t,-t}^\top \mathbf{x}_{-t}^{(i)} + \Omega_{tt} \bar{x}_t^{(i)} \right)^2 \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right)
\end{aligned}$$

$$\stackrel{(a)}{=} \frac{1}{n} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \exp \left( - [\theta_t^{(i)} + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \Theta_{tt} \bar{x}_t^{(i)} \right),$$

where (a) follows from the definitions of  $\Delta_t^{(i)}$  and  $\tilde{\mathbf{x}}^{(i)}$ .

## B.2 Example for Assumption. 2

As seen in (63), Assumption. 2 is used to lower bound  $\mathbb{E}_{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right]$  by  $\|\Omega_t\|_2^2$ . In this section, we show that  $\mathbb{E}_{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right]$  can be lower bounded by  $\|\Omega_t\|_2^2$  without requiring Assumption. 2 if  $\Theta_{tt}^* = 0$  for all  $t \in [p]$  and the row-wise  $\ell_1$  sparsity of  $\Theta$  in Assumption. 1 is assumed to be induced by row-wise  $\ell_0$  sparsity, i.e.,  $\|\Theta_t\|_0 \leq \beta/\alpha$  for all  $t \in [p]$ . To that end, first we claim that the conditional variance of  $x_t^{(i)}$  conditioned on  $\mathbf{x}_{-t} = \mathbf{x}_{-t}^{(i)}$  and  $\mathbf{z} = \mathbf{z}^{(i)}$  is lower bounded by a constant for every  $t \in [p]$  and  $i \in [n]$ . We provide a proof in Appendix. B.2.1.

**Lemma B.6** (Lower bound on the conditional variance). *We have*

$$\text{Var}(x_t^{(i)} | \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}) \geq \frac{2x_{\max}^2}{\pi e C_{2,\tau}^4} \quad \text{for all } t \in [p] \text{ and } i \in [n],$$

where the constant  $C_{2,\tau}$  was defined in (48).

Given this lemma, we proceed. We have

$$\mathbb{E} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right] \stackrel{(a)}{\geq} \text{Var} \left[ [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right] \stackrel{(b)}{=} \text{Var} \left[ \omega_t^{(i)} x_t^{(i)} + 2\Omega_t^\top \mathbf{x}^{(i)} x_t^{(i)} \right], \quad (65)$$

where (a) follows from the fact that for any random variable  $a$ ,  $\mathbb{E}[a^2] \geq \text{Var}[a]$  and (b) follows because we let  $\Omega_{tt} = 0$  since  $\Theta_{tt}^* = 0$ . We define the following set to lower bound  $\text{Var} \left[ \omega_t^{(i)} x_t^{(i)} + 2\Omega_t^\top \mathbf{x}^{(i)} x_t^{(i)} \right]$ :

$$\mathcal{E}(\Theta^*) \triangleq \{ (t, u) \in [p]^2 : t < u, \Theta_{tu}^* \neq 0 \}, \quad (66)$$

and consider the graph  $\mathcal{G}(\Theta^*) = ([p], \mathcal{E}(\Theta^*))$  with  $[p]$  as nodes and  $\mathcal{E}(\Theta^*)$  as edges such that  $f_{\mathbf{x}|\mathbf{z}}(\mathbf{x}|\mathbf{z}; \theta^*(\mathbf{z}), \Theta^*)$  is Markov with respect to  $\mathcal{G}(\Theta^*)$ . We claim that there exists a non-empty set  $\mathcal{R}_t \subset [p] \setminus \{t\}$  such that

- (i)  $\mathcal{R}_t$  is an independent set of  $\mathcal{G}(\Theta^*)$ , i.e., there are no edges between any pair of nodes in  $\mathcal{R}_t$ , and
- (ii) the row vector  $\Omega_t$  satisfies  $\sum_{u \in \mathcal{R}_t} |\Omega_{tu}|^2 \geq \frac{1}{\beta/\alpha + 1} \|\Omega_t\|_2^2$ .

Taking this claim as given at the moment, we continue our proof. Denoting  $\mathcal{R}_t^c \triangleq [p] \setminus \mathcal{R}_t$ , and using the law of total variance, the variance term in (65) can be lower bounded as

$$\begin{aligned} \text{Var} \left[ \omega_t^{(i)} x_t^{(i)} + 2\Omega_t^\top \mathbf{x}^{(i)} x_t^{(i)} \right] &\geq \mathbb{E} \left[ \text{Var} \left[ \omega_t^{(i)} x_t^{(i)} + 2\Omega_t^\top \mathbf{x}^{(i)} x_t^{(i)} \mid \mathbf{x}_{\mathcal{R}_t^c}^{(i)}, \mathbf{z}^{(i)} \right] \right] \\ &\stackrel{(a)}{=} 4\mathbb{E} \left[ (x_t^{(i)})^2 \text{Var} \left( \sum_{u \in \mathcal{R}_t} \Omega_{tu} x_u^{(i)} \mid \mathbf{x}_{\mathcal{R}_t^c}^{(i)}, \mathbf{z}^{(i)} \right) \right] \\ &\stackrel{(b)}{=} 4\mathbb{E} \left[ (x_t^{(i)})^2 \sum_{u \in \mathcal{R}_t} \Omega_{tu}^2 \text{Var} \left( x_u^{(i)} \mid \mathbf{x}_{\mathcal{R}_t^c}^{(i)}, \mathbf{z}^{(i)} \right) \right] \end{aligned}$$

$$\begin{aligned}
&\stackrel{(c)}{=} 4\mathbb{E}\left[(x_t^{(i)})^2 \sum_{u \in \mathcal{R}_t} \Omega_{tu}^2 \text{Var}\left(x_u^{(i)} \mid \mathbf{x}_{-u}^{(i)}, \mathbf{z}^{(i)}\right)\right] \\
&\geq \frac{8x_{\max}^2}{\pi e C_{2,\tau}^4} \sum_{u \in \mathcal{R}_t} \Omega_{tu}^2 \mathbb{E}\left[(x_t^{(i)})^2\right] \\
&\stackrel{(e)}{\geq} \frac{8x_{\max}^2}{\pi e C_{2,\tau}^4} \sum_{u \in \mathcal{R}_t} \Omega_{tu}^2 \text{Var}\left(x_t^{(i)} \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right) \\
&\stackrel{(f)}{\geq} \frac{16x_{\max}^4}{\pi^2 e^2 C_{2,\tau}^8} \sum_{u \in \mathcal{R}_t} \Omega_{tu}^2 \stackrel{(ii)}{\geq} \frac{16x_{\max}^4 \|\Omega_t\|_2^2}{\pi^2 e^2 (\beta/\alpha + 1) C_{2,\tau}^8}, \tag{67}
\end{aligned}$$

where (a) follows because  $(x_u^{(i)})_{u \in \mathcal{R}_t^c}$  are deterministic when conditioned on themselves, and  $t \in \mathcal{R}_t^c$ , (b) follows because  $(x_u^{(i)})_{u \in \mathcal{R}_t}$  are conditionally independent given  $\mathbf{x}_{\mathcal{R}_t^c}^{(i)}$  and  $\mathbf{z}^{(i)}$  which is a direct consequence of (i), (c) follows because of the local Markov property (as the conditioning set includes all the neighbors in  $\mathcal{G}(\Theta^*)$  of each node in  $\mathcal{R}_t$ ), (d) and (f) follow from Lemma B.6, and (e) follows because  $\mathbb{E}\left[(x_t^{(i)})^2\right] = \mathbb{E}\left[\mathbb{E}\left[(x_t^{(i)})^2 \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right]\right] \geq \text{Var}\left(x_t^{(i)} \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right)$ .

Combining (65) and (67), we have

$$\mathbb{E}_{\mathbf{x}^{(i)}, \mathbf{z}^{(i)}} \left[ \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right)^2 \right] \geq \frac{16x_{\max}^4}{\pi^2 e^2 (\beta/\alpha + 1) C_{2,\tau}^8} \cdot \|\Omega_t\|_2^2.$$

It remains to construct the set  $\mathcal{R}_t$  that is an independent set of  $\mathcal{G}(\Theta^*)$  and satisfies (ii).

**Construction of the set  $\mathcal{R}_t$**  For every  $u \in [p]$ , let  $\mathcal{N}(u)$  denote the set of neighbors of  $u$  in  $\mathcal{G}(\Theta^*)$ , i.e.,  $\mathcal{N}(u) \triangleq \{v \in [p] : (u, v) \in \mathcal{E}(\Theta^*)\} \cup \{v \in [p] : (v, u) \in \mathcal{E}(\Theta^*)\}$ . We start by selecting  $r_1 \in [p] \setminus \{t\}$  such that

$$|\Omega_{tr_1}| \geq |\Omega_{tu}| \quad \text{for all } u \in [p] \setminus \{t, r_1\}.$$

Next, we identify  $r_2 \in [p] \setminus \{t, r_1, \mathcal{N}(r_1)\}$  such that

$$|\Omega_{tr_2}| \geq |\Omega_{tu}| \quad \text{for all } u \in [p] \setminus \{t, r_1, \mathcal{N}(r_1), r_2\}.$$

We continue identifying  $r_3, \dots, r_s$  in such a manner till no more nodes are left, where  $s$  denotes the total number of nodes selected. Now we define  $\mathcal{R}_t \triangleq \{r_1, \dots, r_s\}$ . For any  $u \in [p]$ , we have  $|\mathcal{N}(u)| \leq \|\Theta_u^*\|_0 \leq \beta/\alpha$  from (66) and Assumption 1. Using this, we see that  $\mathcal{R}_t$  is an independent set of  $\mathcal{G}(\Theta^*)$  as claimed in (i) such that it satisfies (ii) by construction.

### B.2.1 Proof of Lemma B.6: Lower bound on the conditional variance

For any random variable  $x$ , let  $h(x)$  denote the differential entropy of  $x$ . Fix any  $t \in [p]$  and  $i \in [n]$ . Then, from Shannon's entropy inequality ( $2h(\cdot) \leq \log \sqrt{2\pi e \text{Var}(\cdot)}$ ), we have

$$2\pi e \text{Var}\left(x_t^{(i)} \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right) \stackrel{(a)}{\geq} \exp\left(2h\left(x_t^{(i)} \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right)\right). \tag{68}$$

Therefore, to bound the variance, it suffices to bound the differential entropy. We have

$$\begin{aligned}
& -h(x_t^{(i)} | \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}) \\
&= \int_{\mathcal{X}^p \times \mathcal{Z}^{pz}} f_{\mathbf{x}, \mathbf{z}}(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) \log \left( f_{x_t | \mathbf{x}_{-t}, \mathbf{z}}(x_t^{(i)} | \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}; \theta_t^*(\mathbf{z}^{(i)}), \Theta_t^*) \right) d\mathbf{x}^{(i)} d\mathbf{z}^{(i)} \\
&= \int_{\mathcal{X}^p \times \mathcal{Z}^{pz}} f_{\mathbf{x}, \mathbf{z}}(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) \log \left( \frac{\exp([\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)})}{\int_{\mathcal{X}} \exp([\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)}) dx_t^{(i)}} \right) d\mathbf{x}^{(i)} d\mathbf{z}^{(i)} \\
&\stackrel{(a)}{\geq} \int_{\mathcal{X}^p \times \mathcal{Z}^{pz}} f_{\mathbf{x}, \mathbf{z}}(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) \log \left( \frac{\exp((|\theta_t^*(\mathbf{z}^{(i)})| + 2\|\Theta_t^*\|_1 \|\mathbf{x}^{(i)}\|_\infty) x_{\max})}{\int_{\mathcal{X}} \exp(-(|\theta_t^*(\mathbf{z}^{(i)})| + 2\|\Theta_t^*\|_1 \|\mathbf{x}^{(i)}\|_\infty) x_{\max}) dx_t^{(i)}} \right) d\mathbf{x}^{(i)} d\mathbf{z}^{(i)} \\
&\stackrel{(b)}{\geq} \int_{\mathcal{X}^p \times \mathcal{Z}^{pz}} f_{\mathbf{x}, \mathbf{z}}(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) \log \left( \frac{\exp((\alpha + 2\beta x_{\max}) x_{\max})}{\int_{\mathcal{X}} \exp(-(\alpha + 2\beta x_{\max}) x_{\max}) dx_t^{(i)}} \right) d\mathbf{x}^{(i)} d\mathbf{z}^{(i)} \\
&\stackrel{(c)}{=} \int_{\mathcal{X}^p \times \mathcal{Z}^{pz}} f_{\mathbf{x}, \mathbf{z}}(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) \log \left( \frac{C_{3,\tau}^2}{2x_{\max}} \right) d\mathbf{x}^{(i)} d\mathbf{z}^{(i)} = \log \left( \frac{C_{3,\tau}^2}{2x_{\max}} \right), \tag{69}
\end{aligned}$$

where (a) follows from triangle inequality and Cauchy–Schwarz inequality and because  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , (b) follows because  $\theta^*(\mathbf{z}^{(i)}) \in \Lambda_\theta$  for all  $i \in [n]$ ,  $\Theta^* \in \Lambda_\Theta$ ,  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , and (c) follows because  $\int_{\mathcal{X}} dx_t^{(i)} = 2x_{\max}$ . Combining (68) and (69) completes the proof.

### B.3 Proof of Lemma. B.2: Lipschitzness of the loss function

Consider any direction  $\underline{\Omega} = \tilde{\Theta} - \underline{\Theta}$ . Now, define the function  $q : [0, 1] \rightarrow \mathbb{R}$  as follows

$$q(a) = \mathcal{L}(\underline{\Theta} + a(\tilde{\Theta} - \underline{\Theta})). \tag{70}$$

Then, the desired inequality in (49) is equivalent to

$$|q(1) - q(0)| \leq 2x_{\max}^2 C_{2,\tau} \left( \sum_{t \in [p]} \|\Omega_t\|_1 + \frac{1}{n} \sum_{i \in [n]} \|\omega^{(i)}\|_1 \right).$$

From the mean value theorem, there exists  $a' \in (0, 1)$  such that

$$|q(1) - q(0)| = \left| \frac{dq(a')}{da} \right| \stackrel{(70)}{=} \left| \frac{d\mathcal{L}(\underline{\Theta} + a'(\tilde{\Theta} - \underline{\Theta}))}{da} \right| \stackrel{(57)}{=} \left| \partial_{\underline{\Omega}} \mathcal{L}(\underline{\Theta}) \Big|_{\underline{\Theta} = \underline{\Theta} + a'(\tilde{\Theta} - \underline{\Theta})} \right|. \tag{71}$$

Using (60) in (71), we can write

$$\begin{aligned}
& |q(1) - q(0)| \\
&= \frac{1}{n} \left| \sum_{t \in [p]} \sum_{i \in [n]} \left( [\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)} \right) \times \exp \left( - \left[ (\theta_t^{(i)} + a'(\tilde{\theta}_t^{(i)} - \theta_t^{(i)})) + \right. \right. \right. \\
&\quad \left. \left. \left. 2(\Theta_{t,-t} + a'(\tilde{\Theta}_{t,-t} - \Theta_{t,-t}))^\top \mathbf{x}_{-t}^{(i)} \right] x_t^{(i)} - (\Theta_{tt} + a'(\tilde{\Theta}_{tt} - \Theta_{tt})) \bar{x}_t^{(i)} \right) \right|
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(a)}{\leq} \exp\left(\left([(1-a')\alpha + a'\alpha] + 2[(1-a')\beta + a'\beta]x_{\max}\right)x_{\max}\right) \frac{1}{n} \left| \sum_{t \in [p]} \sum_{i \in [n]} \left([\Delta_t^{(i)}]^\top \tilde{\mathbf{x}}^{(i)}\right) \right| \\
&\stackrel{(b)}{\leq} \frac{2x_{\max}^2 C_{2,\tau}}{n} \sum_{t \in [p]} \sum_{i \in [n]} \|\Delta_t^{(i)}\|_1 \stackrel{(c)}{=} 2x_{\max}^2 C_{2,\tau} \left( \sum_{t \in [p]} \|\Omega_t\|_1 + \frac{1}{n} \sum_{i \in [n]} \|\omega^{(i)}\|_1 \right),
\end{aligned}$$

where (a) follows from triangle inequality, Cauchy–Schwarz inequality,  $\theta^{(i)}, \tilde{\theta}^{(i)} \in \Lambda_\theta$ ,  $\Theta, \tilde{\Theta} \in \Lambda_\Theta$ , and  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , (b) follows from (48), the triangle inequality, and because  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , and (c) follows from the definition of  $\Delta_t^{(i)}$ .

## C Proof of Theorem 1 Part II: Recovering unit-level parameters

To analyze our estimate of the unit-level parameters, we use the estimate  $\hat{\Theta}$  of the population-level parameter  $\Theta^*$  along with the associated guarantee provided in Theorem. 1 Part I. We note that the constraints on the unit-level parameters in (14) are independent across units, i.e.,  $\theta^{(i)} \in \Lambda_\theta$  independently for all  $i \in [n]$ . Therefore, we look at  $n$  independent convex optimization problems by decomposing the loss function  $\mathcal{L}$  in (13) and the estimate  $\hat{\Theta}$  in (14) as follows: For  $i \in [n]$ , we define

$$\begin{aligned}
\mathcal{L}^{(i)}(\theta^{(i)}) &\triangleq \sum_{t \in [p]} \exp\left(-[\theta_t^{(i)} + 2\hat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \hat{\Theta}_{tt}\bar{x}_t^{(i)}\right) \\
\text{and } \hat{\theta}^{(i)} &\triangleq \arg \min_{\theta^{(i)} \in \Lambda_\theta} \mathcal{L}^{(i)}(\theta^{(i)}).
\end{aligned} \tag{72}$$

Now, fix any  $i \in [n]$ . From (72), we have  $\mathcal{L}^{(i)}(\hat{\theta}^{(i)}) \leq \mathcal{L}^{(i)}(\theta^{*(i)})$ . Using contraposition, to prove this part, it is sufficient to show that all points  $\theta^{(i)} \in \Lambda_\theta$  that satisfy  $\|\theta^{(i)} - \theta^{*(i)}\|_2 \geq R(\varepsilon, \delta)$  also uniformly satisfy

$$\mathcal{L}^{(i)}(\theta^{(i)}) \geq \mathcal{L}^{(i)}(\theta^{*(i)}) + R^2(\varepsilon, \delta) \text{ when } n \geq \frac{ce^{c'\beta}p^4}{\varepsilon^4} \left( p \log \frac{p^2}{\delta\varepsilon^2} + \tilde{\mathcal{M}}_{\theta,n}(\varepsilon, \delta) \right), \tag{73}$$

with probability at least  $1 - \delta$  where  $R(\varepsilon, \delta)$  was defined in (17) and  $\tilde{\mathcal{M}}_{\theta,n}(\varepsilon, \delta)$  was defined in (18). Then, the guarantee in Theorem. 1 follows by applying a union bound over all  $i \in [n]$ .

To that end, the lemma below, proven in Appendix. C.1, shows that for any fixed  $\theta^{(i)} \in \Lambda_\theta$ , if  $\theta^{(i)}$  is far from  $\theta^{*(i)}$ , then with high probability  $\mathcal{L}^{(i)}(\theta^{(i)})$  is significantly larger than  $\mathcal{L}^{(i)}(\theta^{*(i)})$ .

**Lemma C.1** (Gap between the loss function for a fixed parameter). *Fix any  $\varepsilon > 0$ ,  $\delta \in (0, 1)$ , and  $i \in [n]$ . Then, for any  $\theta^{(i)} \in \Lambda_\theta$  such that  $\|\theta^{(i)} - \theta^{*(i)}\|_2 \geq \varepsilon\gamma$  (see (17)), we have*

$$\mathcal{L}^{(i)}(\theta^{(i)}) \geq \mathcal{L}^{(i)}(\theta^{*(i)}) + \frac{2^{2.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\theta^{(i)} - \theta^{*(i)}\|_2^2 \text{ for } n \geq \frac{ce^{c'\beta}p^4}{\varepsilon^4} \left( p \log \frac{p^2}{\delta\varepsilon^2} + \mathcal{M}_{\theta,n} \left( \frac{\varepsilon^2}{p} \right) \right),$$

with probability at least  $1 - \delta - c\beta^2 \log p \cdot \exp(-e^{-c'\beta} \|\theta^{(i)} - \theta^{*(i)}\|_2^2)$  where  $C_{2,\tau}$  was defined in (48).

**Note.** When we invoke Lemma. C.1, we ensure that  $c\beta^2 \log p \cdot \exp(-e^{-c'\beta} \|\theta^{(i)} - \theta^{*(i)}\|_2^2)$  is of the same order as  $\delta$ .

Next, we show that the loss function  $\mathcal{L}^{(i)}$  is Lipschitz (see Appendix. C.2 for the proof).

**Lemma C.2** (Lipschitzness of the loss function). *Consider any  $i \in [n]$ . Then, the loss function  $\mathcal{L}^{(i)}$  is Lipschitz with respect to the  $\ell_1$  norm  $\|\cdot\|_1$  and with Lipschitz constant  $x_{\max}C_{2,\tau}$ , i.e.,*

$$|\mathcal{L}^{(i)}(\tilde{\theta}^{(i)}) - \mathcal{L}^{(i)}(\theta^{(i)})| \leq x_{\max}C_{2,\tau}\|\tilde{\theta}^{(i)} - \theta^{(i)}\|_1 \quad \text{for all } \theta^{(i)}, \tilde{\theta}^{(i)} \in \Lambda_\theta, \quad (74)$$

where the constant  $C_{2,\tau}$  was defined in (48).

Given these lemmas, we now proceed with the proof.

**Proof strategy** We want to show that all points  $\theta^{(i)} \in \Lambda_\theta$ , that satisfy  $\|\theta^{(i)} - \theta^{*(i)}\|_2 \geq R(\varepsilon, \delta)$ , uniformly satisfy (73) with probability at least  $1 - \delta$ . To do so, we consider the set of points  $\Lambda_\theta^r \subset \Lambda_\theta$  whose distance from  $\theta^{*(i)}$  is at least  $r > 0$  in  $\ell_2$  norm. Then, using an appropriate covering set of  $\Lambda_\theta^r$  and the Lipschitzness of  $\mathcal{L}^{(i)}$ , we show that the value of  $\mathcal{L}^{(i)}$  at all points in  $\Lambda_\theta^r$  is uniformly  $\Omega(r^2)$  larger than the value of  $\mathcal{L}^{(i)}$  at  $\theta^{*(i)}$  with high probability. Finally, we choose  $r$  small enough to make the failure probability smaller than  $\delta$ .

**Arguments for points in the covering set** Consider any  $r \geq \varepsilon\gamma$  (where  $\gamma$  is defined in (17)) and the set of elements  $\Lambda_\theta^r \triangleq \{\theta^{(i)} \in \Lambda_\theta : \|\theta^{*(i)} - \theta^{(i)}\|_2 \geq r\}$ . Let  $\mathcal{U}(\Lambda_\theta^r, \varepsilon')$  be the  $\varepsilon'$ -cover of the smallest size for the set  $\Lambda_\theta^r$  with respect to  $\|\cdot\|_1$  (see Definition. 2) and let  $\mathcal{C}(\Lambda_\theta^r, \varepsilon')$  be the  $\varepsilon'$ -covering number where

$$\varepsilon' \triangleq \frac{2\sqrt{2}\beta x_{\max}^3 r^2}{\pi e C_{2,\tau}^6}. \quad (75)$$

Now, we argue by a union bound that the value of  $\mathcal{L}^{(i)}$  at all points in  $\mathcal{U}(\Lambda_\theta^r, \varepsilon')$  is uniformly  $\Omega(r^2)$  larger than  $\mathcal{L}^{(i)}(\theta^{*(i)})$  with high probability. For any  $\theta^{(i)} \in \mathcal{U}(\Lambda_\theta^r, \varepsilon')$ , we have

$$\|\theta^{*(i)} - \theta^{(i)}\|_2 \stackrel{(a)}{\geq} r, \quad (76)$$

where (a) follows because  $\mathcal{U}(\Lambda_\theta^r, \varepsilon') \subseteq \Lambda_\theta^r$ . Now, applying Lemma. C.1 with  $\varepsilon \leftarrow \varepsilon$  and  $\delta \leftarrow \delta/2\mathcal{C}(\Lambda_\theta^r, \varepsilon')$ , we have

$$\mathcal{L}^{(i)}(\theta^{(i)}) \geq \mathcal{L}^{(i)}(\theta^{*(i)}) + \frac{4\sqrt{2}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\theta^{*(i)} - \theta^{(i)}\|_2^2 \stackrel{(76)}{\geq} \mathcal{L}^{(i)}(\theta^{*(i)}) + \frac{4\sqrt{2}\beta x_{\max}^4 r^2}{\pi e C_{2,\tau}^5},$$

with probability at least  $1 - \delta/2\mathcal{C}(\Lambda_\theta^r, \varepsilon') - c\beta^2 \log p \cdot \exp(-e^{-c'\beta} \|\theta^{(i)} - \theta^{*(i)}\|_2^2)$  whenever

$$n \geq \frac{ce^{c'\beta} p^4}{\varepsilon^4} \left( p \log \frac{\mathcal{C}(\Lambda_\theta^r, \varepsilon') \cdot p^2}{\delta \varepsilon^2} + \mathcal{M}_{\theta,n} \left( \frac{\varepsilon^2}{p} \right) \right). \quad (77)$$

By applying the union bound over  $\mathcal{U}(\Lambda_\theta^r, \varepsilon')$ , as long as  $n$  satisfies (77), we have

$$\mathcal{L}^{(i)}(\theta^{(i)}) \geq \mathcal{L}^{(i)}(\theta^{*(i)}) + \frac{4\sqrt{2}\beta x_{\max}^4 r^2}{\pi e C_{2,\tau}^5} \quad \text{uniformly for every } \theta^{(i)} \in \mathcal{U}(\Lambda_\theta^r, \varepsilon'), \quad (78)$$

with probability at least  $1 - \delta/2 - c\beta^2 \mathcal{C}(\Lambda_\theta^r, \varepsilon') \log p \cdot \exp(-e^{-c'\beta} \|\theta^{(i)} - \theta^{*(i)}\|_2^2)$  which can lower bounded by  $1 - \delta/2 - c\beta^2 \mathcal{C}(\Lambda_\theta^r, \varepsilon') \log p \cdot \exp(-e^{-c'\beta} r^2)$  using (76).

**Arguments for points outside the covering set** Next, we establish the claim (73) for an arbitrary  $\tilde{\theta}^{(i)} \in \Lambda_\theta^r$  conditional on the event that (78) holds. Given a fixed  $\tilde{\theta}^{(i)} \in \Lambda_\theta^r$ , let  $\theta^{(i)}$  be (one of) the point(s) in the  $\mathcal{U}(\Lambda_\theta^r, \varepsilon')$  that satisfies  $\|\theta^{(i)} - \tilde{\theta}^{(i)}\|_1 \leq \varepsilon'$  (there exists such a point by Definition. 2) Then, the choices (75) and Lemma. C.2 put together imply that

$$\begin{aligned} \mathcal{L}^{(i)}(\tilde{\theta}^{(i)}) &\geq \mathcal{L}^{(i)}(\theta^{(i)}) - x_{\max} C_{2,\tau} \|\theta^{(i)} - \tilde{\theta}^{(i)}\|_1 \geq \mathcal{L}^{(i)}(\theta^{(i)}) - x_{\max} C_{2,\tau} \varepsilon' \\ &\stackrel{(75)}{\geq} \mathcal{L}^{(i)}(\theta^{(i)}) - \frac{2\sqrt{2}\beta x_{\max}^4 r^2}{\pi e C_{2,\tau}^5} \\ &\stackrel{(78)}{\geq} \mathcal{L}^{(i)}(\theta^{*(i)}) + \frac{2\sqrt{2}\beta x_{\max}^4 r^2}{\pi e C_{2,\tau}^5}, \end{aligned}$$

It remains to bound sample size  $n$  and the failure probability  $\delta$ .

**Bounding  $n$**  Using  $\Lambda_\theta^r \subseteq \Lambda_\theta$ , we find that

$$\mathcal{C}(\Lambda_\theta^r, \varepsilon') \stackrel{(a)}{\leq} \mathcal{C}(\Lambda_\theta, \varepsilon'). \quad (79)$$

Putting together (75) and (79), the lower bound (77) can be replaced by

$$n \geq \frac{ce^{c'\beta} p^4}{\varepsilon^4} \left( p \log \frac{p^2}{\delta \varepsilon^2} + p \mathcal{M}_\theta(r^2) + \mathcal{M}_{\theta,n} \left( \frac{\varepsilon^2}{p} \right) \right).$$

**Bounding  $\delta$**  To bound the failure probability by  $\delta$ , it is sufficient to chose  $r$  such that

$$\delta \geq \delta/2 + c\beta^2 \mathcal{C}(\Lambda_\theta^r, \varepsilon') \log p \cdot \exp(-e^{-c'\beta} r^2). \quad (80)$$

From (79) and (80), it is sufficient to chose  $r$  such that

$$\delta \geq \delta/2 + c\beta^2 \mathcal{C}(\Lambda_\theta, \varepsilon') \log p \cdot \exp(-e^{-c'\beta} r^2). \quad (81)$$

Re-arranging and taking logarithm on both sides of (81) and using (75), we have

$$\log \delta \geq c \left[ \log(\beta^2 \log p) + \mathcal{M}_\theta \left( \frac{r^2}{ce^{c'\beta}} \right) - e^{-c'\beta} r^2 \right]. \quad (82)$$

Finally, (82) holds whenever

$$r \geq ce^{c'\beta} \sqrt{\log \frac{\beta^2 \log p}{\delta} + \mathcal{M}_\theta(ce^{-c'\beta})}.$$

Recalling that the choice of  $r$  was such that  $r \geq \varepsilon\gamma$  completes the proof.

### C.1 Proof of Lemma. C.1: Gap between the loss function for a fixed parameter

Fix any  $\varepsilon > 0$ , any  $\delta \in (0, 1)$ , and any  $i \in [n]$ . Consider any direction  $\omega^{(i)} \in \mathbb{R}^p$  along the parameter  $\theta^{(i)}$ , i.e.,

$$\omega^{(i)} = \theta^{(i)} - \theta^{*(i)}. \quad (83)$$

We denote the first-order and the second-order directional derivatives of the loss function  $\mathcal{L}^{(i)}$  in (72) along the direction  $\omega^{(i)}$  evaluated at  $\theta^{(i)}$  by  $\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{(i)}))$  and  $\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)})$ , respectively. Below, we state a lemma (with proof divided across Appendix. C.1.1 and Appendix. C.1.2) that provides us a control on  $\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))$  and  $\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)})$ . The assumptions of Lemma. C.1 remain in force.

**Lemma C.3** (Control on first and second directional derivatives). *For any fixed  $\varepsilon_1, \varepsilon_2 > 0$ ,  $\delta_1 \in (0, 1)$ ,  $i \in [n]$ ,  $\theta^{(i)} \in \Lambda_\theta$  with  $\omega^{(i)}$  defined in (83), we have the following:*

(a) Concentration of first directional derivative: *We have*

$$|\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))| \leq \varepsilon_1 \|\omega^{(i)}\|_1 + \varepsilon_2 \|\omega^{(i)}\|_2^2 \quad \text{for } n \geq \frac{ce^{c'\beta} p^4 (p \log \frac{p^2}{\delta_1 \varepsilon_1^2} + \mathcal{M}_{\theta, n}(\frac{\varepsilon_1}{p}))}{\varepsilon_1^4},$$

$$\text{with probability at least } 1 - \delta_1 - O\left(\beta^2 \log p \exp\left(\frac{-\varepsilon_2^2 \|\omega^{(i)}\|_2^2}{e^{c'\beta}}\right)\right).$$

(b) Anti-concentration of second directional derivative: *We have*

$$\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)}) \geq \frac{32\sqrt{2}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2,$$

$$\text{with probability at least } 1 - O\left(\beta^2 \log p \exp\left(\frac{-\|\omega^{(i)}\|_2^2}{e^{c'\beta}}\right)\right) \text{ where } C_{2,\tau} \text{ was defined in (48).}$$

Given this lemma, we now proceed with the proof. Define a function  $g : [0, 1] \rightarrow \mathbb{R}^p$  as follows

$$g(a) = \theta^{*(i)} + a(\theta^{(i)} - \theta^{*(i)}).$$

Notice that  $g(0) = \theta^{*(i)}$  and  $g(1) = \theta^{(i)}$  as well as

$$\frac{d\mathcal{L}^{(i)}(g(a))}{da} = \partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\tilde{\theta}^{(i)}))|_{\tilde{\theta}^{(i)}=g(a)} \quad \text{and} \quad \frac{d^2\mathcal{L}^{(i)}(g(a))}{da^2} = \partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\tilde{\theta}^{(i)})|_{\tilde{\theta}^{(i)}=g(a)}. \quad (84)$$

By the fundamental theorem of calculus, we have

$$\frac{d\mathcal{L}^{(i)}(g(a))}{da} \geq \frac{d\mathcal{L}^{(i)}(g(a))}{da} \Big|_{a=0} + a \min_{a \in (0,1)} \frac{d^2\mathcal{L}^{(i)}(g(a))}{da^2}. \quad (85)$$

Integrating both sides of (85) with respect to  $a$ , we obtain

$$\begin{aligned} \mathcal{L}^{(i)}(g(a)) - \mathcal{L}^{(i)}(g(0)) &\geq a \frac{d\mathcal{L}^{(i)}(g(a))}{da} \Big|_{a=0} + \frac{a^2}{2} \min_{a \in (0,1)} \frac{d^2\mathcal{L}^{(i)}(g(a))}{da^2} \\ &\stackrel{(84)}{=} a \partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\tilde{\theta}^{(i)}))|_{\tilde{\theta}^{(i)}=g(0)} + \frac{a^2}{2} \min_{a \in (0,1)} \partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\tilde{\theta}^{(i)})|_{\tilde{\theta}^{(i)}=g(a)} \\ &\stackrel{(a)}{=} a \partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)})) + \frac{a^2}{2} \min_{a \in (0,1)} \partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\tilde{\theta}^{(i)})|_{\tilde{\theta}^{(i)}=g(a)} \\ &\stackrel{(b)}{\geq} -a |\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))| + \frac{a^2}{2} \min_{a \in (0,1)} \partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\tilde{\theta}^{(i)})|_{\tilde{\theta}^{(i)}=g(a)}, \end{aligned} \quad (86)$$

where (a) follows because  $g(0) = \theta^{*(i)}$ , and (b) follows by the triangle inequality. Plugging in  $a = 1$  in (86) as well as using  $g(0) = \theta^{*(i)}$  and  $g(1) = \theta^{(i)}$ , we find that

$$\mathcal{L}^{(i)}(\theta^{(i)}) - \mathcal{L}^{(i)}(\theta^{*(i)}) \geq -|\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))| + \frac{1}{2} \min_{a \in (0,1)} \partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\tilde{\theta}^{(i)})|_{\tilde{\theta}^{(i)}=g(a)}.$$

Now, we use Lemma. C.3 with  $\varepsilon_1 \leftarrow 4\sqrt{2}\beta x_{\max}^4 \varepsilon / \pi e C_{2,\tau}^5$ ,  $\varepsilon_2 \leftarrow 8\sqrt{2}\beta x_{\max}^4 / \pi e C_{2,\tau}^5$ , and  $\delta_1 \leftarrow \delta$ . Therefore, with probability at least  $1 - \delta - O\left(\beta^2 \log p \exp\left(\frac{-\|\omega^{(i)}\|_2^2}{e^{c'\beta}}\right)\right)$  and as long as  $n \geq O\left(\frac{e^{c'\beta} p^4 (p \log \frac{p^2}{\delta} + \mathcal{M}_{\theta,n}(\frac{\varepsilon^2}{p}))}{\varepsilon^4}\right)$ , we have

$$\begin{aligned} \mathcal{L}^{(i)}(\theta^{(i)}) - \mathcal{L}^{(i)}(\theta^{*(i)}) &\geq -\frac{2^{2.5}\beta x_{\max}^4 \varepsilon}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_1 - \frac{2^{3.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2 + \frac{2^{4.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2 \\ &= -\frac{2^{2.5}\beta x_{\max}^4 \varepsilon}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_1 + \frac{2^{3.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2 \\ &\stackrel{(17)}{\geq} -\frac{2^{2.5}\beta x_{\max}^4 \varepsilon \gamma}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2 + \frac{2^{3.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2 \\ &\stackrel{(a)}{\geq} -\frac{2^{2.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2 + \frac{2^{3.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2 = \frac{2^{2.5}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2, \end{aligned}$$

where (a) follows because  $\|\omega^{(i)}\|_2 = \|\theta^{(i)} - \theta^{*(i)}\|_2 \geq \varepsilon \gamma$  according to the lemma statement.

### C.1.1 Proof of Lemma. C.3(a): Concentration of first directional derivative

Fix some  $i \in [n]$  and some  $\theta^{(i)} \in \Lambda_\theta$ . Let  $\omega^{(i)}$  be as defined in (83). We claim that the first-order directional derivative of  $\mathcal{L}^{(i)}$  defined in (72) is given by

$$\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{(i)})) = -\sum_{t \in [p]} \omega_t^{(i)} x_t^{(i)} \exp\left(-[\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)}\right). \quad (87)$$

We provide a proof at the end. For now, we assume the claim and proceed.

We note that the pair  $\{\mathbf{x}, \mathbf{z}\}$  corresponds to a  $\tau$ -SGM (see Definition. G.1) with  $\tau \triangleq (\alpha, \beta, x_{\max}, \Theta)$ . To show the concentration, we use Proposition. G.1 (see Appendix. G) with  $\lambda = \frac{1}{4\sqrt{2}x_{\max}^2}$ , decompose  $\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))$  as a sum of  $L = 1024\beta^2 x_{\max}^4 \log 4p$ , and focus on these  $L$  terms. Consider the  $L$  subsets  $S_1, \dots, S_L \in [p]$  obtained from Proposition. G.1 with  $\lambda = \frac{1}{4\sqrt{2}x_{\max}^2}$  and define

$$\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \triangleq \sum_{t \in S_u} \omega_t^{(i)} x_t^{(i)} \exp\left(-[\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)}\right) \text{ for every } u \in L. \quad (88)$$

Now, we decompose  $\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))$  as a sum of the  $L$  terms defined above. More precisely, we have

$$\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)})) \stackrel{(87)}{=} -\sum_{t \in [p]} \omega_t^{(i)} x_t^{(i)} \exp\left(-[\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)}\right)$$

$$\begin{aligned}
&\stackrel{(a)}{=} -\frac{1}{L'} \sum_{u \in [L]} \sum_{t \in S_u} \omega_t^{(i)} x_t^{(i)} \exp \left( -[\theta_t^{\star(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \\
&\stackrel{(88)}{=} -\frac{1}{L'} \sum_{u \in [L]} \psi_u(\mathbf{x}^{(i)}; \omega^{(i)}), \tag{89}
\end{aligned}$$

where (a) follows because each  $t \in [p]$  appears in exactly  $L' = \lceil L/32\sqrt{2}\beta x_{\max}^2 \rceil$  of the sets  $S_1, \dots, S_L$  according to Proposition. G.1(a) (with  $\lambda = \frac{1}{4\sqrt{2}x_{\max}^2}$ ). Now, we focus on the  $L$  terms in (89).

Consider any  $u \in [L]$ . We claim that conditioned on  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , the expected value of  $\psi_u(\mathbf{x}^{(i)}; \omega^{(i)})$  can be upper bounded uniformly across all  $u \in [L]$ . We provide a proof at the end.

**Lemma C.4** (Upper bound on expected  $\psi_u$ ). *Fix  $\varepsilon > 0$ ,  $\delta \in (0, 1)$ ,  $i \in [n]$  and  $\theta^{(i)} \in \Lambda_\theta$ . Then, with  $\omega^{(i)}$  defined in (83) and given  $\mathbf{z}^{(i)}$  and  $\mathbf{x}_{-S_u}^{(i)}$  for all  $u \in [L]$ , we have*

$$\max_{u \in [L]} \mathbb{E} \left[ \psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right] \leq \varepsilon \|\omega^{(i)}\|_1 \quad \text{for } n \geq \frac{ce^{c'\beta} p^4 (p \log \frac{p^2}{\delta \varepsilon^2} + \mathcal{M}_{\theta, n}(\frac{\varepsilon^2}{p}))}{\varepsilon^4},$$

with probability at least  $1 - \delta$ .

Consider again any  $u \in [L]$ . Now, we claim that conditioned on  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ ,  $\psi_u(\mathbf{x}^{(i)}; \omega^{(i)})$  concentrates around its conditional expected value. We provide a proof at the end.

**Lemma C.5** (Concentration of  $\psi_u$ ). *Fix  $\varepsilon > 0$ ,  $i \in [n]$ ,  $u \in [L]$ , and  $\theta^{(i)} \in \Lambda_\theta$ . Then, with  $\omega^{(i)}$  defined in (83) and given  $\mathbf{z}^{(i)}$  and  $\mathbf{x}_{-S_u}^{(i)}$ , we have*

$$\left| \psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) - \mathbb{E}[\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}] \right| \leq \varepsilon,$$

with probability at least  $1 - \exp\left(\frac{-\varepsilon^2}{e^{c'\beta} \|\omega^{(i)}\|_2^2}\right)$ .

Given these lemmas, we proceed to show the concentration of  $\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{\star(i)}))$ . To that end, for any  $u \in [L]$ , given  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , let  $E_u$  denote the event that

$$\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \leq \mathbb{E}[\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}] + \frac{1}{32\sqrt{2}\beta x_{\max}^2} \varepsilon_2 \|\omega^{(i)}\|_2^2. \tag{90}$$

Since  $E_u$  is an indicator event, using the law of total expectation results in

$$\mathbb{P}(E_u) = \mathbb{E} \left[ \mathbb{P}(E_u \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}) \right] \stackrel{(a)}{\geq} 1 - \exp\left(\frac{-\varepsilon_2^2 \|\omega^{(i)}\|_2^2}{e^{c'\beta}}\right).$$

where (a) follows from Lemma. C.5 with  $\varepsilon \leftarrow \frac{\varepsilon_2 \|\omega^{(i)}\|_2^2}{32\sqrt{2}\beta x_{\max}^2}$ . Now, by applying the union bound over all  $u \in [L]$  where  $L = 1024\beta^2 x_{\max}^4 \log 4p$ , we have

$$\mathbb{P}\left(\bigcap_{u \in [L]} E_u\right) \geq 1 - O\left(\beta^2 \log p \exp\left(\frac{-\varepsilon_2^2 \|\omega^{(i)}\|_2^2}{e^{c'\beta}}\right)\right).$$

Now, assume the event  $\cap_{u \in L} E_u$  holds. Whenever this holds, we also have

$$\begin{aligned} |\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))| &\stackrel{(89)}{\leq} \frac{1}{L'} \sum_{u \in [L]} |\psi_u(\mathbf{x}^{(i)}; \omega^{(i)})| \\ &\stackrel{(90)}{\leq} \frac{1}{L'} \sum_{u \in [L]} \left| \mathbb{E}[\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) | \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}] + \frac{1}{32\sqrt{2}\beta x_{\max}^2} \varepsilon_2 \|\omega^{(i)}\|_2^2 \right|, \end{aligned} \quad (91)$$

where  $L' = \lceil L/32\sqrt{2}\beta x_{\max}^2 \rceil$ . Further, using Lemma. C.4 in (91) with  $\varepsilon \leftarrow \frac{\varepsilon_1}{32\sqrt{2}\beta x_{\max}^2}$  and  $\delta \leftarrow \delta_1$ , whenever

$$n \geq \frac{ce^{c'\beta} \cdot p^4 \left( p \log \frac{p^2}{\delta_1 \varepsilon_1^2} + \mathcal{M}_{\theta, n} \left( \frac{\varepsilon_1^2}{p} \right) \right)}{\varepsilon_1^4},$$

with probability at least  $1 - \delta_1$ , we have,

$$\begin{aligned} |\partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{*(i)}))| &\leq \frac{1}{L'} \sum_{u \in [L]} \left( \frac{1}{32\sqrt{2}\beta x_{\max}^2} \varepsilon_1 \|\omega^{(i)}\|_1 + \frac{1}{32\sqrt{2}\beta x_{\max}^2} \varepsilon_2 \|\omega^{(i)}\|_2^2 \right) \\ &= \frac{L}{32\sqrt{2}\beta x_{\max}^2 L'} \left( \varepsilon_1 \|\omega^{(i)}\|_1 + \varepsilon_2 \|\omega^{(i)}\|_2^2 \right) \stackrel{(a)}{\leq} \varepsilon_1 \|\omega^{(i)}\|_1 + \varepsilon_2 \|\omega^{(i)}\|_2^2, \end{aligned}$$

where (a) follows because  $L' = \lceil L/32\sqrt{2}\beta x_{\max}^2 \rceil$ .

**Proof of (87): Expression for first directional derivative** Fix any  $i \in [n]$ . The first-order partial derivatives of  $\mathcal{L}^{(i)}$  (defined in (72)) with respect to the entries of the parameter vector  $\theta^{(i)}$  are given by

$$\frac{\partial \mathcal{L}^{(i)}(\theta^{(i)})}{\partial \theta_t^{(i)}} = -x_t^{(i)} \exp \left( - [\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \quad \text{for all } t \in [p].$$

Now, we can write the first-order directional derivative of  $\mathcal{L}^{(i)}$  as

$$\begin{aligned} \partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{(i)})) &\triangleq \lim_{h \rightarrow 0} \frac{\mathcal{L}^{(i)}(\theta^{(i)} + h\omega^{(i)}) - \mathcal{L}^{(i)}(\theta^{(i)})}{h} = \sum_{t \in [p]} \omega_t^{(i)} \frac{\partial \mathcal{L}^{(i)}(\theta^{(i)})}{\partial \theta_t^{(i)}} \\ &= - \sum_{t \in [p]} \omega_t^{(i)} x_t^{(i)} \exp \left( - [\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right). \end{aligned}$$

**Proof of Lemma. C.4: Upper bound on expected  $\psi_u$**  Fix any  $i \in [n]$ ,  $u \in [L]$ , and  $\theta^{(i)} \in \Lambda_\theta$ . Then, given  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , we have

$$\begin{aligned} &\mathbb{E} \left[ \psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right] \\ &\stackrel{(a)}{=} \mathbb{E} \left[ \sum_{t \in S_u} \omega_t^{(i)} x_t^{(i)} \exp \left( - [\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right] \end{aligned}$$

$$\begin{aligned}
&\stackrel{(b)}{=} \sum_{t \in S_u} \omega_t^{(i)} \mathbb{E} \left[ x_t^{(i)} \exp \left( - [\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right] \\
&\stackrel{(c)}{=} \sum_{t \in S_u} \omega_t^{(i)} \mathbb{E} \left[ \mathbb{E} \left[ x_t^{(i)} \exp \left( - [\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)} \right] \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right], \quad (92)
\end{aligned}$$

where (a) follows from the definition of  $\psi_u(\mathbf{x}^{(i)}; \omega^{(i)})$  in (88), (b) follows from linearity of expectation, and (c) follows from the law of total expectation, i.e.,  $\mathbb{E}[\mathbb{E}[Y|X, Z]|Z] = \mathbb{E}[Y|Z]$  since  $\mathbf{x}_{-S_u}^{(i)} \subseteq \mathbf{x}_{-t}^{(i)}$ . Now, we bound  $\mathbb{E} \left[ x_t^{(i)} \exp \left( - [\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)} \right]$  for every  $t \in S_u$ . We have

$$\begin{aligned}
&\mathbb{E} \left[ x_t^{(i)} \exp \left( - [\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)} \right] \\
&= \int_{\mathcal{X}} x_t^{(i)} \exp \left( - [\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) f_{x_t | \mathbf{x}_{-t}, \mathbf{z}}(x_t^{(i)} | \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}; \theta_t^*(\mathbf{z}^{(i)}), \Theta_t^*) dx_t^{(i)} \\
&\stackrel{(a)}{=} \frac{\int_{\mathcal{X}} x_t^{(i)} \exp \left( 2[\Theta_{t,-t}^* - \widehat{\Theta}_{t,-t}]^\top \mathbf{x}_{-t}^{(i)} x_t^{(i)} + [\Theta_{tt}^* - \widehat{\Theta}_{tt}] \bar{x}_t^{(i)} \right) dx_t^{(i)}}{\int_{\mathcal{X}} \exp \left( [\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)} \right) dx_t^{(i)}} \\
&\stackrel{(b)}{=} \frac{\int_{\mathcal{X}} x_t^{(i)} \left[ 1 + 2[\Theta_{t,-t}^* - \widehat{\Theta}_{t,-t}]^\top \mathbf{x}_{-t}^{(i)} x_t^{(i)} + [\Theta_{tt}^* - \widehat{\Theta}_{tt}] \bar{x}_t^{(i)} \right] dx_t^{(i)}}{\int_{\mathcal{X}} \exp \left( [\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)} \right) dx_t^{(i)}} \\
&\quad + \frac{\int_{\mathcal{X}} x_t^{(i)} \left[ o \left( [\Theta_{t,-t}^* - \widehat{\Theta}_{t,-t}]^\top \mathbf{x}_{-t}^{(i)} x_t^{(i)} + [\Theta_{tt}^* - \widehat{\Theta}_{tt}] \bar{x}_t^{(i)} \right)^2 \right] dx_t^{(i)}}{\int_{\mathcal{X}} \exp \left( [\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)} \right) dx_t^{(i)}} \\
&\stackrel{(c)}{=} \frac{4x_{\max}^3 [\Theta_{t,-t}^* - \widehat{\Theta}_{t,-t}]^\top \mathbf{x}_{-t}^{(i)}}{3 \int_{\mathcal{X}} \exp \left( [\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)} \right) dx_t^{(i)}} \\
&\quad + \frac{x_{\max}^5 ([\Theta_{t,-t}^* - \widehat{\Theta}_{t,-t}]^\top \mathbf{x}_{-t}^{(i)}) (\Theta_{tt}^* - \widehat{\Theta}_{tt}) o(1)}{\int_{\mathcal{X}} \exp \left( [\theta_t^*(\mathbf{z}^{(i)}) + 2\Theta_{t,-t}^{*\top} \mathbf{x}_{-t}^{(i)}] x_t^{(i)} + \Theta_{tt}^* \bar{x}_t^{(i)} \right) dx_t^{(i)}}, \quad (93)
\end{aligned}$$

where (a) follows from (12) and  $\theta^{*(i)} = \theta^*(\mathbf{z}^{(i)}) \forall i \in [n]$ , (b) follows by using the Taylor series expansion  $\exp(y) = 1 + y + o(y^2)$  around zero, (c) follows because  $\int_{\mathcal{X}} x_t^{(i)} dx_t^{(i)} = \int_{\mathcal{X}} x_t^{(i)} \bar{x}_t^{(i)} dx_t^{(i)} = \int_{\mathcal{X}} (x_t^{(i)})^3 dx_t^{(i)} = \int_{\mathcal{X}} x_t^{(i)} (\bar{x}_t^{(i)})^2 dx_t^{(i)} = 0$ ,  $\int_{\mathcal{X}} (x_t^{(i)})^2 dx_t^{(i)} = 2x_{\max}^3/3$ , and  $\int_{\mathcal{X}} (x_t^{(i)})^2 \bar{x}_t^{(i)} dx_t^{(i)} = 8x_{\max}^5/45$ .

Now, we bound the numerators in (93) by using  $\|\Theta_t^* - \widehat{\Theta}_t\|_1 \leq \sqrt{p} \|\Theta_t^* - \widehat{\Theta}_t\|_2$ . Then, we invoke Theorem. 1 to bound  $\|\Theta_t^* - \widehat{\Theta}_t\|_2$  by  $\varepsilon \leftarrow \frac{3\varepsilon}{2C_{2,\tau} x_{\max}^3 \sqrt{p}}$ . Therefore, we subsume the second term by the first term resulting in the following bound:

$$\mathbb{E} \left[ x_t^{(i)} \exp \left( - [\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)} \right] \leq \frac{2C_{2,\tau} x_{\max}^3 \sqrt{p} \|\Theta_t^* - \widehat{\Theta}_t\|_2}{3}, \quad (94)$$

where we have used the triangle inequality,  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$  as well as  $\|\Theta_t^* - \widehat{\Theta}_t\|_1 \leq \sqrt{p} \|\Theta_t^* - \widehat{\Theta}_t\|_2$  to upper bound the numerator, and the arguments used in the proof of Lemma. B.6 as well as  $\int_{\mathcal{X}} dx_t^{(i)} = 2x_{\max}$  to lower bound the denominator.

Using Theorem. 1 in (94) with  $\varepsilon \leftarrow \frac{3\varepsilon}{2C_{2,\tau}x_{\max}^3\sqrt{p}}$  and  $\delta \leftarrow \delta$ , we have

$$\mathbb{E}\left[x_t^{(i)} \exp\left(-[\theta_t^{*(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \widehat{\Theta}_{tt}\bar{x}_t^{(i)}\right) \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right] \leq \varepsilon, \quad (95)$$

with probability at least  $1 - \delta$  as long as

$$n \geq \frac{ce^{c'\beta} \cdot p^4 (p \log \frac{p^2}{\delta\varepsilon^2} + \mathcal{M}_{\theta,n}(\frac{\varepsilon^2}{p}))}{\varepsilon^4}. \quad (96)$$

Using (95) and triangle inequality in (92), we have

$$\mathbb{E}\left[\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right] \leq \varepsilon \sum_{t \in S_u} |\omega_t^{(i)}| \leq \varepsilon \|\omega^{(i)}\|_1,$$

with probability at least  $1 - \delta$  as long as  $n$  satisfies (96).

**Proof of Lemma. C.5: Concentration of  $\psi_u$**  To show this concentration result, we use Corollary. H.1 (187) for the function  $q_2$ . To that end, we note that the pair  $\{\mathbf{x}, \mathbf{z}\}$  corresponds to a  $\tau$ -SGM (Definition. G.1) with  $\tau \triangleq (\alpha, \beta, x_{\max}, \Theta)$ . However, the random vector  $\mathbf{x}$  conditioned on  $\mathbf{z}$  need not satisfy the Dobrushin's uniqueness condition (Definition. F.2). Therefore, we cannot apply Corollary. H.1 (187) as is. To resolve this, we resort to Proposition. G.1 with  $\lambda = \frac{1}{4\sqrt{2}x_{\max}^2}$  to reduce the random vector  $\mathbf{x}$  conditioned on  $\mathbf{z}$  to Dobrushin's regime.

Fix any  $u \in [L]$ . Then, from Proposition. G.1(b), (i) the pair of random vectors  $\{\mathbf{x}_{S_u}, (\mathbf{x}_{-S_u}, \mathbf{z})\}$  corresponds to a  $\tau_1$ -SGM with  $\tau_1 \triangleq (\alpha + 2\beta x_{\max}, \frac{1}{4\sqrt{2}x_{\max}^2}, x_{\max}, \Theta_{S_u})$ , and (ii) the random vector  $\mathbf{x}_{S_u}$  conditioned on  $(\mathbf{x}_{-S_u}, \mathbf{z})$  satisfies the Dobrushin's uniqueness condition (Definition. F.2) with coupling matrix  $2\sqrt{2}x_{\max}^2|\Theta_{S_u}|$  with  $2\sqrt{2}x_{\max}^2\|\Theta_{S_u}\|_{\text{op}} \leq 2\sqrt{2}x_{\max}^2\lambda \leq 1/2$ . Now, for any fixed  $i \in [n]$ , we apply Corollary. H.1 (187) for the function  $q_2$  with  $\varepsilon \leftarrow \varepsilon$  for a given  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , to obtain

$$\mathbb{P}\left(\left|\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) - \mathbb{E}\left[\psi_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right]\right| \geq \varepsilon \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right) \leq \exp\left(\frac{-\varepsilon^2}{e^{c'\beta}\|\omega^{(i)}\|_2^2}\right).$$

### C.1.2 Proof of Lemma. C.3(b): Anti-concentration of second directional derivative

Fix some  $i \in [n]$  and some  $\theta^{(i)} \in \Lambda_\theta$ . Let  $\omega^{(i)}$  be as defined in (83). We claim that the second-order directional derivative of  $\mathcal{L}^{(i)}$  defined in (72) is given by

$$\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)}) = \sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2 \exp\left(-[\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}]x_t^{(i)} - \widehat{\Theta}_{tt}\bar{x}_t^{(i)}\right). \quad (97)$$

We provide a proof at the end. For now, we assume the claim and proceed. Now, we lower bound  $\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)})$  by a quadratic form as follows

$$\begin{aligned} \partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)}) &\stackrel{(a)}{\geq} \sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2 \times \exp\left(-(|\theta_t^{(i)}| + 2\|\widehat{\Theta}_t\|_1 \|\mathbf{x}^{(i)}\|_\infty) x_{\max}\right) \\ &\stackrel{(b)}{\geq} \sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2 \times \exp\left(-(\alpha + 2\beta x_{\max}) x_{\max}\right) \stackrel{(48)}{=} \frac{1}{C_{2,\tau} t \in [p]} \sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2, \end{aligned} \quad (98)$$

where (a) follows from (97) by triangle inequality, Cauchy–Schwarz inequality, and because  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ , and (b) follows because  $\widehat{\Theta} \in \Lambda_\Theta$ ,  $\theta^{(i)} \in \Lambda_\theta$ , and  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ .

Now, to show the anti-concentration of  $\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)})$ , we show the anti-concentration of the quadratic form in (98). To that end, we note that the pair  $\{\mathbf{x}, \mathbf{z}\}$  corresponds to a  $\tau$ -SGM (Definition. G.1) with  $\tau \triangleq (\alpha, \beta, x_{\max}, \Theta)$ . Then, we decompose the quadratic form in (98) as a sum of  $L = 1024\beta^2 x_{\max}^4 \log 4p$  terms using Proposition. G.1 (see Appendix. G) with  $\lambda = \frac{1}{4\sqrt{2}x_{\max}^2}$  and focus on these  $L$  terms. Consider the  $L$  subsets  $S_1, \dots, S_L \in [p]$  obtained from Proposition. G.1 and define

$$\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \triangleq \sum_{t \in S_u} (\omega_t^{(i)} x_t^{(i)})^2 \quad \text{for every } u \in [L]. \quad (99)$$

Then, we have

$$\sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2 \stackrel{(a)}{=} \frac{1}{L'} \sum_{u \in [L]} \sum_{t \in S_u} (\omega_t^{(i)} x_t^{(i)})^2 \stackrel{(99)}{=} \frac{1}{L'} \sum_{u \in [L]} \bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}), \quad (100)$$

where (a) follows because each  $t \in [p]$  appears in exactly  $L' = \lceil L/32\sqrt{2}\beta x_{\max}^2 \rceil$  of the sets  $S_1, \dots, S_L$  according to Proposition. G.1(a) (with  $\lambda = \frac{1}{4\sqrt{2}x_{\max}^2}$ ). Now, we focus on the  $L$  terms in (100).

Consider any  $u \in [L]$ . We claim that conditioned on  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , the expected value of  $\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)})$  can be upper bounded uniformly across all  $u \in [L]$ . We provide a proof at the end.

**Lemma C.6** (Lower bound on expected  $\bar{\psi}_u$ ). *Fix  $i \in [n]$  and  $\theta^{(i)} \in \Lambda_\theta$ . Then, with  $\omega^{(i)}$  defined in (83) and given  $\mathbf{z}^{(i)}$  and  $\mathbf{x}_{-S_u}^{(i)}$ , we have*

$$\min_{u \in [L]} \mathbb{E} \left[ \bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right] \geq \frac{2x_{\max}^2}{\pi e C_{2,\tau}^4} \|\omega^{(i)}\|_2^2,$$

where the constant  $C_{2,\tau}$  was defined in (48).

Consider again any  $u \in [L]$ . Now, we claim that conditioned on  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ ,  $\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)})$  concentrates around its conditional expected value. We provide a proof at the end.

**Lemma C.7** (Concentration of  $\bar{\psi}_u$ ). *Fix  $\varepsilon > 0$ ,  $i \in [n]$ ,  $u \in [L]$ , and  $\theta^{(i)} \in \Lambda_\theta$ . Then, with  $\omega^{(i)}$  defined in (83) and given  $\mathbf{z}^{(i)}$  and  $\mathbf{x}_{-S_u}^{(i)}$ , we have*

$$\left| \bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) - \mathbb{E} \left[ \bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right] \right| \leq \varepsilon,$$

with probability at least  $1 - \exp\left(\frac{-\varepsilon^2}{e^{c'\beta} \|\omega^{(i)}\|_2^2}\right)$ .

Given these lemmas, we proceed to show the anti-concentration of the quadratic form in (98) implying the anti-concentration of  $\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)})$ . To that end, for any  $u \in [L]$ , given  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , let  $E_u$  denote the event that

$$\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \geq \mathbb{E} \left[ \bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)} \right] - \frac{x_{\max}^2}{\pi e C_{2,\tau}^4} \|\omega^{(i)}\|_2^2. \quad (101)$$

Since  $E_u$  in an indicator event, using the law of total expectation results in

$$\mathbb{P}(E_u) = \mathbb{E} \left[ \mathbb{P}(E_u | \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}) \right] \stackrel{(a)}{\geq} 1 - \exp \left( - \frac{\|\omega^{(i)}\|_2^2}{e^{c'\beta}} \right),$$

where (a) follows from Lemma. C.7 with  $\varepsilon \leftarrow \frac{x_{\max}^2}{\pi e C_{2,\tau}^4} \|\omega^{(i)}\|_2^2$ . Now, by applying the union bound over all  $u \in [L]$  where  $L = 1024\beta^2 x_{\max}^4 \log 4p$ , we have

$$\mathbb{P} \left( \bigcap_{u \in [L]} E_u \right) \geq 1 - O \left( \beta^2 \log p \exp \left( - \frac{\|\omega^{(i)}\|_2^2}{e^{c'\beta}} \right) \right).$$

Now, assume the event  $\bigcap_{u \in [L]} E_u$  holds. Whenever this holds, we also have

$$\begin{aligned} \sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2 &\stackrel{(100)}{=} \frac{1}{L'} \sum_{u \in [L]} \bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \\ &\stackrel{(101)}{\geq} \frac{1}{L'} \sum_{u \in [L]} \left( \mathbb{E} [\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) | \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}] - \frac{x_{\max}^2}{\pi e C_{2,\tau}^4} \|\omega^{(i)}\|_2^2 \right) \\ &\stackrel{(a)}{\geq} \frac{1}{L'} \sum_{u \in [L]} \frac{x_{\max}^2}{\pi e C_{2,\tau}^4} \|\omega^{(i)}\|_2^2 = \frac{x_{\max}^2 L}{\pi e L' C_{2,\tau}^4} \|\omega^{(i)}\|_2^2, \end{aligned} \quad (102)$$

where  $L' = \lceil L/32\sqrt{2}\beta x_{\max}^2 \rceil$  and (a) follows from Lemma. C.6. Finally, approximating  $L' = L/32\sqrt{2}\beta x_{\max}^2$  and using (98), we have

$$\partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)}) \geq \frac{1}{C_{2,\tau}} \sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2 \stackrel{(102)}{\geq} \frac{32\sqrt{2}\beta x_{\max}^4}{\pi e C_{2,\tau}^5} \|\omega^{(i)}\|_2^2,$$

which completes the proof.

**Proof of (97): Expression for second directional derivative** Fix any  $i \in [n]$ . The second-order partial derivatives of  $\mathcal{L}^{(i)}$  (defined in (72)) with respect to the entries of the parameter vector  $\theta^{(i)}$  are given by

$$\frac{\partial^2 \mathcal{L}^{(i)}(\theta^{(i)})}{\partial [\theta_t^{(i)}]^2} = [x_t^{(i)}]^2 \exp \left( - [\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \quad \text{for all } t \in [p].$$

Now, we can write the second-order directional derivative of  $\mathcal{L}^{(i)}$  as

$$\begin{aligned} \partial_{[\omega^{(i)}]_2}^2 \mathcal{L}^{(i)}(\theta^{(i)}) &\triangleq \lim_{h \rightarrow 0} \frac{\partial_{\omega^{(i)}} \mathcal{L}^{(i)}(\theta^{(i)} + h\omega^{(i)}) - \partial_{\omega^{(i)}} \mathcal{L}^{(i)}(\theta^{(i)})}{h} = \sum_{t \in [p]} [\omega_t^{(i)}]^2 \frac{\partial^2 \mathcal{L}^{(i)}(\theta^{(i)})}{\partial [\theta_t^{(i)}]^2} \\ &= \sum_{t \in [p]} (\omega_t^{(i)} x_t^{(i)})^2 \exp \left( - [\theta_t^{(i)} + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_{-t}^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right). \end{aligned}$$

**Proof of Lemma. C.6: Lower bound on expected  $\bar{\psi}_u$**  Fix any  $i \in [n]$ ,  $u \in [L]$ , and  $\theta^{(i)} \in \Lambda_\theta$ . Then, given  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , we have

$$\begin{aligned}
\mathbb{E}\left[\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right] &\stackrel{(99)}{=} \mathbb{E}\left[\sum_{t \in S_u} (\omega_t^{(i)} x_t^{(i)})^2 \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right] \\
&\stackrel{(a)}{=} \sum_{t \in S_u} \mathbb{E}\left[(\omega_t^{(i)} x_t^{(i)})^2 \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right] \\
&\stackrel{(b)}{=} \sum_{t \in S_u} \mathbb{E}\left[\mathbb{E}\left[(\omega_t^{(i)} x_t^{(i)})^2 \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right] \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right] \\
&\stackrel{(c)}{\geq} \sum_{t \in S_u} \mathbb{E}\left[\mathbb{V}\text{ar}\left(\omega_t^{(i)} x_t^{(i)} \mid \mathbf{x}_{-t}^{(i)}, \mathbf{z}^{(i)}\right) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right] \\
&\stackrel{(d)}{\geq} \frac{2x_{\max}^2}{\pi e C_{2,\tau}^4} \|\omega^{(i)}\|_2^2,
\end{aligned}$$

where (a) follows from linearity of expectation, (b) follows from the law of total expectation i.e.,  $\mathbb{E}[\mathbb{E}[Y|X, Z]|Z] = \mathbb{E}[Y|Z]$  since  $\mathbf{x}_{-S_u}^{(i)} \subseteq \mathbf{x}_{-t}^{(i)}$ , (c) follows from the fact that for any random variable  $a$ ,  $\mathbb{E}[a^2] \geq \mathbb{V}\text{ar}[a]$ , and (d) follows from Lemma. B.6.

**Proof of Lemma. C.7: Concentration of  $\bar{\psi}_u$**  To show this concentration result, we use Corollary. H.1 (187) for the function  $q_1$ . To that end, we note that the pair  $\{\mathbf{x}, \mathbf{z}\}$  corresponds to a  $\tau$ -SGM (Definition. G.1) with  $\tau \triangleq (\alpha, \beta, x_{\max}, \Theta)$ . However, the random vector  $\mathbf{x}$  conditioned on  $\mathbf{z}$  need not satisfy the Dobrushin's uniqueness condition (Definition. F.2). Therefore, we cannot apply Corollary. H.1 (187) as is. To resolve this, we resort to Proposition. G.1 with  $\lambda = \frac{1}{4\sqrt{2}x_{\max}^2}$  to reduce the random vector  $\mathbf{x}$  conditioned on  $\mathbf{z}$  to Dobrushin's regime.

Fix any  $u \in [L]$ . Then, from Proposition. G.1(b), (i) the pair of random vectors  $\{\mathbf{x}_{S_u}, (\mathbf{x}_{-S_u}, \mathbf{z})\}$  corresponds to a  $\tau_1$ -SGM with  $\tau_1 \triangleq (\alpha + 2\beta x_{\max}, \frac{1}{4\sqrt{2}x_{\max}^2}, x_{\max}, \Theta_{S_u})$ , and (ii) the random vector  $\mathbf{x}_{S_u}$  conditioned on  $(\mathbf{x}_{-S_u}, \mathbf{z})$  satisfies the Dobrushin's uniqueness condition (Definition. F.2) with coupling matrix  $2\sqrt{2}x_{\max}^2|\Theta_{S_u}|$  with  $2\sqrt{2}x_{\max}^2\|\|\Theta_{S_u}\|\|_{\text{op}} \leq 2\sqrt{2}x_{\max}^2\lambda \leq 1/2$ . Now, for any fixed  $i \in [n]$ , we apply Corollary. H.1 (187) for the function  $q_1$  with  $\varepsilon = \varepsilon$  for a given  $\mathbf{x}_{-S_u}^{(i)}$  and  $\mathbf{z}^{(i)}$ , to obtain

$$\mathbb{P}\left(\left|\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) - \mathbb{E}\left[\bar{\psi}_u(\mathbf{x}^{(i)}; \omega^{(i)}) \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right]\right| \geq \varepsilon \mid \mathbf{x}_{-S_u}^{(i)}, \mathbf{z}^{(i)}\right) \leq \exp\left(\frac{-\varepsilon^2}{e^{c'\beta}\|\omega^{(i)}\|_2^2}\right).$$

## C.2 Proof of Lemma. C.2: Lipschitzness of the loss function

Fix any  $i \in [n]$ , any  $\theta^{(i)}, \tilde{\theta}^{(i)} \in \Lambda_\theta$ . Consider the direction  $\omega^{(i)} = \tilde{\theta}^{(i)} - \theta^{(i)}$ , and define the function  $q : [0, 1] \rightarrow \mathbb{R}$  as follows

$$q(a) = \mathcal{L}^{(i)}(\theta^{(i)} + a(\tilde{\theta}^{(i)} - \theta^{(i)})). \tag{103}$$

Then, the desired inequality in (74) is equivalent to

$$|q(1) - q(0)| \leq x_{\max} C_{2,\tau} \|\omega^{(i)}\|_1.$$

From the mean value theorem, there exists  $a' \in (0, 1)$  such that

$$|q(1) - q(0)| = \left| \frac{dq(a')}{da} \right|. \quad (104)$$

Therefore, we have

$$\begin{aligned} |q(1) - q(0)| &\stackrel{(104)}{=} \left| \frac{dq(a')}{da} \right| \stackrel{(103)}{=} \left| \frac{d\mathcal{L}^{(i)}(\theta^{(i)} + a'(\tilde{\theta}^{(i)} - \theta^{(i)}))}{da} \right| \\ &\stackrel{(84)}{=} \left| \partial_{\omega^{(i)}}(\mathcal{L}^{(i)}(\theta^{(i)})) \Big|_{\theta^{(i)} = \theta^{(i)} + a'(\tilde{\theta}^{(i)} - \theta^{(i)})} \right|. \end{aligned} \quad (105)$$

Using (87) in (105), we have

$$\begin{aligned} |q(1) - q(0)| &= \left| \sum_{t \in [p]} \omega_t^{(i)} x_t^{(i)} \exp \left( -[\theta_t^{(i)} + a'(\tilde{\theta}_t^{(i)} - \theta_t^{(i)}) + 2\widehat{\Theta}_{t,-t}^\top \mathbf{x}_t^{(i)}] x_t^{(i)} - \widehat{\Theta}_{tt} \bar{x}_t^{(i)} \right) \right| \\ &\stackrel{(a)}{\leq} x_{\max} \sum_{t \in [p]} |\omega_t^{(i)}| \exp \left( \left[ |(1-a')\theta_t^{(i)}| + |a'\tilde{\theta}_t^{(i)}| + 2\|\widehat{\Theta}_t\|_1 \|\mathbf{x}^{(i)}\|_\infty \right] x_{\max} \right) \\ &\stackrel{(b)}{\leq} x_{\max} \exp \left( ((1-a')\alpha + a'\alpha + 2\beta x_{\max}) x_{\max} \right) \sum_{t \in [p]} |\omega_t^{(i)}| \\ &\stackrel{(48)}{=} x_{\max} C_{2,\tau} \|\omega^{(i)}\|_1, \end{aligned}$$

where (a) follows from triangle inequality, Cauchy–Schwarz inequality, and because  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$  and (b) follows because  $\theta^{(i)}, \tilde{\theta}^{(i)} \in \Lambda_\theta$ ,  $\widehat{\Theta} \in \Lambda_\Theta$ , and  $\|\mathbf{x}^{(i)}\|_\infty \leq x_{\max}$  for all  $i \in [n]$ .

## D Proof of Theorem 2: Guarantee on quality of outcome estimate

Fix any unit  $i \in [n]$  and an alternate intervention  $\tilde{\mathbf{a}}^{(i)} \in \mathcal{A}^{p_a}$ . Then, we have

$$\mu^{(i)}(\tilde{\mathbf{a}}^{(i)}) \stackrel{(8)}{=} \mathbb{E}[\mathbf{y}^{(i)}(\tilde{\mathbf{a}}^{(i)}) | \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}] \stackrel{(a)}{=} \mathbb{E}[\mathbf{y} | \mathbf{a} = \tilde{\mathbf{a}}^{(i)}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}],$$

where (a) follows because the unit-level counterfactual distribution is equivalent to unit-level conditional distribution under the causal framework considered as described in Section 3.1. To obtain a convenient expression for  $\mathbb{E}[\mathbf{y} | \mathbf{a} = \tilde{\mathbf{a}}^{(i)}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}]$ , we identify  $\Phi^{*(u,y)} \in \mathbb{R}^{p_u \times p_y}$  to be the component of  $\Theta^*$  corresponding to  $\mathbf{u}$  and  $\mathbf{y}$  for all  $\mathbf{u} \in \{\mathbf{v}, \mathbf{a}, \mathbf{y}\}$  and  $\theta^{*(i,y)} \in \mathbb{R}^{p_y}$  to be the component of  $\theta^{*(i)}$  corresponding to  $\mathbf{y}$ . Then, the conditional distribution of  $\mathbf{y}$  as a function of the interventions  $\mathbf{a}$ , while keeping  $\mathbf{v}$  and  $\mathbf{z}$  fixed at the corresponding realizations for unit  $i$ , i.e.,  $\mathbf{v}^{(i)}$  and  $\mathbf{z}^{(i)}$ , respectively, can be written as

$$f_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a}) \propto \exp \left( [\theta^{*(i,y)} + 2\mathbf{v}^{(i)\top} \Phi^{*(v,y)} + 2\mathbf{a}^\top \Phi^{*(a,y)}] \mathbf{y} + \mathbf{y}^\top \Phi^{*(y,y)} \mathbf{y} \right). \quad (106)$$

Therefore, we have

$$\mathbb{E}[\mathbf{y} | \mathbf{a} = \tilde{\mathbf{a}}^{(i)}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}] = \mathbb{E}_{f_{\mathbf{y}|\mathbf{a}}^{(i)}}[\mathbf{y} | \mathbf{a} = \tilde{\mathbf{a}}^{(i)}].$$

Now, consider the  $p_u$  dimensional random vector  $\mathbf{u}$  supported on  $\mathcal{X}^{p_u}$  with distribution  $f_{\mathbf{u}}$  parameterized by  $\psi \in \mathbb{R}^{p_y}$  and  $\Psi \in \mathbb{R}^{p_y \times p_y}$  as follows

$$f_{\mathbf{u}}(\mathbf{u}|\psi, \Psi) \propto \exp(\psi^\top \mathbf{u} + \mathbf{u}^\top \Psi \mathbf{u}). \quad (107)$$

Then, note that  $\widehat{f}_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a})$  in (15) and  $f_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a})$  in (106) belong to the set  $\{f_{\mathbf{u}}(\cdot|\psi, \Psi) : \psi \in \mathbb{R}^{p_y}, \Psi \in \mathbb{R}^{p_y \times p_y}\}$  for some  $\psi$  and  $\Psi$ . Now, we consider any two distributions in this set, namely  $f_{\mathbf{u}}(\mathbf{u}|\widehat{\psi}, \widehat{\Psi})$  and  $f_{\mathbf{u}}(\mathbf{u}|\psi^*, \Psi^*)$ . Then, we claim that the two norm of the difference of the mean vectors of these distributions is bounded as below. We provide a proof at the end.

**Lemma D.1** (Perturbation in the mean vector). *For any  $\psi \in \mathbb{R}^{p_y}$  and  $\Psi \in \mathbb{R}^{p_y \times p_y}$ , let  $\mu_{\psi, \Psi}(\mathbf{u}) \in \mathbb{R}^{p_u}$  and  $\text{Cov}_{\psi, \Psi}(\mathbf{u}, \mathbf{u}) \in \mathbb{R}^{p_u \times p_u}$  denote the mean vector and the covariance matrix of  $\mathbf{u}$ , respectively, with respect to  $f_{\mathbf{u}}$  in (107). Then, for any  $\widehat{\psi}, \psi^* \in \mathbb{R}^{p_y}$  and  $\widehat{\Psi}, \Psi^* \in \mathbb{R}^{p_y \times p_y}$ , there exists some  $t \in (0, 1)$ ,  $\widetilde{\psi} \triangleq t\widehat{\psi} + (1-t)\psi^*$  and  $\widetilde{\Psi} \triangleq t\widehat{\Psi} + (1-t)\Psi^*$  such that*

$$\begin{aligned} \|\mu_{\widetilde{\psi}, \widetilde{\Psi}}(\mathbf{u}) - \mu_{\psi^*, \Psi^*}(\mathbf{u})\|_2 &\leq \|\text{Cov}_{\widetilde{\psi}, \widetilde{\Psi}}(\mathbf{u}, \mathbf{u})\|_{\text{op}} \|\widehat{\psi} - \psi^*\|_2 \\ &\quad + \sum_{t_3 \in [p]} \|\text{Cov}_{\widetilde{\psi}, \widetilde{\Psi}}(\mathbf{u}, u_{t_3} \mathbf{u})\|_{\text{op}} \|(\widehat{\Psi}_{t_3} - \Psi_{t_3}^*)\|_2. \end{aligned}$$

Given this lemma, we proceed with the proof. By applying this lemma to  $\widehat{f}_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a})$  in (15) and  $f_{\mathbf{y}|\mathbf{a}}^{(i)}(\mathbf{y}|\mathbf{a})$  in (106), we see that it is sufficient to show the following bound

$$\begin{aligned} \|(\theta^{*(i,y)} - \widehat{\theta}^{(i,y)}) + 2\mathbf{v}^{(i)\top}(\Phi^{*(v,y)} - \widehat{\Phi}^{(v,y)}) + 2\widetilde{\mathbf{a}}^{(i)\top}(\Phi^{*(a,y)} - \widehat{\Phi}^{(a,y)})\|_2 \\ + \sum_{t \in [p_y]} \|\Phi_t^{*(y,y)} - \widehat{\Phi}_t^{(y,y)}\|_2 \leq R(\varepsilon, \delta/n) + p\varepsilon. \end{aligned}$$

To that end, we have

$$\sum_{t \in [p_y]} \|\Phi_t^{*(y,y)} - \widehat{\Phi}_t^{(y,y)}\|_2 \stackrel{(a)}{\leq} \sum_{t \in [p_y]} \|\Theta_t^* - \widehat{\Theta}_t\|_2, \quad (108)$$

where (a) follows because  $\ell_2$  norm of any sub-vector is no more than  $\ell_2$  norm of the vector. Similarly, we have

$$\begin{aligned} &\|(\theta^{*(i,y)} - \widehat{\theta}^{(i,y)}) + 2\mathbf{v}^{(i)\top}(\Phi^{*(v,y)} - \widehat{\Phi}^{(v,y)}) + 2\widetilde{\mathbf{a}}^{(i)\top}(\Phi^{*(a,y)} - \widehat{\Phi}^{(a,y)})\|_2 \\ &\stackrel{(a)}{\leq} \|\theta^{*(i,y)} - \widehat{\theta}^{(i,y)}\|_2 + 2\|\mathbf{v}^{(i)\top}(\Phi^{*(v,y)} - \widehat{\Phi}^{(v,y)})\|_2 + 2\|\widetilde{\mathbf{a}}^{(i)\top}(\Phi^{*(a,y)} - \widehat{\Phi}^{(a,y)})\|_2 \\ &\stackrel{(b)}{\leq} \|\theta^{*(i,y)} - \widehat{\theta}^{(i,y)}\|_2 + 2\|\mathbf{v}^{(i)}\|_2 \|\Phi^{*(v,y)} - \widehat{\Phi}^{(v,y)}\|_{\text{op}} + 2\|\widetilde{\mathbf{a}}^{(i)}\|_2 \|(\Phi^{*(a,y)} - \widehat{\Phi}^{(a,y)})\|_{\text{op}} \\ &\stackrel{(c)}{\leq} \|\theta^{*(i)} - \widehat{\theta}^{(i)}\|_2 + 2\left(\|\mathbf{v}^{(i)}\|_2 + \|\widetilde{\mathbf{a}}^{(i)}\|_2\right) \|\Theta^* - \widehat{\Theta}\|_{\text{op}} \\ &\stackrel{(d)}{\leq} \|\theta^{*(i)} - \widehat{\theta}^{(i)}\|_2 + 2\left(\|\mathbf{v}^{(i)}\|_2 + \|\widetilde{\mathbf{a}}^{(i)}\|_2\right) \|\Theta^* - \widehat{\Theta}\|_1 \\ &\stackrel{(e)}{\leq} \|\theta^{*(i)} - \widehat{\theta}^{(i)}\|_2 + 2x_{\max}(\sqrt{p_v} + \sqrt{p_a}) \|\Theta^* - \widehat{\Theta}\|_1, \end{aligned} \quad (109)$$

where (a) follows from triangle inequality, (b) follows because induced matrix norms are submultiplicative, (c) follows because operator norm of any sub-matrix is no more than operator norm of the matrix and  $\ell_2$  norm of any sub-vector is no more than  $\ell_2$  norm of the vector, (d) follows because  $\Theta^* - \widehat{\Theta}$  is symmetric and because matrix operator norm is bounded by square root of the product of matrix one norm and matrix infinity norm, and (e) follows because  $\max\{\|\mathbf{v}^{(i)}\|_\infty, \|\mathbf{a}^{(i)}\|_\infty\} \leq x_{\max}$  for all  $i \in [n]$ .

Now, combining (108) and (109), we have

$$\begin{aligned} & \|(\boldsymbol{\theta}^{*(i,y)} - \widehat{\boldsymbol{\theta}}^{(i,y)}) + 2\mathbf{v}^{(i)\top}(\Phi^{*(v,y)} - \widehat{\Phi}^{(v,y)}) + 2\widehat{\mathbf{a}}^{(i)\top}(\Phi^{*(a,y)} - \widehat{\Phi}^{(a,y)})\|_2 + \sum_{t \in [p_y]} \|\Phi_t^{*(y,y)} - \widehat{\Phi}_t^{(y,y)}\|_2 \\ & \leq \|\boldsymbol{\theta}^{*(i)} - \widehat{\boldsymbol{\theta}}^{(i)}\|_2 + 2x_{\max}(\sqrt{p_v} + \sqrt{p_a}) \|\Theta^* - \widehat{\Theta}\|_1 + \sum_{t \in [p_y]} \|\Theta_t^* - \widehat{\Theta}_t\|_2 \\ & \stackrel{(a)}{\leq} R(\varepsilon, \delta/n) + 2x_{\max}(\sqrt{p_v} + \sqrt{p_a})\sqrt{p}\varepsilon + p_y\varepsilon, \end{aligned}$$

and (a) follows from Theorem. 1 by using the relationship between vector norms. The proof is complete by rescaling  $\varepsilon$  and absorbing the constants in  $c$ .

**Proof of Lemma. D.1: Perturbation in the mean vector** Let  $Z(\psi, \Psi) \in \mathbb{R}_+$  denote the log-partition function of  $f_{\mathbf{u}}(\cdot|\psi, \Psi)$  in (107). Then, from (Busa-Fekete et al., 2019, Theorem 1), we have

$$\|\mu_{\widehat{\psi}, \widehat{\Psi}}(\mathbf{u}) - \mu_{\psi^*, \Psi^*}(\mathbf{u})\|_2 = \|\nabla_{\widehat{\psi}} Z(\widehat{\psi}, \widehat{\Psi}) - \nabla_{\psi^*} Z(\psi^*, \Psi^*)\|_2. \quad (110)$$

For  $t_1, t_2, t_3 \in [p]$ , consider  $\frac{\partial^2 Z(\psi, \Psi)}{\partial \psi_{t_1} \partial \psi_{t_2}}$  and  $\frac{\partial^2 Z(\psi, \Psi)}{\partial \psi_{t_1} \partial \Psi_{t_2, t_3}}$ . Using the fact that the Hessian of the log partition function of any regular exponential family is the covariance matrix of the associated sufficient statistic, we have

$$\frac{\partial^2 Z(\psi, \Psi)}{\partial \psi_{t_1} \partial \psi_{t_2}} = \text{Cov}_{\psi, \Psi}(\mathbf{u}_{t_1}, \mathbf{u}_{t_2}) \quad \text{and} \quad \frac{\partial^2 Z(\psi, \Psi)}{\partial \psi_{t_1} \partial \Psi_{t_2, t_3}} = \text{Cov}_{\psi, \Psi}(\mathbf{u}_{t_1}, \mathbf{u}_{t_2} \mathbf{u}_{t_3}). \quad (111)$$

Now, for some  $c \in (0, 1)$ ,  $\widetilde{\psi} \triangleq c\widehat{\psi} + (1-c)\psi^*$  and  $\widetilde{\Psi} \triangleq c\widehat{\Psi} + (1-c)\Psi^*$ , we have the following from the mean value theorem

$$\begin{aligned} & \frac{\partial Z(\widehat{\psi}, \widehat{\Psi})}{\partial \widehat{\psi}_{t_1}} - \frac{\partial Z(\psi^*, \Psi^*)}{\partial \psi_{t_1}^*} \\ & = \sum_{t_2 \in [p]} \frac{\partial^2 Z(\widetilde{\psi}, \widetilde{\Psi})}{\partial \widetilde{\psi}_{t_2} \partial \widetilde{\psi}_{t_1}} \cdot (\widehat{\psi}_{t_2} - \psi_{t_2}^*) + \sum_{t_2 \in [p]} \sum_{t_3 \in [p]} \frac{\partial^2 Z(\widetilde{\psi}, \widetilde{\Psi})}{\partial \widetilde{\Psi}_{t_2, t_3} \partial \widetilde{\psi}_{t_1}} \cdot (\widehat{\Psi}_{t_2, t_3} - \Psi_{t_2, t_3}^*) \\ & \stackrel{(111)}{=} \sum_{t_2 \in [p]} \text{Cov}_{\widetilde{\psi}, \widetilde{\Psi}}(\mathbf{u}_{t_1}, \mathbf{u}_{t_2}) \cdot (\widehat{\psi}_{t_2} - \psi_{t_2}^*) + \sum_{t_3 \in [p]} \sum_{t_2 \in [p]} \text{Cov}_{\widetilde{\psi}, \widetilde{\Psi}}(\mathbf{u}_{t_1}, \mathbf{u}_{t_3} \mathbf{u}_{t_2}) \cdot (\widehat{\Psi}_{t_3, t_2} - \Psi_{t_3, t_2}^*). \end{aligned}$$

Now, using the triangle inequality and sub-multiplicativity of induced matrix norms, we have

$$\begin{aligned} \|\nabla_{\widehat{\psi}} Z(\widehat{\psi}, \widehat{\Psi}) - \nabla_{\psi^*} Z(\psi^*, \Psi^*)\|_2 & \leq \|\text{Cov}_{\widetilde{\psi}, \widetilde{\Psi}}(\mathbf{u}, \mathbf{u})\|_{\text{op}} \|\widehat{\psi} - \psi^*\|_2 \\ & \quad + \sum_{t_3 \in [p]} \|\text{Cov}_{\widetilde{\psi}, \widetilde{\Psi}}(\mathbf{u}, \mathbf{u}_{t_3} \mathbf{u})\|_{\text{op}} \|\widehat{\Psi}_{t_3} - \Psi_{t_3}^*\|_2. \end{aligned} \quad (112)$$

Combining (110) and (112) completes the proof.

## D.1 Bounded operator norms for perturbations in the parameters

In Section. 4.2, we assumed the operator norms of (i) the covariance matrix of  $\mathbf{y}$  conditioned on  $\mathbf{a}$ ,  $\mathbf{z}$ , and  $\mathbf{v}$  and (ii) the cross-covariance matrix of  $\mathbf{y}$  and  $y_t \mathbf{y}$  conditioned on  $\mathbf{a}$ ,  $\mathbf{z}$ , and  $\mathbf{v}$  for all  $t \in [p_y]$  to remain bounded for small perturbation in the parameters. In this section, we provide examples where these hold.

Suppose the distribution of  $\mathbf{y}$  conditioned on  $\mathbf{a}$ ,  $\mathbf{z}$ , and  $\mathbf{v}$  is a Gaussian distribution. For simplicity, let the mean of this distribution be zero. Then, for any  $t, u, v \in [p_y]$ ,

$$\text{Cov}_{\theta, \Theta}(y_u, y_t y_v | \mathbf{a}, \mathbf{z}, \mathbf{v}) = \mathbb{E}_{\theta, \Theta}(y_u y_t y_v | \mathbf{a}, \mathbf{z}, \mathbf{v}) \stackrel{(a)}{=} 0.$$

where (a) follows because  $\mathbb{E}_{\theta, \Theta}(y_u y_t y_v | \mathbf{a}, \mathbf{z}, \mathbf{v})$  is the third cumulant of  $y_u y_t y_v | \mathbf{a}, \mathbf{z}, \mathbf{v}$  and the third cumulant for any Gaussian distribution is zero (Holmquist, 1988). Then,

$$\max_{t \in [p_y]} \|\text{Cov}_{\theta, \Theta}(\mathbf{y}, y_t \mathbf{y} | \mathbf{a}, \mathbf{z}, \mathbf{v})\|_{\text{op}} = 0. \quad (113)$$

Further, (113) also holds for small perturbations in  $\theta$  and  $\Theta$  as the distribution of  $\mathbf{y}$  conditioned on  $\mathbf{a}$ ,  $\mathbf{z}$ , and  $\mathbf{v}$  would still be a Gaussian distribution.

Now, we bound  $\|\text{Cov}_{\theta, \Theta}(\mathbf{y}, \mathbf{y} | \mathbf{a}, \mathbf{z}, \mathbf{v})\|_{\text{op}}$  under additional conditions. For simplicity, suppose  $\text{Var}_{\theta, \Theta}(y_t | \mathbf{a}, \mathbf{z}, \mathbf{v}) = 1$  for all  $t \in [p_y]$ . Further, suppose the (undirected) graphical structure associated with elements of  $\mathbf{y}$ , i.e.,  $y_1, \dots, y_{p_y}$ , is a chain (This would be true for the motivating example in Figure. 1(a)). If the correlation between any two elements of  $\mathbf{y}$  connected by an edge in the tree is equal to  $\rho \in [0, 1]$  (This is equivalent to all the off-diagonal non-zero entries of  $\Theta$  being the same), then for any  $u, v \in [p_y]$ ,

$$\text{Cov}_{\theta, \Theta}(y_u, y_v | \mathbf{a}, \mathbf{z}, \mathbf{v}) \stackrel{(a)}{=} \rho^{|u-v|},$$

where (a) follows by the correlation decay property for Gaussian tree models (Tan et al., 2010, Equation. 18). Then, for any  $0 \leq \rho < 1$

$$\|\text{Cov}_{\theta, \Theta}(\mathbf{y}, \mathbf{y} | \mathbf{a}, \mathbf{z}, \mathbf{v})\|_{\text{op}} \stackrel{(a)}{\leq} \frac{1 + \rho}{1 - \rho}, \quad (114)$$

where (a) follows from Trench (1999). Further, (114) holds for small perturbations in  $\theta$  and  $\Theta$  as long as  $\rho < 1$ . Therefore,  $C(\mathbb{B})$  in (22) is a constant (with respect to  $p$ ) for small perturbations in  $\theta$  and  $\Theta$ .

While we showed that  $C(\mathbb{B})$  is a constant for a class of Gaussian distributions, we expect similar results for truncated Gaussian distributions and exponential family distributions in (3).

## E Proof of Proposition 2: Impute missing covariates

We start by decomposing the true covariates  $\mathbf{v}$  into two variables: one to capture the randomness in the noisy observations  $\bar{\mathbf{v}}$  and the other to capture the randomness in the measurement error  $\Delta \mathbf{v}$ , i.e.,  $\mathbf{v} = \bar{\mathbf{v}} - \Delta \mathbf{v}$ . Then, by letting  $\bar{p} \triangleq 2p_v + p_a + p_y$  and using (27), the joint probability distribution  $f_{\bar{\mathbf{w}}}$  of the  $\bar{p}$ -dimensional random vector  $\bar{\mathbf{w}} \triangleq (\Delta \mathbf{v}, \bar{\mathbf{v}}, \mathbf{a}, \mathbf{y})$  can be parameterized by a vector  $\bar{\phi} \in \mathbb{R}^{\bar{p} \times 1}$  and a symmetric matrix  $\bar{\Phi} \in \mathbb{R}^{\bar{p} \times \bar{p}}$  as follows

$$f_{\bar{\mathbf{w}}}(\bar{\mathbf{w}}; \bar{\phi}, \bar{\Phi}) \propto \exp\left(\bar{\phi}^\top \bar{\mathbf{w}} + \bar{\mathbf{w}}^\top \bar{\Phi} \bar{\mathbf{w}}\right), \quad \text{where } \bar{\mathbf{w}} \triangleq (\Delta \mathbf{v}, \bar{\mathbf{v}}, \mathbf{a}, \mathbf{y}),$$

and  $\Delta \mathbf{v}$ ,  $\bar{\mathbf{v}}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  denote realizations of  $\Delta \mathbf{v}$ ,  $\bar{\mathbf{v}}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$ , respectively. More importantly,  $\bar{\phi}$  and  $\bar{\Phi}$  are derived completely from  $\phi$  and  $\Phi$ , respectively, and have special structure:

$$\begin{aligned}\bar{\phi}^{(\bar{v})} &= -\bar{\phi}^{(\Delta v)} = \phi^{(v)}, \\ \bar{\phi}^{(u)} &= \phi^{(u)} \text{ for all } \mathbf{u} \in \{\mathbf{a}, \mathbf{y}\}, \\ \bar{\Phi}^{(\Delta v, \Delta v)} &= \bar{\Phi}^{(\bar{v}, \bar{v})} = -\bar{\Phi}^{(\bar{v}, \Delta v)} = \Phi^{(v, v)} \\ \bar{\Phi}^{(u, \bar{v})} &= -\bar{\Phi}^{(u, \Delta v)} = \Phi^{(u, v)} \text{ for all } \mathbf{u} \in \{\mathbf{a}, \mathbf{y}\}, \text{ and} \\ \bar{\Phi}^{(u_1, u_2)} &= \Phi^{(u_1, u_2)} \text{ for all } \mathbf{u}_1, \mathbf{u}_2 \in \{\mathbf{a}, \mathbf{y}\}.\end{aligned}$$

Now, to learn counterfactuals and measurement errors for units  $i \in \{1, \dots, n/2\}$ , we use the methodology developed in Section. 3 by replacing the role of unobserved covariates  $\mathbf{z}$  by  $\Delta \mathbf{v}$ . In particular, we consider learning  $f_{\mathbf{y}|\mathbf{a}, \Delta \mathbf{v}, \bar{\mathbf{v}}}(\mathbf{y} = \cdot | \mathbf{a} = \cdot, \Delta \mathbf{v}, \bar{\mathbf{v}})$  as a function of  $\mathbf{a}$ . From (4) and the structure on  $\bar{\phi}$  and  $\bar{\Phi}$  described above, this reduces to learning

$$\begin{aligned}\text{(i)} \quad & \bar{\phi}^{(y)} + 2\bar{\Phi}^{(\Delta v, y)\top} \Delta \mathbf{v} + 2\bar{\Phi}^{(\bar{v}, y)\top} \bar{\mathbf{v}} = \phi^{(y)} - 2\Phi^{(v, y)\top} \Delta \mathbf{v} + 2\Phi^{(v, y)\top} \bar{\mathbf{v}}, \\ \text{(ii)} \quad & \bar{\Phi}^{(a, y)} = \Phi^{(a, y)}, \text{ and} \\ \text{(iii)} \quad & \bar{\Phi}^{(y, y)} = \Phi^{(y, y)}.\end{aligned}$$

To learn these, we consider the distribution of  $\mathbf{x} \triangleq (\bar{\mathbf{v}}, \mathbf{a}, \mathbf{y})$  conditioned on  $\Delta \mathbf{v} = \Delta \mathbf{v}$ . From (6), we have

$$f_{\mathbf{x}|\Delta \mathbf{v}}(\mathbf{x}|\Delta \mathbf{v}; \theta(\Delta \mathbf{v}), \Theta) \propto \exp\left([\theta(\Delta \mathbf{v})]^\top \mathbf{x} + \mathbf{x}^\top \Theta \mathbf{x}\right) \text{ with } \theta(\Delta \mathbf{v}) \triangleq \begin{bmatrix} \phi^{(v)} - 2\Phi^{(v, v)\top} \Delta \mathbf{v} \\ \phi^{(a)} - 2\Phi^{(v, a)\top} \Delta \mathbf{v} \\ \phi^{(y)} - 2\Phi^{(v, y)\top} \Delta \mathbf{v} \end{bmatrix}, \quad (115)$$

$\mathbf{x} \triangleq (\bar{\mathbf{v}}, \mathbf{a}, \mathbf{y})$ ,  $\Theta \triangleq \bar{\Phi}$ , and  $\bar{\mathbf{v}}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$  denoting realizations of  $\bar{\mathbf{v}}$ ,  $\mathbf{a}$ , and  $\mathbf{y}$ , respectively. The special structure on  $\bar{\Phi}$  discussed above implies that  $\Phi^{(v, v)}$ ,  $\Phi^{(v, a)}$ , and  $\Phi^{(v, y)}$  affect both  $\theta(\Delta \mathbf{v})$  and  $\Theta$  which we exploit. As mentioned in Section. 6.1, we denote the true distribution of  $\mathbf{x}$  conditioned on  $\Delta \mathbf{v} = \Delta \mathbf{v}$  by  $f_{\mathbf{x}|\Delta \mathbf{v}}(\cdot | \Delta \mathbf{v}; \theta^*(\Delta \mathbf{v}), \Theta^*)$ .

**Proof idea** First, we use units  $i \in \{n/2 + 1, \dots, n\}$  without any measurement error to estimate  $\phi^*$  and  $\Phi^* = \Theta^*$ , i.e., the parameters corresponding to the distribution of  $(\mathbf{v}, \mathbf{a}, \mathbf{y})$  (see Section. 6.1). Next, for units  $i \in \{1, \dots, n/2\}$  with measurement error, we estimate  $\theta^*(\Delta \mathbf{v}^{(i)})$  by expressing it as a linear combination of the estimates of  $\phi^*$  and  $\Phi^*$  (enabling the use of Example. 1). The coefficients of this linear combination turn out to be our estimates of the measurement error  $\Delta \mathbf{v}^{(i)}$ .

**Estimate  $\phi^*$  and  $\Phi^*$**  For units  $i \in \{n/2 + 1, \dots, n\}$ , under our assumption  $\Delta \mathbf{v}^{(i)} = 0$  implying  $\theta^*(\Delta \mathbf{v}^{(i)}) = \phi^*$ . Therefore, in addition to the population-level parameter  $\Theta^* = \Phi^*$ , the unit-level parameter  $\theta^*(\Delta \mathbf{v}) = \phi^*$  is also shared for these units. As a result, the set of distributions  $\{f_{\mathbf{x}|\Delta \mathbf{v}}(\cdot | \Delta \mathbf{v}; \theta^*(\Delta \mathbf{v}), \Theta^*)\}_{i=1}^n$  all coincide. Thus, learning  $\phi^*$  and  $\Phi^*$  boils down to learning parameters of a sparse graphical model (because of the assumptions in Section. 6.1) from  $n/2$  samples. We use the methodology and analysis from Shah et al. (2023) (which is closely related to the one in this work) to obtain estimates  $\hat{\phi}$  and  $\hat{\Phi}$  such that with probability at least  $1 - \delta$ , we have

$$\max \left\{ \|\phi^* - \hat{\phi}\|_2, \|\Phi^* - \hat{\Phi}\|_{2, \infty} \right\} \leq \varepsilon_1 \quad \text{whenever} \quad n \geq \frac{ce^{c\beta} \log \frac{p}{\sqrt{\delta}}}{\varepsilon_1^2}. \quad (116)$$

**Recover the unit-level parameters** Now, for units  $i \in \{1, \dots, n/2\}$ , we express the true unit-level parameters  $\theta^\star(\Delta \mathbf{v}^{(i)})$  as a linear combination of known vectors. To that end, fix any  $i \in [n/2]$ . Then, using (115), we can write  $\theta^\star(i) \triangleq \theta^\star(\Delta \mathbf{v}^{(i)})$  as a linear combination of  $p_v + 1$  vectors, i.e.,

$$\theta^\star(i) = \mathbf{B} \mathbf{a}^{(i)}, \quad (117)$$

where

$$\mathbf{B} \triangleq [\phi^\star, -2\Phi_1^\star, \dots, -2\Phi_{p_v}^\star] \in \mathbb{R}^{p \times (p_v+1)} \quad \text{and} \quad \mathbf{a}^{(i)} \triangleq \begin{bmatrix} 1 \\ \Delta \mathbf{v}^{(i)} \end{bmatrix} \in \mathbb{R}^{(p_v+1) \times 1}. \quad (118)$$

While we do not know the matrix  $\mathbf{B}$ , we can produce an estimate  $\widehat{\mathbf{B}}$  using  $\widehat{\phi}$  and  $\widehat{\Phi}$  such that, with probability at least  $1 - \delta$ ,

$$\|\widehat{\mathbf{B}} - \mathbf{B}\|_{2,\infty} \leq \varepsilon_1 \quad \text{whenever} \quad n \geq \frac{ce^{c'\beta} \log \frac{p}{\sqrt{\delta}}}{\varepsilon_1^2}. \quad (119)$$

This guarantee follows directly from (116) and the definition of  $\mathbf{B}$  in (118). Then, we can write

$$\theta^\star(i) = \widehat{\mathbf{B}} \widetilde{\mathbf{a}}^{(i)} \quad \text{where} \quad \widetilde{\mathbf{a}}^{(i)} \triangleq \mathbf{a}^{(i)} + \zeta, \quad (120)$$

for some error term  $\zeta$ . Conditioned on the event (119),  $\zeta$  can be controlled in following manner

$$\begin{aligned} \|\widehat{\mathbf{B}}\zeta\|_2 &\stackrel{(120)}{=} \|\theta^\star(i) - \widehat{\mathbf{B}}\mathbf{a}^{(i)}\|_2 \stackrel{(117)}{=} \|\mathbf{B}\mathbf{a}^{(i)} - \widehat{\mathbf{B}}\mathbf{a}^{(i)}\|_2 \\ &\stackrel{(a)}{\leq} \|\mathbf{B} - \widehat{\mathbf{B}}\|_{\text{op}} \|\mathbf{a}^{(i)}\|_2 \\ &\stackrel{(b)}{\leq} (\sqrt{p} \|\mathbf{B} - \widehat{\mathbf{B}}\|_{2,\infty}) \cdot (\sqrt{p_v + 1} \|\mathbf{a}^{(i)}\|_\infty) \stackrel{(c)}{\leq} \alpha \varepsilon_1 \sqrt{(p_v + 1)p}, \end{aligned} \quad (121)$$

where (a) follows from sub-multiplicativity of induced matrix norms, (b) follows from standard matrix norm inequalities, and (c) follows from (119) and because the measurement errors are bounded by  $\alpha$ .

Then, performing an analysis similar to one in Appendix. C while using the bound on  $n$  in (116) instead of the one in (19), and using Example. 1, we obtain estimates  $\widehat{\theta}^{(1)}, \dots, \widehat{\theta}^{(n/2)}$  such that (see Corollary. 1(a) for reference), with probability at least  $1 - \delta$ , we have

$$\max_{i \in [n/2]} \text{MSE}(\widehat{\theta}^{(i)}, \theta^\star(i)) \leq \max \left\{ \varepsilon_1^2, \frac{ce^{c'\beta} (p_v + \log(\log \frac{np}{\delta}))}{p} \right\}, \quad (122)$$

whenever  $n \geq ce^{c'\beta} \varepsilon_1^{-2} (\log \frac{\sqrt{np}}{\sqrt{\delta}} + p_v)$ .

**Recover the measurement error** We condition on the event (122) happening and note that the above estimate  $\widehat{\theta}^{(i)}$  of the unit-level parameter  $\theta^\star(i)$  is of the form  $\widehat{\theta}^{(i)} = \widehat{\mathbf{B}} \widehat{\mathbf{a}}^{(i)}$  for  $i \in [n/2]$ . We declare  $\widehat{\mathbf{a}}^{(i)}$  as our estimate of the measurement error for unit  $i \in [n/2]$  and prove the corresponding guarantee below.

Fix any  $i \in [n/2]$ . From (120) and triangle inequality, we find that

$$\|\theta^\star(i) - \widehat{\theta}^{(i)}\|_2 = \|\widehat{\mathbf{B}}\mathbf{a}^{(i)} + \widehat{\mathbf{B}}\zeta - \widehat{\mathbf{B}}\widehat{\mathbf{a}}^{(i)}\|_2 \geq \|\widehat{\mathbf{B}}\mathbf{a}^{(i)} - \widehat{\mathbf{B}}\widehat{\mathbf{a}}^{(i)}\|_2 - \|\widehat{\mathbf{B}}\zeta\|_2. \quad (123)$$

Then, doing standard algebra with (123) yields that

$$\text{MSE}(\hat{\theta}^{(i)}, \theta^{*(i)}) + \frac{\|\widehat{\mathbf{B}}\zeta\|_2^2}{p} \geq \frac{\|\widehat{\mathbf{B}}\mathbf{a}^{(i)} - \widehat{\mathbf{B}}\widehat{\mathbf{a}}^{(i)}\|_2^2}{2p} = \frac{(\mathbf{a}^{(i)} - \widehat{\mathbf{a}}^{(i)})^\top \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}} (\mathbf{a}^{(i)} - \widehat{\mathbf{a}}^{(i)})}{2p}. \quad (124)$$

Combining (121), (122), and (124) with the choice  $\varepsilon_1 = \kappa\varepsilon_2/\alpha\sqrt{p_v+1}$ , we have

$$\frac{(\mathbf{a}^{(i)} - \widehat{\mathbf{a}}^{(i)})^\top \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}} (\mathbf{a}^{(i)} - \widehat{\mathbf{a}}^{(i)})}{2p} \leq \max \left\{ \frac{\varepsilon_2^2 \kappa^2}{\alpha^2(p_v+1)}, \frac{ce^{c'\beta}(p_v + \log(\log \frac{np}{\delta}))}{p} \right\} + \varepsilon_2^2 \kappa^2, \quad (125)$$

uniformly for all  $i \in [n/2]$ , with probability at least  $1 - \delta$ , whenever  $n \geq ce^{c'\beta} \kappa^{-2} \varepsilon_2^{-2} (p_v+1) (\log \frac{\sqrt{np}}{\sqrt{\delta}} + p_v)$ . Next, we claim that the eigenvalues of  $\widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}$  can be lower bounded by  $\kappa p/2$  whenever  $\varepsilon_2 \leq \sqrt{p/(p_v+1)}/8$ . Taking this claim as given at the moment, we continue our proof. We have

$$\frac{\kappa}{4} \|\mathbf{a}^{(i)} - \widehat{\mathbf{a}}^{(i)}\|_2^2 \leq \frac{(\mathbf{a}^{(i)} - \widehat{\mathbf{a}}^{(i)})^\top \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}} (\mathbf{a}^{(i)} - \widehat{\mathbf{a}}^{(i)})}{2p} \quad \text{whenever} \quad \varepsilon_2 \leq \frac{1}{8} \sqrt{\frac{p}{p_v+1}}, \quad (126)$$

uniformly for all  $i \in [n/2]$ . Combining (125) and (126) completes the proof.

It remains to show that the eigenvalues of  $\widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}$  can be lower bounded by  $\kappa p/2$  conditioned on (116). For any matrix  $\mathbf{M}$ , let  $\lambda_{\max}(\mathbf{M})$  and  $\lambda_{\min}(\mathbf{M})$  denote the largest and the smallest eigenvalues of  $\mathbf{M}$ , respectively. Then from Weyl's inequality (Bhatia, 2007, Theorem. 8.2), we have

$$\lambda_{\min}(\widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}) \geq \lambda_{\min}(\mathbf{B}^\top \mathbf{B}) - \lambda_{\max}(\mathbf{B}^\top \mathbf{B} - \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}) \stackrel{(a)}{\geq} \kappa p - \lambda_{\max}(\mathbf{B}^\top \mathbf{B} - \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}),$$

where (a) follows from the assumption on the eigenvalues of  $\mathbf{B}^\top \mathbf{B}$ . Now, it suffices to upper bound  $\lambda_{\max}(\mathbf{B}^\top \mathbf{B} - \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}})$  by  $\kappa p/2$ . We have

$$\begin{aligned} |\lambda_{\max}(\mathbf{B}^\top \mathbf{B} - \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}})| &\stackrel{(a)}{=} \|\mathbf{B}^\top \mathbf{B} - \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}\|_{\text{op}} \\ &\stackrel{(b)}{\leq} (p_v+1) \|\mathbf{B}^\top \mathbf{B} - \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}\|_{\text{max}} \\ &\stackrel{(c)}{\leq} (p_v+1) \left( \|\mathbf{B}^\top (\mathbf{B} - \widehat{\mathbf{B}})\|_{\text{max}} + \|(\mathbf{B} - \widehat{\mathbf{B}})^\top \widehat{\mathbf{B}}\|_{\text{max}} \right) \\ &\stackrel{(d)}{\leq} (p_v+1) (\|\mathbf{B}^\top\|_{2,\infty} + \|\widehat{\mathbf{B}}^\top\|_{2,\infty}) \|(\mathbf{B} - \widehat{\mathbf{B}})^\top\|_{2,\infty} \\ &\stackrel{(e)}{\leq} (p_v+1) (2\alpha\sqrt{p} + 2\alpha\sqrt{p}) \cdot \varepsilon_1 \stackrel{(f)}{\leq} 4\kappa\varepsilon_2\sqrt{p_v+1}\sqrt{p} \stackrel{(g)}{\leq} \frac{\kappa p}{2}, \end{aligned}$$

where (a) follows because  $\mathbf{B}^\top \mathbf{B} - \widehat{\mathbf{B}}^\top \widehat{\mathbf{B}}$  is symmetric, (b) follows from because  $\|\mathbf{M}\|_{\text{op}} \leq \|\mathbf{M}\|_{\text{F}} \leq d\|\mathbf{M}\|_{\text{max}}$  for any square matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$ , (c) follows from the triangle inequality, (d) follows by Cauchy-Schwarz inequality, (e) follows because  $\|\widehat{\mathbf{B}}\|_{\text{max}} \leq 2\alpha$ ,  $\|\mathbf{B}\|_{\text{max}} \leq 2\alpha$  (because of the assumptions in Section. 6.1), and from (116) and (118), (f) follows from the choice of  $\varepsilon_1$ , and (g) follows whenever  $\varepsilon_2 \leq \frac{1}{8} \sqrt{\frac{p}{p_v+1}}$ .

## F Logarithmic Sobolev inequality and tail bounds

In this section, we present two results which may be of independent interest. First, we show that a random vector supported on a compact set satisfies the logarithmic Sobolev inequality (to be defined) if it satisfies the Dobrushin's uniqueness condition (to be defined). This result is a generalization of the result in [Marton \(2015\)](#) for discrete random vectors to continuous random vectors supported on a compact set. Next, we show that if a random vector satisfies the logarithmic Sobolev inequality, then any arbitrary function of the random vector concentrates around its mean. This result is a generalization of the result in [Dagan et al. \(2021\)](#) for discrete random vectors to continuous random vectors.

Throughout this section, we consider a  $p$ -dimensional random vector  $\mathbf{x}$  supported on  $\mathcal{X}^p$  with distribution  $f_{\mathbf{x}}$  where  $p \geq 1$ . We start by defining the logarithmic Sobolev inequality (LSI). We use the convention  $0 \log 0 = 0$ .

**Definition F.1** (Logarithmic Sobolev inequality). *A random vector  $\mathbf{x}$  satisfies the logarithmic Sobolev inequality with constant  $\sigma^2 > 0$  (abbreviated as  $\text{LSI}_{\mathbf{x}}(\sigma^2)$ ) if*

$$\text{Ent}_{\mathbf{x}}(q^2) \leq \sigma^2 \mathbb{E}_{\mathbf{x}} \left[ \|\nabla_{\mathbf{x}} q(\mathbf{x})\|_2^2 \right] \quad \text{for all } q : \mathcal{X}^p \rightarrow \mathbb{R}, \quad (127)$$

where  $\text{Ent}_{\mathbf{x}}(g) \triangleq \mathbb{E}_{\mathbf{x}}[g(\mathbf{x}) \log g(\mathbf{x})] - \mathbb{E}_{\mathbf{x}}[g(\mathbf{x})] \log \mathbb{E}_{\mathbf{x}}[g(\mathbf{x})]$  denotes the entropy of the function  $g : \mathcal{X}^p \rightarrow \mathbb{R}_+$ .

Next, we state the Dobrushin's uniqueness condition. For any distributions  $f$  and  $g$ , let  $\|f - g\|_{\text{TV}}$  denote the total variation distance between  $f$  and  $g$ .

**Definition F.2** (Dobrushin's uniqueness condition). *A random vector  $\mathbf{x}$  satisfies the Dobrushin's uniqueness condition with coupling matrix  $\Theta \in \mathbb{R}_+^{p \times p}$  if  $\|\Theta\|_{\text{op}} < 1$ , and for every  $t \in [p], u \in [p] \setminus \{t\}$ , and  $\mathbf{x}_{-t}, \tilde{\mathbf{x}}_{-t} \in \mathcal{X}^{p-1}$  differing only in the  $u^{\text{th}}$  coordinate,*

$$\|f_{x_t | \mathbf{x}_{-t} = \mathbf{x}_{-t}} - f_{x_t | \mathbf{x}_{-t} = \tilde{\mathbf{x}}_{-t}}\|_{\text{TV}} \leq \Theta_{tu}. \quad (128)$$

We note that the Dobrushin's uniqueness condition, as originally stated (see [Marton \(2015\)](#)) for Ising model, also requires  $\Theta_{tt} = 0$  for all  $t \in [p]$ . This condition makes sense for Ising model where  $x_t^2 = 1$  for all  $t \in [p]$ . However, this is not true for continuous random vectors necessitating a need for modification in the condition.

From hereon, we let  $\mathcal{X}^p$  be compact unless otherwise specified. Moreover, we define

$$f_{\min} \triangleq \min_{t \in [p], \mathbf{x} \in \mathcal{X}^p} f_{x_t | \mathbf{x}_{-t}}(x_t | \mathbf{x}_{-t}). \quad (129)$$

Now, we provide the first main result of this section with a proof in [Appendix F.1](#).

**Proposition F.1** (Logarithmic Sobolev inequality). *If a random vector  $\mathbf{x}$  with  $f_{\min} > 0$  (see [\(129\)](#)) satisfies (a) the Dobrushin's uniqueness condition ([Definition F.2](#)) with coupling matrix  $\Theta \in \mathbb{R}_+^{p \times p}$ , and (b)  $x_t | \mathbf{x}_{-t}$  satisfies  $\text{LSI}_{x_t | \mathbf{x}_{-t} = \mathbf{x}_{-t}}(\sigma^2)$  for all  $t \in [p]$  and  $\mathbf{x}_{-t} \in \mathcal{X}^{p-1}$  (see [Definition F.1](#)), then it satisfies  $\text{LSI}_{\mathbf{x}}(2\sigma^2 / (f_{\min}(1 - \|\Theta\|_{\text{op}})^2))$ .*

Next, we define the notion of pseudo derivative and pseudo Hessian that come in handy in our proofs for providing upper bounds on the norm of the derivative and the Hessian.

**Definition F.3** (Pseudo derivative and Hessian). *For a function  $q : \mathcal{X}^p \rightarrow \mathbb{R}$ , the functions  $\tilde{\nabla}q : \mathcal{X}^p \rightarrow \mathbb{R}^{p_1}$  and  $\tilde{\nabla}^2q : \mathcal{X}^p \rightarrow \mathbb{R}^{p_1 \times p_2}$  ( $p_1, p_2 \geq 1$ ) are, respectively, called a pseudo derivative and a pseudo Hessian for  $q$  if for all  $\mathbf{y} \in \mathcal{X}^p$  and  $\rho \in \mathbb{R}^{p_1 \times 1}$ , we have*

$$\|\tilde{\nabla}q(\mathbf{y})\|_2 \geq \|\nabla q(\mathbf{y})\|_2 \quad \text{and} \quad \|\rho^\top \tilde{\nabla}^2q(\mathbf{y})\|_2 \geq \|\nabla[\rho^\top \tilde{\nabla}q(\mathbf{y})]\|_2. \quad (130)$$

Finally, we provide the second main result of this section with a proof in Appendix. F.2.

**Proposition F.2** (Tail bounds for arbitrary functions under LSI). *Given a random vector  $\mathbf{x}$  satisfying  $\text{LSI}_{\mathbf{x}}(\sigma^2)$ , any function  $q : \mathcal{X}^p \rightarrow \mathbb{R}$  with a pseudo derivative  $\tilde{\nabla}q$  and pseudo Hessian  $\tilde{\nabla}^2q$  (see Definition. F.3),  $\mathbf{x}$  satisfies a tail bound, namely for any fixed  $\varepsilon > 0$ , we have*

$$\mathbb{P}\left[|q_c(\mathbf{x})| \geq \varepsilon\right] \leq \exp\left(\frac{-c}{\sigma^4} \min\left(\frac{\varepsilon^2}{\mathbb{E}[\|\tilde{\nabla}q(\mathbf{x})\|_2]^2 + \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{F}}^2}, \frac{\varepsilon}{\max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}}}\right)\right),$$

where  $q_c(\mathbf{x}) = q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]$  and  $c$  is a universal constant.

### F.1 Proof of Proposition. F.1: Logarithmic Sobolev inequality

We start by defining the notion of  $W_2$  distance (Marton, 2015) which is useful in the proof. We note that  $W_2$  distance is a metric on the space of probability measures and satisfies triangle inequality.

**Definition F.4.** (Marton, 2015,  $W_2$  distance) *For random vectors  $\mathbf{x}$  and  $\mathbf{y}$  supported on  $\mathcal{X}^p$  with distributions  $f$  and  $g$ , respectively, the  $W_2$  distance is given by  $W_2^2(g_{\mathbf{y}}, f_{\mathbf{x}}) \triangleq \inf_{\pi} \sum_{t \in [p]} \left[\mathbb{P}_{\pi}(\mathbf{x}_t \neq \mathbf{y}_t)\right]^2$ , where the infimum is taken over all couplings  $\pi(\mathbf{x}, \mathbf{y})$  such that  $\pi(\mathbf{x}) = f(\mathbf{x})$  and  $\pi(\mathbf{y}) = g(\mathbf{y})$ .*

Given Definition. F.4, our next lemma states that if appropriate  $W_2$  distances are bounded, then the KL divergence (denoted by  $\text{KL}(\cdot \|\cdot)$ ) and the entropy approximately tensorize. We provide a proof in Appendix. F.1.1.

**Lemma F.1** (Approximate tensorization of KL divergence and entropy). *Given random vectors  $\mathbf{x}$  and  $\mathbf{y}$  supported on  $\mathcal{X}^p$  with distributions  $f$  and  $g$ , respectively, such that  $f_{\min} > 0$  (see (129)), if for all subsets  $S \subseteq [p]$  (with  $S^C \triangleq [p] \setminus S$ ) and all  $\mathbf{y}_{S^C} \in \mathcal{X}^{p-|S|}$ ,*

$$W_2^2(g_{\mathbf{y}_S | \mathbf{y}_{S^C} = \mathbf{y}_{S^C}}, f_{\mathbf{x}_S | \mathbf{x}_{S^C} = \mathbf{y}_{S^C}}) \leq C \sum_{t \in S} \mathbb{E} \left[ \|g_{\mathbf{y}_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} - f_{\mathbf{x}_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}\|_{\text{TV}}^2 \mid \mathbf{y}_{S^C} = \mathbf{y}_{S^C} \right], \quad (131)$$

almost surely for some constant  $C \geq 1$ , then

$$\text{KL}(g_{\mathbf{y}} \| f_{\mathbf{x}}) \leq \frac{2C}{f_{\min}} \sum_{t \in [p]} \mathbb{E}[\text{KL}(g_{\mathbf{y}_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} \| f_{\mathbf{x}_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}})], \quad \text{and} \quad (132)$$

$$\text{Ent}_{\mathbf{x}}(q) \leq \frac{2C}{f_{\min}} \sum_{t \in [p]} \mathbb{E}_{\mathbf{x}_{-t}}[\text{Ent}_{\mathbf{x}_t | \mathbf{x}_{-t}}(q)] \quad \text{for any function } q : \mathcal{X}^p \rightarrow \mathbb{R}_+. \quad (133)$$

Next, we claim that if the random vector  $\mathbf{x}$  satisfies Dobrushin's uniqueness condition, then the condition (131) of Lemma. F.1 is naturally satisfied. We provide a proof in Appendix. F.1.2.

**Lemma F.2** (Dobrushin’s uniqueness implies approximate tensorization). *Given random vectors  $\mathbf{x}$  and  $\mathbf{y}$  supported on  $\mathcal{X}^p$  with distributions  $f$  and  $g$ , respectively, if  $\mathbf{x}$  satisfies Dobrushin’s uniqueness condition (see Definition. F.2) with coupling matrix  $\Theta \in \mathbb{R}^{p \times p}$ , then for all subsets  $S \subseteq [p]$  (with  $S^C \triangleq [p] \setminus S$ ) and all  $\mathbf{y}_{S^C} \in \mathcal{X}^{p-|S|}$ ,*

$$W_2^2(g_{\mathbf{y}_S | \mathbf{y}_{S^C} = \mathbf{y}_{S^C}}, f_{\mathbf{x}_S | \mathbf{x}_{S^C} = \mathbf{y}_{S^C}}) \leq C \sum_{t \in S} \mathbb{E} \left[ \|g_{y_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} - f_{x_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}\|_{\text{TV}}^2 \Big| \mathbf{y}_{S^C} = \mathbf{y}_{S^C} \right], \quad (134)$$

almost surely where  $C = (1 - \|\Theta\|_{\text{op}})^2$ .

Now to prove Proposition. F.1, applying Lemmas. F.1 and F.2 for an arbitrary function  $f : \mathcal{X}^p \rightarrow \mathbb{R}$ , we find that

$$\begin{aligned} \text{Ent}_{\mathbf{x}}(q^2) &\leq \frac{2}{f_{\min}(1 - \|\Theta\|_{\text{op}})^2} \sum_{t \in [p]} \mathbb{E}_{\mathbf{x}_{-t}} \left[ \text{Ent}_{x_t | \mathbf{x}_{-t}}(q^2) \right] \\ &\stackrel{(a)}{\leq} \frac{2\sigma^2}{f_{\min}(1 - \|\Theta\|_{\text{op}})^2} \sum_{t \in [p]} \mathbb{E}_{\mathbf{x}_{-t}} \left[ \mathbb{E}_{x_t | \mathbf{x}_{-t}} \left[ \|\nabla_{x_t} q(x_t; \mathbf{x}_{-t})\|_2^2 \right] \right] \\ &\stackrel{(b)}{=} \frac{2\sigma^2}{f_{\min}(1 - \|\Theta\|_{\text{op}})^2} \mathbb{E}_{\mathbf{x}_{-t}} \left[ \mathbb{E}_{x_t | \mathbf{x}_{-t}} \left[ \sum_{t \in [p]} \|\nabla_{x_t} q(x_t; \mathbf{x}_{-t})\|_2^2 \right] \right] \\ &\stackrel{(c)}{=} \frac{2\sigma^2}{f_{\min}(1 - \|\Theta\|_{\text{op}})^2} \mathbb{E}_{\mathbf{x}} \left[ \|\nabla_{\mathbf{x}} q(\mathbf{x})\|_2^2 \right], \end{aligned}$$

where (a) follows because  $x_t | \mathbf{x}_{-t}$  satisfies  $\text{LSI}_{x_t | \mathbf{x}_{-t} = \mathbf{x}_{-t}}(\sigma^2)$  for all  $t \in [p]$  and  $\mathbf{x}_{-t} \in \mathcal{X}^{p-1}$ , (b) follows by the linearity of expectation and (c) follows by the law of total expectation. The claim follows.

### F.1.1 Proof of Lemma. F.1: Approximate tensorization of KL divergence and entropy

We start by establishing a reverse-Pinsker style inequality for distributions with compact support to bound their KL divergence by their total variation distance. We provide a proof at the end.

**Lemma F.3** (Reverse-Pinsker inequality). *For any distributions  $f$  and  $g$  supported on  $\mathcal{X} \subset \mathbb{R}$  such that  $\min_{x \in \mathcal{X}} f(x) > 0$ , we have  $\text{KL}(g \| f) \leq \frac{4}{\min_{x \in \mathcal{X}} f(x)} \|g - f\|_{\text{TV}}^2$ .*

Given Lemma. F.3, we proceed to prove Lemma. F.1.

**Proof of bound (132)** To prove (132), we show that the following inequality holds using the technique of mathematical induction on  $p$ :

$$\text{KL}(g_{\mathbf{y}} \| f_{\mathbf{x}}) \leq \frac{4C}{f_{\min}} \sum_{t \in [p]} \mathbb{E} \left[ \|g_{y_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} - f_{x_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}\|_{\text{TV}}^2 \right]. \quad (135)$$

Then, (132) follows by using Pinsker’s inequality to bound the right hand side of (135).

**Base case:  $p = 1$**  For the base case, we need to establish that the claim holds for all distributions supported on  $\mathcal{X}$  that satisfy the required conditions. In other words, we need to show that

$$\text{KL}(g_y \| f_x) \leq \frac{4C}{f_{\min}} \|g_y - f_x\|_{\text{TV}}^2 \quad \text{for every } t \in [p],$$

for all random variables  $x$  and  $y$  supported on  $\mathcal{X}$  such that  $f_{\min} = \min_{x \in \mathcal{X}} f_x(x) > 0$ . This follows from Lemma. F.3 by observing that  $C \geq 1$ .

**Inductive step** Now, we assume that the claim holds for all distributions supported on  $\mathcal{X}^{p-1}$  that satisfy the required conditions, and establish it for distributions supported on  $\mathcal{X}^p$ . From the chain rule of KL divergence, we have

$$\text{KL}(g_{\mathbf{y}} \| f_{\mathbf{x}}) = \text{KL}(g_{y_t} \| f_{x_t}) + \mathbb{E}[\text{KL}(g_{\mathbf{y}_{-t}|y_t} \| f_{\mathbf{x}_{-t}|x_t})] \quad \text{for every } t \in [p].$$

Taking an average over all  $t \in [p]$ , we have

$$\text{KL}(g_{\mathbf{y}} \| f_{\mathbf{x}}) = \frac{1}{p} \sum_{t \in [p]} \text{KL}(g_{y_t} \| f_{x_t}) + \frac{1}{p} \sum_{t \in [p]} \mathbb{E}[\text{KL}(g_{\mathbf{y}_{-t}|y_t} \| f_{\mathbf{x}_{-t}|x_t})]. \quad (136)$$

Now, we bound the first term in (136). Let  $\pi^*$  be the coupling between  $\mathbf{x}$  and  $\mathbf{y}$  that achieves  $W_2(g_{\mathbf{y}}, f_{\mathbf{x}})$  i.e.,<sup>10</sup>

$$\pi^* = \arg \min_{\pi: \pi(\mathbf{x})=f(\mathbf{x}), \pi(\mathbf{y})=g(\mathbf{y})} \sum_{t \in [p]} \left[ \mathbb{P}_{\pi}(x_t \neq y_t) \right]^2. \quad (137)$$

Then, we have

$$\begin{aligned} \frac{1}{p} \sum_{t \in [p]} \text{KL}(g_{y_t} \| f_{x_t}) &\stackrel{(a)}{\leq} \frac{1}{p} \sum_{t \in [p]} \frac{4}{f_{\min}} \|g_{y_t} - f_{x_t}\|_{\text{TV}}^2 \\ &\stackrel{(b)}{\leq} \frac{4}{pf_{\min}} \sum_{t \in [p]} \left[ \mathbb{P}_{\pi^*}(x_t \neq y_t) \right]^2 \\ &\stackrel{(c)}{=} \frac{4}{pf_{\min}} W_2^2(g_{\mathbf{y}}, f_{\mathbf{x}}) \\ &\stackrel{(131)}{\leq} \frac{4C}{pf_{\min}} \sum_{t \in [p]} \mathbb{E} \left[ \|g_{y_t|y_{-t}=\mathbf{y}_{-t}} - f_{x_t|x_{-t}=\mathbf{y}_{-t}}\|_{\text{TV}}^2 \right], \end{aligned} \quad (138)$$

where (a) follows from Lemma. F.3 because lower bound on conditional implies lower bound on marginals, i.e.,  $\min_{t \in [p], x_t \in \mathcal{X}} f_{x_t}(x_t) = \min_{t \in [p], x_t \in \mathcal{X}} \int_{\mathbf{x}_{-t} \in \mathcal{X}^{p-1}} f_{x_t|x_{-t}}(x_t|\mathbf{x}_{-t}) f_{\mathbf{x}_{-t}}(\mathbf{x}_{-t}) d\mathbf{x}_{-t} > f_{\min}$ , (b) follows from the connections of total variation distance to optimal transportation cost, i.e.,  $\|g_{\mathbf{y}} - f_{\mathbf{x}}\|_{\text{TV}} = \inf_{\pi: \pi(x)=f(x), \pi(y)=g(y)} \mathbb{P}_{\pi}(x \neq y)$ , and (c) follows from Definition. F.4 and (137).

Next, we bound the second term in (136). We have

$$\begin{aligned} &\frac{1}{p} \sum_{t \in [p]} \mathbb{E}[\text{KL}(g_{\mathbf{y}_{-t}|y_t} \| f_{\mathbf{x}_{-t}|x_t})] \\ &\leq \frac{1}{p} \sum_{t \in [p]} \mathbb{E} \left[ \frac{4C}{f_{\min}} \sum_{u \in [p] \setminus \{t\}} \mathbb{E} \left[ \|g_{y_u|y_{-u}=\mathbf{y}_{-u}} - f_{x_u|x_{-u}=\mathbf{y}_{-u}}\|_{\text{TV}}^2 \middle| y_t = y_t \right] \right] \\ &\stackrel{(b)}{=} \frac{4C}{pf_{\min}} \sum_{t \in [p]} \sum_{u \in [p] \setminus \{t\}} \mathbb{E} \left[ \|g_{y_u|y_{-u}=\mathbf{y}_{-u}} - f_{x_u|x_{-u}=\mathbf{y}_{-u}}\|_{\text{TV}}^2 \right] \end{aligned}$$

<sup>10</sup>The minimum is achieved by using arguments similar to the ones used to show that the Wasserstein distance attains its minimum (Villani, 2009, Chapter 4).

$$= \frac{4C(p-1)}{pf_{\min}} \sum_{u \in [p]} \mathbb{E} \left[ \|g_{y_u | \mathbf{y}_{-u} = \mathbf{y}_{-u}} - f_{x_u | \mathbf{x}_{-u} = \mathbf{y}_{-u}}\|_{\text{TV}}^2 \right], \quad (139)$$

where (a) follows from the inductive hypothesis and (b) follows from the law of total expectation. Then, (135) follows by putting (136), (138), and (139) together.

**Proof of bound (133)** To prove (133), we note that (132) holds for any random vector  $\mathbf{y}$  supported on  $\mathcal{X}^p$ . Consider  $\mathbf{y}$  to be such that  $q(\mathbf{x})/\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]$  is the Radon-Nikodym derivative of  $g_{\mathbf{y}}$  with respect to  $f_{\mathbf{x}}$ . For any  $\mathcal{A}^p \subseteq \mathcal{X}^p$ , we have

$$\int_{\mathbf{y} \in \mathcal{A}^p} g_{\mathbf{y}} d\mathbf{y} = \int_{\mathbf{x} \in \mathcal{A}^p} \frac{q(\mathbf{x})}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]} f_{\mathbf{x}} d\mathbf{x}.$$

Integrating out  $y_t$  and  $x_t$  for  $t \in [p]$ , we have

$$\int_{\mathbf{y}_{-t} \in \mathcal{A}^{p-1}} g_{\mathbf{y}_{-t}} d\mathbf{y}_{-t} = \int_{\mathbf{x}_{-t} \in \mathcal{A}^{p-1}} \frac{\mathbb{E}_{x_t | \mathbf{x}_{-t}}[q(\mathbf{x})]}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]} f_{\mathbf{x}_{-t}} d\mathbf{x}_{-t},$$

implying

$$\frac{dg_{\mathbf{y}_{-t}}}{df_{\mathbf{x}_{-t}}} = \frac{\mathbb{E}_{x_t | \mathbf{x}_{-t}}[q(\mathbf{x})]}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]} \quad \text{and} \quad \frac{dg_{y_t | \mathbf{y}_{-t}}}{df_{x_t | \mathbf{x}_{-t}}} = \frac{q(\mathbf{x})}{\mathbb{E}_{x_t | \mathbf{x}_{-t}}[q(\mathbf{x})]} \quad \text{for all } t \in [p]. \quad (140)$$

We have

$$\begin{aligned} \text{KL}(g_{\mathbf{y}} \| f_{\mathbf{x}}) &\stackrel{(a)}{=} \mathbb{E}_{\mathbf{x}} \left[ \frac{dg_{\mathbf{y}}}{df_{\mathbf{x}}} \log \frac{dg_{\mathbf{y}}}{df_{\mathbf{x}}} \right] \\ &\stackrel{(b)}{=} \mathbb{E}_{\mathbf{x}} \left[ \frac{q(\mathbf{x})}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]} \log \frac{q(\mathbf{x})}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]} \right] \\ &= \frac{1}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]} \left( \mathbb{E}_{\mathbf{x}}[q(\mathbf{x}) \log q(\mathbf{x})] - \mathbb{E}_{\mathbf{x}}[q(\mathbf{x})] \log \mathbb{E}_{\mathbf{x}}[q(\mathbf{x})] \right) = \frac{\text{Ent}_{\mathbf{x}}(q)}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]}, \end{aligned} \quad (141)$$

where (a) follows from the definition of KL divergence and (b) follows from the choice of  $\mathbf{y}$ . Similarly, for every  $t \in [p]$ , we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{y}_{-t}} \left[ \text{KL}(g_{y_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} \| f_{x_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}) \right] \\ &\stackrel{(a)}{=} \mathbb{E}_{\mathbf{y}_{-t}} \left[ \mathbb{E}_{y_t | \mathbf{y}_{-t}} \left[ \log \frac{dg_{y_t | \mathbf{y}_{-t}}}{df_{x_t | \mathbf{x}_{-t}}} \right] \right] \\ &\stackrel{(b)}{=} \mathbb{E}_{\mathbf{y}} \left[ \log \frac{dg_{y_t | \mathbf{y}_{-t}}}{df_{x_t | \mathbf{x}_{-t}}} \right] \\ &\stackrel{(c)}{=} \mathbb{E}_{\mathbf{x}} \left[ \frac{dg_{\mathbf{y}}}{df_{\mathbf{x}}} \log \frac{dg_{y_t | \mathbf{y}_{-t}}}{df_{x_t | \mathbf{x}_{-t}}} \right] \\ &\stackrel{(d)}{=} \mathbb{E}_{\mathbf{x}} \left[ \frac{q(\mathbf{x})}{\mathbb{E}_{\mathbf{x}}[q(\mathbf{x})]} \log \frac{q(\mathbf{x})}{\mathbb{E}_{x_t | \mathbf{x}_{-t}}[q(\mathbf{x})]} \right] \end{aligned}$$

$$\begin{aligned}
& \underline{(e)} \frac{\mathbb{E}_{\mathbf{x}_{-t}} \left[ \mathbb{E}_{x_t | \mathbf{x}_{-t}} [q(\mathbf{x}) \log q(\mathbf{x})] - \mathbb{E}_{x_t | \mathbf{x}_{-t}} [q(\mathbf{x}) \log \mathbb{E}_{x_t | \mathbf{x}_{-t}} [q(\mathbf{x})]] \right]}{\mathbb{E}_{\mathbf{x}} [q(\mathbf{x})]} \\
& \underline{(f)} \frac{\mathbb{E}_{\mathbf{x}_{-t}} [\text{Ent}_{x_t | \mathbf{x}_{-t}}(q)]}{\mathbb{E}[q(\mathbf{x})]}, \tag{142}
\end{aligned}$$

where (a) follows from the definition of KL divergence, (b) follows from the law of total expectation, (c) follows from the definition of Radon-Nikodym derivative, (d) follows from the choice of  $\mathbf{y}$  and (140), (e) follows from the law of total expectation, (f) follows from the definition of entropy. Then, (133) follows by putting (132), (141), and (142) together.

**Proof of Lemma. F.3: Reverse-Pinsker inequality** Using the facts (a)  $\log a \geq 1 - \frac{1}{a}$  for all  $a > 0$ , and (b)  $\min_{x \in \mathcal{X}} f(x) > 0$ , we find that

$$\log \frac{f(x)}{g(x)} \geq 1 - \frac{g(x)}{f(x)} \quad \text{for every } x \in \mathcal{X}. \tag{143}$$

Multiplying both sides of (143) by  $g(x) \geq 0$  and rearranging terms yields that

$$g(x) \log \frac{g(x)}{f(x)} \leq \frac{g^2(x)}{f(x)} - g(x) \quad \text{for every } x \in \mathcal{X}. \tag{144}$$

Now, we have

$$\begin{aligned}
\text{KL}(g \| f) &= \int_{x \in \mathcal{X}} g(x) \log \frac{g(x)}{f(x)} dx \stackrel{(144)}{\leq} \int_{x \in \mathcal{X}} \left( \frac{g^2(x)}{f(x)} - g(x) \right) dx \\
&\stackrel{(a)}{=} \int_{x \in \mathcal{X}} \frac{(g(x) - f(x))^2}{f(x)} dx \\
&\leq \frac{1}{\min_{x \in \mathcal{X}} f(x)} \int_{x \in \mathcal{X}} (g(x) - f(x))^2 dx \\
&\stackrel{(b)}{\leq} \frac{1}{\min_{x \in \mathcal{X}} f(x)} \left( \int_{x \in \mathcal{X}} |g(x) - f(x)| dx \right)^2 \\
&\stackrel{(c)}{=} \frac{1}{\min_{x \in \mathcal{X}} f(x)} \left( 2 \|g - f\|_{\text{TV}} \right)^2 \\
&= \frac{4}{\min_{x \in \mathcal{X}} f(x)} \|g - f\|_{\text{TV}}^2,
\end{aligned}$$

where (a) follows by simple manipulations, (b) follows by using the order of norms on Euclidean space, and (c) follows by the definition of the total variation distance.

### F.1.2 Proof of Lemma. F.2: Dobrushin's uniqueness implies approximate tensorization

We start by defining the notion of Gibbs sampler which is useful in the proof.

**Definition F.5.** (*Marton, 2015, Gibbs Sampler*) For a random vector  $\mathbf{x}$  with distribution  $f$ , define the Markov kernels and the Gibbs sampler as follows

$$\Gamma_t(\mathbf{x} | \mathbf{x}') \triangleq \mathbf{1}(\mathbf{x}_{-t} = \mathbf{x}'_{-t}) f_{x_t | \mathbf{x}_{-t}}(x_t | \mathbf{x}'_{-t}) \quad \text{and} \quad \Gamma(\mathbf{x} | \mathbf{x}') \triangleq p^{-1} \sum_{t \in [p]} \Gamma_t(\mathbf{x} | \mathbf{x}'), \tag{145}$$

for all  $t \in [p]$  and  $x, x' \in \mathcal{X}^p$ . That is, the kernel  $\Gamma_t$  leaves all but the  $t^{\text{th}}$  coordinate unchanged, and updates the  $t^{\text{th}}$  coordinate according to  $f_{x_t | \mathbf{x}_{-t}}$ , and the sampler  $\Gamma$  selects an index  $t \in [p]$  at random, and applies  $\Gamma_t$ . Further, for a random vector  $\mathbf{y}$  with distribution  $g$  supported on  $\mathcal{X}^p$ , we also define

$$\begin{aligned} g_{\mathbf{y}}\Gamma_t(\mathbf{y}) &\triangleq \int g_{\mathbf{y}}(\mathbf{y}')\Gamma_t(\mathbf{y}|\mathbf{y}')d\mathbf{y}' \text{ for } t \in [p], \text{ and} \\ g_{\mathbf{y}}\Gamma(\mathbf{y}) &\triangleq \int g_{\mathbf{y}}(\mathbf{y}')\Gamma(\mathbf{y}|\mathbf{y}')d\mathbf{y}' \text{ for all } \mathbf{y} \in \mathcal{X}^p. \end{aligned} \quad (146)$$

We now proceed to prove Lemma. F.2 and split it in two cases: (i)  $S = [p]$ , and (ii)  $S \subset [p]$ .

**Case (i)** ( $S = [p]$ ) Let  $\Gamma$  be the Gibbs sampler associated with the distribution  $f$ . Then,

$$W_2(g_{\mathbf{y}_S | \mathbf{y}_{S^c}}, f_{\mathbf{x}_S | \mathbf{x}_{S^c}}) = W_2(g_{\mathbf{y}}, f_{\mathbf{x}}) \stackrel{(a)}{\leq} W_2(g_{\mathbf{y}}, g_{\mathbf{y}}\Gamma) + W_2(g_{\mathbf{y}}\Gamma, f_{\mathbf{x}}), \quad (147)$$

where (a) follows from the triangle inequality. We claim that

$$W_2(g_{\mathbf{y}}, g_{\mathbf{y}}\Gamma) \leq \frac{1}{p} \sqrt{\sum_{t \in [p]} \mathbb{E}_{\mathbf{y}_{-t}} \left[ \|g_{y_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} - f_{x_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}\|_{\text{TV}}^2 \right]}, \quad \text{and} \quad (148)$$

$$W_2(g_{\mathbf{y}}\Gamma, f_{\mathbf{x}}) \leq \left( 1 - \frac{(1 - \|\Theta\|_{\text{op}})}{p} \right) W_2(g_{\mathbf{y}}, f_{\mathbf{x}}). \quad (149)$$

Putting (147) to (149) together, we have

$$\begin{aligned} W_2(g_{\mathbf{y}}, f_{\mathbf{x}}) &\leq \frac{1}{p} \sqrt{\sum_{t \in [p]} \mathbb{E}_{\mathbf{y}_{-t}} \left[ \|g_{y_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} - f_{x_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}\|_{\text{TV}}^2 \right]} \\ &\quad + \left( 1 - \frac{(1 - \|\Theta\|_{\text{op}})}{p} \right) W_2(g_{\mathbf{y}}, f_{\mathbf{x}}). \end{aligned} \quad (150)$$

Rearranging (150) results in (134) for  $S = [p]$  as desired. It remains to prove our earlier claims (148) and (149) which we now do one-by-one.

**Proof of bound (148) on  $W_2(g_{\mathbf{y}}, g_{\mathbf{y}}\Gamma)$**  To bound  $W_2(g_{\mathbf{y}}, g_{\mathbf{y}}\Gamma)$ , we construct a random vector  $\mathbf{y}^\Gamma$  such that it is coupled with the random vector  $\mathbf{y}$ . We select an index  $b \in [p]$  at random, and define

$$y_v^\Gamma \triangleq y_v \quad \text{for all } v \in [p] \setminus \{b\}.$$

Then, given  $b$  and  $\mathbf{y}_{-b} = \mathbf{y}_{-b}$ , we define the joint distribution of  $(y_b, y_b^\Gamma)$  to be the maximal coupling of  $g_{y_b | \mathbf{y}_{-b} = \mathbf{y}_{-b}}$  and  $f_{x_b | \mathbf{x}_{-b} = \mathbf{y}_{-b}}$  that achieves  $\|g_{y_b | \mathbf{y}_{-b} = \mathbf{y}_{-b}} - f_{x_b | \mathbf{x}_{-b} = \mathbf{y}_{-b}}\|_{\text{TV}}$ . It is easy to see that the marginal distribution of  $\mathbf{y}$  is  $g_{\mathbf{y}}$  and the marginal distribution of  $\mathbf{y}^\Gamma$  is  $g_{\mathbf{y}}\Gamma$  (see Definition. F.5). Then, we have

$$\begin{aligned} W_2^2(g_{\mathbf{y}}, g_{\mathbf{y}}\Gamma) &\stackrel{(a)}{\leq} \sum_{t \in [p]} \left[ \mathbb{P}(b = t) \mathbb{P}(y_t \neq y_t^\Gamma | b = t) + \mathbb{P}(b \neq t) \mathbb{P}(y_t \neq y_t^\Gamma | b \neq t) \right]^2 \\ &\stackrel{(b)}{=} \sum_{t \in [p]} \left[ \frac{1}{p} \mathbb{P}(y_t \neq y_t^\Gamma | b = t) \right]^2 \end{aligned}$$

$$\begin{aligned}
&\stackrel{(c)}{=} \frac{1}{p^2} \sum_{t \in [p]} \left[ \int_{\mathbf{y}_{-t} \in \mathcal{X}^{p-1}} \mathbb{P}(y_t \neq y_t^\Gamma | b = t, \mathbf{y}_{-t} = \mathbf{y}_{-t}) g_{\mathbf{y}_{-t}|b=t}(\mathbf{y}_{-t}|b = t) d\mathbf{y}_{-t} \right]^2 \\
&\stackrel{(d)}{=} \frac{1}{p^2} \sum_{t \in [p]} \left[ \int_{\mathbf{y}_{-t} \in \mathcal{X}^{p-1}} \|g_{y_t|\mathbf{y}_{-t}=\mathbf{y}_{-t}} - f_{x_t|\mathbf{x}_{-t}=\mathbf{y}_{-t}}\|_{\text{TV}} g_{\mathbf{y}_{-t}}(\mathbf{y}_{-t}) d\mathbf{y}_{-t} \right]^2 \\
&= \frac{1}{p^2} \sum_{t \in [p]} \left[ \mathbb{E}_{\mathbf{y}_{-t}} \left[ \|g_{y_t|\mathbf{y}_{-t}=\mathbf{y}_{-t}} - f_{x_t|\mathbf{x}_{-t}=\mathbf{y}_{-t}}\|_{\text{TV}} \right] \right]^2, \tag{151}
\end{aligned}$$

where (a) follows from Definition. F.4 and the Bayes rule, (b) follows because  $\mathbb{P}(b = t) = \frac{1}{p}$  and  $\mathbb{P}(y_t \neq y_t^\Gamma | b \neq t) = 0$ , (c) follows by the law of total probability, and (d) follows because  $g_{\mathbf{y}_{-t}|b=t}(\mathbf{y}_{-t}|b = t) = g_{\mathbf{y}_{-t}}(\mathbf{y}_{-t})$  and by the construction of the coupling between  $\mathbf{y}$  and  $\mathbf{y}^\Gamma$ . Then, (148) follows by using Jensen's inequality in (151).

**Proof of bound (149) on  $W_2(g_{\mathbf{y}}\Gamma, f_{\mathbf{x}})$**  We first show that  $f_{\mathbf{x}}$  is an invariant measure for  $\Gamma$ , i.e.,  $f_{\mathbf{x}} = f_{\mathbf{x}}\Gamma$ , implying  $W_2(g_{\mathbf{y}}\Gamma, f_{\mathbf{x}}) = W_2(g_{\mathbf{y}}\Gamma, f_{\mathbf{x}}\Gamma)$ , and then  $\Gamma$  is a contraction with respect to the  $W_2$  distance with rate  $1 - \frac{(1-\|\Theta\|_{\text{op}})}{p}$ , i.e.,  $W_2(g_{\mathbf{y}}\Gamma, f_{\mathbf{x}}\Gamma) \leq \left(1 - \frac{(1-\|\Theta\|_{\text{op}})}{p}\right) W_2(g_{\mathbf{y}}, f_{\mathbf{x}})$ , implying (149).

**Proof of  $f_{\mathbf{x}}$  being an invariant measure for  $\Gamma$**  We have

$$\begin{aligned}
f_{\mathbf{x}}\Gamma(\mathbf{x}) &\stackrel{(146)}{=} \int_{\mathbf{x}' \in \mathcal{X}^p} f_{\mathbf{x}}(\mathbf{x}') \Gamma(\mathbf{x}|\mathbf{x}') d\mathbf{x}' \\
&\stackrel{(145)}{=} \int_{\mathbf{x}' \in \mathcal{X}^p} f_{\mathbf{x}}(\mathbf{x}') \left( \frac{1}{p} \sum_{t \in [p]} \Gamma_t(\mathbf{x}|\mathbf{x}') \right) d\mathbf{x}' \\
&\stackrel{(145)}{=} \frac{1}{p} \sum_{t \in [p]} \int_{\mathbf{x}' \in \mathcal{X}^p} f_{\mathbf{x}}(\mathbf{x}') \mathbb{1}(\mathbf{x}_{-t} = \mathbf{x}'_{-t}) f_{x_t|\mathbf{x}_{-t}}(x_t|\mathbf{x}'_{-t}) d\mathbf{x}' \\
&= \frac{1}{p} \sum_{t \in [p]} f_{x_t|\mathbf{x}_{-t}}(x_t|\mathbf{x}_{-t}) \int_{x'_t \in \mathcal{X}} f_{\mathbf{x}}(\mathbf{x}_{-t}, x'_t) dx'_t \\
&= \frac{1}{p} \sum_{t \in [p]} f_{x_t|\mathbf{x}_{-t}}(x_t|\mathbf{x}_{-t}) f_{\mathbf{x}_{-t}}(\mathbf{x}_{-t}) = f_{\mathbf{x}}(\mathbf{x}).
\end{aligned}$$

**Proof of  $\Gamma$  being a contraction w.r.t the  $W_2$  distance** Let  $\pi^*$  be the coupling between  $\mathbf{x}$  and  $\mathbf{y}$  that achieves  $W_2(g_{\mathbf{y}}, f_{\mathbf{x}})$  i.e.,<sup>11</sup>

$$\pi^* = \arg \min_{\pi: \pi(\mathbf{x})=f(\mathbf{x}), \pi(\mathbf{y})=g(\mathbf{y})} \sqrt{\sum_{t \in [p]} \left[ \mathbb{P}_{\pi}(x_t \neq y_t) \right]^2}. \tag{152}$$

We construct random variables  $\mathbf{x}'$  and  $\mathbf{y}'$  as well as a coupling  $\pi'$  between them such that the marginal distribution of  $\mathbf{x}'$  is  $f_{\mathbf{x}}\Gamma$  and the marginal distribution of  $\mathbf{y}'$  is  $g_{\mathbf{y}}\Gamma$ . We start by selecting an index

<sup>11</sup>The minimum is achieved by using arguments similar to the ones used to show that the Wasserstein distance attains its minimum (Villani, 2009, Chapter 4).

$b \in [p]$  at random, and defining

$$y'_v \triangleq y_v \quad \text{and} \quad x'_v \triangleq x_v \quad \text{for all } v \neq b. \quad (153)$$

Then, given  $b$ ,  $\mathbf{y}'_{-b} = \mathbf{y}_{-b}$ , and  $\mathbf{x}'_{-b} = \mathbf{x}_{-b}$ , we define the joint distribution of  $(y'_b, x'_b)$  to be the maximal coupling of  $f_{x_b|\mathbf{x}_{-b}}(\cdot|\mathbf{y}_{-b})$  and  $f_{x_b|\mathbf{x}_{-b}}(\cdot|\mathbf{x}_{-b})$  that achieves  $\|f_{x_b|\mathbf{x}_{-b}=\mathbf{y}_{-b}} - f_{x_b|\mathbf{x}_{-b}=\mathbf{x}_{-b}}\|_{\text{TV}}$ .

Now, for every  $t \in [p]$ , we bound  $\mathbb{P}_{\pi'}(y'_t \neq x'_t)$  in terms of  $\mathbb{P}_{\pi^*}(y_t \neq x_t)$ . To that end, we have

$$\begin{aligned} \mathbb{P}_{\pi'}(y'_t \neq x'_t) &\stackrel{(a)}{=} \mathbb{P}(b=t)\mathbb{P}_{\pi'}(y'_t \neq x'_t|b=t) + \mathbb{P}(b \neq t)\mathbb{P}_{\pi'}(y'_t \neq x'_t|b \neq t) \\ &\stackrel{(b)}{=} \frac{1}{p}\mathbb{P}_{\pi'}(y'_t \neq x'_t|b=t) + \left(1 - \frac{1}{p}\right)\mathbb{P}_{\pi^*}(y_t \neq x_t), \end{aligned} \quad (154)$$

where (a) follows from the Bayes rule and (b) follows because  $\mathbb{P}(b=t) = \frac{1}{p}$  and (153). Focusing on  $\mathbb{P}_{\pi'}(y'_t \neq x'_t|b=t)$  and using the law of total probability, we have

$$\begin{aligned} &\mathbb{P}_{\pi'}(y'_t \neq x'_t|b=t) \\ &= \int_{\mathbf{y}_{-t}, \mathbf{x}_{-t} \in \mathcal{X}^{p-1}} \mathbb{P}_{\pi'}(y'_t \neq x'_t|b=t, \mathbf{y}'_{-t} = \mathbf{y}_{-t}, \mathbf{x}'_{-t} = \mathbf{x}_{-t}) \pi'_{\mathbf{y}'_{-t}, \mathbf{x}'_{-t}|b=t}(\mathbf{y}_{-t}, \mathbf{x}_{-t}|b=t) d\mathbf{y}_{-t} d\mathbf{x}_{-t} \\ &\stackrel{(a)}{=} \int_{\mathbf{y}_{-t}, \mathbf{x}_{-t} \in \mathcal{X}^{p-1}} \|f_{x_t|\mathbf{x}_{-t}=\mathbf{y}_{-t}} - f_{x_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}}\|_{\text{TV}} \pi_{\mathbf{y}_{-t}, \mathbf{x}_{-t}}^*(\mathbf{y}_{-t}, \mathbf{x}_{-t}) d\mathbf{y}_{-t} d\mathbf{x}_{-t} \\ &= \mathbb{E}_{\pi_{\mathbf{y}_{-t}, \mathbf{x}_{-t}}^*} \left[ \|f_{x_t|\mathbf{x}_{-t}=\mathbf{y}_{-t}} - f_{x_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}}\|_{\text{TV}} \right] \end{aligned} \quad (155)$$

where (a) follows by the construction of the coupling between  $\mathbf{y}'$  and  $\mathbf{x}'$ . Now, using the triangle inequality in (155), we have

$$\begin{aligned} \mathbb{P}_{\pi'}(y'_t \neq x'_t|b=t) &\leq \mathbb{E}_{\pi_{\mathbf{y}_{-t}, \mathbf{x}_{-t}}^*} \left[ \sum_{u \in [p] \setminus \{t\}} \mathbf{1}(r_v = s_v = y_v \forall v < u) \mathbf{1}(r_v = s_v = x_v \forall v > u) \times \right. \\ &\quad \left. \mathbf{1}(r_u = y_u, x_u = s_u) \|f_{x_t|\mathbf{x}_{-t}=r_{-t}} - f_{x_t|\mathbf{x}_{-t}=s_{-t}}\|_{\text{TV}} \right] \\ &\stackrel{(128)}{\leq} \mathbb{E}_{\pi_{\mathbf{y}_{-t}, \mathbf{x}_{-t}}^*} \left[ \sum_{u \in [p] \setminus \{t\}} \Theta_{tu} \mathbf{1}(y_u \neq x_u) \right] = \sum_{u \in [p] \setminus \{t\}} \Theta_{tu} \mathbb{P}_{\pi^*}(y_u \neq x_u). \end{aligned} \quad (156)$$

Putting together (154) and (156), we have

$$\mathbb{P}_{\pi'}(y'_t \neq x'_t) \leq \frac{1}{p} \sum_{u \in [p] \setminus \{t\}} \Theta_{tu} \mathbb{P}_{\pi^*}(y_u \neq x_u) + \left(1 - \frac{1}{p}\right) \mathbb{P}_{\pi^*}(y_t \neq x_t). \quad (157)$$

Next, we use (157) to show contraction of  $\Gamma$ . To that end, we define  $\text{diag}(\Theta) \in \mathbb{R}^{p \times p}$  to be the matrix with diagonal same as  $\Theta$  and all non-diagonal entries equal to zeros. Then, we have

$$W_2^2(g_{\mathbf{y}}\Gamma, f_{\mathbf{x}}\Gamma) \stackrel{(a)}{\leq} \sum_{t \in [p]} \left[ \mathbb{P}_{\pi'}(y'_t \neq x'_t) \right]^2$$

$$\begin{aligned}
&\stackrel{(157)}{\leq} \sum_{t \in [p]} \left[ \frac{1}{p} \sum_{u \in [p] \setminus \{t\}} \Theta_{tu} \mathbb{P}_{\pi^*}(y_u \neq x_u) + \left(1 - \frac{1}{p}\right) \mathbb{P}_{\pi^*}(y_t \neq x_t) \right]^2 \\
&\stackrel{(b)}{\leq} \left\| \left(1 - \frac{1}{p}\right) I + \frac{1}{p} (\Theta - \text{diag}(\Theta)) \right\|_{\text{op}}^2 \sum_{t \in [p]} \left[ \mathbb{P}_{\pi^*}(y_t \neq x_t) \right]^2 \\
&\stackrel{(c)}{=} \left\| \left(1 - \frac{1}{p}\right) I + \frac{1}{p} (\Theta - \text{diag}(\Theta)) \right\|_{\text{op}}^2 W_2^2(g_{\mathbf{y}}, f_{\mathbf{x}}) \\
&\stackrel{(d)}{\leq} \left( \left(1 - \frac{1}{p}\right) + \frac{1}{p} \|\Theta - \text{diag}(\Theta)\|_{\text{op}} \right)^2 W_2^2(g_{\mathbf{y}}, f_{\mathbf{x}}) \\
&\stackrel{(e)}{\leq} \left( \left(1 - \frac{1}{p}\right) + \frac{1}{p} \|\Theta\|_{\text{op}} \right)^2 W_2^2(g_{\mathbf{y}}, f_{\mathbf{x}}), \tag{158}
\end{aligned}$$

where (a) follows from Definition. F.4, (b) follows by some linear algebraic manipulations, (c) follows from Definition. F.4 and (152), (d) follows from the triangle inequality, and (e) follows because  $\|\mathbf{M}_1\|_{\text{op}} \leq \|\mathbf{M}_2\|_{\text{op}}$  for any matrices  $\mathbf{M}_1$  and  $\mathbf{M}_2$  such that  $0 \leq \mathbf{M}_1 \leq \mathbf{M}_2$  (component-wise). Then, contraction of  $\Gamma$  follows by taking square root on both sides of (158).

**Case (ii)** ( $S \subset [p]$ ) We can directly verify that the matrix  $\Theta_S \triangleq \{\Theta_{tu}\}_{t,u \in S}$  is such that  $\|\Theta_S\|_{\text{op}} \leq \|\Theta\|_{\text{op}}$ . This is true because the operator norm of any sub-matrix is no more than the operator norm of the matrix. Further, we note that for any  $\mathbf{y}_{S^c} \in \mathcal{X}^{p-|S|}$ , the random vector  $\mathbf{x}_S | \mathbf{x}_{S^c} = \mathbf{y}_{S^c}$  with distribution  $f_{\mathbf{x}_S | \mathbf{x}_{S^c} = \mathbf{y}_{S^c}}$  satisfies the Dobrushin's uniqueness condition (Definition. F.2) with coupling matrix  $\Theta_S$ . Then, by performing an analysis similar to the one above, we have

$$\begin{aligned}
W_2(g_{\mathbf{y}_S | \mathbf{y}_{S^c}}, f_{\mathbf{x}_S | \mathbf{x}_{S^c}}) &\leq \frac{1}{(1 - \|\Theta_S\|_{\text{op}})} \sqrt{\sum_{t \in S} \mathbb{E} \left[ \|g_{y_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} - f_{x_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}\|_{\text{TV}}^2 \mid \mathbf{y}_{S^c} = \mathbf{y}_{S^c} \right]} \\
&\stackrel{(a)}{\leq} \frac{1}{(1 - \|\Theta\|_{\text{op}})} \sqrt{\sum_{t \in S} \mathbb{E} \left[ \|g_{y_t | \mathbf{y}_{-t} = \mathbf{y}_{-t}} - f_{x_t | \mathbf{x}_{-t} = \mathbf{y}_{-t}}\|_{\text{TV}}^2 \mid \mathbf{y}_{S^c} = \mathbf{y}_{S^c} \right]},
\end{aligned}$$

where (a) follows because  $\frac{1}{(1 - \|\Theta_S\|_{\text{op}})} \leq \frac{1}{(1 - \|\Theta\|_{\text{op}})}$ . This completes the proof.

## F.2 Proof of Proposition. F.2: Tail bounds for arbitrary functions under LSI

Fix a function  $q : \mathcal{X}^p \rightarrow \mathbb{R}$ . Fix any pseudo derivative  $\tilde{\nabla}q$  for  $q$  and any pseudo Hessian  $\tilde{\nabla}^2q$  for  $q$ . To prove Proposition. F.2, we bound the  $p$ -th moment of  $q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]$  by certain norms of  $\tilde{\nabla}^2q$  and  $\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]$ . To that end, first, we claim that in order to control the  $p$ -th moment of  $q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]$ , it is sufficient to control the  $p$ -th moment of  $\|\nabla q(\mathbf{x})\|_2$ . Then, using (130), we note that the  $p$ -th moment of  $\|\nabla q(\mathbf{x})\|_2$  is bounded by the  $p$ -th moment of  $\|\tilde{\nabla}q(\mathbf{x})\|_2$ . Next, we claim that the  $p$ -th moment of  $\|\tilde{\nabla}q(\mathbf{x})\|_2$  is bounded by a linear combination of appropriate norms of  $\tilde{\nabla}^2q$  and  $\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]$ . We formalize the claims below and divide the proof across Appendix. F.2.1 and Appendix. F.2.2.

**Lemma F.4** (Bounded  $p$ -th moments of  $q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]$  and  $\|\tilde{\nabla}q(\mathbf{x})\|_2$ ). *If a random vector  $\mathbf{x}$  satisfies LSI $_{\mathbf{x}}(\sigma^2)$ , then for any arbitrary function  $q : \mathcal{X}^p \rightarrow \mathbb{R}$ ,*

$$\|q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]\|_{L_p} \leq \sigma \sqrt{2p} \|\nabla q(\mathbf{x})\|_2 \|L_p \quad \text{for any } p \geq 2. \tag{159}$$

Further, for any pseudo derivative  $\tilde{\nabla}q(\mathbf{x})$  and any pseudo Hessian  $\tilde{\nabla}^2q(\mathbf{x})$  for  $q$ , and even  $p \geq 2$ ,

$$\|\|\tilde{\nabla}q(\mathbf{x})\|_2\|_{L_p} \leq 2c\sigma \left( \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_F + \sqrt{p} \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}} \right) + 4\|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2, \quad (160)$$

where  $c \geq 0$  is a universal constant.

Given these lemmas, we proceed to prove Proposition. F.2. We let  $q_c(\mathbf{x}) = q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]$ . Combining (159) and (160) for any even  $p \geq 2$ , there exists a universal constant  $c'$  such that

$$\|q_c(\mathbf{x})\|_{L_p} \leq c'\sigma^2 \left( \sqrt{p} \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_F + p \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}} + \sqrt{p} \|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2 \right). \quad (161)$$

Now, we complete the proof by using (161) along with Markov's inequality for a specific choice of  $p$ . For any even  $p \geq 2$ , we have

$$\begin{aligned} & \mathbb{P} \left[ |q_c(\mathbf{x})| > ec'\sigma^2 \left( \sqrt{p} \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_F + p \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}} + \sqrt{p} \|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2 \right) \right] \\ &= \mathbb{P} \left[ |q_c(\mathbf{x})|^p > (ec'\sigma^2)^p \left( \sqrt{p} \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_F + p \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}} + \sqrt{p} \|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2 \right)^p \right] \\ &\stackrel{(a)}{\leq} \frac{\mathbb{E}[|q_c(\mathbf{x})|^p]}{(ec'\sigma^2)^p \left( \sqrt{p} \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_F + p \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}} + \sqrt{p} \|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2 \right)^p} \\ &\stackrel{(161)}{\leq} e^{-p}, \end{aligned}$$

where (a) follows from Markov's inequality. The proof is complete by choosing an appropriate universal constant  $c''$ , and performing basic algebraic manipulations after letting

$$p = \frac{1}{c''\sigma^2} \min \left( \frac{\varepsilon^2}{\mathbb{E}[\|\tilde{\nabla}q(\mathbf{x})\|_2]^2 + \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_F^2}, \frac{\varepsilon}{\max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}}} \right).$$

We note that an even  $p \geq 2$  can be ensured by choosing appropriate  $c''$ .

### F.2.1 Proof of Lemma. F.4(159): Bounded $p$ -th moment of $q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]$

Fix any  $p \geq 2$ . We start by using the following result from (Aida and Stroock, 1994, Theorem 3.4) since  $\mathbf{x}$  satisfies  $\text{LSI}_{\mathbf{x}}(\sigma^2)$ :

$$\|q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]\|_{L_p}^2 \leq \|q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]\|_{L_2}^2 + 2\sigma^2(p-2) \|\|\nabla q(\mathbf{x})\|_2\|_{L_p}^2. \quad (162)$$

Then, we bound the first term in (162) by using the fact that logarithmic Sobolev inequality implies Poincare inequality with the same constant:

$$\|q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]\|_{L_2}^2 = \text{Var}(q(\mathbf{x})) \leq \sigma^2 \mathbb{E}_{\mathbf{x}} \left[ \|\nabla q(\mathbf{x})\|_2^2 \right]. \quad (163)$$

Putting together (162) and (163), we have

$$\begin{aligned} \|q(\mathbf{x}) - \mathbb{E}[q(\mathbf{x})]\|_{L_p}^2 &\leq \sigma^2 \mathbb{E}_{\mathbf{x}} \left[ \|\nabla q(\mathbf{x})\|_2^2 \right] + 2\sigma^2(p-2) \|\|\nabla q(\mathbf{x})\|_2\|_{L_p}^2 \\ &\stackrel{(a)}{\leq} \sigma^2 \left( \mathbb{E}_{\mathbf{x}} \left[ \|\nabla q(\mathbf{x})\|_2^p \right] \right)^{2/p} + 2\sigma^2(p-2) \|\|\nabla q(\mathbf{x})\|_2\|_{L_p}^2 \end{aligned}$$

$$\begin{aligned}
&\stackrel{(b)}{=} \sigma^2 \|\|\nabla q(\mathbf{x})\|_2\|_{L_p}^2 + 2\sigma^2(p-2) \|\|\nabla q(\mathbf{x})\|_2\|_{L_p}^2 \\
&\leq 2\sigma^2 p \|\|\nabla q(\mathbf{x})\|_2\|_{L_p}^2,
\end{aligned} \tag{164}$$

where (a) follows by Jensen's inequality and (b) follows by the definition of  $p$ -th moment. Taking square root on both sides of (164) completes the proof.

### F.2.2 Proof of Lemma. F.4(160): Bounded $p$ -th moment of $\|\tilde{\nabla}q(\mathbf{x})\|_2$

Fix any even  $p \geq 2$ . Fix any pseudo derivative  $\tilde{\nabla}q$  and any pseudo Hessian  $\tilde{\nabla}^2q$ . We start by obtaining a convenient bound on  $\|\tilde{\nabla}q(\mathbf{x})\|_2$  for every  $\mathbf{x} \in \mathcal{X}^p$  and then proceed to bound the  $p$ -th moment of  $\|\tilde{\nabla}q(\mathbf{x})\|_2$ .

Consider a  $p$ -dimensional standard normal random vector  $\mathbf{g}$  independent of  $\mathbf{x}$ . For a given  $\mathbf{x} = \mathbf{x} \in \mathcal{X}^p$ , the random variable  $\frac{\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}}{\|\tilde{\nabla}q(\mathbf{x})\|_2}$  is a standard normal random variable. Then, for every  $\mathbf{x} \in \mathcal{X}^p$ , we have

$$\left\| \frac{\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}}{\|\tilde{\nabla}q(\mathbf{x})\|_2} \right\|_{L_p} \stackrel{(a)}{=} \left( \mathbb{E}_{\mathbf{g}|\mathbf{x}=\mathbf{x}} \left[ \left( \frac{\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}}{\|\tilde{\nabla}q(\mathbf{x})\|_2} \right)^p \right] \right)^{1/p} \stackrel{(b)}{\geq} \frac{\sqrt{p}}{2}, \tag{165}$$

where (a) follows from the definition of  $p$ -th moment, and (b) follows since  $\|\mathbf{g}\|_{L_p} \geq \frac{\sqrt{p}}{2}$  for any standard normal random variable  $\mathbf{g}$  and even  $p \geq 2$ . Rearranging (165), we have

$$\|\tilde{\nabla}q(\mathbf{x})\|_2 \leq \frac{2}{\sqrt{p}} \left( \mathbb{E}_{\mathbf{g}|\mathbf{x}=\mathbf{x}} \left[ (\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g})^p \right] \right)^{1/p}. \tag{166}$$

Now, we proceed to bound the  $p$ -th moment of  $\|\tilde{\nabla}q(\mathbf{x})\|_2$  as follows

$$\begin{aligned}
\|\|\tilde{\nabla}q(\mathbf{x})\|_2\|_{L_p} &\stackrel{(a)}{=} \left( \mathbb{E}_{\mathbf{x}} [\|\tilde{\nabla}q(\mathbf{x})\|_2^p] \right)^{1/p} \\
&\stackrel{(166)}{\leq} \frac{2}{\sqrt{p}} \left( \mathbb{E}_{\mathbf{x}, \mathbf{g}} \left[ (\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g})^p \right] \right)^{1/p} \\
&\stackrel{(b)}{=} \frac{2}{\sqrt{p}} \left\| \tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} \right\|_{L_p} \\
&\stackrel{(c)}{\leq} \frac{2}{\sqrt{p}} \left( \left\| \tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}} [\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}] \right\|_{L_p} + \left\| \mathbb{E}_{\mathbf{x}} [\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}] \right\|_{L_p} \right),
\end{aligned} \tag{167}$$

where (a) and (b) follow from the definition of  $p$ -th moment and (c) follows by Minkowski's inequality. We claim that

$$\left\| \tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}} [\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}] \right\|_{L_p} \leq c\sigma \left( \sqrt{p} \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{F}} + p \max_{\mathbf{x} \in \mathcal{X}^p} \|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}} \right), \quad \& \tag{168}$$

$$\left\| \mathbb{E}_{\mathbf{x}} [\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}] \right\|_{L_p} \leq 2\sqrt{p} \left\| \mathbb{E}_{\mathbf{x}} [\tilde{\nabla}q(\mathbf{x})] \right\|_2, \tag{169}$$

where  $c \geq 0$  is a universal constant. Putting together (167) to (169) completes the proof. It remains to prove our claims (168) and (169) which we now do one-by-one.

**Proof of bound (168)** To start, we bound  $(\mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}[(\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}])^p])^{1/p}$  for every  $\mathbf{g} = \mathbf{g}$ , and then proceed to bound  $\|\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}]\|_{L_p}$ .

To that end, we define  $h_{\mathbf{g}}(\mathbf{x}) \triangleq \tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}]$  and observe that  $\mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}[h_{\mathbf{g}}(\mathbf{x})] = 0$ . Now, applying Lemma. F.4 (159) to  $h_{\mathbf{g}}(\cdot)$ , we have

$$\begin{aligned} \|h_{\mathbf{g}}(\mathbf{x})\|_{L_p} &\leq \sigma\sqrt{2p}\left(\mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}\left[\|\nabla h_{\mathbf{g}}(\mathbf{x})\|_2^p\right]\right)^{1/p} \stackrel{(a)}{\leq} \sigma\sqrt{2p}\left(\mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}\left[\|\nabla[\mathbf{g}^\top \tilde{\nabla}q(\mathbf{x})]\|_2^p\right]\right)^{1/p} \\ &\stackrel{(130)}{\leq} \sigma\sqrt{2p}\left(\mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}\left[\|\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\|_2^p\right]\right)^{1/p}, \end{aligned} \quad (170)$$

where (a) follows from the definition of  $h_{\mathbf{g}}(\mathbf{x})$ . Now, to obtain a bound on the RHS of (170), we further fix  $\mathbf{x} = \mathbf{x}$ . Then, we let  $\mathbf{g}'$  be another  $p$ -dimensional standard normal vector and apply an inequality similar to (166) to  $\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})$  obtaining

$$\left\|\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\right\|_2 \leq \frac{2}{\sqrt{p}}\left(\mathbb{E}_{\mathbf{g}'|\mathbf{x}=\mathbf{x}, \mathbf{g}=\mathbf{g}}\left[\left(\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\mathbf{g}'\right)^p\right]\right)^{1/p},$$

which implies

$$\left(\mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}\left[\left\|\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\right\|_2^p\right]\right)^{1/p} \leq \frac{2}{\sqrt{p}}\left(\mathbb{E}_{\mathbf{x}, \mathbf{g}'|\mathbf{g}=\mathbf{g}}\left[\left(\nabla\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\mathbf{g}'\right)^p\right]\right)^{1/p}. \quad (171)$$

Putting together (170) and (171), and using the definition of  $h_{\mathbf{g}}(\mathbf{x})$ , we have

$$\mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}\left[\left(\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}|\mathbf{g}=\mathbf{g}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}]\right)^p\right] \leq (2\sqrt{2}\sigma)^p \mathbb{E}_{\mathbf{x}, \mathbf{g}'|\mathbf{g}=\mathbf{g}}\left[\left(\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\mathbf{g}'\right)^p\right]. \quad (172)$$

Now, we proceed to bound  $\|\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}]\|_{L_p}$  as follows

$$\begin{aligned} \left\|\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}]\right\|_{L_p} &\stackrel{(a)}{=} \left(\mathbb{E}_{\mathbf{x}, \mathbf{g}}\left[\left(\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g} - \mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}]\right)^p\right]\right)^{1/p} \\ &\stackrel{(172)}{\leq} 2\sqrt{2}\sigma\left(\mathbb{E}_{\mathbf{g}, \mathbf{x}, \mathbf{g}'}\left[\left(\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\mathbf{g}'\right)^p\right]\right)^{1/p}, \end{aligned} \quad (173)$$

where (a) follows from the definition of  $p$ -th moment. Finally, to bound the RHS of (173), we fix  $\mathbf{x} = \mathbf{x}$  and bound the  $p$ -th norm of the quadratic form  $\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\mathbf{g}'$  by the Hanson-Wright inequality resulting in

$$\begin{aligned} \left(\mathbb{E}_{\mathbf{g}, \mathbf{g}'|\mathbf{x}=\mathbf{x}}\left[\left(\mathbf{g}^\top \tilde{\nabla}^2q(\mathbf{x})\mathbf{g}'\right)^p\right]\right)^{1/p} &\leq c\left(\sqrt{p}\|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{F}} + p\|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}}\right) \\ &\leq c\left(\sqrt{p}\max_{\mathbf{x}\in\mathcal{X}^p}\|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{F}} + p\max_{\mathbf{x}\in\mathcal{X}^p}\|\tilde{\nabla}^2q(\mathbf{x})\|_{\text{op}}\right), \end{aligned} \quad (174)$$

where  $c \geq 0$  is a universal constant. Then, (168) follows by putting together (173) and (174).

**Proof of bound (169)** By linearity of expectation, we have

$$\|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})^\top \mathbf{g}]\|_{L_p} = \|(\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})])^\top \mathbf{g}\|_{L_p}. \quad (175)$$

We note that the random variable  $\frac{(\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})])^\top \mathbf{g}}{\|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2}$  is a standard normal random variable. Therefore,

$$\left\| \frac{(\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})])^\top \mathbf{g}}{\|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2} \right\|_{L_p} \stackrel{(a)}{=} \left( \mathbb{E}_{\mathbf{g}} \left[ \left( \frac{(\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})])^\top \mathbf{g}}{\|\mathbb{E}_{\mathbf{x}}[\tilde{\nabla}q(\mathbf{x})]\|_2} \right)^p \right] \right)^{1/p} \stackrel{(b)}{\leq} 2\sqrt{p}, \quad (176)$$

where (a) follows from the definition of  $p$ -th moment, and (b) follows since  $\|\mathbf{g}\|_{L_p} \leq 2\sqrt{p}$  for any standard normal variable  $\mathbf{g}$ . Then, (169) follows by using (176) in (175).

## G Identifying weakly dependent random variables

In Appendix. F, we derived (in Proposition. F.1) that a random vector (supported on a compact set) satisfies the logarithmic Sobolev inequality if it satisfies the Dobrushin's uniqueness condition (in Definition. F.2). Further, we also derived (Proposition. F.2) tail bounds for a random vector satisfying the logarithmic Sobolev inequality. Combining the two, we see that in order to use the tail bound, the random vector needs to satisfy the Dobrushin's uniqueness condition, i.e, the elements of the random vector should be weakly dependent. In this section, we show that any random vector (outside Dobrushin's regime) that is a  $\tau$ -Sparse Graphical Model (to be defined) can be reduced to satisfy the Dobrushin's uniqueness condition. In particular, we show that by conditioning on a subset of the random vector, the unconditioned subset of the random vector (in the conditional distribution) are only weakly dependent. We exploit this trick in Lemma. C.5 and Lemma. C.7 to enable application of the tail bound in Appendix. F. The result below is a generalization of the result in Dagan et al. (2021) for discrete random vectors to continuous random vectors.

We start by defining the notion of  $\tau$ -Sparse Graphical Model.

**Definition G.1** ( $\tau$ -Sparse Graphical Model). *A pair of random vectors  $\{\mathbf{x}, \mathbf{z}\}$  supported on  $\mathcal{X}^p \times \mathcal{Z}^{p_z}$  is a  $\tau$ -Sparse Graphical Model for model-parameters  $\tau \triangleq (\alpha, \zeta, x_{\max}, \Theta)$  and denoted by  $\tau$ -SGM if  $\mathcal{X} = \{-x_{\max}, x_{\max}\}$ , and*

1. *for any realization  $\mathbf{z} \in \mathcal{Z}^{p_z}$ , the conditional probability distribution of  $\mathbf{x}$  given  $\mathbf{z} = \mathbf{z}$  is given by  $f_{\mathbf{x}|\mathbf{z}}(\cdot | \mathbf{z}; \theta(\mathbf{z}), \Theta)$  in (6) for a vector  $\theta(\mathbf{z}) \in \mathbb{R}^p$  depending on  $\mathbf{z}$  and a symmetric matrix  $\Theta \in \mathbb{R}^{p \times p}$  (independent of  $\mathbf{z}$ ),*
2.  $\max \{ \max_{\mathbf{z} \in \mathcal{Z}^{p_z}} \|\theta(\mathbf{z})\|_\infty, \|\Theta\|_{\max} \} \leq \alpha$ , and
3.  $\|\Theta\|_\infty \leq \zeta$ .

Now, we provide the main result of this section.

**Proposition G.1** (Identifying weakly dependent random variables). *Given a pair of random vectors  $\{\mathbf{x}, \mathbf{z}\}$  supported on  $\mathcal{X}^p \times \mathcal{Z}^{p_z}$  that is a  $\tau$ -SGM (Definition. G.1) with  $\tau \triangleq (\alpha, \zeta, x_{\max}, \Theta)$ , and a scalar  $\lambda \in (0, \zeta]$ , there exists  $L \triangleq 32\zeta^2 \log 4p/\lambda^2$  subsets  $S_1, \dots, S_L \subseteq [p]$  that satisfy the following properties:*

- (a) *For any  $t \in [p]$ , we have  $\sum_{u=1}^L \mathbf{1}(t \in S_u) = \lceil \lambda L / (8\zeta) \rceil$ .*
- (b) *For any  $u \in [L]$ ,*

- (i) the pair of random vectors  $\{\mathbf{x}_{S_u}, (\mathbf{x}_{-S_u}, \mathbf{z})\}$  correspond to a  $\tau_1$ -SGM with  $\tau_1 \triangleq (\alpha + 2x_{\max}\zeta, \lambda, x_{\max}, \Theta_{S_u})$  where  $\Theta_{S_u} \triangleq \{\Theta_{tv}\}_{t,v \in S_u}$ , and
- (ii) the random vector  $\mathbf{x}_{S_u}$  conditioned on  $(\mathbf{x}_{-S_u}, \mathbf{z})$  satisfies the Dobrushin's uniqueness condition (Definition. F.2) with coupling matrix  $2\sqrt{2}x_{\max}^2|\Theta_{S_u}|$  whenever  $\lambda \in \left(0, \frac{1}{2\sqrt{2}x_{\max}^2}\right]$  with  $\|\Theta_{S_u}\|_{\text{op}} \leq \lambda$ .

*Proof of Proposition. G.1: Identifying weakly dependent random variables.* We prove each part one-by-one using a generalization of Dagan et al. (2021, Lemma. 12).

Recall Dagan et al. (2021, Lemma. 12): Let  $A \in \mathbb{R}^{p \times p}$  be a matrix with zeros on the diagonal and  $\|A\|_{\infty} \leq 1$ . Let  $0 < \eta < 1$ . Then, there exists subsets  $\bar{S}_1, \dots, \bar{S}_{\bar{L}} \subseteq [p]$  with  $\bar{L} \triangleq 32 \log 4p/\eta^2$  such that

- (a) For any  $t \in [p]$ , we have  $\sum_{u=1}^{\bar{L}} \mathbf{1}(t \in \bar{S}_u) = \lceil \eta \bar{L} / 8 \rceil$ , and
- (b) For any  $u \in [\bar{L}]$  and  $t \in \bar{S}_u$ ,  $\sum_{v \in \bar{S}_u} |A_{tv}| \leq \eta$ .

We claim that Dagan et al. (2021, Lemma. 12) holds even when  $A$  does not have zeros on the diagonal. The proof is exactly the same as the proof of Dagan et al. (2021, Lemma. 12).

**Proof of part (a)** From Definition. G.1, for any realization  $\mathbf{z} \in \mathcal{Z}^{p_z}$ , the conditional probability distribution of  $\mathbf{x}$  given  $\mathbf{z} = \mathbf{z}$  is given by  $f_{\mathbf{x}|\mathbf{z}}(\cdot | \mathbf{z}; \theta(\mathbf{z}), \Theta)$  in (6) where  $\theta(\mathbf{z}) \in \mathbb{R}^p$  is a vector and  $\Theta \in \mathbb{R}^{p \times p}$  is a symmetric matrix with  $\|\Theta\|_{\infty} \leq \zeta$ . Consider the matrix  $A \triangleq \frac{1}{\zeta} \Theta$ . Since  $\|A\|_{\infty} \leq 1$ , we can apply the generalization of Dagan et al. (2021, Lemma. 12) on  $A$  with  $\eta = \frac{\lambda}{\zeta}$ . Then part (a) follows directly from Dagan et al. (2021, Lemma. 12.1).

**Proof of part (b)(i)** To prove this part, consider the distribution of  $\mathbf{x}_{S_u}$  conditioned on  $\mathbf{x}_{-S_u} = \mathbf{x}_{-S_u}$  and  $\mathbf{z} = \mathbf{z}$  for any  $u \in [L]$ , i.e.,  $f_{\mathbf{x}_{S_u}|\mathbf{x}_{-S_u}, \mathbf{z}}(\mathbf{x}_{S_u} | \mathbf{x}_{-S_u}, \mathbf{z}; \theta(\mathbf{z}), \Theta) \triangleq f(\mathbf{x}_{S_u} | \mathbf{x}_{-S_u}, \mathbf{z}; \theta(\mathbf{z}), \Theta)$  as follows

$$f(\mathbf{x}_{S_u} | \mathbf{x}_{-S_u}, \mathbf{z}; \theta(\mathbf{z}), \Theta) \propto \exp \left( \sum_{t \in S_u} (\theta_t(\mathbf{z}) + 2 \sum_{v \notin S_u} \Theta_{tv} x_v) x_t + \sum_{t \in S_u} \sum_{v \in S_u} \Theta_{tv} x_t x_v \right). \quad (177)$$

We can re-parameterize  $f(\mathbf{x}_{S_u} | \mathbf{x}_{-S_u}, \mathbf{z}; \theta(\mathbf{z}), \Theta)$  in (177) as follows

$$f_{\mathbf{x}_{S_u}|\mathbf{x}_{-S_u}, \mathbf{z}}(\mathbf{x}_{S_u} | \mathbf{x}_{-S_u}, \mathbf{z}; v(\mathbf{z}, \mathbf{x}_{-S_u}), \Upsilon) \propto \exp \left( [v(\mathbf{z}, \mathbf{x}_{-S_u})]^\top \mathbf{x}_{S_u} + \mathbf{x}_{S_u}^\top \Upsilon \mathbf{x}_{S_u} \right)$$

where

$$v(\mathbf{z}, \mathbf{x}_{-S_u}) \in \mathbb{R}^{|S_u| \times 1}, \text{ with } v_t(\mathbf{z}, \mathbf{x}_{-S_u}) \triangleq \theta_t(\mathbf{z}) + 2 \sum_{k \notin S_u} \Theta_{tk} x_k \text{ for } t \in S_u, \text{ and} \quad (178)$$

$$\Upsilon = \Upsilon^\top \in \mathbb{R}^{|S_u| \times |S_u|} \text{ with } \Upsilon_{tv} \triangleq \Theta_{tv}, \text{ for all } t, v \in S_u. \quad (179)$$

Now, to show that the random vector  $\mathbf{x}_{S_u}$  conditioned on  $\mathbf{x}_{-S_u}$  and  $\mathbf{z}$  corresponds to an  $\tau_1$ -SGM with  $\tau_1 \triangleq (\alpha + 2x_{\max}\zeta, \lambda, x_{\max}, \Theta_{S_u})$ , it suffices to establish that

$$\max \left\{ \max_{\mathbf{z} \in \mathcal{Z}^{p_z}} \|v(\mathbf{z}, \mathbf{x}_{-S_u})\|_{\infty}, \|\Upsilon\|_{\max} \right\} \stackrel{(i)}{\leq} \alpha + 2x_{\max}\zeta \quad \text{and} \quad \|\Upsilon\|_{\infty} \stackrel{(ii)}{\leq} \lambda. \quad (180)$$

To establish (i) in (180), we note that

$$\|\Upsilon\|_{\max} \stackrel{(179)}{\leq} \|\Theta\|_{\max} \stackrel{(a)}{\leq} \alpha \quad \text{and} \quad (181)$$

$$\begin{aligned} \|v(\mathbf{z}, \mathbf{x}_{-S_u})\|_{\infty} &\stackrel{(b)}{\leq} \|\theta(\mathbf{z})\|_{\infty} + 2 \max_{t \in S_u} \|\Theta_t\|_1 \|\mathbf{x}\|_{\infty} \stackrel{(c)}{\leq} \|\theta(\mathbf{z})\|_{\infty} + 2x_{\max} \|\Theta\|_{\infty} \\ &\stackrel{(d)}{\leq} \alpha + 2x_{\max} \zeta, \end{aligned} \quad (182)$$

where (a) and (d) follow from Definition. G.1, (b) follows from (178) and the triangle inequality, and (c) follows from the definition of  $\|\cdot\|_{\infty}$  and Definition. G.1. Then, from (181) and (182), we have

$$\max \left\{ \max_{\mathbf{z} \in \mathcal{Z}^{p_z}} \|v(\mathbf{z}, \mathbf{x}_{-S_u})\|_{\infty}, \|\Upsilon\|_{\max} \right\} \leq \alpha + 2x_{\max} \zeta,$$

as claimed. Next, to establish (ii) in (180), we again apply the generalization of Dagan et al. (2021, Lemma. 12) on the matrix  $A = \frac{1}{\zeta} \Theta$  with  $\eta = \frac{\lambda}{\zeta}$ . Then, we have

$$\sum_{v \in S_u} \left| \frac{\Theta_{tv}}{\zeta} \right| \leq \frac{\lambda}{\zeta} \quad \text{for all } t \in S_u, u \in [L]. \quad (183)$$

Therefore, we have

$$\|\Upsilon\|_{\infty} = \max_{t \in S_u} \left( \sum_{v \in S_u} |\Upsilon_{tv}| \right) \stackrel{(179)}{=} \max_{t \in S_u} \left( \sum_{v \in S_u} |\Theta_{tv}| \right) \stackrel{(183)}{\leq} \lambda, \quad (184)$$

as desired. The proof for this part is now complete.

**Proof of part (b)(ii)** We start by noting that the operator norm of a symmetric matrix is bounded by the infinity norm of the matrix. Then, from the analysis in part (b) (i), for any  $u \in S_u$ , we have

$$\|\|\Theta_{S_u}\|_{\text{op}} \leq \|\|\Theta_{S_u}\|_{\infty} \stackrel{(179)}{=} \|\|\Upsilon\|_{\infty} \stackrel{(184)}{\leq} \lambda.$$

Therefore,  $\|2\sqrt{2}x_{\max}^2|\Theta_{S_u}\|_{\infty} \leq 1$  whenever  $\lambda \leq 1/2\sqrt{2}x_{\max}^2$ . It remains to show that for every  $u \in [L]$ ,  $t \in S_u$ ,  $v \in S_u \setminus \{t\}$ ,  $\mathbf{z} = \mathbf{z}$ , and  $\mathbf{x}_{-t}, \tilde{\mathbf{x}}_{-t} \in \mathcal{X}^{p-1}$  differing only in the  $v^{\text{th}}$  coordinate,

$$\|f_{x_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}, \mathbf{z}=\mathbf{z}} - f_{x_t|\mathbf{x}_{-t}=\tilde{\mathbf{x}}_{-t}, \mathbf{z}=\mathbf{z}}\|_{\text{TV}} \leq 2\sqrt{2}x_{\max}^2|\Theta_{tv}|.$$

To that end, fix any  $u \in [L]$ , any  $t \in S_u$ , any  $v \in S_u \setminus \{t\}$ , any  $\mathbf{z} = \mathbf{z}$ , and any  $\mathbf{x}_{-t}, \tilde{\mathbf{x}}_{-t} \in \mathcal{X}^{p-1}$  differing only in the  $v^{\text{th}}$  coordinate. We have

$$\begin{aligned} \|f_{x_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}, \mathbf{z}=\mathbf{z}} - f_{x_t|\mathbf{x}_{-t}=\tilde{\mathbf{x}}_{-t}, \mathbf{z}=\mathbf{z}}\|_{\text{TV}}^2 &\stackrel{(a)}{\leq} \frac{1}{2} \text{KL} (f_{x_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}, \mathbf{z}=\mathbf{z}} \| f_{x_t|\mathbf{x}_{-t}=\tilde{\mathbf{x}}_{-t}, \mathbf{z}=\mathbf{z}}) \\ &\stackrel{(b)}{=} \frac{1}{2} (2\Theta_{tv}x_v - 2\Theta_{tv}\tilde{x}_v)^2 x_{\max}^2 \stackrel{(c)}{\leq} 8x_{\max}^4 \Theta_{tv}^2, \end{aligned}$$

where (a) follows from Pinsker's inequality, (b) follows by (i) applying (Busa-Fekete et al., 2019, Theorem 1) to the exponential family parameterized as per  $f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}$  in (12), (ii) noting that  $f_{x_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}, \mathbf{z}=\mathbf{z}} \propto \exp([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^{\top} \mathbf{x}_{-t}]x_t + \Theta_{tt}\bar{x}_t)$  and  $f_{x_t|\mathbf{x}_{-t}=\tilde{\mathbf{x}}_{-t}, \mathbf{z}=\mathbf{z}} \propto \exp([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^{\top} \tilde{\mathbf{x}}_{-t}]x_t + \Theta_{tt}\bar{x}_t)$  where  $\bar{x}_t \triangleq x_t^2 - x_{\max}^2/3$ , and (iii) noting that the Hessian of the log partition function for any regular exponential family is the covariance matrix of the associated sufficient statistic which is bounded by  $x_{\max}^2$  when  $\mathcal{X} = \{-x_{\max}, x_{\max}\}$ , and (c) follows because  $x_v, \tilde{x}_v \in \{-x_{\max}, x_{\max}\}$ . This completes the proof.  $\square$

## H Supporting concentration results

In this section, we provide a corollary of Proposition. F.2 that is used to prove the concentration results in Lemma. C.5 and Lemma. C.7. To show any concentration result for the random vector  $\mathbf{x}$  conditioned on  $\mathbf{z}$  via Proposition. F.2, we need  $\mathbf{x}|\mathbf{z}$  to satisfy the logarithmic Sobolev inequality (defined in (127)). From Proposition. F.1, for this to be true, we need the random vector  $\mathbf{x}_t$  conditioned on  $(\mathbf{x}_{-t}, \mathbf{z})$  to satisfy the logarithmic Sobolev inequality for all  $t \in [p]$ . In the result below, we show this holds with a proof in Appendix. H.1. We define a  $\tau \triangleq (\alpha, \zeta, x_{\max}, \Theta)$ -dependent constant:

$$C_{3,\tau} \triangleq \exp(x_{\max}(\alpha + 2\zeta x_{\max})). \quad (185)$$

**Lemma H.1** (Logarithmic Sobolev inequality for  $\mathbf{x}_t|\mathbf{x}_{-t}, \mathbf{z}$ ). *Given a pair of random vectors  $\{\mathbf{x}, \mathbf{z}\}$  supported on  $\mathcal{X}^p \times \mathcal{Z}^{p_z}$  that is a  $\tau$ -SGM (Definition. G.1) with  $\tau \triangleq (\alpha, \zeta, x_{\max}, \Theta)$ ,  $\mathbf{x}_t|\mathbf{x}_{-t}, \mathbf{z}$  satisfies  $\text{LSI}_{\mathbf{x}_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}, \mathbf{z}=\mathbf{z}} \left( \frac{8x_{\max}^2}{\pi^2} C_{3,\tau}^2 \right)$  for all  $t \in [p]$ ,  $\mathbf{x}_{-t} \in \mathcal{X}^{p-1}$ , and  $\mathbf{z} \in \mathcal{Z}^{p_z}$ .*

Now, we state the desired corollary of Proposition. F.2 with a proof in Appendix. H.2. The corollary makes use of some  $\tau \triangleq (\alpha, \zeta, x_{\max}, \Theta)$ -dependent constants:

$$C_{4,\tau} \triangleq 1 + \alpha x_{\max} + 4x_{\max}^2 \zeta \quad \text{and} \quad C_{5,\tau} \triangleq \frac{32x_{\max}^3 C_{3,\tau}^4}{\pi^2}. \quad (186)$$

**Corollary H.1** (Supporting concentration bounds). *Suppose a pair of random vectors  $\{\mathbf{x}, \mathbf{z}\}$  supported on  $\mathcal{X}^p \times \mathcal{Z}^{p_z}$  corresponds to a  $\tau$ -SGM (Definition. G.1) with  $\tau \triangleq (\alpha, \zeta, x_{\max}, \Theta)$ , and  $\mathbf{x}$  conditioned on  $\mathbf{z}$  satisfies the Dobrushin's uniqueness condition (Definition. F.2) with coupling matrix  $\bar{\Theta}$ . For any  $\theta, \bar{\theta} \in \Lambda_\theta$  and  $\Theta \in \Lambda_\Theta$ , define the functions  $q_1$  and  $q_2$  as*

$$q_1(\mathbf{x}) \triangleq \sum_{t \in [p]} (\omega_t \mathbf{x}_t)^2 \quad \text{and} \quad q_2(\mathbf{x}) \triangleq \sum_{t \in [p]} \omega_t \mathbf{x}_t \exp \left( -[\theta_t + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}] \mathbf{x}_t - \Theta_{tt} \bar{\mathbf{x}}_t \right),$$

where  $\omega = \bar{\theta} - \theta$  and  $\bar{\mathbf{x}}_t \triangleq \mathbf{x}_t^2 - x_{\max}^2/3$ . Then, for any  $\varepsilon > 0$

$$\mathbb{P} \left[ |q_i(\mathbf{x}) - \mathbb{E}[q_i(\mathbf{x})|\mathbf{z}]| \geq \varepsilon \mid \mathbf{z} \right] \leq \exp \left( \frac{-c(1 - \|\bar{\Theta}\|_{\text{op}})^4 \varepsilon^2}{c_i \|\omega\|_2^2} \right) \quad \text{for } i = 1, 2, \quad (187)$$

where  $c$  is a universal constant,  $c_1 \triangleq 16\alpha^2 x_{\max}^2 C_{5,\tau}^2$ , and  $c_2 \triangleq C_{3,\tau}^2 C_{4,\tau}^2 C_{5,\tau}^2$  with  $C_{3,\tau}$  defined in (185) and  $C_{4,\tau}$  and  $C_{5,\tau}$  defined in (186).

### H.1 Proof of Lemma. H.1: Logarithmic Sobolev inequality for $\mathbf{x}_t|\mathbf{x}_{-t}, \mathbf{z}$

Let  $u$  be the uniform distribution on  $\mathcal{X}$ . Then,  $u$  satisfies  $\text{LSI}_u \left( \frac{8x_{\max}^2}{\pi^2} \right)$  (see Ghang et al. (2014, Corollary. 2.4)). Then, using the Holley-Stroock perturbation principle (see Holley and Stroock (1987, Page. 31), Ledoux (2001, Lemma. 1.2)), for every  $t \in [p]$ ,  $\mathbf{x}_{-t} \in \mathcal{X}^{p-1}$ , and  $\mathbf{z} \in \mathcal{Z}^{p_z}$ ,  $\mathbf{x}_t|\mathbf{x}_{-t} = \mathbf{x}_{-t}, \mathbf{z} = \mathbf{z}$  satisfies the logarithmic Sobolev inequality with a constant bounded by

$$\frac{8x_{\max}^2 \exp(\sup_{\mathbf{x}_t \in \mathcal{X}} \psi(\mathbf{x}_t; \mathbf{x}_{-t}, \mathbf{z}) - \inf_{\mathbf{x}_t \in \mathcal{X}} \psi(\mathbf{x}_t; \mathbf{x}_{-t}, \mathbf{z}))}{\pi^2},$$

where  $\psi(x_t; \mathbf{x}_{-t}, \mathbf{z}) \triangleq -[\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t$  where  $\bar{x}_t = x_t^2 - x_{\max}^2/3$ . We have

$$\begin{aligned} \exp\left(\sup_{x_t \in \mathcal{X}} \psi(x_t; \mathbf{x}_{-t}, \mathbf{z}) - \inf_{x_t \in \mathcal{X}} \psi(x_t; \mathbf{x}_{-t}, \mathbf{z})\right) &\stackrel{(a)}{\leq} \exp\left(2|\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}|x_{\max} + \Theta_{tt}x_{\max}^2\right) \\ &\stackrel{(b)}{\leq} \exp\left((2\alpha + 4\zeta x_{\max})x_{\max}\right) \stackrel{(185)}{=} C_{3,\tau}^2, \end{aligned}$$

where (a) follows from Definition. G.1 and (b) follows by using Definition. G.1 along with triangle inequality and Cauchy–Schwarz inequality.

## H.2 Proof of Corollary. H.1: Supporting concentration bounds

To apply Proposition. F.2 to the random vector  $\mathbf{x}$  conditioned on  $\mathbf{z}$ , we need  $\mathbf{x}|\mathbf{z}$  to satisfy the logarithmic Sobolev inequality. From Proposition. F.1, this is true if (i)  $f_{\min} = \min_{t \in [p], \mathbf{x} \in \mathcal{X}^p, \mathbf{z} \in \mathcal{X}^{pz}} f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}(x_t|\mathbf{x}_{-t}, \mathbf{z}) > 0$  (see (129)), (ii)  $\mathbf{x}|\mathbf{z}$  satisfies the Dobrushin’s uniqueness condition, and (iii)  $x_t|\mathbf{x}_{-t}, \mathbf{z}$  satisfies the logarithmic Sobolev inequality for all  $t \in [p]$ . By assumption,  $\mathbf{x}|\mathbf{z}$  satisfies the Dobrushin’s uniqueness condition with coupling matrix  $\bar{\Theta}$ . From Lemma. H.1,  $x_t|\mathbf{x}_{-t}, \mathbf{z}$  satisfies  $\text{LSI}_{x_t|\mathbf{x}_{-t}=\mathbf{x}_{-t}, \mathbf{z}=\mathbf{z}}\left(\frac{8x_{\max}^2 C_{3,\tau}^2}{\pi^2}\right)$ . It remains to show that  $f_{\min} > 0$ . Consider any  $t \in [p]$ , any  $\mathbf{x} \in \mathcal{X}^p$ , and any  $\mathbf{z} \in \mathcal{X}^{pz}$ . Let  $\bar{x}_t = x_t^2 - x_{\max}^2/3$ . We have

$$\begin{aligned} f_{x_t|\mathbf{x}_{-t}, \mathbf{z}}(x_t|\mathbf{x}_{-t}, \mathbf{z}) &\stackrel{(a)}{=} \frac{\exp\left([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t + \Theta_{tt}\bar{x}_t\right)}{\int_{\mathcal{X}} \exp\left([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t + \Theta_{tt}\bar{x}_t\right) dx_t} \\ &\stackrel{(b)}{\geq} \frac{\exp\left(-|\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}|x_{\max} - \Theta_{tt}x_{\max}^2\right)}{\int_{\mathcal{X}} \exp\left([\theta_t(\mathbf{z}) + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_{\max} + \Theta_{tt}x_{\max}^2\right) dx_t} \\ &\stackrel{(c)}{\geq} \frac{\exp\left(-(|\theta(\mathbf{z})| + 2\|\Theta_{t,-t}\|_1\|\mathbf{x}\|_\infty)x_{\max} - \Theta_{tt}x_{\max}^2\right)}{\int_{\mathcal{X}} \exp\left((|\theta(\mathbf{z})| + 2\|\Theta_{t,-t}\|_1\|\mathbf{x}\|_\infty)x_{\max} + \Theta_{tt}x_{\max}^2\right) dx_t} \\ &\stackrel{(d)}{\geq} \frac{\exp\left(-(\alpha + 2\zeta x_{\max})x_{\max}\right)}{\int_{\mathcal{X}} \exp\left((\alpha + 2\zeta x_{\max})x_{\max}\right) dx_t} \stackrel{(e)}{=} \frac{1}{2x_{\max}C_{3,\tau}^2}, \end{aligned}$$

where (a) follows from (12), (b) and (d) follow from Definition. G.1, (c) follows by triangle inequality and Cauchy–Schwarz inequality, and (e) follows because  $\int_{\mathcal{X}} dx_t = 2x_{\max}$ . Therefore,  $f_{\min} = \frac{1}{2x_{\max}C_{3,\tau}^2}$ . Putting (i), (ii), and (iii) together, and using Proposition. F.1, we see that  $\mathbf{x}|\mathbf{z}$  satisfies  $\text{LSI}_{\mathbf{x}}\left(\frac{C_{5,\tau}}{(1-\|\bar{\Theta}\|_{\text{op}})^2}\right)$  where  $C_{5,\tau}$  was defined in (186).

Now, we apply Proposition. F.2 to  $q_1$  and  $q_2$  one-by-one. The general strategy is to choose appropriate pseudo derivatives and pseudo Hessians for both  $q_1$  and  $q_2$ , and evaluate the corresponding terms appearing in Proposition. F.2.

**Concentration for  $q_1$**  Fix any  $\mathbf{x} \in \mathcal{X}^p$ . We start by decomposing  $q_1(\mathbf{x})$  as follows

$$q_1(\mathbf{x}) = \bar{w}^\top r(\mathbf{x}), \tag{188}$$

where  $\bar{\omega} \triangleq (\omega_1^2, \dots, \omega_p^2)$  and  $r(\mathbf{x}) \triangleq (r_1(\mathbf{x}), \dots, r_p(\mathbf{x}))$  with  $r_t(\mathbf{x}) = x_t^2$  for every  $t \in [p]$ . Next, we define  $H : \mathcal{X}^p \rightarrow \mathbb{R}^{p \times p}$  such that

$$H_{tu}(\mathbf{x}) = \frac{dr_u(\mathbf{x})}{dx_t} \quad \text{for every } t, u \in [p]. \quad (189)$$

**Pseudo derivative** We bound the  $\ell_2$  norm of the gradient of  $q_1(\mathbf{x})$  as follows

$$\begin{aligned} \|\nabla q_1(\mathbf{x})\|_2^2 &= \sum_{t \in [p]} \left( \frac{dq_1(\mathbf{x})}{dx_t} \right)^2 \stackrel{(188)}{=} \sum_{t \in [p]} \left( \bar{\omega}^\top \frac{dr(\mathbf{x})}{dx_t} \right)^2 \\ &\stackrel{(189)}{=} \|H(\mathbf{x})\bar{\omega}\|_2^2 \\ &\stackrel{(a)}{\leq} \|H(\mathbf{x})\|_{\text{op}}^2 \|\bar{\omega}\|_2^2 \stackrel{(b)}{\leq} \|H(\mathbf{x})\|_1 \|H(\mathbf{x})\|_\infty \|\bar{\omega}\|_2^2, \end{aligned} \quad (190)$$

where (a) follows because induced matrix norms are submultiplicative and (b) follows because the matrix operator norm is bounded by square root of the product of matrix one norm and matrix infinity norm. Now, we claim that the one norm and the infinity norm of  $H(\mathbf{x})$  are bounded as follows

$$\max \left\{ \max_{\mathbf{x} \in \mathcal{X}^p} \|H(\mathbf{x})\|_1, \max_{\mathbf{x} \in \mathcal{X}^p} \|H(\mathbf{x})\|_\infty \right\} \leq 2x_{\max}. \quad (191)$$

Taking this claim as given at the moment, we continue with our proof. Combining (190) and (191), we have

$$\max_{\mathbf{x} \in \mathcal{X}^p} \|\nabla q_1(\mathbf{x})\|_2^2 \leq 4x_{\max}^2 \|\bar{\omega}\|_2^2 = 4x_{\max}^2 \sum_{t \in [p]} \omega_t^4 \leq 4x_{\max}^2 \max_{u \in [p]} \omega_u^2 \sum_{t \in [p]} \omega_t^2 \stackrel{(a)}{\leq} 16x_{\max}^2 \alpha^2 \|\omega\|_2^2,$$

where (a) follows because  $\omega \in 2\Lambda_\theta$ . Therefore, we choose the pseudo derivative (see Definition. F.3) as follows

$$\tilde{\nabla} q_1(\mathbf{x}) = 4x_{\max} \alpha \|\omega\|_2. \quad (192)$$

**Pseudo Hessian** Fix any  $\rho \in \mathbb{R}$ . We bound  $\|\nabla(\rho^\top \tilde{\nabla} q_1(\mathbf{x}))\|_2^2$  (see Definition. F.3) as follows

$$\|\nabla(\rho^\top \tilde{\nabla} q_1(\mathbf{x}))\|_2^2 = \sum_{u \in [p]} \left( \frac{d\rho^\top \tilde{\nabla} q_1(\mathbf{x})}{dx_u} \right)^2 \stackrel{(192)}{=} 0.$$

Therefore, we choose the pseudo Hessian (see Definition. F.3) as follows

$$\tilde{\nabla}^2 q_1(\mathbf{x}) = 0. \quad (193)$$

The concentration result in (187) for  $q_1$  follows by applying Proposition. F.2 with the pseudo discrete derivative defined in (192) and the pseudo discrete Hessian defined in (193).

It remains to show that the one-norm and the infinity-norm of  $H(\mathbf{x})$  are bounded as in (191).

**Bounds on the one-norm and the infinity-norm of  $H(\mathbf{x})$**  We have

$$H_{tu}(\mathbf{x}) = \begin{cases} 2x_t & \text{if } t = u, \\ 0 & \text{otherwise.} \end{cases} \quad (194)$$

Therefore,

$$\begin{aligned} \|H(\mathbf{x})\|_1 &= \max_{u \in [p]} \sum_{t \in [p]} |H_{tu}(\mathbf{x})| \stackrel{(194)}{\leq} \max_{u \in [p]} 2|x_u| \stackrel{(a)}{\leq} 2x_{\max} \quad \text{and} \\ \|H(\mathbf{x})\|_\infty &= \max_{t \in [p]} \sum_{u \in [p]} |H_{tu}(\mathbf{x})| \stackrel{(194)}{\leq} \max_{t \in [p]} 2|x_t| \stackrel{(a)}{\leq} 2x_{\max}, \end{aligned}$$

where (a) follows from Definition. G.1.

**Concentration for  $q_2$**  Fix any  $\mathbf{x} \in \mathcal{X}^p$ . We start by decomposing  $q_2(\mathbf{x})$  as follows

$$q_2(\mathbf{x}) = \omega^\top r(\mathbf{x}), \quad (195)$$

where  $r(\mathbf{x}) \triangleq (r_1(\mathbf{x}), \dots, r_p(\mathbf{x}))$  with  $r_t(\mathbf{x}) = x_t \exp(-[\theta_t + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t)$  for every  $t \in [p]$ . Next, we define  $H : \mathcal{X}^p \rightarrow \mathbb{R}^{p \times p}$  such that

$$H_{tu}(\mathbf{x}) = \frac{dr_u(\mathbf{x})}{dx_t} \quad \text{for every } t, u \in [p]. \quad (196)$$

**Pseudo derivative** We bound the  $\ell_2$  norm of the gradient of  $q_2(\mathbf{x})$  as follows

$$\begin{aligned} \|\nabla q_2(\mathbf{x})\|_2^2 &= \sum_{t \in [p]} \left( \frac{dq_2(\mathbf{x})}{dx_t} \right)^2 \stackrel{(195)}{=} \sum_{t \in [p]} \left( \frac{\omega^\top dr(\mathbf{x})}{dx_t} \right)^2 \\ &\stackrel{(196)}{=} \|H(\mathbf{x})\omega\|_2^2 \\ &\stackrel{(a)}{\leq} \|H(\mathbf{x})\|_{\text{op}}^2 \|\omega\|_2^2 \stackrel{(b)}{\leq} \|H(\mathbf{x})\|_1 \|H(\mathbf{x})\|_\infty \|\omega\|_2^2, \end{aligned} \quad (197)$$

where (a) follows because induced matrix norms are submultiplicative and (b) follows because the matrix operator norm is bounded by square root of the product of matrix one norm and matrix infinity norm. Now, we claim that the one norm and the infinity norm of  $H(\mathbf{x})$  are bounded as follows

$$\max \left\{ \max_{\mathbf{x} \in \mathcal{X}^p} \|H(\mathbf{x})\|_1, \max_{\mathbf{x} \in \mathcal{X}^p} \|H(\mathbf{x})\|_\infty \right\} \leq C_{3,\tau} C_{4,\tau}. \quad (198)$$

where  $C_{3,\tau}$  and  $C_{4,\tau}$  were defined in (185) and (186) respectively. Taking this claim as given at the moment, we continue with our proof. Combining (197) and (198), we have

$$\max_{\mathbf{x} \in \mathcal{X}^p} \|\nabla q_2(\mathbf{x})\|_2^2 \leq C_{3,\tau}^2 C_{4,\tau}^2 \|\omega\|_2^2.$$

Therefore, we choose the pseudo derivative (see Definition. F.3) as follows

$$\tilde{\nabla} q_2(\mathbf{x}) = C_{3,\tau} C_{4,\tau} \|\omega\|_2. \quad (199)$$

**Pseudo Hessian** Fix any  $\rho \in \mathbb{R}$ . We bound  $\|\nabla(\rho^\top \tilde{\nabla} q_2(\mathbf{x}))\|_2^2$  (see Definition. F.3) as follows

$$\|\nabla(\rho^\top \tilde{\nabla} q_2(\mathbf{x}))\|_2^2 = \sum_{u \in [p]} \left( \frac{d\rho^\top \tilde{\nabla} q_2(\mathbf{x})}{dx_u} \right)^2 \stackrel{(199)}{=} 0.$$

Therefore, we choose the pseudo Hessian (see Definition. F.3) as follows

$$\tilde{\nabla}^2 q_2(\mathbf{x}) = 0. \quad (200)$$

The concentration result in (187) for  $q_1$  follows by applying Proposition. F.2 with the pseudo discrete derivative defined in (199) and the pseudo discrete Hessian defined in (200).

It remains to show that the one-norm and the infinity-norm of  $H(\mathbf{x})$  are bounded as in (198).

**Bounds on the one-norm and the infinity-norm of  $H$**  We have

$$H_{tu}(\mathbf{x}) = \begin{cases} [1 - [\theta_u + 2\Theta_u^\top \mathbf{x}]x_u] \exp(-[\theta_u + 2\Theta_{u,-u}^\top \mathbf{x}_{-u}]x_u - \Theta_{uu}\bar{x}_u) & \text{if } t = u, \\ -2\Theta_{tu}x_u^2 \exp(-[\theta_u + 2\Theta_{u,-u}^\top \mathbf{x}_{-u}]x_u - \Theta_{uu}\bar{x}_u) & \text{otherwise.} \end{cases} \quad (201)$$

Therefore,

$$\begin{aligned} \|H(\mathbf{x})\|_1 &= \max_{u \in [p]} \sum_{t \in [p]} |H_{tu}(\mathbf{x})| \\ &\stackrel{(201)}{=} \max_{u \in [p]} |1 - [\theta_u + 2\Theta_u^\top \mathbf{x}]x_u| \exp(-[\theta_u + 2\Theta_{u,-u}^\top \mathbf{x}_{-u}]x_u - \Theta_{uu}\bar{x}_u) \\ &\quad + 2 \max_{u \in [p]} x_u^2 \exp(-[\theta_u + 2\Theta_{u,-u}^\top \mathbf{x}_{-u}]x_u - \Theta_{uu}\bar{x}_u) \sum_{t \neq u} |\Theta_{tu}| \\ &\stackrel{(a)}{\leq} (1 + \alpha x_{\max} + 4x_{\max}^2 \zeta) \exp(x_{\max}(\alpha + 2\zeta x_{\max})) \stackrel{(b)}{=} C_{3,\tau} C_{4,\tau}, \end{aligned}$$

where (a) follows from Definition. G.1 along with triangle inequality and Cauchy–Schwarz inequality and (b) follows from (185) and (186). Similarly, we have

$$\begin{aligned} \|H(\mathbf{x})\|_\infty &= \max_{t \in [p]} \sum_{u \in [p]} |H_{tu}(\mathbf{x})| \\ &\stackrel{(201)}{=} \max_{t \in [p]} |1 - [\theta_t + 2\Theta_t^\top \mathbf{x}]x_t| \exp(-[\theta_t + 2\Theta_{t,-t}^\top \mathbf{x}_{-t}]x_t - \Theta_{tt}\bar{x}_t) \\ &\quad + 2 \max_{t \in [p]} \sum_{u \neq t} |\Theta_{tu}| x_u^2 \exp(-[\theta_u + 2\Theta_{u,-u}^\top \mathbf{x}_{-u}]x_u - \Theta_{uu}\bar{x}_u) \\ &\stackrel{(a)}{\leq} (1 + \alpha x_{\max} + 4x_{\max}^2 \zeta) \exp(x_{\max}(\alpha + 2\zeta x_{\max})) \stackrel{(b)}{=} C_{3,\tau} C_{4,\tau}, \end{aligned}$$

where (a) follows from Definition. G.1 along with triangle inequality and Cauchy–Schwarz inequality and (b) follows from (185) and (186).

## References

- A. Abadie and J. Gardeazabal. The economic costs of conflict: A case study of the Basque country. *American economic review*, 93(1):113–132, 2003. (Cited on page 6.)
- A. Abadie, A. Diamond, and J. Hainmueller. Synthetic control methods for comparative case studies: Estimating the effect of California’s tobacco control program. *Journal of the American Statistical Association*, 105(490):493–505, 2010. (Cited on page 6.)
- A. Agarwal, D. Shah, and D. Shen. Synthetic A/B testing using synthetic interventions. *arXiv preprint arXiv:2006.07691*, 2020. (Cited on page 6.)
- S. Aida and D. Stroock. Moment estimates derived from Poincaré and logarithmic Sobolev inequalities. *Mathematical Research Letters*, 1(1):75–86, 1994. (Cited on page 69.)
- J. D. Angrist and J.-S. Pischke. *Mostly harmless econometrics: An empiricist’s companion*. Princeton university press, 2009. (Cited on page 6.)
- D. Arkhangelsky and G. Imbens. The role of the propensity score in fixed effect models. Technical report, National Bureau of Economic Research, 2018. (Cited on pages 4 and 6.)
- D. Arkhangelsky, S. Athey, D. A. Hirshberg, G. W. Imbens, and S. Wager. Synthetic difference-in-differences. *American Economic Review*, 111(12):4088–4118, 2021. (Cited on page 6.)
- S. Athey, J. Tibshirani, and S. Wager. Generalized random forests. *The Annals of Statistics*, 47(2): 1148–1178, 2019. (Cited on page 5.)
- M. Bertrand, E. Duflo, and S. Mullainathan. How much should we trust differences-in-differences estimates? *The Quarterly journal of economics*, 119(1):249–275, 2004. (Cited on page 6.)
- J. Besag. Statistical analysis of non-lattice data. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 24(3):179–195, 1975. (Cited on page 5.)
- R. Bhatia. *Perturbation bounds for matrix eigenvalues*. SIAM, 2007. (Cited on page 58.)
- B. B. Bhattacharya and S. Mukherjee. Inference in Ising models. *Bernoulli*, 24(1):493–525, 2018. (Cited on page 5.)
- G. Bresler. Efficiently learning Ising models on arbitrary graphs. In *Proceedings of the 47th Annual ACM Symposium on Theory of Computing*, pages 771–782, 2015. (Cited on pages 4 and 9.)
- G. Bresler and R.-D. Buhai. Learning restricted Boltzmann machines with sparse latent variables. *Advances in Neural Information Processing Systems*, 33:7020–7030, 2020. (Cited on page 5.)
- G. Bresler, F. Koehler, and A. Moitra. Learning restricted Boltzmann machines via influence maximization. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 828–839, 2019. (Cited on page 5.)
- S. Bubeck et al. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015. (Cited on page 11.)

- R. Busa-Fekete, D. Fotakis, B. Szörényi, and M. Zampetakis. Optimal learning of Mallows block model. In *Conference on Learning Theory*, pages 529–532. PMLR, 2019. (Cited on pages 54 and 74.)
- R. Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997. (Cited on page 25.)
- V. Chandrasekaran, P. A. Parrilo, and A. S. Willsky. Latent variable graphical model selection via convex optimization. *The Annals of Statistics*, 40(4):1935–1967, 2012. (Cited on page 5.)
- S. Chatterjee. Estimation in spin glasses: A first step. *The Annals of Statistics*, 35(5):1931–1946, 2007. (Cited on page 5.)
- Y. Dagan, C. Daskalakis, N. Dikkala, and A. V. Kandiros. Learning Ising models from one or multiple samples. In *Proceedings of the 53rd Annual ACM SIGACT Symposium on Theory of Computing*, pages 161–168, 2021. (Cited on pages 3, 4, 5, 12, 23, 25, 59, 72, 73, and 74.)
- C. Daskalakis, N. Dikkala, and I. Panageas. Regression from dependent observations. In *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, pages 881–889, 2019. (Cited on page 5.)
- R. Dwivedi, K. Tian, S. Tomkins, P. Klasnja, S. Murphy, and D. Shah. Counterfactual inference for sequential experiments. *arXiv preprint arXiv:2202.06891*, 2022a. (Cited on page 6.)
- R. Dwivedi, K. Tian, S. Tomkins, P. Klasnja, S. Murphy, and D. Shah. Doubly robust nearest neighbors in factor models. *arXiv preprint arXiv:2211.14297*, 2022b. (Cited on page 6.)
- I. Fernández-Val and M. Weidner. Fixed effects estimation of large-T panel data models. *Annual Review of Economics*, 10(1):109–138, 2018. doi: 10.1146/annurev-economics-080217-053542. (Cited on page 6.)
- W. Ghang, Z. Martin, and S. Waruhiu. The sharp log-Sobolev inequality on a compact interval. *Involve*, 7:181–186, 2014. (Cited on page 75.)
- P. Ghosal and S. Mukherjee. Joint estimation of parameters in Ising model. *The Annals of Statistics*, 48(2):785–810, 2020. (Cited on page 5.)
- S. Goel. Learning Ising and Potts models with latent variables. In *International Conference on Artificial Intelligence and Statistics*, pages 3557–3566. PMLR, 2020. (Cited on page 5.)
- J. Hartford, G. Lewis, K. Leyton-Brown, and M. Taddy. Deep IV: A flexible approach for counterfactual prediction. In *International Conference on Machine Learning*, pages 1414–1423. PMLR, 2017. (Cited on page 5.)
- M. Hernán and J. Robins. *Causal Inference: What If*. Boca Raton: Chapman & Hall/CRC, 2020. (Cited on page 4.)
- R. Holley and D. Stroock. Logarithmic Sobolev inequalities and stochastic Ising models. *Journal of Statistical Physics*, 46:1159–1194, 1987. (Cited on page 75.)
- B. Holmquist. Moments and cumulants of the multivariate normal distribution. *Stochastic Analysis and Applications*, 6(3):273–278, 1988. (Cited on page 55.)

- G. W. Imbens and J. D. Angrist. Identification and estimation of local average treatment effects. *Econometrica: journal of the Econometric Society*, pages 467–475, 1994. (Cited on page 5.)
- G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015. (Cited on pages 4 and 8.)
- E. T. Jaynes. Information theory and statistical mechanics. *Physical review*, 106(4):620, 1957. (Cited on page 3.)
- A. Jesson, S. Mindermann, Y. Gal, and U. Shalit. Quantifying ignorance in individual-level causal-effect estimates under hidden confounding. In *International Conference on Machine Learning*, pages 4829–4838. PMLR, 2021. (Cited on page 6.)
- Y. Jin, Z. Ren, and E. J. Candès. Sensitivity analysis of individual treatment effects: A robust conformal inference approach. *Proceedings of the National Academy of Sciences*, 120(6):e2214889120, 2023. (Cited on page 6.)
- N. Kallus, X. Mao, and A. Zhou. Interval estimation of individual-level causal effects under unobserved confounding. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 2281–2290. PMLR, 2019. (Cited on page 6.)
- V. Kandiros, Y. Dagan, N. Dikkala, S. Goel, and C. Daskalakis. Statistical estimation from dependent data. In *International Conference on Machine Learning*, pages 5269–5278. PMLR, 2021. (Cited on pages 3, 4, 5, 14, and 25.)
- A. Klivans and R. Meka. Learning graphical models using multiplicative weights. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 343–354. IEEE, 2017. (Cited on pages 4 and 9.)
- M. Ledoux. Logarithmic Sobolev inequalities for unbounded spin systems revisited. In *Séminaire de Probabilités XXXV*, pages 167–194. Springer, 2001. (Cited on page 75.)
- J. Ma and G. Michailidis. Joint structural estimation of multiple graphical models. *The Journal of Machine Learning Research*, 17(1):5777–5824, 2016. (Cited on page 14.)
- S. Ma, L. Xue, and H. Zou. Alternating direction methods for latent variable Gaussian graphical model selection. *Neural computation*, 25(8):2172–2198, 2013. (Cited on page 5.)
- K. Marton. Logarithmic Sobolev inequalities in discrete product spaces: a proof by a transportation cost distance. *arXiv preprint arXiv:1507.02803*, 2015. (Cited on pages 23, 59, 60, and 64.)
- S. Mukherjee, S. Halder, B. B. Bhattacharya, and G. Michailidis. High dimensional logistic regression under network dependence. *arXiv preprint arXiv:2110.03200*, 2021. (Cited on pages 3 and 5.)
- S. N. Negahban, P. Ravikumar, M. J. Wainwright, B. Yu, et al. A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. *Statistical Science*, 27(4):538–557, 2012. (Cited on page 14.)
- J. Neyman. Sur les applications de la théorie des probabilités aux expériences agricoles: Essai des principes. *Roczniki Nauk Rolniczych*, 10(1):1–51, 1923. (Cited on page 2.)

- J. Pearl. *Causality*. Cambridge university press, 2009. (Cited on pages 4, 8, and 24.)
- J. Pearl, M. Glymour, and N. P. Jewell. *Causal inference in statistics: A primer*. John Wiley & Sons, 2016. (Cited on page 4.)
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, et al. Scikit-learn: Machine learning in python. *The Journal of Machine Learning Research*, 12:2825–2830, 2011. (Cited on page 20.)
- C. X. Ren, S. Misra, M. Vuffray, and A. Y. Lokhov. Learning continuous exponential families beyond Gaussian. *arXiv preprint arXiv:2102.09198*, 2021. (Cited on page 5.)
- P. R. Rosenbaum and D. B. Rubin. Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome. *Journal of the Royal Statistical Society: Series B (Methodological)*, 45(2):212–218, 1983. (Cited on page 5.)
- D. B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5):688, 1974. (Cited on page 2.)
- D. B. Rubin. Randomization analysis of experimental data: The Fisher randomization test comment. *Journal of the American Statistical Association*, 75(371):591–593, 1980. (Cited on page 7.)
- N. P. Santhanam and M. J. Wainwright. Information-theoretic limits of selecting binary graphical models in high dimensions. *IEEE Transactions on Information Theory*, 58(7):4117–4134, 2012. (Cited on pages 9 and 14.)
- V. Semenova and V. Chernozhukov. Debiased machine learning of conditional average treatment effects and other causal functions. *The Econometrics Journal*, 24(2):264–289, 2021. (Cited on page 5.)
- A. Shah, D. Shah, and G. Wornell. On learning continuous pairwise markov random fields. In *International Conference on Artificial Intelligence and Statistics*, pages 1153–1161. PMLR, 2021a. (Cited on pages 4, 5, 9, 10, 11, 14, and 23.)
- A. Shah, D. Shah, and G. Wornell. A computationally efficient method for learning exponential family distributions. *Advances in Neural Information Processing Systems*, 34:15841–15854, 2021b. (Cited on pages 3, 4, 10, 12, and 14.)
- A. Shah, D. Shah, and G. Wornell. On computationally efficient learning of exponential family distributions. *arXiv preprint arXiv:2309.06413*, 2023. (Cited on pages 14 and 56.)
- R. Singh, M. Sahani, and A. Gretton. Kernel instrumental variable regression. *Advances in Neural Information Processing Systems*, 32:4593–4605, 2019. (Cited on page 5.)
- V. Syrgkanis, V. Lei, M. Oprescu, M. Hei, K. Battocchi, and G. Lewis. Machine learning estimation of heterogeneous treatment effects with instruments. *Advances in Neural Information Processing Systems*, 32, 2019. (Cited on page 5.)
- A. Taeb, P. Shah, and V. Chandrasekaran. Learning exponential family graphical models with latent variables using regularized conditional likelihood. *arXiv preprint arXiv:2010.09386*, 2020. (Cited on pages 5 and 8.)

- V. Y. Tan, A. Anandkumar, and A. S. Willsky. Learning Gaussian tree models: Analysis of error exponents and extremal structures. *IEEE Transactions on Signal Processing*, 58(5):2701–2714, 2010. (Cited on page 55.)
- W. F. Trench. Asymptotic distribution of the spectra of a class of generalized Kac–Murdock–Szegő matrices. *Linear algebra and its applications*, 294(1-3):181–192, 1999. (Cited on page 55.)
- C. Villani. *Optimal transport: old and new*, volume 338. Springer, 2009. (Cited on pages 62 and 66.)
- M. Vinyes and G. Obozinski. Learning the effect of latent variables in Gaussian graphical models with unobserved variables. *arXiv preprint arXiv:1807.07754*, 2018. (Cited on page 5.)
- M. Vuffray, S. Misra, A. Lokhov, and M. Chertkov. Interaction screening: Efficient and sample-optimal learning of Ising models. *Advances in Neural Information Processing Systems*, 29, 2016. (Cited on pages 4, 5, 9, 10, 11, and 14.)
- M. Vuffray, S. Misra, and A. Y. Lokhov. Efficient learning of discrete graphical models. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124017, 2022. (Cited on pages 4, 5, 9, 10, 11, and 14.)
- M. J. Wainwright, M. I. Jordan, et al. Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1–2):1–305, 2008. (Cited on pages 4, 10, 11, and 15.)
- G. Wang, J. Li, and W. J. Hopp. An instrumental variable forest approach for detecting heterogeneous treatment effects in observational studies. *Management Science*, 68(5):3399–3418, 2022. (Cited on page 5.)
- K. Wang, A. Franks, and S.-Y. Oh. Learning Gaussian graphical models with latent confounders. *Journal of Multivariate Analysis*, 198:105213, 2023. (Cited on page 5.)
- S. Wilhelm and B. Manjunath. tmvtnorm: A package for the truncated multivariate normal distribution. *SIGMA*, 2(2):1–25, 2010. (Cited on page 20.)
- J. H. Won and S.-J. Kim. Maximum likelihood covariance estimation with a condition number constraint. In *2006 Fortieth Asilomar Conference on Signals, Systems and Computers*, pages 1445–1449. IEEE, 2006. (Cited on page 14.)
- L. Xu, Y. Chen, S. Srinivasan, N. de Freitas, A. Doucet, and A. Gretton. Learning deep features in instrumental variable regression. *arXiv preprint arXiv:2010.07154*, 2020. (Cited on page 5.)
- S. Yadlowsky, H. Namkoong, S. Basu, J. Duchi, and L. Tian. Bounds on the conditional and average treatment effect with unobserved confounding factors. *The Annals of Statistics*, 50(5):2587–2615, 2022. (Cited on page 6.)
- M. Yin, C. Shi, Y. Wang, and D. M. Blei. Conformal sensitivity analysis for individual treatment effects. *Journal of the American Statistical Association*, pages 1–14, 2022. (Cited on page 6.)
- S. Zhou, P. Rütimann, M. Xu, and P. Bühlmann. High-dimensional covariance estimation based on Gaussian graphical models. *The Journal of Machine Learning Research*, 12:2975–3026, 2011. (Cited on page 14.)