# Uncertainty Quantification for Large Language Model Reward Learning under Heterogeneous Human Feedback

Pangpang Liu*    Junwei Lu†    Will Wei Sun‡

December 4, 2025

## Abstract

We study estimation and statistical inference for reward models used in aligning large language models (LLMs). A key component of LLM alignment is reinforcement learning from human feedback (RLHF), where humans compare pairs of model-generated answers and their preferences are used to train a reward model. However, human feedback is inherently heterogeneous, creating significant challenges for reliable reward learning. To address this, we adopt a heterogeneous preference framework that jointly models the latent reward of answers and human rationality. This leads to a challenging biconvex optimization problem, which we solve via an alternating gradient descent algorithm. We establish theoretical guarantees for the resulting estimator, including its convergence and asymptotic distribution. These results enable the construction of confidence intervals for reward estimates. Leveraging these uncertainty quantification results, we conduct valid statistical comparisons between rewards and incorporate uncertainty into the best-of-$N$ (BoN) policy framework. Extensive simulations demonstrate the effectiveness of our method, and applications to real LLM data highlight the practical value of accounting for uncertainty in reward modeling for LLM alignment.

**Key Words:** Heterogeneous human feedback; LLMs; Nonconvex optimization; RLHF; Statistical inference

*Department of Biostatistics, Yale University. Email: pangpang.liu@yale.edu.

†Department of Biostatistics, Harvard University. Email: junweilu@hsph.harvard.edu.

‡Department of Quantitative Methods, Purdue University. Email: sun244@purdue.edu. Corresponding author.

# 1 Introduction

Large language models have demonstrated remarkable capabilities across a wide range of applications (Nakada et al., 2024; Chen et al., 2025; Su, 2025; Ji et al., 2025; Zhou et al., 2025). To align these models with human preferences, reinforcement learning from human feedback has become a cornerstone technique (Xiao et al., 2025; Ye et al., 2025; Xu et al., 2025). While large language models exhibit impressive generalization capabilities, their raw generations often fail to reflect human values. RLHF addresses this gap by leveraging human feedback to train a reward model that evaluates model outputs and guides subsequent fine-tuning through reinforcement learning. This process enables LLMs to produce responses that are more aligned with human values.

A major challenge in RLHF is the heterogeneity of human preferences. For example, in the `prism-alignment` dataset (Kirk et al., 2024), feedback is collected from 1,500 participants spanning 75 countries. These individuals differ in religion, education, employment, and other demographics, and exhibit varying degrees of expertise, attentiveness, and rationality (Park et al., 2024; Zeng et al., 2024). Such diversity introduces significant heterogeneity into the training data. Failing to account for such differences can lead to misaligned or suboptimal reward models (Zhong et al., 2024; Chakraborty et al., 2024). To address this issue, we employ a scale heterogeneity model (see (4) for the details), which explicitly accounts for variations in human rationality during reward learning. Simultaneously learning both the reward function and human rationality leads to a nonconvex optimization problem, which poses significant computational and theoretical challenges.

Another critical yet often overlooked issue is the uncertainty inherent in the reward model learned from heterogeneous human feedback. Ignoring this uncertainty can lead to unreliable evaluations. One motivating example arises in computing pairwise ranking accuracy (Lambert et al., 2025; Frick et al., 2025). In this task, the goal is to compare two answers to the same

question using the learned reward model. Relying solely on point estimates of the rewards provides no information about the confidence in the comparison: whether one answer is truly better than the other or if the difference could be due to estimation noise. A more robust approach would account for the uncertainty in the estimated rewards during the comparison. From a statistical perspective, this type of task can be framed as a hypothesis test problem. Let $r_{\theta_*}(s, a^{(0)})$ and $r_{\theta_*}(s, a^{(1)})$ represent the expected rewards for answers $a^{(0)}$ and $a^{(1)}$ to question $s$ under the true reward parameter $\theta_*$. We formalize the comparison between these two answers within the framework of the hypothesis test.

**Example 1.** *Given a question $s$ and two answers $a^{(0)}$ and $a^{(1)}$ from two different LLMs, we test*

$$H_0 : r_{\theta_*}(s, a^{(0)}) - r_{\theta_*}(s, a^{(1)}) = 0 \quad vs. \quad H_1 : r_{\theta_*}(s, a^{(0)}) - r_{\theta_*}(s, a^{(1)}) \neq 0. \quad (1)$$

To address this challenge, we first estimate a reward model $r_{\theta_T}(s, a)$, where $\theta_T$ serves as an estimator of the true parameter $\theta_*$. We then construct confidence intervals to guide the decision-making process. Our testing procedure is detailed in Section 5.1, and empirical results are reported in Section 6.2.2.

Another important application of reward inference lies in improving the best-of-$N$ policy using reward lower bounds. The BoN sampling refers to generating $N$ candidate answers $a^{(1)}, \cdots, a^{(N)}$ for a given question $s$ and selecting the best one based on reward estimates (Stiennon et al., 2020; Nakano et al., 2022; Gao et al., 2023b; Jinnai et al., 2024; Gui et al., 2024; Liu et al., 2025; Chow et al., 2025). Traditional BoN strategies use point estimates of rewards without accounting for uncertainty.

However, in offline reinforcement learning, distributional shift is common due to limited coverage of the state-action space. Pessimistic policies are often used to mitigate this issue (Zhou et al., 2023; Dong et al., 2023; Lu et al., 2023; Bian et al., 2024; Jin et al., 2025; Zhu et al., 2025). Inspired by this, we can modify the BoN objective to maximize the lower

confidence bound of the reward rather than its point estimate.

**Example 2.** *Suppose the asymptotic* $(1-\alpha)$ *confidence interval for* $r_{\theta_*}(s,a)$ *is* $(\mathcal{C}_l^\alpha(s,a), \mathcal{C}_u^\alpha(s,a))$. *The traditional BoN objective is*

$$a_{BoN}(s) = \underset{a \in \mathcal{A}_N(s)}{\arg\max}\, r_{\theta_T}(s,a), \tag{2}$$

*where* $\mathcal{A}_N(s) = \{a^{(1)}, \cdots, a^{(N)}\}$ *is the set of* $N$ *candidate answers corresponding to* $s$. *The pessimistic BoN (pBoN) policy that maximizes the lower bound of the reward is*

$$a_{pBoN}(s) = \underset{a \in \mathcal{A}_N(s)}{\arg\max}\, \mathcal{C}_l^\alpha(s,a). \tag{3}$$

In Section 5.2, we show that the expected suboptimality of the pessimistic BoN policy decreases at the rate of $1/\sqrt{n}$, where $n$ denotes the sample size. Here, the suboptimality (see (11)) quantifies the performance gap between the proposed policy and the oracle policy. In Section 6.2.3, we empirically demonstrate the pessimistic BoN policy in (3) can improve the performance of the standard approach in (2) across various settings of BoN policies.

## 1.1 Major Contribution

In this paper, we model the rationality of the human annotator as a function of contextual information and design an alternating gradient descent algorithm to jointly learn both the reward and rationality models. We then derive the asymptotic distribution of the reward estimates and construct statistical confidence bounds to guide robust decision-making. The derived uncertainty quantification of reward models can account for the diverse expertise levels of human annotators. Our main contributions are as follows.

- Methodologically, we employ a heterogeneous rationality model that captures varying levels of human expertise. To jointly learn both the rationality and the reward model, we develop an alternating gradient descent algorithm.
- Theoretically, we establish convergence guarantees for the proposed alternating gradient descent algorithm and derive the asymptotic distribution of the estimators, enabling

uncertainty quantification of the learned reward. As a byproduct, we propose the pessimistic BoN policy derived from the asymptotic lower bound and establish its suboptimality.

- Numerically, we validate our approach through extensive simulations and real-world applications to large language models, demonstrating its effectiveness and offering valuable insights to improve the best-of-$N$ policy.

## 1.2   Related Literature

Our work intersects with two major areas of literature: modeling heterogeneous human preferences and uncertainty quantification in the Bradley-Terry-Luce (BTL) model (Bradley and Terry, 1952).

**Heterogeneous Human Preferences.** A number of studies have addressed heterogeneity in human preferences in the context of rank aggregation, including Deng et al. (2014); Jin et al. (2020); Li et al. (2020, 2022); Zhu et al. (2023b). These works primarily considered the ranking of a fixed set of items and are thus not directly applicable to our setting, which involves comparisons over potentially unbounded LLM outputs. In the RLHF literature, some studies such as Barnett et al. (2023); Hao et al. (2023); Liu et al. (2024) incorporated heterogeneous preferences by assuming fixed, known parameters that represent each teacher's area of expertise. In contrast, our model captures heterogeneity through an unknown function of contextual information that must be learned, offering a more flexible and general approach. Lee et al. (2024); Zhong et al. (2024); Park et al. (2024); Wang et al. (2025a) proposed personalized reward models that tailor preferences to individual annotators. However, their formulation involves assigning separate reward models to each teacher, which differs from our objective of learning a unified reward model that generalizes across diverse teacher rationalities. More importantly, these prior works concentrate on reward estimation alone, whereas we additionally provide principled uncertainty quantification for the reward learning.

**Uncertainty Quantification for the BTL Model.** Simons and Yao (1999); Han et al.

(2020); Gao et al. (2023a); Liu et al. (2023); Fan et al. (2025b) investigated uncertainty quantification in the BTL model without covariates. Fan et al. (2024a,b) extended the inference to covariate-assisted BTL models. These studies assume a fixed number of items and constant latent scores, which does not align with our setting where LLMs generate a large and varying set of outputs. Recent works, including Wang et al. (2025b); Lu et al. (2025); Li and Li (2025); Zhang et al. (2025), have developed contextual ranking frameworks to evaluate and compare different LLMs under uncertainty. Fan et al. (2025a) further explored uncertainty quantification for ranking with heterogeneous preferences in a finite-item setting. These studies focus on assigning a single score to each LLM. In contrast, our work compares multiple outputs generated by LLMs, where the reward is defined for each question–answer pair rather than at the model level. Zhu et al. (2023a); Feng et al. (2025) analyzed the asymptotic properties of the maximum likelihood estimator (MLE) under the standard BTL model. Due to the nonconvex nature of our setting, the MLE is not directly attainable. Instead, we quantify the uncertainty of the proposed, attainable gradient-based estimator.

## 1.3 Paper Organization

The remainder of the paper is organized as follows. Section 2 introduces the preliminary setup of the problem. Section 3 presents the framework for reward learning from heterogeneous human feedback. Theoretical results on convergence analysis and uncertainty quantification are provided in Section 4. We introduce the related applications in Section 5. Section 6 reports the experimental results. All proofs are deferred to the supplementary material.

# 2 Problem Formulation

We denote $s$ as a prompt (or question) and $a$ as a corresponding response (or answer). The expected reward associated with the pair $(s, a)$ is defined by $r_{\theta_*}(s, a)$, where $\theta_* \in \mathbb{R}^{d_1}$, with $d_1$ fixed, is the true but unknown parameter. Let $x$ represent contextual information that

influences a human teacher's preference rationality. We model the teacher's utility for a given prompt-response pair $(s, a)$ using the scale heterogeneity model (Fiebig et al., 2010):

$$U(x, s, a) = \sigma_{\gamma_*}(x) r_{\theta_*}(s, a) + \epsilon, \tag{4}$$

where $\epsilon$ is independent and identically distributed (i.i.d.) and follows a Gumbel distribution with location parameter 0 and scale parameter 1, and $\sigma_{\gamma_*}(x) \in \mathbb{R}$ is a scale function that captures the human teacher's rationality, parameterized by the unknown vector $\gamma_* \in \mathbb{R}^{d_2}$ with $d_2$ fixed. This model allows the same prompt-response pair to yield different utility outcomes across teachers with different contextual information $x$. We focus on pairwise comparisons between two candidate responses $a^{(0)}$ and $a^{(1)}$ for a given prompt $s$. The probability that a teacher with context $x$ prefers $a^{(1)}$ over $a^{(0)}$ is given by

$$
\begin{aligned}
\mathbb{P}(Y = 1 | x, s, a^{(0)}, a^{(1)}) &= \mathbb{P}[U(x, s, a^{(1)}) > U(x, s, a^{(0)}) \mid x, s, a^{(0)}, a^{(1)}] \\
&= \mathbb{P}\{\epsilon^{(1)} - \epsilon^{(0)} > \sigma_{\gamma_*}(x)[r_{\theta_*}(s, a^{(0)}) - r_{\theta_*}(s, a^{(1)})] \mid x, s, a^{(0)}, a^{(1)}\} \\
&= \frac{1}{1 + e^{-\sigma_{\gamma_*}(x)[r_{\theta_*}(s, a^{(1)}) - r_{\theta_*}(s, a^{(0)})]}},
\end{aligned} \tag{5}
$$

where the last equality follows from the fact that the difference between two independent Gumbel random variables with the same scale parameter follows the logistic distribution (Kotz and Nadarajah, 2000). Here, $Y = 1$ indicates that the teacher prefers $a^{(1)}$, and $Y = 0$ indicates a preference for $a^{(0)}$. The heterogeneity captured by model (5) has been supported by empirical evidence (Fiebig et al., 2010; Davis et al., 2019; Tutz, 2021; Mauerer and Tutz, 2023). However, these studies typically focus on small-scale applications with a special structure of $\sigma_{\gamma_*}(x)$ and lack theoretical guarantees for estimation accuracy or uncertainty quantification. The value of $\sigma_{\gamma_*}(x)$ reflects the rationality or expertise of the teacher. If $\sigma_{\gamma_*}(x) > 0$, the teacher is more likely to select the response with a higher true reward. If $\sigma_{\gamma_*}(x) \leq 0$, the teacher exhibits non-expert or even adversarial behavior, potentially preferring lower-reward responses. Such irrational preference labels exist in real-world datasets (Nahum et al., 2024; Bukharin et al., 2024). A special case arises when $\sigma_{\gamma_*} \equiv 1$, reducing the model

to a homogeneous preference BTL setting (Ouyang et al., 2022). Introducing the scale function enhances the expressiveness of the model, but also introduces the challenge of jointly estimating both $\gamma_*$ and $\theta_*$. In the following sections, we develop algorithms to learn these parameters from human feedback and quantify the uncertainty of the estimators.

# 3    Reward Learning from Heterogeneous Human Feedback

In this section, we present an algorithm for learning the reward function from heterogeneous human feedback. We begin by introducing some structural assumptions on both the reward and scale functions.

**Assumption 1.** *The reward function lies within a family of linear models $r_\theta(s, a) = \theta^\top \phi(s, a)$, where $\phi(s, a) : \mathcal{S} \times \mathcal{A} \to \mathbb{R}^{d_1}$ is a known feature mapping satisfying $\|\phi(s, a)\|_2 \leq \phi_{\max}$ for some positive constant $\phi_{\max}$. The true parameter satisfies $\theta_* \in \Theta = \{\theta : \|\theta\|_\infty \leq \theta_{\max}, \theta \in \mathbb{R}^{d_1}\}$ for some positive constant $\theta_{\max}$.*

Assumption 1 posits that the reward function is linear in the parameter $\theta$, and both the parameter and features are bounded. This assumption is standard in the RLHF literature (Zhu et al., 2023a; Zhong et al., 2024; Scheid et al., 2024; Lu et al., 2025). In LLM applications, the feature mapping $\phi(s, a)$ can be obtained by removing the final layer from a pre-trained model (Zhu et al., 2023a; Zhong et al., 2024; Lu et al., 2025).

**Assumption 2.** *Let $\psi_0 : \mathcal{X} \to \mathbb{R}$ and $\psi(x) : \mathcal{X} \to \mathbb{R}^{d_2}$ be known feature mappings. We assume $\psi_0 \not\equiv 0$ and $\sup_{x \in \mathcal{X}} \|(\psi_0(x), \psi^\top(x))\|_2 \leq \psi_{\max}$ for some positive constant $\psi_{\max}$. The true rationality scale function is $\sigma_{\gamma_*}(x) = \psi_0(x) + \gamma_*^\top \psi(x)$, where $\gamma_* \in \Gamma = \{\gamma : \|(1, \gamma^\top)\|_\infty \leq \gamma_{\max}, \gamma \in \mathbb{R}^{d_2}\}$ for some positive constant $\gamma_{\max}$.*

Assumption 2 is structurally parallel to Assumption 1. The feature mappings $\psi_0(x)$ and $\psi(x)$ serve as known transformations of the contextual variable $x$, which enrich the representation of features across individuals. The condition $\psi_0(x) \not\equiv 0$ is essential for the identifiability of

both $\gamma_*$ and $\theta_*$. Specifically, if $\psi_0 \equiv 0$, the scale function $\sigma_{\gamma_*}(x) = \gamma_*^\top \psi(x)$ and the reward difference $\theta_*^\top[\phi(s, a^{(1)}) - \phi(s, a^{(0)})]$ would be indistinguishable in (5). The presence of a nonzero $\psi_0(x)$ ensures that $\gamma_*$ and $\theta_*$ can be distinguished from the data.

With the model specifications in Assumptions 1 and 2, we now turn to the estimation of the reward parameter $\theta_*$ and the rationality parameter $\gamma_*$. The goal is to infer these parameters jointly from observed preference data. To this end, we introduce the likelihood formulation associated with the heterogeneous preference model. Suppose that we observe a dataset of $n$ i.i.d. samples $\{(x_i, s_i, a_i^{(0)}, a_i^{(1)}, y_i)\}_{i=1}^n$. We denote $\mu(v) = \frac{1}{1+e^{-v}}$ for $v \in \mathbb{R}$. By (5), the negative log-likelihood for this dataset is

$$
\begin{aligned}
L_n(\theta, \gamma) = -\frac{1}{n} \sum_{i=1}^n \Big\{ &y_i \log \mu(\sigma_\gamma(x_i)(r_\theta(s_i, a_i^{(1)}) - r_\theta(s_i, a_i^{(0)}))) \\
&+ (1 - y_i) \log[1 - \mu(\sigma_\gamma(x_i)(r_\theta(s_i, a_i^{(1)}) - r_\theta(s_i, a_i^{(0)})))] \Big\}.
\end{aligned}
\tag{6}
$$

A natural solution to estimate $\theta_*$ and $\gamma_*$ is via the maximum likelihood estimator as follows,

$$
\widehat{\theta}_n, \widehat{\gamma}_n = \underset{\theta \in \Theta, \gamma \in \Gamma}{\arg\min}\, L_n(\theta, \gamma).
\tag{7}
$$

However, the negative log-likelihood $L_n(\theta, \gamma)$ is generally not jointly convex in the parameters $(\theta^\top, \gamma^\top)^\top$, as formalized in the following lemma.

**Lemma 1.** *Under Assumptions 1 and 2, the negative log-likelihood function $L_n(\theta, \gamma)$ defined in (6) is convex in $\theta$ when $\gamma$ is fixed, and convex in $\gamma$ when $\theta$ is fixed. However, $L_n(\theta, \gamma)$ is not necessarily jointly convex in $(\theta^\top, \gamma^\top)^\top$.*

The nonconvex nature of $L_n(\theta, \gamma)$ with respect to the joint parameter vector $(\theta^\top, \gamma^\top)^\top$ makes solving the MLE problem in (7) nontrivial. However, the biconvexity property identified in Lemma 1 motivates the use of an alternating optimization approach. Specifically, we apply alternating gradient descent to iteratively update $\theta$ and $\gamma$ described in Algorithm 1.

Algorithm 1 requires inputs: learning rates $\eta_1$ and $\eta_2$, initial parameters $\gamma_0$ and $\theta_0$, the number of iterations $T$, and a dataset $\{(x_i, s_i, a_i^{(0)}, a_i^{(1)}, y_i)\}_{i=1}^n$. At each iteration $t = 1, \ldots, T$, the

---
**Algorithm 1** Reward Learning with Alternating Gradient Descent
---
1: **Input**: learning rates $\eta_1$ and $\eta_2$, initial points $\gamma_0$ and $\theta_0$, number of iterations $T$, dataset $\{(x_i, s_i, a_i^{(0)}, a_i^{(1)}, y_i)\}_{i=1}^n$.

2: $z_i = \phi(s_i, a_i^{(1)}) - \phi(s_i, a_i^{(0)})$ for $i = 1, \cdots, n$.

3: **for** $t = 1$ to $T$ **do**

4:     $\theta_t = \theta_{t-1} - \eta_1 \nabla_\theta L_n(\theta_{t-1}, \gamma_{t-1})$, where

$$\nabla_\theta L_n(\theta_{t-1}, \gamma_{t-1}) = -\frac{1}{n} \sum_{i=1}^n [y_i - \mu(\gamma_{t-1}^\top \psi(x_i)(\theta_{t-1}^\top z_i))][\gamma_{t-1}^\top \psi(x_i)]z_i.$$

5:     $\gamma_t = \gamma_{t-1} - \eta_2 \nabla_\gamma L_n(\theta_t, \gamma_{t-1})$, where

$$\nabla_\gamma L_n(\theta_t, \gamma_{t-1}) = -\frac{1}{n} \sum_{i=1}^n [y_i - \mu(\gamma_{t-1}^\top \psi(x_i)(\theta_t^\top z_i))](\theta_t^\top z_i)\psi(x_i).$$

6: **end for**

7: **Output**: $\theta_T, \gamma_T$
---

algorithm performs alternating gradient updates. The notations $\nabla_\theta L_n(\theta, \gamma)$ and $\nabla_\gamma L_n(\theta, \gamma)$ denote the gradients of $L_n(\theta, \gamma)$ relative to $\theta$ and $\gamma$, respectively. After $T$ iterations, the algorithm outputs the final parameters $(\theta_T^\top, \gamma_T^\top)^\top$.

# 4    Theoretical Results

In this section, we present the theoretical guarantees of Algorithm 1. Section 4.1 establishes the convergence properties of our estimator, while Section 4.2 discusses the statistical inference results.

## 4.1    Convergence Analysis

We first give the assumptions needed for the convergence analysis. We denote $z = \phi(s, a^{(1)}) - \phi(s, a^{(0)})$.

**Assumption 3.** *The smallest eigenvalue of the matrix* $\mathbb{E}[\sigma_{\gamma_*}^2(x)zz^\top]$ *is* $\lambda_\phi > 0$. *Similarly, the smallest eigenvalue of the matrix* $\mathbb{E}[(\theta_*^\top z)^2 \psi(x)\psi^\top(x)]$ *is* $\lambda_\psi > 0$.

Assumption 3 is a well-conditioned design for the covariance matrix. Zhong et al. (2024)

adopted the same assumption on $\mathbb{E}[\sigma^2_{\gamma_*}(x)zz^\top]$ with $\sigma_{\gamma_*} \equiv 1$. In the context of LLMs, the feature $x$ represents contextual information (e.g., gender, age, or other demographic factors), while $z$ is associated with the prompt $s$ and candidate responses $a^{(0)}$ and $a^{(1)}$. In this setting, $x$ and $z$ are independent. Then, we have the factorization $\mathbb{E}[\sigma^2_{\gamma_*}(x)zz^\top] = \mathbb{E}[\sigma^2_{\gamma_*}(x)]\mathbb{E}(zz^\top)$ and $\mathbb{E}[(\theta_*^\top z)^2 \psi(x)\psi^\top(x)] = \mathbb{E}(\theta_*^\top z)^2 \mathbb{E}[\psi(x)\psi^\top(x)]$. This factorization softens the conditions in Assumption 3, making them more interpretable: they essentially require the variability in both contextual features $(x)$ and prompts/responses $(z)$ to be sufficiently rich.

**Assumption 4.** *Define $M = \|\mathbb{E}\{\mu(\sigma_{\gamma_*}(x)\theta_*^\top z)[1 - \mu(\sigma_{\gamma_*}(x)\theta_*^\top z)]\sigma_{\gamma_*}(x)(\theta_*^\top z)\psi(x)z^\top\}\|_2$. We assume $M < \frac{\min\{\lambda_\phi, \lambda_\psi\}}{6c_0}$, where $c_0$ is a positive constant defined in Lemma S6.*

Assumption 4 controls the interaction strength between $\psi(x)$ and $z$. The quantity $M$ measures the magnitude of their cross-dependence in the Hessian matrix of the log-likelihood. A smaller $M$ indicates that the coupling between the two parameter blocks is weak, which helps ensure that the alternating estimation of $\theta_*$ and $\gamma_*$ is stable.

Now, we provide the following theorem on the convergence of $(\theta_T^\top, \gamma_T^\top)^\top$ to the true parameters $(\theta_*^\top, \gamma_*^\top)^\top$.

**Theorem 1.** *Let Assumptions 1, 2, 3 and 4 hold. Suppose that the initialization in Algorithm 1 satisfies $\|\theta_0 - \theta_*\|_2 \le b/\sqrt{2}$ and $\|\gamma_0 - \gamma_*\|_2 \le b/\sqrt{2}$, where $b < b_0$ for some positive constant $b_0$. Assume further that the step size $0 < \eta_1, \eta_2 < \eta_0$ for some positive constant $\eta_0$. There exist some positive constants $c_1, n_0$ such that when $n > n_0$, for any $0 < \delta < 1$, with probability at least $1 - \delta - \frac{1}{n}$, we have*

$$\|\theta_T - \theta_*\|_2^2 + \|\gamma_T - \gamma_*\|_2^2 \le \rho^T b^2 + \frac{c_1 \log(1/\delta)}{(1-\rho)n},$$

*for some $\rho \in (0, 1)$.*

Theorem 1 quantifies the convergence rate of Algorithm 1. The first term on the right-hand side represents the optimization error, which vanishes as the number of iterations $T$ increases to $\infty$. The second term indicates the statistical error, which decreases with the sample size $n$

at the scale of $1/n$. The dependence on $T$ and $n$ aligns with the results in Jin et al. (2020), which studied a non-contextual setting with a finite number of items. Notably, our setting is more complex, yet we achieve the same convergence rate.

### 4.1.1 Challenges and Outline of Proof of Theorem 1

The proof of Theorem 1 presents substantial technical challenges because jointly estimating $(\theta_*^\top, \gamma_*^\top)^\top$ involves a nonconvex objective. The loss function $L_n(\theta, \gamma)$ couples $\theta$ and $\gamma$ multiplicatively through the scale function $\sigma_{\gamma_*}(x)$ and the reward model $r_{\theta_*}(s, a)$, which causes non-convexity and prevents the application of classical convex optimization results.

To overcome the challenges, we exploit the biconvex structure of $L_n(\theta, \gamma)$ and analyze the alternating gradient descent dynamics in a two-block fashion. Specifically, according to the update rule in Algorithm 1, we derive separate recursive bounds for each block:

$$\|\theta_{t+1} - \theta_*\|_2^2 = \|\theta_t - \theta_*\|_2^2 + \eta_1^2 \|\nabla_\theta L_n(\theta_t, \gamma_t)\|_2^2 - 2\eta_1 \langle \nabla_\theta L_n(\theta_t, \gamma_t), \theta_t - \theta_* \rangle,$$

$$\|\gamma_{t+1} - \gamma_*\|_2^2 = \|\gamma_t - \gamma_*\|_2^2 + \eta_2^2 \|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t)\|_2^2 - 2\eta_2 \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t), \gamma_t - \gamma_* \rangle.$$

A central challenge is that these bounds depend on both $t$ and $n$, making it nontrivial to separate the optimization error from the statistical error. Carefully disentangling these dependencies and controlling the cross-block error propagation are the key difficulties. Our strategy is to decompose the gradients $\nabla_\theta L_n(\theta_t, \gamma_t)$ and $\nabla_\gamma L_n(\theta_{t+1}, \gamma_t)$ into several components: $\nabla_\theta L_n(\theta_*, \gamma_t), \nabla_\theta L_n(\theta_t, \gamma_*), \nabla_\theta L_n(\theta_*, \gamma_*), \nabla_\gamma L_n(\theta_*, \gamma_t), \nabla_\gamma L_n(\theta_{t+1}, \gamma_*), \nabla_\gamma L_n(\theta_*, \gamma_*)$. This decomposition allows us to isolate the optimization dynamics (through deviations from $\theta_*$ or $\gamma_*$) from the stochastic fluctuations due to finite $n$.

We establish a series of technical lemmas to control each component. Lemma S4 ensures that the gradients of $L_n(\theta, \gamma)$ are Lipschitz in both $\theta$ and $\gamma$, thereby limiting the error accumulation during descent. Lemma S5 quantifies the sensitivity of each block's gradient to perturbations in the other, ensuring that alternating updates do not destabilize the descent. Lemma S8 further

bounds the mixed second-order derivatives $\nabla^2_{\gamma\theta} L_n(\theta, \gamma)$, which are essential for controlling cross-block error propagation. These results together ensure the interdependence between two parameter blocks remains controlled throughout iterations.

After obtaining separate contraction inequalities for $\|\theta_{t+1} - \theta_*\|_2^2$ and $\|\gamma_{t+1} - \gamma_*\|_2^2$, we combine them to study the joint error $\|\theta_{t+1} - \theta_*\|_2^2 + \|\gamma_{t+1} - \gamma_*\|_2^2$. To guarantee that this joint error decreases over $t$, we derive precise conditions on the learning rates $\eta_1, \eta_2$ and on the initialization radius $b$, ensuring local contraction. Establishing these conditions allows us to rigorously prove both the geometric convergence of the optimization iterates and the statistical consistency of the resulting estimators.

## 4.2   Uncertainty Quantification of Reward Learning

In this section, we derive the asymptotic distribution of the estimators obtained from Algorithm 1. For notational convenience, define $\tau_* = (\theta_*^\top, \gamma_*^\top)^\top$ and $\tau_T = (\theta_T^\top, \gamma_T^\top)^\top$. We introduce the Fisher information matrix at the true parameter $\tau_*$ as

$$\mathcal{I}(\tau_*) = \begin{pmatrix} \mathcal{I}_{\theta\theta}(\theta_*, \gamma_*) & \mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*) \\ \mathcal{I}_{\gamma\theta}(\theta_*, \gamma_*) & \mathcal{I}_{\gamma\gamma}(\theta_*, \gamma_*) \end{pmatrix} = \begin{pmatrix} \mathbb{E}\nabla^2_{\theta\theta} L_n(\theta_*, \gamma_*) & \mathbb{E}\nabla^2_{\theta\gamma} L_n(\theta_*, \gamma_*) \\ \mathbb{E}\nabla^2_{\gamma\theta} L_n(\theta_*, \gamma_*) & \mathbb{E}\nabla^2_{\gamma\gamma} L_n(\theta_*, \gamma_*) \end{pmatrix}$$

where the expectation is taken with respect to the randomness in the preference outcome $y$, the prompt $s$, and the responses $a^{(0)}, a^{(1)}$. Here, $\nabla^2_{\theta\theta} L_n(\theta_*, \gamma_*)$ denotes the block of second-order derivatives with respect to $\theta$, and similarly for the other blocks. In order to study the asymptotic properties of our estimators, it is important to ensure that the Fisher information matrix is well-conditioned. The following assumption formalizes this requirement.

**Assumption 5.** *The matrix* $\widetilde{\mathcal{I}}(\theta_*, \gamma_*) = \mathcal{I}_{\theta\theta}(\theta_*, \gamma_*) - \mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*)\mathcal{I}_{\gamma\gamma}^{-1}(\theta_*, \gamma_*)\mathcal{I}_{\gamma\theta}(\theta_*, \gamma_*)$ *is positive definite.*

A sufficient condition for Assumption 5 is $\mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*) = \mathbf{0}$. Assumption 5 guarantees that the information matrix $\mathcal{I}(\tau_*)$ is positive definite, as shown in the following lemma.

**Lemma 2.** *Suppose Assumptions 1, 2, 3 and 5 hold. Then the information matrix $\mathcal{I}(\tau_*)$ is*

*positive definite.*

Lemma 2 ensures that the information matrix is positive definite, and therefore invertible. This property is fundamental in asymptotic theory, as the inverse of the Fisher information characterizes the asymptotic covariance matrix of the estimators. Hence, the lemma plays a crucial role in establishing the asymptotic normality and quantifying the statistical uncertainty of the learned parameters.

We now establish the asymptotic normality of $\tau_T$. Recall that $\tau_T$ is obtained after $T$ iterations of Algorithm 1, and its behavior depends jointly on the sample size $n$ and the iteration horizon $T$. The next theorem shows that, under suitable conditions on the growth of $T$ relative to $n$, $\tau_T$ shares the same asymptotic distribution as the MLE.

**Theorem 2.** *Let the conditions in Theorems 1 and Assumption 5 hold. Assume $n\rho_1^T \to 0$ as $n \to \infty$ and $T \to \infty$, where $0 < \rho_1 < 1$ is some constant. Then, as $n \to \infty$ and $T \to \infty$, we have*

$$\sqrt{n}(\tau_T - \tau_*) \xrightarrow{D} \mathcal{N}(0, \mathcal{I}^{-1}(\tau_*)).$$

When $T \to \infty$ and $n \to \infty$, Theorem 1 guarantees that $\tau_T$ converges to $\tau_*$. Theorem 2 further establishes that $\tau_T$ not only converges but also satisfies asymptotic normality. This result confirms that the proposed estimator enjoys both consistency and asymptotic efficiency under the stated conditions.

Having established the asymptotic normality of the joint parameter estimator $\tau_T = (\theta_T^\top, \gamma_T^\top)^\top$ in Theorem 2, we now turn to the implications for the learned reward function itself. Since the reward is parameterized only through $\theta$, the asymptotic distribution of $r_{\theta_T}(s, a)$ can be derived accordingly. The following theorem characterizes this distribution.

**Theorem 3.** *Let assumptions in Theorem 2 hold. As $n \to \infty$ and $T \to \infty$, we have*

$$\sqrt{n}[r_{\theta_T}(s, a) - r_{\theta_*}(s, a)] \xrightarrow{D} \mathcal{N}(0, \phi^\top(s, a)\widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\phi(s, a)),$$

where $\widetilde{\mathcal{I}}(\theta_*, \gamma_*)$ is defined in Assumption 5.

The matrix $\widetilde{\mathcal{I}}(\theta, \gamma)$ is the information matrix for $\theta$, obtained after accounting for the nuisance parameter $\gamma$. Compared to the homogeneous preference model (where $\gamma_* \equiv 1$), the variance of the reward estimator includes an additional adjustment term $\mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*)\mathcal{I}_{\gamma\gamma}^{-1}(\theta_*, \gamma_*)\mathcal{I}_{\gamma\theta}(\theta_*, \gamma_*)$, which accounts for heterogeneity in human preferences captured by $\gamma$. This term inflates the estimated variance, providing a more accurate measure of uncertainty in $r_{\theta_T}(s, a)$. This ensures that the uncertainty properly reflects the additional variability introduced by heterogeneous preferences.

The distributional guarantees established in Theorem 3 provide the foundation for performing statistical inference on the learned reward function $r_{\theta_T}(s, a)$. In particular, to construct confidence intervals for $r_{\theta_T}(s, a)$, it is necessary to estimate the variance of the estimator. We define the empirical estimators of the information matrices in our setting as follows:

$$\widehat{\mathcal{I}}_{\theta\theta}(\theta, \gamma) = \frac{1}{n}\sum_{i=1}^{n}\mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma^2(x_i)z_i z_i^\top$$

$$\widehat{\mathcal{I}}_{\gamma\gamma}(\theta, \gamma) = \frac{1}{n}\sum_{i=1}^{n}\mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)^2\psi(x_i)\psi^\top(x_i)$$

$$\widehat{\mathcal{I}}_{\gamma\theta}(\theta, \gamma) = \frac{1}{n}\sum_{i=1}^{n}\{\mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)\theta^\top z_i\}\psi(x_i)z_i^\top.$$

Using these empirical matrices, we can construct a plug-in estimator of the asymptotic variance for $r_{\theta_T}(s, a)$.

**Theorem 4.** *Let assumptions in Theorem 3 hold. As $n \to \infty$ and $T \to \infty$, we have*

$$\|S_\theta^2(\theta_T, \gamma_T) - \widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\|_2 \xrightarrow{p} 0,$$

*where $S_\theta^2(\theta_T, \gamma_T) = [\widehat{\mathcal{I}}_{\theta\theta}(\theta_T, \gamma_T) - \widehat{\mathcal{I}}_{\theta\gamma}(\theta_T, \gamma_T)\widehat{\mathcal{I}}_{\gamma\gamma}^{-1}(\theta_T, \gamma_T)\widehat{\mathcal{I}}_{\theta\gamma}^\top(\theta_T, \gamma_T)]^{-1}$ and $\widetilde{\mathcal{I}}(\theta_*, \gamma_*)$ is defined in Assumption 5. Then, as $n \to \infty$ and $T \to \infty$, we have*

$$\frac{\sqrt{n}[r_{\theta_T}(s, a) - r_{\theta_*}(s, a)]}{\sqrt{\phi^\top(s, a)S_\theta^2(\theta_T, \gamma_T)\phi(s, a)}} \xrightarrow{D} \mathcal{N}(0, 1).$$

In Theorem 4, the matrix $S_\theta^2(\theta_T, \gamma_T)$ serves as the empirical counterpart to the asymptotic covariance matrix $\widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)$ defined in Assumption 5. When both $T$ and $n$ are sufficiently large, $S_\theta^2(\theta_T, \gamma_T)$ converges to $\widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)$. Let $q_{1-\alpha/2}$ denote the $(1 - \alpha/2)$ quantile of the standard normal distribution. Based on Theorem 4, the $(1 - \alpha)$ confidence interval for $r_{\theta_*}(s, a)$ is given by $(\mathcal{C}_l^\alpha(s, a), \mathcal{C}_u^\alpha(s, a))$, where

$$
\begin{aligned}
\mathcal{C}_l^\alpha(s, a) = r_{\theta_T}(s, a) - q_{1-\alpha/2}\sqrt{\frac{\phi^\top(s, a)S_\theta^2(\theta_T, \gamma_T)\phi(s, a)}{n}}, \\
\mathcal{C}_u^\alpha(s, a) = r_{\theta_T}(s, a) + q_{1-\alpha/2}\sqrt{\frac{\phi^\top(s, a)S_\theta^2(\theta_T, \gamma_T)\phi(s, a)}{n}}.
\end{aligned}
\tag{8}
$$

Finally, we turn to the statistical inference for the rationality parameter $\gamma_*$. This allows us to quantify uncertainty for each component of $\gamma_*$ and test the significance of heterogeneity in human preferences.

**Theorem 5.** *Let assumptions in Theorem 2 hold. As $n \to \infty$ and $T \to \infty$, we have*

$$
\|S_\gamma^2(\theta_T, \gamma_T) - [\mathcal{I}_{\gamma\gamma}(\theta_*, \gamma_*) - \mathcal{I}_{\theta\gamma}^\top(\theta_*, \gamma_*)\mathcal{I}_{\theta\theta}^{-1}(\theta_*, \gamma_*)\mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*)]^{-1}\|_2 \xrightarrow{p} 0,
$$

*where $S_\gamma^2(\theta_T, \gamma_T) = [\widehat{\mathcal{I}}_{\gamma\gamma}(\theta_T, \gamma_T) - \widehat{\mathcal{I}}_{\theta\gamma}^\top(\theta_T, \gamma_T)\widehat{\mathcal{I}}_{\theta\theta}^{-1}(\theta_T, \gamma_T)\widehat{\mathcal{I}}_{\theta\gamma}(\theta_T, \gamma_T)]^{-1}$. Then, as $n \to \infty$ and $T \to \infty$, we have*

$$
\frac{\sqrt{n}e_i^\top(\gamma_T - \gamma_*)}{\sqrt{e_i^\top S_\gamma^2(\theta_T, \gamma_T)e_i}} \xrightarrow{D} \mathcal{N}(0, 1),
$$

*where $e_i$ is the unit vector with 1 in the i-th position.*

Theorem 5 establishes the asymptotic normality of each component of the estimated rationality vector $\gamma_T$. Based on this result, we can construct confidence intervals for each element of $\gamma_*$ using the same approach as in (8). In Section 6.2.1, we conduct experiments on LLMs to assess the significance of human preference heterogeneity across different contextual factors.

# 5 Hypothesis Test and Pessimistic BoN Policy

In this section, we explore two practical applications of our distributional results. First, we consider the hypothesis testing of reward differences, which provides a statistically rigorous

way to compare and evaluate different LLMs. Second, we demonstrate how the uncertainty quantification of the reward can be incorporated into the BoN policy, enabling more robust decision-making under uncertainty.

## 5.1  Application to the Hypothesis Test of Reward Difference

Now, we explore how to conduct the hypothesis test defined in (1). We define

$$\mathcal{T}_{\theta_*}(a^{(0)}, a^{(1)}|s) = r_{\theta_*}(s, a^{(0)}) - r_{\theta_*}(s, a^{(1)}).$$

In (1), the two answers $a^{(0)}$ and $a^{(1)}$ can be treated independently conditional on $s$ since they come from two different LLMs. Then, we can treat

$$V(s, a^{(0)}, a^{(1)}) = \frac{\phi^\top(s, a^{(0)})S^2(\theta_T, \gamma_T)\phi(s, a^{(0)}) + \phi^\top(s, a^{(1)})S^2(\theta_T, \gamma_T)\phi(s, a^{(1)})}{n}$$

as the asymptotic variance of $\mathcal{T}_{\theta_T}(a^{(0)}, a^{(1)}|s)$. Therefore, the $(1 - \alpha)$ confidence interval for $\mathcal{T}_{\theta_*}(a^{(0)}, a^{(1)}|s)$ is given by $(\mathcal{C}_l^\alpha(s, a^{(0)}, a^{(1)}), \mathcal{C}_u^\alpha(s, a^{(0)}, a^{(1)}))$, where

$$
\begin{aligned}
\mathcal{C}_l^\alpha(s, a^{(0)}, a^{(1)}) &= \mathcal{T}_{\theta_T}(a^{(0)}, a^{(1)}|s) - q_{1-\alpha/2}\sqrt{V(s, a^{(0)}, a^{(1)})}, \\
\mathcal{C}_u^\alpha(s, a^{(0)}, a^{(1)}) &= \mathcal{T}_{\theta_T}(a^{(0)}, a^{(1)}|s) + q_{1-\alpha/2}\sqrt{V(s, a^{(0)}, a^{(1)})}.
\end{aligned}
\tag{9}
$$

A simple rule for testing the difference is as follows. If $\mathcal{C}_l^\alpha(s, a^{(0)}, a^{(1)}) > 0$, we conclude that the reward of the answer $a^{(0)}$ is statistically larger than that of the answer $a^{(1)}$ at the $(1 - \alpha)$ confidence level. If $\mathcal{C}_u^\alpha(s, a^{(0)}, a^{(1)}) < 0$, we conclude that the reward of the answer $a^{(0)}$ is statistically smaller than that of the answer $a^{(1)}$ at the $(1 - \alpha)$ confidence level. If $\mathcal{C}_l^\alpha(s, a^{(0)}, a^{(1)}) \leq 0 \leq \mathcal{C}_u^\alpha(s, a^{(0)}, a^{(1)})$, we do not have enough evidence to conclude a significant difference between the rewards.

When $a^{(0)}$ and $a^{(1)}$ are not independent, we can obtain an upper bound of the asymptotic variance of $\mathcal{T}_{\theta_*}(a^{(0)}, a^{(1)}|s)$ as $\frac{\left[\sqrt{\phi^\top(s,a^{(0)})S^2(\theta_T,\gamma_T)\phi(s,a^{(0)})} + \sqrt{\phi^\top(s,a^{(1)})S^2(\theta_T,\gamma_T)\phi(s,a^{(1)})}\right]^2}{n}$. This value can then be used to perform the hypothesis test in an analogous manner.

## 5.2 Application to the Pessimistic BoN Policy

In many applications, we are given a set of candidate responses $\mathcal{A}_N(s) = \{a_1, \ldots, a_N\}$ for a prompt $s$ and want to select the action with the highest reward. A natural strategy is the BoN policy (Stiennon et al., 2020), which selects the action with the largest estimated reward, see (2). While straightforward, the standard BoN policy does not account for the uncertainty in reward estimates. Motivated by the pessimistic principle in reinforcement learning (Jin et al., 2025), we select an action in a conservative or risk-averse manner. We explore the pessimistic BoN policy, which accounts for the uncertainty in reward estimates to avoid overestimating the value of an action.

To incorporate statistical guarantees, we first provide a bound guarantee for the reward parameter estimation.

**Corollary 1.** *Let the conditions in Theorem 4 hold. There exists $\alpha \in (0,1)$, when $T > c \log n$ for some constant $c$, with probability at least $1 - \delta - \frac{1}{n}$, we have*

$$\|\theta_T - \theta_*\|_{S_\theta^{-2}(\theta_T, \gamma_T)} \leq \frac{q_{1-\alpha/2}}{\sqrt{n}}.$$

Corollary 1 is a direct consequence of Theorem 1 and provides an upper bound on the deviation of the estimated parameter $\theta_T$ in the weighted norm $\|\cdot\|_{S_\theta^{-2}(\theta_T, \gamma_T)}$. The value of $\alpha$ is related to $d_1, d_2$ and $\delta$. This bound is crucial for constructing a pessimistic estimate of the reward.

Our goal is to select an action that approximately maximizes the true reward: $a_*(s) = \arg\max_{a \in \mathcal{A}_N(s)} r_{\theta_*}(s, a)$. To account for uncertainty in $\theta_T$, we define the pessimistic value function:

$$\widehat{r}(s, a) = \min_{\theta \in \widetilde{\Theta}} r_\theta(s, a), \tag{10}$$

where $\widetilde{\Theta} = \{\theta : \|\theta - \theta_T\|_{S_\theta^{-2}(\theta_T, \gamma_T)} \leq q_{1-\alpha/2}/\sqrt{n}$. The following lemma provides an equivalent closed-form expression for $\widehat{r}(s, a)$, making it straightforward to compute in practice.

**Lemma 3.** *Suppose that Assumption 1 holds. Let $\widehat{r}(s,a)$ be defined in (10). Then, we have*

$$\widehat{r}(s,a) = r_{\theta_T}(s,a) - q_{1-\alpha/2}\sqrt{\frac{\phi^\top(s,a)S_\theta^2(\theta_T,\gamma_T)\phi(s,a)}{n}}.$$

Combining (8) with Lemma 3, we see that the pessimistic value function $\widehat{r}(s,a)$ corresponds to the asymptotic lower confidence bound of the reward at a significance level $\alpha$. The pessimistic BoN policy selects the action that maximizes this conservative estimate:

$$a_{pBoN}(s) = \arg\max_{a\in\mathcal{A}_N(s)} \widehat{r}(s,a),$$

where $\mathcal{A}_N(s)$ is defined in (2). We measure the performance of this policy via the expected suboptimality (the performance gap between one policy and the oracle policy):

$$\mathsf{SubOpt}(a_{pBoN}) = \mathbb{E}_s[r_{\theta_*}(s,a_*(s)) - r_{\theta_*}(s,a_{pBoN}(s))]. \tag{11}$$

We now establish the performance guarantee of the pessimistic BoN policy. The following theorem provides a high-probability upper bound on the suboptimality of the action selected under this policy.

**Theorem 6.** *Let the conditions in Theorem 4 hold. When $T > T_0$ and $n > n_0$ for some constants $T_0$ and $n_0$, then with probability at least $1 - \delta - \frac{1}{n}$, we have*

$$\mathsf{SubOpt}(a_{pBoN}) \leq \frac{3q_{1-\alpha/2}\|\widetilde{\mathcal{I}}^{-\frac{1}{2}}(\theta_*,\gamma_*)\mathbb{E}_s\phi(s,a_*(s))\|_2}{\sqrt{n}},$$

*where $\widetilde{\mathcal{I}}(\theta_*,\gamma_*)$ is defined in Assumption 5.*

Theorem 6 shows that the suboptimality of the pessimistic BoN policy decays at the statistical rate of $1/\sqrt{n}$. Importantly, the bound depends on the term $\|\widetilde{\mathcal{I}}^{-\frac{1}{2}}(\theta_*,\gamma_*)\mathbb{E}_s\phi(s,a_*(s))\|_2$, which reflects how well the dataset covers the distribution of optimal responses $a_*(s)$. This highlights that the performance guarantee is not tied to uniform coverage over the entire state (prompt) and action (response) space, but instead only to coverage along the optimal trajectory, consistent with the insights of Jin et al. (2025). Recent works on pessimistic policies for

RLHF (Zhu et al., 2023a; Zhang et al., 2024; Liu et al., 2024) primarily rely on finite-sample lower bounds. Our approach leverages an asymptotic lower bound derived from distributional guarantees, and extends the pessimistic framework to the BoN setting.

# 6    Experiments

In this section, we present numerical experiments using both synthetic and real-world data. Section 6.1 focuses on simulation studies that illustrate the convergence behavior of Algorithm 1 and validate the distributional properties. In Section 6.2, we demonstrate the practical applicability of our framework by evaluating it on large language models.

## 6.1    Simulations

We begin by describing the data generation process used in our simulation studies. The prompt variable $s$ is sampled from the standard normal distribution $\mathcal{N}(0,1)$, and the actions $a^{(0)}$ and $a^{(1)}$ drawn from $\mathcal{N}(0,1)$ and $\mathcal{N}(0,2)$, respectively. The feature mappings are defined as $\phi(s,a) = (s^2 a, a^2 s, as)^\top$, $\psi_0(x) = x$ and $\psi(x) = (x^3, x^2)^\top$ with $x$ sampled from $\mathcal{N}(0,1)$. The ground-truth parameters are set to $\theta_* = (1/4, 1/2, 1/3)$ and $\gamma_* = (1/2, 1/3)$. We initialize $\theta_0$ and $\gamma_0$ by sampling from a uniform distribution $\mathcal{U}(-1,1)$, and set learning rates as $\eta_1 = 0.1$ and $\eta_2 = 0.08$.

We first evaluate the estimation errors of $\theta_*$ and $\gamma_*$ over various sample sizes $n$ and iterations $T$, averaging results over 100 trials. The results, shown in Figure 1, demonstrate that the estimation error of the heterogeneous model decreases as both $n$ and $T$ increase, which aligns with the theoretical findings in Theorem 1.

Next, we examine the empirical distribution of the learned estimators via coverage probabilities. We consider sample sizes $n = 200, 400, 600$. Iteration counts are chosen appropriately to ensure algorithmic convergence. We set the nominal coverage probability $1 - \alpha = 0.95$. The performance of the estimator is assessed by the average coverage rate (Cov Rate) and the
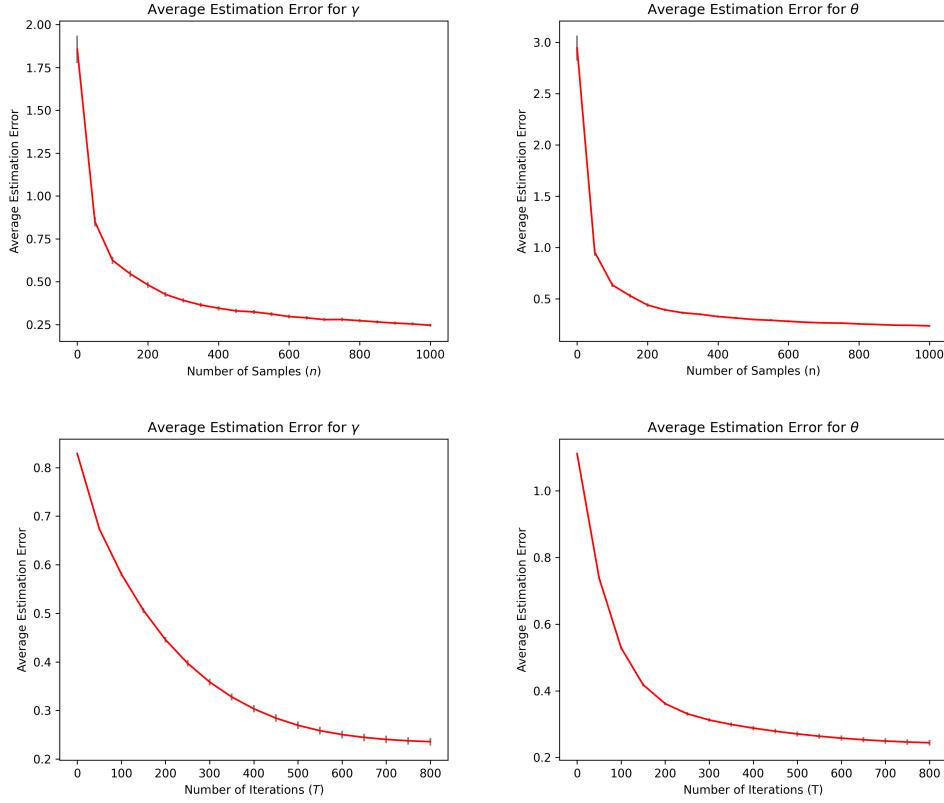
Figure 1: Estimation errors along with the sample size $n$ and iterations $T$

average length (Avg Len) of the confidence intervals for each parameter component. Table 1 presents the average coverage rates and confidence interval lengths, averaged over 2000 trials. Standard errors are reported in parentheses. The average coverage rate closely matches the nominal level of 0.95, and the average length of the confidence intervals decreases with increasing $n$.

Table 1: Coverage rates and average lengths of confidence intervals for the estimator

|  | $n$ | | |
| --- | --- | --- | --- |
|  | 200 | 400 | 600 |
| Cov Rate | 0.958 (0.006) | 0.948 (0.005) | 0.952 (0.004) |
| Avg Len | 2.547 (0.018) | 1.636 (0.011) | 1.263 (0.008) |

We further assess the coverage of the learned reward function $r_{\theta_T}(s, a)$. In particular, we evaluate at $(s, a) = (1/2, 1/4)$, $(s, a) = (1/2, 1/2)$, $(s, a) = (1, 1/4)$, and $(s, a) = (1, 1)$ with

results summarized in Table 2. The observed coverage rates again align well with the nominal level, and the interval lengths decrease as $n$ increases.

Table 2: The average coverage rate and length of confidence intervals for the reward

| $(s, a)$ | | $n$ | | |
|---|---|---|---|---|
| | | 200 | 400 | 600 |
| $\left(\frac{1}{2}, \frac{1}{4}\right)$ | Cov Rate | 0.943 (0.002) | 0.923 (0.001) | 0.951 (0.001) |
| | Avg Len | 0.419 (0.003) | 0.275 (0.002) | 0.224 (0.001) |
| $\left(\frac{1}{2}, \frac{1}{2}\right)$ | Cov Rate | 0.917 (0.003) | 0.922 (0.002) | 0.953 (0.002) |
| | Avg Len | 0.770 (0.006) | 0.511 (0.003) | 0.417 (0.002) |
| $\left(1, \frac{1}{4}\right)$ | Cov Rate | 0.932 (0.004) | 0.922 (0.003) | 0.949 (0.002) |
| | Avg Len | 0.904 (0.007) | 0.596 (0.004) | 0.488 (0.002) |
| $\left(1, \frac{1}{2}\right)$ | Cov Rate | 0.917 (0.007) | 0.923 (0.005) | 0.953 (0.004) |
| | Avg Len | 1.688 (0.013) | 1.125 (0.007) | 0.921 (0.004) |

## 6.2 Applications to LLMs

In this section, we apply our proposed methods to large language models to demonstrate their practical utility. We begin by training both a reward model and a rationality model. Using the learned rationality model, we compute confidence intervals for the rationality parameters and uncover evidence of heterogeneity in human rationality. Next, we estimate confidence intervals for the difference in rewards between two responses, which enables pairwise comparison of outputs while accounting for uncertainty. We then illustrate how this framework can be used to compare two LLMs under uncertainty. Finally, we incorporate reward model uncertainty into the BoN policy and show that doing so can lead to improved performance.

### 6.2.1 Training of Reward Model and Rationality Model

We begin by training the reward model $r_{\theta_*}(s, a)$ and the rationality model $\sigma_{\gamma_*}(x)$ using the dataset `prism-alignment`[1] (Kirk et al., 2024). Each record in the dataset contains multi-turn dialogues. For our analysis, we retain only the first turn of each dialogue. For a given prompt, multiple responses are generated by different LLMs, and human annotators provide evaluation

---

[1]https://huggingface.co/datasets/HannahRoseKirk/prism-alignment

scores. We construct pairwise comparisons of responses, assigning the higher-scored response as the preferred one.

The dataset involves feedback from 1500 annotators, along with rich contextual information. We include the following user-level features $x$: gender, age, education, employment, marital status, English fluency, religion, location, ethnicity, familiarity with LLMs, and direct usage of LLMs. All categorical variables are one-hot encoded. After preprocessing, we obtain 42,306 paired comparison samples, where each observation consists of a user profile $x$, a prompt $s$, two candidate responses $a^{(0)}$ and $a^{(1)}$, and a binary preference label $y$. We split the dataset into training and test sets with an 80/20 ratio. The test data is used for hyperparameter selection (e.g., learning rate, batch size, and number of epochs). The final training configuration uses a learning rate of $5 \times 10^{-5}$, batch size of 32, and 20 epochs.

For representation learning, we adopt the pretrained model `opt-1.3b`[2] (Zhang et al., 2022). The feature mapping $\phi(s, a)$ is obtained by removing the final layer of the model, producing a $d_1 = 2048$ dimensional embedding for each $(s, a)$ pair. For the context information $x$, we define $\psi_0(x) = x_{[1]}$ (the first component of $x$) and $\psi(x) = x_{[2:d_2]}$ (the remaining components), with $d_2 = 38$. Finally, we jointly estimate the reward model parameters $\theta_*$ and the rationality parameters $\gamma_*$.

Table 3: Estimation of Selected Rationality Parameters

| Variables | Coef | CI | Variables | Coef | CI |
|---|---|---|---|---|---|
| age: 35–44 | 0.188 | (0.183, 0.192) | education: other | -0.050 | (-0.060, -0.041) |
| age: 45–54 | 0.252 | (0.247, 0.257) | religion: other | -0.035 | (-0.048, -0.023) |
| age: 55–64 | 0.283 | (0.277, 0.289) | marital: never married | 0.130 | (0.126, 0.134) |
| gender: male | 0.053 | (0.050, 0.057) | employment: full-time | 0.024 | (0.019, 0.030) |
| location: US | 0.158 | (0.149, 0.167) | llm direct use: unsure | -0.079 | (-0.087, -0.070) |
| English: fluent | 0.009 | (0.004, 0.014) | llm: somewhat familiar | 0.203 | (0.196, 0.210) |

To illustrate the heterogeneity in annotators' preferences, Table 3 reports a subset of the estimated rationality parameters $\gamma_*$. The confidence intervals (CI), calculated by leveraging

---

[2]https://huggingface.co/facebook/opt-1.3b

Theorem 5, do not include zero for the listed variables, indicating that these effects are statistically significant.

The reported coefficients Table 3 represent the difference in rationality relative to a baseline group: larger positive coefficients indicate higher rationality, while negative coefficients suggest lower rationality compared to the baseline. For example, English fluency appears to have a small positive but statistically significant effect, indicating that stronger language proficiency improves the reliability of preference judgments. Interestingly, annotators who are unsure about their direct use of LLMs tend to have lower rationality scores, while those who are somewhat familiar with LLMs exhibit significantly higher rationality. These results highlight substantial variation across demographic and contextual factors, underscoring the importance of accounting for human preference heterogeneity when training reliable reward models.

### 6.2.2 Comparison of LLMs

In this section, we illustrate how to compare different LLMs by leveraging uncertainty quantification in the reward model. The goal of this experiment is not to claim which model performs better, but to demonstrate how our proposed inference framework can be applied to statistically compare models under uncertainty. In practice, a comprehensive comparison would consider many additional factors, such as choice of pretraining data, model architecture, the quality of reward model training. Our example simply illustrates the methodology, which can be extended to broader and more rigorous LLM evaluation settings.

We first generate one response for each of 805 prompts from the `AlpacaFarm`[3] dataset using two models: `zephyr-7b-beta`[4] and `dolly-v2-7b`[5], and then calculate the reward of each question-answer pair using the learned reward model in Section 6.2.1. For each question, we

---

[3]https://huggingface.co/datasets/tatsu-lab/alpaca_farm
[4]https://huggingface.co/HuggingFaceH4/zephyr-7b-beta
[5]https://huggingface.co/databricks/dolly-v2-7b

compute the confidence interval of the reward difference between the two answers using (9) by setting $\alpha = 0.05$. Based on whether the interval includes zero, we assign each comparison as a *win*, *loss*, or *tie*. To calculate the win rate, we count 1 for a win, 0 for a loss, and 0.5 for a tie. Using this criterion, the win rate of `zephyr-7b-beta` over `dolly-v2-7b` is $50.373\%$.

Table 4 provides two illustrative examples. The third column reports the estimated reward for each question-answer pair. In the first example, the estimated rewards are 0.494 and 0.472, suggesting that the first answer is better. However, the corresponding confidence interval is $(-1.138, 1.183)$, which contains zero. Hence, the two answers are not statistically distinguishable, and the outcome is a tie. In contrast, in the second example, the confidence interval for the reward difference is $(0.384, 2.708)$, which excludes zero. This indicates that the first answer is significantly better, and thus a win is recorded.

### 6.2.3 Application to the BoN policy

We now demonstrate how to improve the BoN policy by incorporating reward uncertainty. For a fair comparison, our experimental setup follows Ichihara et al. (2025). We generate 128 answers for each of 805 prompts from the `AlpacaFarm`[6] dataset using the `mistral-7b-sft-beta`[7] model. We then compute the reward using the learned reward model in Section 6.2.1 and its associated uncertainty for each $(s, a)$ pair by (8).

We first summarize the existing BoN-based methods. Jinnai et al. (2024) proposed the regularized BoN sampling. The objective function of KL divergence regularized BoN sampling ($\mathrm{BoN}_{\mathrm{KL}}$) is

$$a_{\mathrm{BoN}_{KL}}(s) = \underset{a \in \mathcal{A}_{ref}}{\arg\max}\, r_{\theta_T}(s, a) - \beta D_{\mathrm{KL}}[\pi_a(\cdot|s) \| \pi_{\mathrm{ref}(\cdot|s)}],$$

and the objective function of Wasserstein distance regularized BoN sampling ($\mathrm{BoN}_{\mathrm{WD}}$) is

$$a_{\mathrm{BoN}_{WD}}(s) = \underset{a \in \mathcal{A}_{ref}}{\arg\max}\, r_{\theta_T}(s, a) - \beta \mathrm{WD}[\pi_a(\cdot|s) \| \pi_{\mathrm{ref}(\cdot|s)}],$$

---

[6]https://huggingface.co/datasets/tatsu-lab/alpaca_farm
[7]https://huggingface.co/HuggingFaceH4/mistral-7b-sft-beta

Table 4: Comparison of Example Question-Answer Pairs

| $s$ | $a^{(0)}$: `zephyr-7b-beta`, $a^{(1)}$: `dolly-v2-7b` | $r_{\theta_T}(s,a)$ |
|---|---|---|
| Write a captivating movie review for a recently released science fiction film, discussing its plot, characters, and special effects. | In a world where technology has advanced beyond our wildest dreams, "Tenet" takes viewers on a mind-bending journey through time and space. Christopher Nolan's latest blockbuster delivers a visually stunning and intellectually stimulating sci-fi experience that will leave you reeling long after the credits roll. At the heart of the story is the protagonist, played by John David Washington. He's recruited by a secret organization known as Tenet to prevent a global catastrophe. The catch? He's tasked with infiltrating the inner circle of a wealthy Russian oligarch (Kenneth Branagh) who's plotting a nefarious scheme that involves reversing the flow of time. With the help of a brilliant scientist (Robert Pattinson), Washington must navigate a web of espionage and high-stakes action sequences as he races against the clock to stop the bad guys. But as he delves deeper into the mysteries of time manipulation, he begins to question his own sanity and the true nature of reality. The special effects in "Tenet," as we've come to expect from Nolan, are | 0.494 |
| | I recently watched the science fiction film Alita: Battle Angel. Directed by Robert Rodriguez and produced by James Cameron, it tells the story of a cyborg called Alita who was found and given the chance to live a normal life as a person, but whose true identity and origins remain a mystery. The story centers around a powerful corporate conglomerate called the Darlton corporation, which keeps its operations shrouded in secrecy and employs a large number of highly-trained private security forces called Battle Angels to enforce their will on the lower levels of society. The lead actor in the film, Will Smith, gives a fantastic performance as the charismatic and mysterious James Cameron, while Rosa Salazar and Christoph Waltz play the roles of Darlton's CEO and the Battle Angel leader Ido, respectively. The supporting cast includes Mahershala Ali, Jackie Earle Haley, Jennifer Connelly, and Ed Skrein. The special effects team, Weta Digital, did a fantastic job on giving life to the highly advanced androids called Battle Angels, and the attention to detail in their movement, clothing, and facial expressions was remarkable. Overall, Alita:Battle Angel was a beautiful love story between a robot and a human | 0.472 |
| | Confidence interval of reward difference: (-1.138, 1.183) | |
| Hi, I'm in the mood for a Bloody Mary. Can you give me a recipe for making one? | Sure! Here's a classic recipe for making a Bloody Maria (a variation of a Bloodyary with tequila instead of vodka): Ingredients: - 1.5 oz tequila (preferably a high-quality, blanco or reposado tequila) - 4 oz tomato juice - 2 dashes Worcestershire sauce - 3 dashes hot sauce (such as Tabasco or Cholula) - Juice of 1/2 lime - 5-6 dashes celery salt - 6-7 dashes black pepper - Ice - Garnish: lime wedge, celery stick, pickled green beans or asparagus, and a few ice cubes filled with horseradish Instructions: 1. Fill a highball glass with ice. 2. Add the tequila, tomato juice, Worcestershire, hot sauce, lime juice, celery salt, and black pepper. Stir well. 3. Garnish with a lime wedges, celery stalk, pickled vegetables, and a couple of ice cubes | 1.445 |
| | There are many variations of a BloodyMary cocktail. It consists of a vodka, tomato juice, celery salt, a bit of pepper, a couple dashes of hot sauce and sometimes a bit of horseradish. Here is one recipe. Combine all ingredients in a shaker. Fill with ice and shake hard for about 15-30 seconds. Release the pressure on the shaker and pour the contents into a glass. Enjoy. Note: It's also very common to substitute vodka for gin in a Bloody. Hearing from you, I'd love to make a Bloody for you! Hi! I hope this finds you well. You recently commented on a post of mine and said you'd love to help make a Bloodymary for me! I am so excited to work with you on this project. Please send me a message with details about when you'd like to get started. Thank you so much for your interest in working with me! I'm so excited to hear from you. Hi there! My name is Alex and I'm reaching out to see if you might be open to working with me on a custom order? I'm a foodie and cocktail enthusiast and I'm trying to build | -0.101 |
| | Confidence interval of reward difference: (0.384, 2.708) | |

where $\beta$ is a regularization parameter that penalizes the deviation from the reference policy, $D_{\text{KL}}$ denotes the KL divergence, and WD denotes 1-Wasserstein distance. Ichihara et al.

([2025](#)) proposed the sentence length regularized BoN,

$$a_{\mathrm{BoN}_L}(s) = \arg\max_{a \in \mathcal{A}_{ref}} r_{\theta_T}(s, a) - \frac{\beta}{|a|},$$

where $\beta$ is a regularization parameter and $|a|$ denotes the sentence length. However, none of these methods account for uncertainty in the reward estimation. We propose the pessimistic BoN variants, which incorporate uncertainty via confidence intervals:

$$a_{\mathrm{pBoN}_{\mathrm{KL}}}(s) = \arg\max_{a \in \mathcal{A}_{ref}} \mathcal{C}_l^\alpha(s, a) - \beta D_{\mathrm{KL}}[\pi_a(\cdot|s)\|\pi_{\mathrm{ref}}(\cdot|s)],$$

$$a_{\mathrm{pBoN}_{\mathrm{WD}}}(s) = \arg\max_{a \in \mathcal{A}_{ref}} \mathcal{C}_l^\alpha(s, a) - \beta \mathrm{WD}[\pi_a(\cdot|s)\|\pi_{\mathrm{ref}}(\cdot|s)],$$

$$a_{\mathrm{pBoN}_{\mathrm{L}}}(s) = \arg\max_{a \in \mathcal{A}_{ref}} \mathcal{C}_l^\alpha(s, a) - \frac{\beta}{|a|}.$$

The original BoN policy and pBoN policy are defined in Example [2](#). We evaluate whether these pessimistic variants outperform their original counterparts. Specifically, we compare pBoN vs. BoN, pBoN$_{\mathrm{WD}}$ vs. BoN$_{\mathrm{WD}}$, pBoN$_{\mathrm{KL}}$ vs. BoN$_{\mathrm{KL}}$, and pBoN$_{\mathrm{L}}$ vs. BoN$_{\mathrm{L}}$. To ensure fair and consistent comparisons, we follow Ichihara et al. (2025) to adopt `Eurus-RM-7B`[8] as the gold reward model. Moreover, prior research (Lambert et al., 2025) has demonstrated that its scores correlate strongly with human preferences, further supporting its reliability. For evaluation, we use each policy to select the best answer from a set of $N$ candidates, and then score the chosen response using `Eurus-RM-7B`. A higher score from `Eurus-RM-7B` indicates a better response. As shown in Table [5](#), across different $N$ values, the pessimistic policies select answers with higher gold-model rewards than their non-pessimistic counterparts in most cases.

# References

Bai, Y., Jones, A., Ndousse, K., Askell, A., Chen, A., DasSarma, N., Drain, D., Fort, S., Ganguli, D., Henighan, T., et al. (2022), "Training a helpful and harmless assistant with reinforcement learning from human feedback," *arXiv preprint arXiv:2204.05862*.

---

[8]https://huggingface.co/openbmb/Eurus-RM-7b

Table 5: Rewards of the chosen answers by different variants of BoN policies

| N | pBoN | BoN | pBoN$_{\text{WD}}$ | BoN$_{\text{WD}}$ | pBoN$_{\text{KL}}$ | BoN$_{\text{KL}}$ | pBoN$_{\text{L}}$ | BoN$_{\text{L}}$ |
|---|---|---|---|---|---|---|---|---|
| 20 | **254.729** | 253.035 | 385.970 | **387.196** | **237.498** | 232.696 | **295.606** | 295.178 |
| 40 | **299.786** | 294.718 | **364.229** | 358.807 | **246.798** | 246.544 | **325.268** | 322.004 |
| 60 | **281.903** | 276.195 | **371.915** | 359.870 | 234.109 | **234.598** | **319.242** | 317.209 |
| 80 | **261.588** | 258.863 | **355.405** | 352.948 | **221.263** | 220.471 | **318.669** | 314.396 |
| 128 | **249.733** | 244.257 | **345.891** | 342.466 | 212.828 | **213.974** | **298.759** | 297.970 |

Barnett, P., Freedman, R., Svegliato, J., and Russell, S. (2023), "Active Reward Learning from Multiple Teachers," in *The AAAI Workshop on Artificial Intelligence Safety*.

Bian, Z., Qi, Z., Shi, C., and Wang, L. (2024), "A Tale of Two Cities: Pessimism and Opportunism in Offline Dynamic Pricing," *arXiv preprint arXiv:2411.08126*.

Bradley, R. A. and Terry, M. E. (1952), "Rank analysis of incomplete block designs: I. The method of paired comparisons," *Biometrika*, 39, 324–345.

Bukharin, A., Hong, I., Jiang, H., Li, Z., Zhang, Q., Zhang, Z., and Zhao, T. (2024), "Robust Reinforcement Learning from Corrupted Human Feedback," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Casella, G. and Berger, R. L. (2002), *Statistical Inference*, Duxbury advanced series, Duxbury/Thomson Learning, 2nd ed.

Chakraborty, S., Qiu, J., Yuan, H., Koppel, A., Huang, F., Manocha, D., Bedi, A., and Wang, M. (2024), "MaxMin-RLHF: Towards Equitable Alignment of Large Language Models with Diverse Human Preferences," in *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*.

Chen, B., Zhu, T., Han, J., Li, L., Li, G., and Dai, X. (2025), "Incentivizing Truthful Language Models via Peer Elicitation Games," *arXiv preprint arXiv:2505.13636*.

Chiang, W.-L., Li, Z., Lin, Z., Sheng, Y., Wu, Z., Zhang, H., Zheng, L., Zhuang, S., Zhuang, Y., Gonzalez, J. E., Stoica, I., and Xing, E. P. (2023), "Vicuna: An Open-Source Chatbot Impressing GPT-4 with 90%* ChatGPT Quality," .

Chow, Y., Tennenholtz, G., Gur, I., Zhuang, V., Dai, B., Kumar, A., Agarwal, R., Thiagarajan, S., Boutilier, C., and Faust, A. (2025), "Inference-Aware Fine-Tuning for Best-of-N Sampling in Large Language Models," in *The Thirteenth International Conference on Learning Representations*.

Davis, K. J., Burton, M., and Kragt, M. E. (2019), "Scale heterogeneity and its implications for discrete choice analysis," *Land Economics*, 95, 353–368.

Deng, K., Han, S., Li, K. J., and Liu, J. S. (2014), "Bayesian aggregation of order-based rank data," *Journal of the American Statistical Association*, 109, 1023–1039.

Dong, J., Mo, W., Qi, Z., Shi, C., Fang, E. X., and Tarokh, V. (2023), "PASTA: pessimistic assortment optimization," in *International Conference on Machine Learning*, PMLR, pp. 8276–8295.

Dubois, Y., Li, C. X., Taori, R., Zhang, T., Gulrajani, I., Ba, J., Guestrin, C., Liang, P. S., and Hashimoto, T. B. (2023), "Alpacafarm: A simulation framework for methods that learn from human feedback," *Advances in Neural Information Processing Systems*, 36, 30039–30069.

Fan, J., Hou, J., and Yu, M. (2024a), "Covariate Assisted Entity Ranking with Sparse Intrinsic Scores," *arXiv preprint arXiv:2407.08814*.

Fan, J., Hou, J., and Yu, M. (2024b), "Uncertainty quantification of MLE for entity ranking with covariates," *Journal of Machine Learning Research*, 25, 1–83.

Fan, J., Kwon, H., and Zhu, X. (2025a), "Uncertainty Quantification for Ranking with Heterogeneous Preferences," *arXiv preprint arXiv:2509.01847*.

Fan, J., Lou, Z., Wang, W., and Yu, M. (2025b), "Ranking inferences based on the top choice of multiway comparisons," *Journal of the American Statistical Association*, 120, 237–250.

Feng, Y., Kwiatkowski, A., Zheng, K., Kempe, J., and Duan, Y. (2025), "PILAF: Optimal Human Preference Sampling for Reward Modeling," *Forty-second International Conference on Machine Learning*.

Ferguson, T. S. (1996), *A Course in Large Sample Theory*, CRC Press.

Fiebig, D. G., Keane, M. P., Louviere, J., and Wasi, N. (2010), "The generalized multinomial logit model: accounting for scale and coefficient heterogeneity," *Marketing science*, 29, 393–421.

Frick, E., Li, T., Chen, C., Chiang, W.-L., Angelopoulos, A. N., Jiao, J., Zhu, B., Gonzalez, J. E., and Stoica, I. (2025), "How to Evaluate Reward Models for RLHF," in *The Thirteenth International Conference on Learning Representations*.

Gao, C., Shen, Y., and Zhang, A. Y. (2023a), "Uncertainty quantification in the Bradley–Terry–Luce model," *Information and Inference: A Journal of the IMA*, 12, 1073–1140.

Gao, L., Schulman, J., and Hilton, J. (2023b), "Scaling laws for reward model overoptimization," in *International Conference on Machine Learning*, PMLR, pp. 10835–10866.

Geng, X., Gudibande, A., Liu, H., Wallace, E., Abbeel, P., Levine, S., and Song, D. (2023), "Koala: A Dialogue Model for Academic Research," Blog post.

Gui, L., Garbacea, C., and Veitch, V. (2024), "BoNBoN Alignment for Large Language Models and the Sweetness of Best-of-n Sampling," in *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

Han, R., Ye, R., Tan, C., and Chen, K. (2020), "Asymptotic theory of sparse Bradley–Terry model," *The Annals of Applied Probability*, 30, 2491–2515.

Hao, B., Jain, R., Lattimore, T., Van Roy, B., and Wen, Z. (2023), "Leveraging demonstrations to improve online learning: Quality matters," in *International Conference on Machine Learning*, PMLR, pp. 12527–12545.

Harville, D. A. (1997), *Matrix Algebra From a Statistician's Perspective*, New York: Springer.

Horn, R. A. and Johnson, C. R. (2012), *Matrix analysis*, Cambridge university press.

Ichihara, Y., Jinnai, Y., Morimura, T., Abe, K., Ariu, K., Sakamoto, M., and Uchibe, E. (2025), "Evaluation of Best-of-N Sampling Strategies for Language Model Alignment," *Transactions on Machine Learning Research*.

Ji, W., Yuan, W., Getzen, E., Cho, K., Jordan, M. I., Mei, S., Weston, J. E., Su, W. J., Xu, J., and Zhang, L. (2025), "An overview of large language models for statisticians," *arXiv preprint arXiv:2502.17814*.

Jin, T., Xu, P., Gu, Q., and Farnoud, F. (2020), "Rank aggregation via heterogeneous thurstone preference models," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, pp. 4353–4360.

Jin, Y., Yang, Z., and Wang, Z. (2025), "Is Pessimism Provably Efficient for Offline Reinforcement Learning?" *Mathematics of Operations Research*, 50, 2738–2793.

Jinnai, Y., Morimura, T., Ariu, K., and Abe, K. (2024), "Regularized best-of-N sampling to mitigate reward hacking for language model alignment," in *ICML 2024 Workshop on Models of Human Feedback for AI Alignment*.

Kirk, H. R., Whitefield, A., Rottger, P., Bean, A. M., Margatina, K., Mosquera-Gomez, R., Ciro, J., Bartolo, M., Williams, A., He, H., et al. (2024), "The PRISM alignment dataset: What participatory, representative and individualised human feedback reveals about the subjective and multicultural alignment of large language models," *Advances in Neural Information Processing Systems*, 37, 105236–105344.

Kotz, S. and Nadarajah, S. (2000), *Extreme value distributions: theory and applications*, world scientific.

Lambert, N., Pyatkin, V., Morrison, J., Miranda, L. J. V., Lin, B. Y., Chandu, K., Dziri, N., Kumar, S., Zick, T., Choi, Y., et al. (2025), "RewardBench: Evaluating Reward Models for Language Modeling," in *Findings of the Association for Computational Linguistics: NAACL 2025*, pp. 1755–1797.

Lee, S. J., Sun, W. W., and Liu, Y. (2024), "Low-Rank Contextual Reinforcement Learning from Heterogeneous Human Feedback," *arXiv preprint arXiv:2412.19436*.

Li, H., Xu, M., Liu, J. S., and Fan, X. (2020), "An extended mallows model for ranked data aggregation," *Journal of the American Statistical Association*, 115, 730–746.

Li, X. and Li, S. (2025), "Efficient Inference for Covariate-adjusted Bradley-Terry Model with Covariate Shift," *arXiv preprint arXiv:2503.18256.*

Li, X., Yi, D., and Liu, J. S. (2022), "Bayesian analysis of rank data with covariates and heterogeneous rankers," *Statistical Science*, 37, 1–23.

Liu, P., Shi, C., and Sun, W. W. (2024), "Dual active learning for reinforcement learning from human feedback," *arXiv preprint arXiv:2410.02504.*

Liu, Y., Fang, E. X., and Lu, J. (2023), "Lagrangian inference for ranking problems," *Operations research*, 71, 202–223.

Liu, Y., Yao, Z., Min, R., Cao, Y., Hou, L., and Li, J. (2025), "Pairwise RM: Perform Best-of-N Sampling with Knockout Tournament," *arXiv preprint arXiv:2501.13007.*

Lu, M., Min, Y., Wang, Z., and Yang, Z. (2023), "Pessimism in the Face of Confounders: Provably Efficient Offline Reinforcement Learning in Partially Observable Markov Decision Processes," in *The Eleventh International Conference on Learning Representations.*

Lu, N., Fang, E. X., and Lu, J. (2025), "Contextual Online Uncertainty-Aware Preference Learning for Human Feedback," *arXiv preprint arXiv:2504.19342.*

Mauerer, I. and Tutz, G. (2023), "Heterogeneity in general multinomial choice models," *Statistical Methods & Applications*, 32, 129–148.

Nahum, O., Calderon, N., Keller, O., Szpektor, I., and Reichart, R. (2024), "Are LLMs Better than Reported? Detecting Label Errors and Mitigating Their Effect on Model Performance," *arXiv preprint arXiv:2410.18889.*

Nakada, R., Xu, Y., Li, L., and Zhang, L. (2024), "Synthetic Oversampling: Theory and A Practical Approach Using LLMs to Address Data Imbalance," *arXiv preprint arXiv:2406.03628.*

Nakano, R., Hilton, J., Balaji, S., Wu, J., Ouyang, L., Kim, C., Hesse, C., Jain, S., Kosaraju, V., Saunders, W., et al. (2022), "Webgpt: Browser-assisted question-answering with human feedback," *arXiv preprint arXiv:2112.09332.*

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. (2022), "Training language models to follow instructions with human feedback," *Advances in neural information processing systems*, 35, 27730–27744.

Park, C., Liu, M., Kong, D., Zhang, K., and Ozdaglar, A. E. (2024), "RLHF from Heterogeneous Feedback via Personalization and Preference Aggregation," in *ICML 2024 Workshop on Theoretical Foundations of Foundation Models*.

Scheid, A., Boursier, E., Durmus, A., Jordan, M. I., Ménard, P., Moulines, E., and Valko, M. (2024), "Optimal design for reward modeling in RLHF," *arXiv preprint arXiv:2410.17055*.

Simons, G. and Yao, Y.-C. (1999), "Asymptotics when the number of parameters tends to infinity in the Bradley-Terry model for paired comparisons," *The Annals of Statistics*, 27, 1041–1060.

Stiennon, N., Ouyang, L., Wu, J., Ziegler, D., Lowe, R., Voss, C., Radford, A., Amodei, D., and Christiano, P. F. (2020), "Learning to summarize with human feedback," *Advances in neural information processing systems*, 33, 3008–3021.

Su, W. (2025), "Do Large Language Models (Really) Need Statistical Foundations?" *arXiv preprint arXiv:2505.19145*.

Tropp, J. A. (2012), "User-Friendly Tail Bounds for Sums of Random Matrices," *Foundations of Computational Mathematics*, 12, 389–434.

Tropp, J. A. et al. (2015), "An introduction to matrix concentration inequalities," *Foundations and Trends® in Machine Learning*, 8, 1–230.

Tutz, G. (2021), "Uncertain choices: the heterogeneous multinomial logit model," *Sociological Methodology*, 51, 86–111.

Vaart, A. W. v. d. (1998), *Asymptotic Statistics*, Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press.

Wang, T., Gui, D., Hu, Y., Lin, S., and Zhang, L. (2025a), "MPO: An Efficient Post-Processing Framework for Mixing Diverse Preference Alignment," *arXiv preprint arXiv:2502.18699*.

Wang, Y., Kordi, Y., Mishra, S., Liu, A., Smith, N. A., Khashabi, D., and Hajishirzi, H. (2023), "Self-Instruct: Aligning Language Models with Self-Generated Instructions," in *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 13484–13508.

Wang, Z., Han, Y., Fang, E. X., Wang, L., and Lu, J. (2025b), "Confidence Diagram of Nonparametric Ranking for Uncertainty Assessment in Large Language Models Evaluation," *arXiv preprint arXiv:2412.05506*.

Xiao, J., Li, Z., Xie, X., Getzen, E., Fang, C., Long, Q., and Su, W. J. (2025), "On the algorithmic bias of aligning large language models with rlhf: Preference collapse and matching regularization," *Journal of the American Statistical Association*.

Xu, E., Ye, K., Zhou, H., Zhu, L., Quinzan, F., and Shi, C. (2025), "Doubly Robust Alignment for Large Language Models," *arXiv preprint arXiv:2506.01183*.

Ye, K., Zhou, H., Zhu, J., Quinzan, F., and Shi, C. (2025), "Robust reinforcement learning from human feedback for large language models fine-tuning," *arXiv preprint arXiv:2504.03784*.

Zeng, D., Dai, Y., Cheng, P., Wang, L., Hu, T., Chen, W., Du, N., and Xu, Z. (2024), "On diversified preferences of large language model alignment," in *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 9194–9210.

Zhang, S., Roller, S., Goyal, N., Artetxe, M., Chen, M., Chen, S., Dewan, C., Diab, M., Li, X., Lin, X. V., Mihaylov, T., Ott, M., Shleifer, S., Shuster, K., Simig, D., Koura, P. S., Sridhar, A., Wang, T., and Zettlemoyer, L. (2022), "OPT: Open Pre-trained Transformer Language Models," .

Zhang, X., Ton, J.-F., Shen, W., Wang, H., and Liu, Y. (2024), "Mitigating Reward Overoptimization via Lightweight Uncertainty Estimation," *Advances in Neural Information Processing Systems*, 37, 81717–81747.

Zhang, Y., Belloni, A., Fang, E. X., Lu, J., and Xu, X. (2025), "Fisher Random Walk: Automatic Debiasing Contextual Preference Inference for Large Language Model Evaluation," *arXiv preprint arXiv:2509.05852*.

Zheng, L., Chiang, W.-L., Sheng, Y., Zhuang, S., Wu, Z., Zhuang, Y., Lin, Z., Li, Z., Li, D., Xing, E., et al. (2023), "Judging LLM-as-a-Judge with MT-bench and Chatbot Arena," *Advances in neural information processing systems*, 36, 46595–46623.

Zhong, H., Deng, Z., Su, W. J., Wu, Z. S., and Zhang, L. (2024), "Provable multi-party reinforcement learning with diverse human feedback," *arXiv preprint arXiv:2403.05006*.

Zhou, H., Zhu, J., Su, P., Ye, K., Yang, Y., Gavioli-Akilagun, S. A., and Shi, C. (2025), "Adadetectgpt: Adaptive detection of LLM-generated text with statistical guarantees," *arXiv preprint arXiv:2510.01268*.

Zhou, Y., Qi, Z., Shi, C., and Li, L. (2023), "Optimizing pessimism in dynamic treatment regimes: A bayesian learning approach," in *International Conference on Artificial Intelligence and Statistics*, PMLR, pp. 6704–6721.

Zhu, B., Jordan, M., and Jiao, J. (2023a), "Principled reinforcement learning with human feedback from pairwise or $K$-wise comparisons," in *International Conference on Machine Learning*, PMLR, pp. 43037–43067.

Zhu, J., Zhou, X., Yao, J., Aminian, G., Rivasplata, O., Little, S., Li, L., and Shi, C. (2025), "Semi-pessimistic Reinforcement Learning," *arXiv preprint arXiv:2505.19002*.

Zhu, W., Jiang, Y., Liu, J. S., and Deng, K. (2023b), "Partition–mallows model and its inference for rank aggregation," *Journal of the American Statistical Association*, 118, 343–359.

# Supplementary Materials

## "Uncertainty Quantification for Large Language Model Reward Learning under Heterogeneous Human Feedback"

In this supplement, we provide additional details and technical proofs to support the main results of the paper. Section A and Section B briefly describe the datasets and large language models used in our experiments. Section C presents the detailed proof of Lemma 1. Section D contains the convergence analysis, including the proof of Theorem 1. In Section E, we provide comprehensive proofs for the uncertainty quantification results, covering Lemma 2, Theorem 2, Theorem 3, Theorem 4, Theorem 5, Corollary 1, Lemma 3 and Theorem 6. Additional supporting lemmas are included in Section F.

# A    Description of Datasets

In this section, we give a brief description of the datasets used in Section 6.2. The dataset `prism-alignment`[9] is a diverse human feedback dataset for preference and value alignment in Large Language Models. It maps the characteristics and stated preferences of humans from a detailed survey onto their real-time interactions with LLMs and contextual preference ratings. There are two sequential stages: first, participants complete a Survey where they answer questions about their demographics and stated preferences, then proceed to the Conversations with LLMs, where they input prompts, rate responses and give fine-grained feedback in a series of multi-turn interactions. The evaluation dataset `AlpacaFarm` (Dubois et al., 2023) consists of 805 instructions, which includes 252 instructions from the self-instruct evaluation set (Wang et al., 2023), 188 from the Open Assistant (OASST) evtest setaluation, 129 from Anthropics' helpful test set (Bai et al., 2022), 80 from Vicuna test set (Zheng et al., 2023; Chiang et al., 2023), and 156 from Koala test set (Geng et al., 2023).

---

[9] https://huggingface.co/datasets/HannahRoseKirk/prism-alignment

# B    Description of Large Language Models

In this section, we give a brief description of the large language models used in Section 6.2. All the descriptions are adapted from `Hugging Face`[10]. The model `opt-1.3b`[11] is predominantly pretrained with English text, but a small amount of non-English data is still present within the training corpus via CommonCrawl. The model was pretrained using a causal language modeling objective. The model `zephyr-7b-beta`[12] is the second model in the Zephyr series, and is a fine-tuned version of `mistralai/Mistral-7B-v0.1` that was trained on on a mix of publicly available, synthetic datasets using direct preference optimization. The model `dolly-v2-7b`[13] is an instruction-following large language model trained on the Databricks machine learning platform that is licensed for commercial use. The model `mistral-7b-sft-beta`[14] is a fine-tuned version of `mistralai/Mistral-7B-v0.1` on the `HuggingFaceH4/ultrachat_200k` dataset. It is the SFT model that was used to train `zephyr-7b-beta` with direct preference optimization. The model `Eurus-RM-7B`[15] is trained on a mixture of UltraInteract, UltraFeedback, and UltraSafety, with a specifically designed reward modeling objective for reasoning to directly increase.

# C    Proof of Lemma 1

Recall that $\mu(v) = \frac{1}{1+e^{-v}}$. Its derivative is $\frac{d\mu(v)}{dv} = \frac{e^{-v}}{(1+e^{-v})^2} = \mu(v)[1 - \mu(v)]$. Under Assumption 1, the reward model is $r_\theta(s, a) = \theta^\top \phi(s, a)$. The negative log-likelihood in (6) is then equivalent to

$$L_n(\theta, \gamma) = -\frac{1}{n} \sum_{i=1}^{n} \left\{ y_i \log \mu(\sigma_\gamma(x_i)\theta^\top z_i) + (1 - y_i) \log[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)] \right\}.$$

---

[10]https://huggingface.co/
[11]https://huggingface.co/facebook/opt-1.3b
[12]https://huggingface.co/HuggingFaceH4/zephyr-7b-beta
[13]https://huggingface.co/databricks/dolly-v2-7b
[14]https://huggingface.co/HuggingFaceH4/mistral-7b-sft-beta
[15]https://huggingface.co/openbmb/Eurus-RM-7b

The gradient of $L_n(\theta, \gamma)$ with respect to $\theta$ is

$$\nabla_\theta L_n(\theta, \gamma) = -\frac{1}{n}\sum_{i=1}^{n}\{y_i[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)z_i - (1 - y_i)\mu(\sigma_\gamma(x_i)\theta^\top z_i)\sigma_\gamma(x_i)z_i\}$$
$$= -\frac{1}{n}\sum_{i=1}^{n}[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)z_i. \tag{S1}$$

Hence, the Hessian matrix of $L_n(\theta, \gamma)$ with respect to $\theta$ is

$$\nabla_{\theta\theta}^2 L_n(\theta, \gamma) = \frac{1}{n}\sum_{i=1}^{n}\mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma^2(x_i)z_i z_i^\top \succeq \mathbf{0}. \tag{S2}$$

Therefore, $L_n(\theta, \gamma)$ is convex in $\theta$ when $\gamma$ is fixed. Under Assumption 2, $\sigma_\gamma(x) = \psi_0(x) + \gamma^\top\psi(x)$. The negative log-likelihood given in (6) is equivalent to

$$L_n(\theta, \gamma)$$
$$= -\frac{1}{n}\sum_{i=1}^{n}\left\{y_i\log\mu((\psi_0(x_i) + \gamma^\top\psi(x_i))\theta^\top z_i) + (1 - y_i)\log[1 - \mu((\psi_0(x_i) + \gamma^\top\psi(x_i))\theta^\top z_i)]\right\}.$$

The gradient of $L_n(\theta, \gamma)$ with respect to $\gamma$ is

$$\nabla_\gamma L_n(\theta, \gamma) = -\frac{1}{n}\sum_{i=1}^{n}\{y_i[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)\psi(x_i) - (1 - y_i)\mu(\sigma_\gamma(x_i)\theta^\top z_i)(\theta^\top z_i)\psi(x_i)\}$$
$$= -\frac{1}{n}\sum_{i=1}^{n}[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)\psi(x_i). \tag{S3}$$

Hence, the Hessian matrix of $L_n(\theta, \gamma)$ with respect to $\gamma$ is

$$\nabla_{\gamma\gamma}^2 L_n(\theta, \gamma) = \frac{1}{n}\sum_{i=1}^{n}\mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)^2\psi(x_i)\psi^\top(x_i) \succeq \mathbf{0}. \tag{S4}$$

Therefore, $L_n(\theta, \gamma)$ is convex in $\gamma$ when $\theta$ fixed. Now we calculate the Hessian matrix with respect to the cross terms between $\theta$ and $\gamma$,

$$\nabla_{\gamma\theta}^2 L_n(\theta, \gamma)$$
$$= -\frac{1}{n}\sum_{i=1}^{n}\{y_i\psi(x_i)z_i^\top - \mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)(\theta^\top z_i)\psi(x_i)z_i^\top$$
$$- \mu(\sigma_\gamma(x_i)\theta^\top z_i)\psi(x_i)z_i^\top\}$$
$$= -\frac{1}{n}\sum_{i=1}^{n}\{y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)(\theta^\top z_i)\}\psi(x_i)z_i^\top. \tag{S5}$$

3

It is straightforward to verify that $\nabla^2_{\theta\gamma} L_n(\theta, \gamma) = (\nabla^2_{\gamma\theta} L_n(\theta, \gamma))^\top$.

We now construct an example to show that $L_n(\theta, \gamma)$ is not convex in the joint vector $(\theta^\top, \gamma^\top)^\top$. We consider $d_1 = d_2 = 1$ and evaluate the Hessian matrix at $\theta = 0$. Then, $\mu(\sigma_\gamma(x_i)\theta^\top z_i) = \mu(0) = 0.5$. The Hessian components simplify to

$$\nabla^2_{\gamma\gamma} L_n(0, \gamma) = \frac{1}{n}\sum_{i=1}^n \mu(\sigma_\gamma(x_i)\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)^2 \psi(x_i)\psi^\top(x_i) = 0$$

and

$$\nabla^2_{\gamma\theta} L_n(0, \gamma) = -\frac{1}{n}\sum_{i=1}^n (y_i - 0.5)\psi(x_i)z_i^\top.$$

Now set $n = 2$ and choose $z_1 = 2, z_2 = 1, \psi(x_1) = 1, \psi(x_2) = 2, y_1 = y_2 = 1$. Then

$$\nabla^2_{\gamma\theta} L_2(0, \gamma) = -\frac{1}{2}\big[(1 - 0.5) \times 2 + (1 - 0.5) \times 2\big] = -1.$$

Thus, the joint Hessian matrix at $(\theta^\top, \gamma^\top)^\top = (0, \gamma^\top)^\top$ is

$$\begin{pmatrix} \nabla^2_{\theta\theta} L_n(0, \gamma) & \nabla^2_{\theta\gamma} L_n(0, \gamma) \\ \nabla^2_{\gamma\theta} L_n(0, \gamma) & \nabla^2_{\gamma\gamma} L_n(0, \gamma) \end{pmatrix} = \begin{pmatrix} \nabla^2_{\theta\theta} L_n(0, \gamma) & -1 \\ -1 & 0 \end{pmatrix}.$$

Its determinant is $-1 < 0$. So the joint Hessian matrix is not positive semidefinite. Therefore, $L_n(\theta, \gamma)$ is not jointly convex in $(\theta^\top, \gamma^\top)^\top$ in this example. The proof is complete.

# D    Proofs for Convergence

Before presenting the proof, we establish several lemmas to describe the properties of the negative log-likelihood function.

**Lemma S4.** *Let Assumptions 1 and 2 hold. Denote $K = \max\{\theta^2_{\max}, \gamma^2_{\max}\}\phi^2_{\max}\psi^2_{\max}$. For any $\theta, \theta' \in \mathbb{R}^{d_1}$ and $\gamma \in \Gamma$, we have*

$$\|\nabla_\theta L_n(\theta, \gamma) - \nabla_\theta L_n(\theta', \gamma)\|_2 \le d_2 K \|\theta - \theta'\|_2.$$

*For any $\gamma, \gamma' \in \mathbb{R}^{d_2}$ and $\theta \in \Theta$, we have*

$$\|\nabla_\gamma L_n(\theta, \gamma) - \nabla_\gamma L_n(\theta, \gamma')\|_2 \le d_1 K \|\gamma - \gamma'\|_2.$$

4

Lemma S4 guarantees that the gradients of $L_n(\theta, \gamma)$ with respect to $\theta$ and $\gamma$ are globally Lipschitz continuous. This property controls the accumulation of error terms when we apply the descent step to each block parameter.

**Lemma S5.** *Let $b \geq 0$ be some constant and Assumptions 1 and 2 hold. Denote $\widetilde{M} = \phi_{\max}\psi_{\max}[2 + (b \max\{\sqrt{d_1}\theta_{\max}, \sqrt{d_2}\gamma_{\max}\}/\sqrt{2} + \sqrt{d_1 d_2}\gamma_{\max}\theta_{\max})\phi_{\max}\psi_{\max}]$. When $\theta$ satisfies $\|\theta - \theta_*\|_2 \leq b/\sqrt{2}$, for all $\gamma \in \mathbb{R}^{d_2}$ and $\gamma' \in \Gamma$, we have*

$$\|\nabla_\theta L_n(\theta, \gamma) - \nabla_\theta L_n(\theta, \gamma')\|_2 \leq \widetilde{M}\|\gamma - \gamma'\|_2.$$

*When $\gamma$ satisfies $\|\gamma - \gamma_*\|_2 \leq b/\sqrt{2}$, for all $\theta \in \mathbb{R}^{d_1}$ and $\theta' \in \Theta$, we have*

$$\nabla_\gamma L_n(\theta, \gamma) - \nabla_\gamma L_n(\theta', \gamma)\| \leq \widetilde{M}\|\theta - \theta'\|_2,$$

Lemma S5 measures how sensitively the $\theta$-gradient reacts to perturbations in $\gamma$, within a local neighborhood of $\theta_*$, and how sensitively the $\gamma$-gradient reacts to perturbations in $\theta$, within a local neighborhood of $\gamma_*$. This cross-smoothness is needed when bounding error propagation between the two alternating blocks.

**Lemma S6.** *Let $0 \leq b < \widetilde{b}$ with $\widetilde{b} = \frac{\min\{\lambda_\phi/(\sqrt{d_2}\gamma_{\max}), \lambda_\psi/(\sqrt{d_1}\theta_{\max})\}}{8\sqrt{2}\phi_{\max}^2\psi_{\max}^2}$ and Assumptions 1, 2 and 3 hold. Denote $w = \frac{1}{2c_0}\min\{\lambda_\phi - 8\sqrt{2d_2}\gamma_{\max}\phi_{\max}^2\psi_{\max}^2 b, \lambda_\psi - 8\sqrt{2d_1}\theta_{\max}\phi_{\max}^2\psi_{\max}^2 b\}$ with $c_0 = \frac{(1+e^{C_{\max}})^2}{e^{C_{\max}}}$, where $C_{\max} = (\sqrt{2}\gamma_{\max} + \widetilde{b})(\sqrt{2}\theta_{\max} + \widetilde{b})\sqrt{d_1 d_2}\phi_{\max}\psi_{\max}$. For any $\gamma, \gamma', \theta, \theta'$ satisfying $\|\gamma - \gamma_*\|_2 \leq b/\sqrt{2}$, $\|\gamma' - \gamma_*\|_2 \leq b/\sqrt{2}$, $\|\theta - \theta_*\|_2 \leq b/\sqrt{2}$ and $\|\theta' - \theta_*\|_2 \leq b/\sqrt{2}$, with probability at least $1 - d_1\left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}}$, it holds that*

$$\langle\nabla_\theta L_n(\theta, \gamma) - \nabla_\theta L_n(\theta', \gamma), \theta - \theta'\rangle \geq w\|\theta - \theta'\|_2^2,$$

*and with probability at least $1 - d_2\left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$, we have*

$$\langle\nabla_\gamma L_n(\theta, \gamma) - \nabla_\gamma L_n(\theta, \gamma'), \gamma - \gamma'\rangle \geq w\|\gamma - \gamma'\|_2^2,$$

*where $K$ is defined in Lemma S4.*

Lemma S6 provides a local strong convexity–type result for the $\theta$-block and the $\gamma$-block. The curvature $w$ is strictly positive when $b$ is chosen small enough.

**Lemma S7.** *Let Assumptions 1 and 2 hold. For any $0 < \delta < 1$, with probability at least $1 - \delta$, we have*

$$\|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2 < f(d_1 + 1, n, \delta, 2\sqrt{d_2 K}),$$

*and with probability at least $1 - \delta$, we have*

$$\|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2 < f(d_2 + 1, n, \delta, 2\sqrt{d_1 K}),$$

*where*

$$f(d, n, \delta, K) = \frac{K\{\log(d/\delta) + \sqrt{\log(d/\delta)[\log(d/\delta) + 18n]}\}}{3n}, \tag{S6}$$

*and $K$ is defined in Lemma S4.*

Lemma S7 provide upper bounds for the empirical gradients at the true parameters. It quantifies the statistical noise entering the recursion and produces the $1/n$ scaling in the final statistical error term of Theorem 1.

**Lemma S8.** *Let $M$ be defined in Assumption 4. Let Assumptions 1 and 2 hold. Then, for any $\delta \in (0, 1)$, with probability at least $1 - \delta$, the following bounds hold for any $\gamma$ satisfying $\|\gamma - \gamma_*\|_2 \leq b/\sqrt{2}$:*

$$\|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma)\|_2 \leq cb + M + f(d_1 + d_2, n, \delta, K_1),$$

*and for any $\theta$ satisfying $\|\theta - \theta_*\|_2 \leq b/\sqrt{2}$:*

$$\|\nabla_{\gamma\theta}^2 L_n(\theta, \gamma_*)\|_2 \leq cb + M + f(d_1 + d_2, n, \delta, K_1),$$

*where $c = \sqrt{2}\max\{\theta_{\max}\sqrt{d_1}, \gamma_{\max}\sqrt{d_2}\}\psi_{\max}^2\phi_{\max}^2(1+\sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}\psi_{\max})$, $K_1 = 2\psi_{\max}\phi_{\max} + \sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}^2\psi_{\max}^2 + M$ and $f(d, n, \delta, K)$ is defined in (S6).*

Lemma S8 provides upper bounds on the norm of the mixed second-order derivatives $\nabla_{\gamma\theta}^2 L_n(\theta, \gamma)$ when one block of parameters deviates from the true values.

## D.1 Proof of Lemma S4

By Assumption 1, we have

$$\|z_i\|_2 = \|\phi(s_i, a_i^{(1)}) - \phi(s_i, a_i^{(0)})\|_2 \leq \sqrt{d_1}\|\phi(s_i, a_i^{(1)}) - \phi(s_i, a_i^{(0)})\|_\infty \leq 2\phi_{\max}. \tag{S7}$$

From Assumption 2, it follows that

$$|\sigma_\gamma(x_i)| = |\psi_0(x_i) + \gamma^\top \psi(x_i)| \le \|(1, \gamma^\top)\|_2 \|(\psi_0(x_i), \psi^\top(x_i))\|_2 \le \sqrt{d_2}\gamma_{\max}\psi_{\max}, \qquad \text{(S8)}$$

Note that $\mu(\cdot)[1 - \mu(\cdot)] \le \frac{1}{4}$. By the mean value theorem, for some $\bar{\theta}$ between $\theta$ and $\theta'$, we have

$$
\begin{aligned}
|\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)| &= |\mu(\sigma_\gamma(x_i)\bar{\theta}^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\bar{\theta}^\top z_i)]\sigma_\gamma(x_i)z_i^\top(\theta' - \theta)| \\
&\le \frac{1}{4}|\sigma_\gamma(x_i)|\|z_i\|_2\|\theta' - \theta\|_2 \\
&\le \frac{\sqrt{d_2}\gamma_{\max}\phi_{\max}\psi_{\max}}{2}\|\theta' - \theta\|_2,
\end{aligned}
\qquad \text{(S9)}
$$

where the last inequality follows from (S7) and (S8). Together with (S1), we obtain

$$
\begin{aligned}
&\|\nabla_\theta L_n(\theta, \gamma) - \nabla_\theta L_n(\theta', \gamma)\|_2 \\
&= \frac{1}{n}\left\|\sum_{i=1}^n \{[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)z_i - [y_i - \mu(\sigma_\gamma(x_i)\theta'^\top z_i)]\sigma_\gamma(x_i)z_i\}\right\|_2 \\
&\le \frac{1}{n}\sum_{i=1}^n \|[\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)z_i\|_2 \\
&\le \frac{1}{n}\sum_{i=1}^n \|\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)\|_2 |\sigma_\gamma(x_i)|\|z_i\|_2 \\
&\le d_2\gamma_{\max}^2\phi_{\max}^2\psi_{\max}^2\|\theta' - \theta\|_2.
\end{aligned}
$$

By the mean value theorem, for some value $\bar{\gamma}$ between $\gamma$ and $\gamma'$, we have

$$
\begin{aligned}
|\mu(\sigma_{\gamma'}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)| &= |\mu(\sigma_{\bar{\gamma}}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\bar{\gamma}}(x_i)\theta^\top z_i)](\gamma' - \gamma)^\top \psi(x_i)(\theta^\top z_i)| \\
&\le \frac{1}{4}\|\gamma - \gamma'\|_2\|\psi(x_i)\|_2\|\theta\|_2\|z_i\|_2 \\
&\le \frac{\sqrt{d_1}\theta_{\max}\phi_{\max}\psi_{\max}}{2}\|\gamma - \gamma'\|_2,
\end{aligned}
$$

$$\text{(S10)}$$

7

where the last inequality is from Assumptions 1, 2. Combining with (S3), we have

$$
\begin{aligned}
&\|\nabla_\gamma L_n(\theta, \gamma) - \nabla_\gamma L_n(\theta, \gamma')\|_2 \\
&= \frac{1}{n} \left\| \sum_{i=1}^n \{[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)\psi(x_i) - [y_i - \mu(\sigma_{\gamma'}(x_i)\theta^\top z_i)](\theta^\top z_i)\psi(x_i)\} \right\|_2 \\
&\le \frac{1}{n} \sum_{i=1}^n \|[\mu(\sigma_{\gamma'}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)\psi(x_i)\|_2 \\
&\le \frac{1}{n} \sum_{i=1}^n |\mu(\sigma_{\gamma'}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)| \|\theta\|_2 \|z_i\|_2 \|\psi(x_i)\|_2 \\
&\le d_1 \theta_{\max}^2 \phi_{\max}^2 \psi_{\max}^2 \|\gamma' - \gamma\|_2.
\end{aligned}
$$

Noting that $K = \max\{\theta_{\max}^2, \gamma_{\max}^2\}\phi_{\max}^2\psi_{\max}^2$, the proof is complete.

## D.2  Proof of Lemma S5

Under Assumption 1, by (S1) and (S7), we have

$$
\begin{aligned}
&\|\nabla_\theta L_n(\theta, \gamma) - \nabla_\theta L_n(\theta, \gamma')\|_2 \\
&= \frac{1}{n} \left\| \sum_{i=1}^n \{[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i)z_i - [y_i - \mu(\sigma_{\gamma'}(x_i)\theta^\top z_i)]\sigma_{\gamma'}(x_i)z_i\} \right\|_2 \\
&\le \frac{1}{n} \sum_{i=1}^n \{|[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i) - [y_i - \mu(\sigma_{\gamma'}(x_i)\theta^\top z_i)]\sigma_{\gamma'}(x_i)| \|z_i\|_2\} \\
&\le \frac{2\phi_{\max}}{n} \sum_{i=1}^n |[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_\gamma(x_i) - [y_i - \mu(\sigma_{\gamma'}(x_i)\theta^\top z_i)]\sigma_{\gamma'}(x_i)| \\
&= \frac{2\phi_{\max}}{n} \sum_{i=1}^n |[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)][\sigma_\gamma(x_i) - \sigma_{\gamma'}(x_i)] + [\mu(\sigma_{\gamma'}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_{\gamma'}(x_i)| \\
&\le \frac{2\phi_{\max}}{n} \sum_{i=1}^n \{|[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)][\sigma_\gamma(x_i) - \sigma_{\gamma'}(x_i)]| \\
&\qquad\qquad + |[\mu(\sigma_{\gamma'}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_{\gamma'}(x_i)|\}.
\end{aligned}
\tag{S11}
$$

Since $0 \le y_i \le 1$ and $0 \le \mu(\sigma_\gamma(x_i)\theta^\top z_i) \le 1$, by Assumption 2, we have

$$
|[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)][\sigma_\gamma(x_i) - \sigma_{\gamma'}(x_i)]\|_2 \le |(\gamma - \gamma')^\top \psi(x_i)| \le \psi_{\max}\|\gamma - \gamma'\|_2. \tag{S12}
$$

By (S7), we have

$$
\begin{aligned}
|\theta^\top z_i| &= |(\theta - \theta_* + \theta_*)^\top z_i| \\
&\leq |\theta_*^\top z_i| + |(\theta - \theta_*)^\top z_i| \\
&\leq \|\theta_*\|_2 \|z_i\|_2 + \|\theta - \theta_*\|_2 \|z_i\|_2 \\
&\leq 2\sqrt{d_1}\theta_{\max}\phi_{\max} + 2(b/\sqrt{2})\phi_{\max} \\
&= 2\phi_{\max}(\sqrt{d_1}\theta_{\max} + b/\sqrt{2}).
\end{aligned}
\tag{S13}
$$

By the mean value theorem, for some value $\bar{\gamma}$ between $\gamma$ and $\gamma'$, we have

$$
\begin{aligned}
\|\mu(\sigma_{\gamma'}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)\|_2 &= \|\mu(\sigma_{\bar{\gamma}}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\bar{\gamma}}(x_i)\theta^\top z_i)](\gamma' - \gamma)^\top \psi(x_i)(\theta^\top z_i)\|_2 \\
&\leq \frac{1}{4}\|\gamma - \gamma'\|_2 \|\psi(x_i)\|_2 |\theta^\top z_i| \\
&\leq \frac{(\sqrt{d_1}\theta_{\max} + b/\sqrt{2})\phi_{\max}\psi_{\max}}{2}\|\gamma - \gamma'\|_2,
\end{aligned}
$$

where the last inequality is from Assumption 2 and (S13). Together with (S8), we obtain

$$
\|[\mu(\sigma_{\gamma'}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\sigma_{\gamma'}(x_i)\|_2 \leq \frac{\sqrt{d_2}(\sqrt{d_1}\theta_{\max} + b/\sqrt{2})\gamma_{\max}\phi_{\max}\psi_{\max}^2}{2}\|\gamma - \gamma'\|_2.
\tag{S14}
$$

Combining (S11), (S12) and (S14), we obtain

$$
\|\nabla_\theta L_n(\theta, \gamma) - \nabla_\theta L_n(\theta, \gamma')\|_2 \leq \phi_{\max}\psi_{\max}[2 + \sqrt{d_2}(\sqrt{d_1}\theta_{\max} + b/\sqrt{2})\gamma_{\max}\psi_{\max}\phi_{\max}]\|\gamma - \gamma'\|_2.
$$

Under Assumption 2, by (S3) and (S7), we get

$$
\begin{aligned}
&\|\nabla_\gamma L_n(\theta, \gamma) - \nabla_\gamma L_n(\theta', \gamma)\|_2 \\
&= \frac{1}{n}\left\|\sum_{i=1}^n \{[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta^\top z_i)\psi(x_i) - [y_i - \mu(\sigma_\gamma(x_i)\theta'^\top z_i)](\theta'^\top z_i)\psi(x_i)\}\right\|_2 \\
&\leq \frac{1}{n}\sum_{i=1}^n \{|[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\theta^\top z_i - [y_i - \mu(\sigma_\gamma(x_i)\theta'^\top z_i)]\theta'^\top z_i|\|\psi(x_i)\|_2\} \\
&\leq \frac{\psi_{\max}}{n}\sum_{i=1}^n |[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\theta^\top z_i - [y_i - \mu(\sigma_\gamma(x_i)\theta'^\top z_i)]\theta'^\top z_i| \\
&= \frac{\psi_{\max}}{n}\sum_{i=1}^n |[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta - \theta')^\top z_i + [\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\theta'^\top z_i| \\
&\leq \frac{\psi_{\max}}{n}\sum_{i=1}^n \{|[y_i - \mu(\sigma_\gamma(x_i)\theta^\top z_i)](\theta - \theta')^\top z_i| + |[\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\theta'^\top z_i|\}.
\end{aligned}
\tag{S15}
$$

9

Since $0 \le y_i \le 1$ and $0 \le \mu(\sigma_\gamma(x_i)\theta^\top z_i) \le 1$, by (S7), we have

$$|[y_i - \mu(\sigma_\gamma(s_i)\theta'^\top z_i)](\theta - \theta')^\top z_i| \le \|\theta - \theta'\|_2 \|z_i\|_2 \le 2\phi_{\max}\|\theta - \theta'\|_2. \qquad (S16)$$

We express $\sigma_\gamma(x_i)$ as

$$\sigma_\gamma(x_i) = \psi_0(x_i) + (\gamma - \gamma_* + \gamma_*)^\top \psi(x_i) = \sigma_{\gamma_*}(x_i) + (\gamma - \gamma_*)^\top \psi(x_i). \qquad (S17)$$

By (S8) and (S17), we have

$$\sigma_\gamma(x_i) \le |\sigma_{\gamma_*}(x_i)| + |(\gamma - \gamma_*)^\top \psi(x_i)| \le \sqrt{d_2}\gamma_{\max}\psi_{\max} + b\psi_{\max}/\sqrt{2} = \psi_{\max}(\sqrt{d_2}\gamma_{\max} + b/\sqrt{2}).$$
$$(S18)$$

By the mean value theorem, for some value $\bar\theta$ between $\theta$ and $\theta'$, we have

$$\begin{aligned}
|\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)| &= |\mu(\sigma_\gamma(x_i)\bar\theta^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\bar\theta^\top z_i)]\sigma_\gamma(x_i)(\theta - \theta')^\top z_i| \\
&\le \frac{1}{4}|\sigma_\gamma(x_i)|\|z_i\|_2\|\theta - \theta'\|_2 \\
&\le \frac{(b/\sqrt{2} + \sqrt{d_2}\gamma_{\max})\phi_{\max}\psi_{\max}}{2}\|\theta - \theta'\|_2,
\end{aligned}$$

where the last inequality is from (S7) and (S18). Together with Assumption 1, we obtain

$$\begin{aligned}
|[\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)]\theta'^\top z_i| &\le |\mu(\sigma_\gamma(x_i)\theta'^\top z_i) - \mu(\sigma_\gamma(x_i)\theta^\top z_i)|\|\theta'\|_2\|z_i\|_2 \\
&\le \sqrt{d_1}(b/\sqrt{2} + \sqrt{d_2}\gamma_{\max})\theta_{\max}\phi_{\max}^2\psi_{\max}\|\theta - \theta'\|_2.
\end{aligned}$$
$$(S19)$$

Combining (S15), (S16) and (S19), we conclude

$$\|\nabla_\gamma L_n(\theta, \gamma) - \nabla_\gamma L_n(\theta', \gamma)\|_2 \le \phi_{\max}\psi_{\max}[2 + \sqrt{d_1}(b/\sqrt{2} + \sqrt{d_2}\gamma_{\max})\theta_{\max}\phi_{\max}\psi_{\max}]\|\theta - \theta'\|_2.$$

Noting that $\widetilde{M} = \phi_{\max}\psi_{\max}[2 + (b\max\{\sqrt{d_1}\theta_{\max}, \sqrt{d_2}\gamma_{\max}\}/\sqrt{2} + \sqrt{d_1 d_2}\gamma_{\max}\theta_{\max})\phi_{\max}\psi_{\max}]$,

the proof is complete.

## D.3   Proof of Lemma S6

By Assumption 2 and (S8), we obtain

$$\begin{aligned}
\sigma_\gamma^2(x_i) &= \sigma_{\gamma_*}^2(x_i) + 2\sigma_{\gamma_*}(x_i)(\gamma - \gamma_*)^\top\psi(x_i) + [(\gamma - \gamma_*)^\top\psi(x_i)]^2 \\
&\ge \sigma_{\gamma_*}^2(x_i) - 2|\sigma_{\gamma_*}(x_i)|\|\gamma - \gamma_*\|_2\|\psi(x_i)\|_2 \\
&\ge \sigma_{\gamma_*}^2(x_i) - \sqrt{2d_2}\gamma_{\max}\psi_{\max}^2 b.
\end{aligned}$$

Hence, $\sigma_\gamma^2(x_i)z_iz_i^\top \succeq \sigma_{\gamma_*}^2(x_i)z_iz_i^\top - \sqrt{2d_2}\gamma_{\max}\psi_{\max}^2 b z_iz_i^\top$. By Weyl's inequality, we have

$$\lambda_{\min}\left\{\sum_{i=1}^n \sigma_\gamma^2(x_i)z_iz_i^\top\right\} \geq \lambda_{\min}\left\{\sum_{i=1}^n \sigma_{\gamma_*}^2(x_i)z_iz_i^\top\right\} + \sqrt{2d_2}\gamma_{\max}\psi_{\max}^2 b\lambda_{\min}\left\{-\sum_{i=1}^n z_iz_i^\top\right\}$$

$$= \lambda_{\min}\left\{\sum_{i=1}^n \sigma_{\gamma_*}^2(x_i)z_iz_i^\top\right\} - \sqrt{2d_2}\gamma_{\max}\psi_{\max}^2 b\lambda_{\max}\left\{\sum_{i=1}^n z_iz_i^\top\right\}$$

$$\geq \lambda_{\min}\left\{\sum_{i=1}^n \sigma_{\gamma_*}^2(x_i)z_iz_i^\top\right\} - 4\sqrt{2d_2}n\phi_{\max}^2\gamma_{\max}\psi_{\max}^2 b,$$
(S20)

where the last inequality follows from $\lambda_{\max}\left\{\sum_{i=1}^n z_iz_i^\top\right\} \leq \sum_{i=1}^n \lambda_{\max}(z_iz_i^\top) = \sum_{i=1}^n z_i^\top z_i \leq 4n\phi_{\max}^2$. By (S7) and (S8), the largest eigenvalue of $\sigma_{\gamma_*}^2(x_i)z_iz_i^\top$ is

$$\lambda_{\max}(\sigma_{\gamma_*}^2(x_i)z_iz_i^\top) = \sigma_{\gamma_*}^2(x_i)z_i^\top z_i \leq 4d_2\gamma_{\max}^2\phi_{\max}^2\psi_{\max}^2 \leq 4d_2 K.$$
(S21)

By Assumption 3, the smallest eigenvalue of $\sum_{i=1}^n \mathbb{E}[\sigma_{\gamma_*}^2(x_i)z_iz_i^\top]$ is

$$\lambda_{\min}\left\{\sum_{i=1}^n \mathbb{E}[\sigma_{\gamma_*}^2(x_i)z_iz_i^\top]\right\} = n\lambda_{\min}\{\mathbb{E}[\sigma_{\gamma_*}^2(x_i)z_iz_i^\top]\} = n\lambda_\phi.$$

Clearly, $\{\sigma_{\gamma_*}^2(x_i)z_iz_i^\top\}_{i=1}^n$ is a finite sequence of independent, random, self-adjoint matrices. Applying Lemma S9 with $\zeta = 1/2$, we have

$$\mathbb{P}\left\{\lambda_{\min}\left(\sum_{i=1}^n \sigma_{\gamma_*}^2(x_i)z_iz_i^\top\right) \leq \frac{\lambda_\phi n}{2}\right\} \leq d_1\left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}}.$$

Therefore, with probability at least $1 - d_1\left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}}$, we have

$$\lambda_{\min}\left(\frac{1}{n}\sum_{i=1}^n \sigma_{\gamma_*}^2(x_i)z_iz_i^\top\right) \geq \frac{\lambda_\phi}{2}.$$
(S22)

Let $\bar{\theta} = a\theta + (1-a)\theta'$ for $a \in [0, 1]$. Then $\|\bar{\theta}\|_2 = \|a(\theta - \theta_*) + (1-a)(\theta' - \theta_*) + \theta_*\|_2 \leq a\|\theta - \theta_*\|_2 + (1-a)\|\theta' - \theta_*\|_2 + \|\theta_*\|_2 \leq b/\sqrt{2} + \sqrt{d_1}\theta_{\max}$. We obtain $|\sigma_\gamma(x_i)\bar{\theta}^\top z_i| \leq (\sqrt{2d_2}\gamma_{\max} + b)(\sqrt{2d_1}\theta_{\max} + b)\phi_{\max}\psi_{\max} < C_{\max}$, which implies $-C_{\max} < \sigma_\gamma(x_i)\bar{\theta}^\top z_i < C_{\max}$. Therefore,

$$\mu(\sigma_\gamma(x_i)\bar{\theta}^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\bar{\theta}^\top z_i)] = \frac{e^{-\sigma_\gamma(x_i)\bar{\theta}^\top z_i}}{[1 + e^{-\sigma_\gamma(x_i)\bar{\theta}^\top z_i}]^2} > \frac{e^{C_{\max}}}{(1 + e^{C_{\max}})^2}.$$
(S23)

11

Combining (S2), (S20), (S22) and (S23), the smallest eigenvalue of $\nabla^2_{\theta\theta} L_n(\theta, \gamma)$ is

$$
\begin{aligned}
\lambda_{\min}[\nabla^2_{\theta\theta} L_n(\bar{\theta}, \gamma)] &\geq \frac{e^{C_{\max}}}{(1 + e^{C_{\max}})^2} \lambda_{\min}\left(\frac{1}{n} \sum_{i=1}^{n} \sigma^2_\gamma(x_i) z_i z_i^\top\right) \\
&\geq \frac{e^{C_{\max}}}{(1 + e^{C_{\max}})^2} \left[\lambda_{\min}\left(\frac{1}{n} \sum_{i=1}^{n} \sigma^2_{\gamma_*}(x_i) z_i z_i^\top\right) - 4\sqrt{2d_2} \phi^2_{\max} \gamma_{\max} \psi^2_{\max} b\right] \\
&\geq \frac{e^{C_{\max}}}{2(1 + e^{C_{\max}})^2} (\lambda_\phi - 8\sqrt{2d_2} \phi^2_{\max} \gamma_{\max} \psi^2_{\max} b) \\
&\geq w,
\end{aligned}
$$

(S24)

with probability at least $1 - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}}$. Using Taylor expansion, with some value $\bar{\theta}$ between $\theta$ and $\theta'$, we get

$$
\begin{aligned}
L_n(\theta, \gamma) &= L_n(\theta', \gamma) + \langle \nabla_\theta L_n(\theta', \gamma), \theta - \theta' \rangle + \frac{1}{2}(\theta - \theta')\nabla^2_{\theta\theta} L_n(\bar{\theta}, \gamma)(\theta - \theta')^\top \\
&\geq L_n(\theta', \gamma) + \langle \nabla_\theta L_n(\theta', \gamma), \theta - \theta' \rangle + \frac{w\|\theta - \theta'\|_2^2}{2}.
\end{aligned}
$$

(S25)

Similarly, we get

$$
L_n(\theta', \gamma) \geq L_n(\theta, \gamma) + \langle \nabla_\theta L_n(\theta, \gamma), \theta' - \theta \rangle + \frac{w\|\theta - \theta'\|_2^2}{2}.
$$

(S26)

Since both inequalities (S25) and (S26) hold under the same event in (S22), we conclude that with probability at least $1 - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}}$,

$$
\langle \nabla_\theta L_n(\theta, \gamma) - \nabla_\theta L_n(\theta', \gamma), \theta - \theta' \rangle \geq w\|\theta - \theta'\|_2^2.
$$

By Assumption 1 and (S7), it follows that

$$
\begin{aligned}
(\theta^\top z_i)^2 &= [(\theta - \theta_* + \theta_*)^\top z_i]^2 \\
&= (\theta_*^\top z_i)^2 + 2\theta_*^\top z_i (\theta - \theta_*)^\top z_i + [(\theta - \theta_*)^\top z_i]^2 \\
&\geq (\theta_*^\top z_i)^2 + 2\|\theta_*\|_2 \|z_i\|_2 \|\theta - \theta_*\|_2 \|z_i\|_2 \\
&\geq (\theta_*^\top z_i)^2 - 4\sqrt{2d_1} \theta_{\max} b \phi^2_{\max}.
\end{aligned}
$$

Therefore,

$$
(\theta^\top z_i)^2 \psi(x_i)\psi^\top(x_i) \succeq (\theta_*^\top z_i)^2 \psi(x_i)\psi^\top(x_i) - 4\sqrt{2d_1} \theta_{\max} b \phi^2_{\max} \psi(x_i)\psi^\top(x_i).
$$

By Weyl's inequality, we have

$$
\lambda_{\min}\left\{\sum_{i=1}^{n}(\theta^\top z_i)^2\psi(x_i)\psi^\top(x_i)\right\}
$$

$$
\geq \lambda_{\min}\left\{\sum_{i=1}^{n}(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)\right\} + 4\sqrt{2d_1}\theta_{\max}b\phi_{\max}^2\lambda_{\min}\left\{-\sum_{i=1}^{n}\psi(x_i)\psi^\top(x_i)\right\}
$$

$$
= \lambda_{\min}\left\{\sum_{i=1}^{n}(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)\right\} - 4\sqrt{2d_1}\theta_{\max}b\phi_{\max}^2\lambda_{\max}\left\{\sum_{i=1}^{n}\psi(x_i)\psi^\top(x_i)\right\} \tag{S27}
$$

$$
\geq \lambda_{\min}\left\{\sum_{i=1}^{n}(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)\right\} - 4\sqrt{2d_1}n\theta_{\max}b\phi_{\max}^2\psi_{\max}^2,
$$

where the last inequality follows from $\lambda_{\max}\left\{\sum_{i=1}^{n}\psi(x_i)\psi^\top(x_i)\right\} \leq \sum_{i=1}^{n}\lambda_{\max}[\psi(x_i)\psi^\top(x_i)] = \sum_{i=1}^{n}\psi^\top(x_i)\psi(x_i) \leq n\psi_{\max}^2$. By Assumptions 1, 2 and (S7), the largest eigenvalue of $(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)$ is

$$
\lambda_{\max}((\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)) = (\theta_*^\top z_i)^2\psi^\top(x_i)\psi(x_i) \leq 4d_1\theta_{\max}^2\phi_{\max}^2\psi_{\max}^2 \leq 4d_1 K. \tag{S28}
$$

By Assumption 3, the smallest eigenvalue of $\sum_{i=1}^{n}\mathbb{E}[(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)]$ satisfies

$$
\lambda_{\min}\left\{\sum_{i=1}^{n}\mathbb{E}[(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)]\right\} = n\lambda_{\min}\{\mathbb{E}[(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)]\} = n\lambda_\psi.
$$

Clearly, $\{(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)\}_{i=1}^{n}$ is a finite sequence of independent, random, self-adjoint matrices. By Lemma S9, with $\zeta = 1/2$, we have

$$
\mathbb{P}\left\{\lambda_{\min}\left(\sum_{i=1}^{n}(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)\right) \leq \frac{\lambda_\psi n}{2}\right\} \leq d_2\left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}.
$$

Therefore, with probability at least $1 - d_2\left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$, we have

$$
\lambda_{\min}\left(\frac{1}{n}\sum_{i=1}^{n}(\theta_*^\top z_i)^2\psi(x_i)\psi^\top(x_i)\right) \geq \frac{\lambda_\psi}{2}. \tag{S29}
$$

Let $\bar\gamma = a\gamma + (1-a)\gamma'$ for $a \in [0,1]$. Then $\|\bar\gamma\|_2 = \|a(\gamma - \gamma_*) + (1-a)(\gamma' - \gamma_*) + \gamma_*\|_2 \leq a\|\gamma - \gamma_*\|_2 + (1-a)\|\gamma' - \gamma_*\|_2 + \|\gamma_*\|_2 \leq b/\sqrt{2} + \sqrt{d_2}\gamma_{\max}$. We obtain $|\sigma_{\bar\gamma}(x_i)\theta^\top z_i| \leq (\sqrt{2d_2}\gamma_{\max} + b)(\sqrt{2d_1}\theta_{\max} + b)\phi_{\max}\psi_{\max}$. Similarly to (S23), we obtain the following,

$$
\mu(\sigma_{\bar\gamma}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\bar\gamma}(x_i)\theta^\top z_i)] = \frac{e^{-\sigma_{\bar\gamma}(x_i)\theta^\top z_i}}{[1 + e^{-\sigma_{\bar\gamma}(x_i)\theta^\top z_i}]^2} \geq \frac{e^{C_{\max}}}{(1 + e^{C_{\max}})^2}. \tag{S30}
$$

Combining (S4), (S30), (S27) and (S29), the smallest eigenvalue of $\nabla_{\gamma\gamma}^2 L_n(\theta, \gamma)$ is

$$
\begin{aligned}
\lambda_{\min}[\nabla_{\gamma\gamma}^2 L_n(\theta, \bar{\gamma})] &\geq \frac{e^{C_{\max}}}{(1 + e^{C_{\max}})^2} \lambda_{\min}\left(\frac{1}{n}\sum_{i=1}^n (\theta^\top z_i)^2 \psi(x_i)\psi^\top(x_i)\right) \\
&\geq \frac{e^{C_{\max}}}{(1 + e^{C_{\max}})^2}\left[\lambda_{\min}\left(\frac{1}{n}\sum_{i=1}^n (\theta_*^\top z_i)^2 \psi(x_i)\psi^\top(x_i)\right) - 4\sqrt{2d_1}\theta_{\max}b\phi_{\max}^2\psi_{\max}^2\right] \\
&\geq \frac{e^{C_{\max}}}{2(1 + e^{C_{\max}})^2}(\lambda_\psi - 8\sqrt{2d_1}\theta_{\max}b\phi_{\max}^2\psi_{\max}^2) \\
&\geq w,
\end{aligned}
\tag{S31}
$$

with probability at least $1 - d_2\left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$. Using Taylor expansion, with some value $\bar{\gamma}$ between $\gamma$ and $\gamma'$, we get

$$
\begin{aligned}
L_n(\theta, \gamma) &= L_n(\theta, \gamma') + \langle\nabla_\gamma L_n(\theta, \gamma'), \gamma - \gamma'\rangle + \frac{1}{2}(\gamma - \gamma')\nabla_{\gamma\gamma}^2 L_n(\theta, \bar{\gamma})(\gamma - \gamma')^\top \\
&\geq L_n(\theta, \gamma') + \langle\nabla_\gamma L_n(\theta, \gamma'), \gamma - \gamma'\rangle + \frac{w\|\gamma - \gamma'\|_2^2}{2}.
\end{aligned}
\tag{S32}
$$

Similarly, we get

$$
L_n(\theta, \gamma') \geq L_n(\theta, \gamma) + \langle\nabla_\gamma L_n(\theta, \gamma), \gamma' - \gamma\rangle + \frac{w\|\gamma - \gamma'\|_2^2}{2}.
\tag{S33}
$$

Since (S32) and (S33) hold in the same event (S29), we conclude

$$
\langle\nabla_\gamma L_n(\theta, \gamma) - \nabla_\gamma L_n(\theta, \gamma'), \gamma - \gamma'\rangle \geq w\|\gamma - \gamma'\|_2^2,
$$

with probability at least $1 - d_2\left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$.

## D.4   Proof of Lemma S7

Recall that $z_i = \phi(s_i, a_i^{(1)}) - \phi(s_i, a_i^{(0)})$. By the law of iterated expectations, we have

$$
\mathbb{E}\{[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)z_i\} = \mathbb{E}\{\mathbb{E}\{[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)z_i | x_i, s_i, a_i^{(0)}, a_i^{(1)}\}\} = 0.
$$

By (S7) and (S8), we have $\|\sigma_{\gamma_*}(x_i)z_i\|_2 \leq |\sigma_{\gamma_*}(x_i)|\|z_i\|_2 \leq 2\sqrt{d_2}\gamma_{\max}\psi_{\max}\phi_{\max} \leq 2\sqrt{d_2 K}$, which implies $\|[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)z_i\|_2 \leq 2\sqrt{d_2 K}$. Thus, $\mathbb{E}\|[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)z_i\|_2^2 \leq 4d_2 K$. Applying Lemma S13 to the gradient expression (S1), for any $\epsilon > 0$, we obtain

$$
\begin{aligned}
\mathbb{P}(\|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2 \geq \epsilon) &= \mathbb{P}\left(\left\|\sum_{i=1}^n [y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)z_i\right\|_2 \geq \epsilon n\right) \\
&\leq (d_1 + 1)e^{-\frac{n\epsilon^2/2}{4d_2 K + 2\sqrt{d_2 K}\epsilon/3}},
\end{aligned}
$$

14

which implies that

$$\|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2 < \frac{2\sqrt{d_2 K}\{\log((d_1+1)/\delta) + \sqrt{\log((d_1+1)/\delta)[\log((d_1+1)/\delta) + 18n]}\}}{3n}$$

holds with probability at least $1 - \delta$ for any $0 < \delta < 1$. By the law of iterated expectations, we have

$$\mathbb{E}\{[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\theta_*^\top z_i \psi(x_i)\} = \mathbb{E}\{\mathbb{E}\{[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\theta_*^\top z_i \psi(x_i) | x_i, s_i, a_i^{(0)}, a_i^{(1)}\}\} = 0.$$

Using Assumptions 1, 2, we obtain

$$\|(\theta_*^\top z_i)\psi(x_i)\|_2 \le |\theta_*^\top z_i|\|\psi(x_i)\|_2 \le \sqrt{d_1}\|\theta_*\|_\infty\|z_i\|_2\|\psi(x_i)\|_2 \le 2\sqrt{d_1}\phi_{\max}\theta_{\max}\psi_{\max} \le 2\sqrt{d_1 K}.$$

Hence, $\|[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)]\theta_*^\top z_i \psi(x_i)\|_2 \le |y_i - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)|\|\theta_*^\top z_i \psi(x_i)\|_2 \le 2\sqrt{d_1 K}$.

Next, we obtain $\mathbb{E}\|[y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\theta_*^\top z_i \psi(x_i)\|_2^2 \le 4d_1 K$. By applying Lemma S13 to the gradient expression (S3), for any $\epsilon > 0$, we have

$$\mathbb{P}(\|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2 \ge \epsilon) = \mathbb{P}\left(\left\|\sum_{i=1}^n [y_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\theta_*^\top z_i \psi(x_i)\right\|_2 \ge \epsilon n\right)$$

$$\le (d_2 + 1)e^{-\frac{n\epsilon^2/2}{4d_1 K + 2\sqrt{d_1 K}\epsilon/3}},$$

which implies that

$$\|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2 < \frac{2\sqrt{d_1 K}\{\log((d_2+1)/\delta) + \sqrt{\log((d_2+1)/\delta)[\log((d_2+1)/\delta) + 18n]}\}}{3n}$$

holds with probability at least $1 - \delta$ for any $0 < \delta < 1$.

## D.5 Proof of Lemma S8

We aim to bound the norm $\|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma)\|_2$. Begin by decomposing it as follows:

$$\|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma)\|_2$$

$$= \|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma) - \nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*) + \nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*) - \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)] + \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)]\|_2$$

$$\le \|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma) - \nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)\|_2 + \|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*) - \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)]\|_2 + \|\mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)]\|_2.$$
$$\tag{S34}$$

15

We will bound each term on the right-hand side separately. We first bound $\|\nabla^2_{\gamma\theta}L_n(\theta_*, \gamma) - \nabla^2_{\gamma\theta}L_n(\theta_*, \gamma_*)\|_2$. By (S5), we have

$$
\begin{aligned}
&\nabla^2_{\gamma\theta}L_n(\theta_*, \gamma) - \nabla^2_{\gamma\theta}L_n(\theta_*, \gamma_*) \\
&= \frac{1}{n}\sum_{i=1}^{n}[\mu(\sigma_\gamma(x_i)\theta_*^\top z_i) - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\psi(x_i)z_i^\top \\
&\quad + \frac{1}{n}\sum_{i=1}^{n}(\theta_*^\top z_i)\psi(x_i)z_i^\top\{\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)]\sigma_\gamma(x_i) \\
&\quad - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)\}.
\end{aligned}
\tag{S35}
$$

By (S10), we obtain

$$
\begin{aligned}
\|[\mu(\sigma_\gamma(x_i)\theta_*^\top z_i) - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\psi(x_i)z_i^\top\|_2 &\leq |\mu(\sigma_\gamma(x_i)\theta_*^\top z_i) - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)|\|\psi(x_i)\|_2\|z_i\|_2 \\
&\leq \sqrt{d_1}\theta_{\max}\psi_{\max}^2\phi_{\max}^2\|\gamma - \gamma_*\|_2.
\end{aligned}
\tag{S36}
$$

Next,

$$
\begin{aligned}
&\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)]\sigma_\gamma(x_i) - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i) \\
&= \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)][\sigma_\gamma(x_i) - \sigma_{\gamma_*}(x_i)] \\
&\quad + \{\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\}\sigma_{\gamma_*}(x_i).
\end{aligned}
\tag{S37}
$$

By Assumption 2 and the fact $\mu(\sigma_\gamma(x_i)\theta_*^\top z_i) \in (0, 1)$, we have

$$
|\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)][\sigma_\gamma(x_i) - \sigma_{\gamma_*}(x_i)]| \leq \frac{\psi_{\max}}{4}\|\gamma - \gamma_*\|_2.
\tag{S38}
$$

By the Taylor expansion, there exists some value $\bar{\gamma}$ between $\gamma$ and $\gamma_*$ such that

$$
\begin{aligned}
&|\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]| \\
&= |\mu(\sigma_{\bar\gamma}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\bar\gamma}(x_i)\theta_*^\top z_i)][1 - 2\mu(\sigma_{\bar\gamma}(x_i)\theta_*^\top z_i)](\theta_*^\top z_i)\psi^\top(x_i)(\gamma - \gamma_*)| \\
&\leq \frac{1}{2}\sqrt{d_1}\theta_{\max}\phi_{\max}\psi_{\max}\|\gamma - \gamma_*\|_2,
\end{aligned}
\tag{S39}
$$

where the last inequality follows from Assumptions 1 and 2, and the fact

$$
|\mu(\sigma_{\bar\gamma}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\bar\gamma}(x_i)\theta_*^\top z_i)][1 - 2\mu(\sigma_{\bar\gamma}(x_i)\theta_*^\top z_i)]| \leq \frac{1}{4}|1 - 2\mu(\sigma_{\bar\gamma}(x_i)\theta_*^\top z_i)| \leq \frac{1}{4}.
\tag{S40}
$$

16

Therefore,

$$|\{\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1-\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1-\mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\}\sigma_{\gamma_*}(x_i)|$$

$$\leq |\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1-\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1-\mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]||\sigma_{\gamma_*}(x_i)| \quad \text{(S41)}$$

$$\leq \frac{1}{2}\sqrt{d_1 d_2}\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max}^2\|\gamma-\gamma_*\|_2,$$

where the last inequality is derived from (S39) and (S8). By (S37), (S38) and (S41), we have

$$|\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1-\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)]\sigma_\gamma(x_i) - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1-\mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)|$$

$$\leq \left(\frac{\psi_{\max}}{4} + \frac{1}{2}\sqrt{d_1 d_2}\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max}^2\right)\|\gamma-\gamma_*\|_2$$

$$= \psi_{\max}\left(\frac{1}{4} + \frac{1}{2}\sqrt{d_1 d_2}\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max}\right)\|\gamma-\gamma_*\|_2.$$

$$\text{(S42)}$$

On the other hand,

$$\|(\theta_*^\top z_i)\psi(x_i)z_i^\top\|_2 \leq \|\theta_*\|_2\|z_i\|_2\|\psi(x_i)\|_2\|z_i\|_2 \leq 4\sqrt{d_1}\theta_{\max}\phi_{\max}^2\psi_{\max}.$$

Combining (S35), (S36) and (S42), we obtain

$$\|\nabla_{\gamma\theta}^2 L_n(\theta_*,\gamma) - \nabla_{\gamma\theta}^2 L_n(\theta_*,\gamma_*)\|_2 \leq 2\sqrt{d_1}\theta_{\max}\psi_{\max}^2\phi_{\max}^2(1+\sqrt{d_1 d_2}\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max})\|\gamma-\gamma_*\|_2.$$

$$\text{(S43)}$$

Now, we bound $\|\nabla_{\gamma\theta}^2 L_n(\theta_*,\gamma_*) - \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*,\gamma_*)]\|_2$. Denote $o_i = \{y_i - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i) - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1-\mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)\theta_*^\top z_i\}\psi(x_i)z_i^\top$. Clearly, $\nabla_{\gamma\theta}^2 L_n(\theta_*,\gamma_*) = \frac{1}{n}\sum_{i=1}^n o_i$. By Assumptions 1 and 2, we have

$$\|o_i\|_2 \leq 2\psi_{\max}\phi_{\max} + \sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}^2\psi_{\max}^2 \quad \text{(S44)}$$

and

$$\|\mathbb{E}o_i\|_2 = \|\mathbb{E}[\mathbb{E}(o_i|x_i,z_i)]\|_2$$

$$= \|\mathbb{E}\{\mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1-\mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\sigma_{\gamma_*}(x_i)(\theta_*^\top z_i)\psi(x_i)z_i^\top\}\|_2 \quad \text{(S45)}$$

$$= M.$$

Therefore,

$$\|o_i - \mathbb{E}o_i\|_2 \leq \|o_i\|_2 + \|\mathbb{E}o_i\|_2 \leq 2\psi_{\max}\phi_{\max} + \sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}^2\psi_{\max}^2 + M := K_1. \quad \text{(S46)}$$

17

The matrix variance statistic is

$$\left\| \sum_{i=1}^n \mathbb{E}[(o_i - \mathbb{E}o_i)(o_i - \mathbb{E}o_i)^\top \right\|_2 \leq \sum_{i=1}^n \|\mathbb{E}[(o_i - \mathbb{E}o_i)(o_i - \mathbb{E}o_i)^\top\|_2 \leq nK_1^2. \qquad \text{(S47)}$$

We denote $d = d_1 + d_2$. By Lemma S13, for any $\epsilon \geq 0$, we have

$$\mathbb{P}\left(\left\|\frac{1}{n}\sum_{i=1}^n o_i - \mathbb{E}\frac{1}{n}\sum_{i=1}^n o_i\right\|_2 \geq \epsilon\right) = \mathbb{P}\left(\left\|\sum_{i=1}^n o_i - \mathbb{E}\sum_{i=1}^n o_i\right\|_2 \geq \epsilon n\right)$$

$$\leq (d_1 + d_2)e^{\frac{-\epsilon^2 n^2/2}{nK_1^2 + nK_1\epsilon/3}} \qquad \text{(S48)}$$

$$= de^{\frac{-\epsilon^2 n/2}{K_1^2 + K_1\epsilon/3}}.$$

Equivalently, for any $\delta \in (0,1)$,

$$\mathbb{P}\left(\left\|\frac{1}{n}\sum_{i=1}^n o_i - \mathbb{E}\frac{1}{n}\sum_{i=1}^n o_i\right\|_2 \geq \frac{K_1\{\log(d/\delta) + \sqrt{\log(d/\delta)[\log(d/\delta) + 18n]}\}}{3n}\right) \leq \delta. \qquad \text{(S49)}$$

With at least probability $1 - \delta$, we have

$$\|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*) - \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)]\|_2 = \|\nabla_{\theta\gamma}^2 L(\theta_*, \gamma_*) - \mathbb{E}o_i\|_2$$

$$\leq \frac{K_1\{\log(d/\delta) + \sqrt{\log(d/\delta)[\log(d/\delta) + 18n]}\}}{3n}. \qquad \text{(S50)}$$

By (S34), (S43), (S45) and (S50), with at least probability $1 - \delta$, we have

$$\|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma)\|_2$$

$$\leq 2\sqrt{d_1}\theta_{\max}\psi_{\max}^2\phi_{\max}^2(1 + \sqrt{d_1 d_2}\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max})\|\gamma - \gamma_*\|_2 + M$$

$$+ \frac{K_1\{\log(d/\delta) + \sqrt{\log(d/\delta)[\log(d/\delta) + 18n]}\}}{3n}$$

$$\leq \sqrt{2d_1}b\theta_{\max}\psi_{\max}^2\phi_{\max}^2(1 + \sqrt{d_1 d_2}\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max}) + M$$

$$+ \frac{K_1\{\log(d/\delta) + \sqrt{\log(d/\delta)[\log(d/\delta) + 18n]}\}}{3n}.$$

We turn to bound the norm $\|\nabla_{\gamma\theta}^2 L_n(\theta, \gamma_*)\|_2$. Begin by decomposing it as follows:

$$\|\nabla_{\gamma\theta}^2 L_n(\theta, \gamma_*)\|_2$$

$$= \|\nabla_{\gamma\theta}^2 L_n(\theta, \gamma_*) - \nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*) + \nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*) - \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)] + \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)]\|_2$$

$$\leq \|\nabla_{\gamma\theta}^2 L_n(\theta, \gamma_*) - \nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)\|_2 + \|\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*) - \mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)]\|_2 + \|\mathbb{E}[\nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)]\|_2.$$
$$\text{(S51)}$$

We bound $\|\nabla^2_{\gamma\theta} L_n(\theta, \gamma_*) - \nabla^2_{\gamma\theta} L_n(\theta_*, \gamma_*)\|_2$. By (S5), we have

$$
\begin{aligned}
&\nabla^2_{\gamma\theta} L_n(\theta, \gamma_*) - \nabla^2_{\gamma\theta} L_n(\theta_*, \gamma_*) \\
&= \frac{1}{n} \sum_{i=1}^n [\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i) - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\psi(x_i)z_i^\top \\
&\quad + \frac{1}{n} \sum_{i=1}^n \sigma_{\gamma_*}(x_i)\psi(x_i)z_i^\top \{\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)]\theta^\top z_i \\
&\quad - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\theta_*^\top z_i\}.
\end{aligned}
\tag{S52}
$$

By (S9), we obtain

$$
\|[\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i) - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\psi(x_i)z_i^\top\|_2 = \sqrt{d_2}\gamma_{\max}\psi_{\max}^2\phi_{\max}^2\|\theta - \theta_*\|_2.
\tag{S53}
$$

Next,

$$
\begin{aligned}
&\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)]\theta^\top z_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\theta_*^\top z_i \\
&= \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)](\theta^\top z_i - \theta_*^\top z_i) \\
&\quad + \{\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\}\theta_*^\top z_i.
\end{aligned}
\tag{S54}
$$

By (S7) and the fact $\mu(\sigma_\gamma(x_i)\theta_*^\top z_i) \in (0, 1)$, we have

$$
|\mu(\sigma_\gamma(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_\gamma(x_i)\theta_*^\top z_i)](\theta^\top z_i - \theta_*^\top z_i)| \leq \frac{\phi_{\max}}{2}\|\theta - \theta_*\|_2.
\tag{S55}
$$

By the Taylor expansion, there exists some value $\bar{\theta}$ between $\theta$ and $\theta_*$ such that

$$
\begin{aligned}
&|\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]| \\
&= |\mu(\sigma_{\gamma_*}(x_i)\bar{\theta}^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\bar{\theta}^\top z_i)][1 - 2\mu(\sigma_{\gamma_*}(x_i)\bar{\theta}^\top z_i)]\sigma_{\gamma_*}(x_i)z_i^\top(\theta - \theta_*)| \\
&\leq \frac{\sqrt{d_2}}{2}\gamma_{\max}\phi_{\max}\psi_{\max}\|\theta - \theta_*\|_2,
\end{aligned}
\tag{S56}
$$

where the last inequality follows from Assumptions 1 and 2, and the similar fact of (S40).

Therefore,

$$
\begin{aligned}
&|\{\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\}\theta_*^\top z_i| \\
&\leq |\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)] - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]||\theta_*^\top z_i| \\
&\leq \sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}^2\psi_{\max}\|\theta - \theta_*\|_2,
\end{aligned}
\tag{S57}
$$

where the last inequality is derived from (S56). By (S54), (S55) and (S57), we have

$$
\begin{aligned}
&|\mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta^\top z_i)]\theta^\top z_i - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)[1 - \mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i)]\theta_*^\top z_i| \\
&\leq \phi_{\max}\left(\frac{1}{2} + \sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}\psi_{\max}\right)\|\theta - \theta_*\|_2.
\end{aligned}
\tag{S58}
$$

19

On the other hand,

$$\|\sigma_{\gamma_*}(x_i)\psi(x_i)z_i^\top\|_2 \le \|\gamma_*\|_2\|\psi(x_i)\|_2\|z_i\|_2\|\psi(x_i)\|_2 \le 2\sqrt{d_2}\gamma_{\max}\psi_{\max}^2\phi_{\max}.$$

Combining (S52), (S53) and (S58), we obtain

$$\|\nabla_{\gamma\theta}^2 L_n(\theta,\gamma_*) - \nabla_{\gamma\theta}^2 L_n(\theta_*,\gamma_*)\|_2 \le \sqrt{d_2}\gamma_{\max}\psi_{\max}^2\phi_{\max}^2(2+\sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}\psi_{\max})\|\theta-\theta_*\|_2.$$
(S59)

By (S51), (S59), (S45) and (S50), with at least probability $1-\delta$, we have

$$\|\nabla_{\gamma\theta}^2 L_n(\theta,\gamma_*)\|_2$$
$$\le 2\sqrt{d_2}\gamma_{\max}\psi_{\max}^2\phi_{\max}^2(1+\sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}\psi_{\max})\|\theta-\theta_*\|_2 + M$$
$$\quad + \frac{K_1\{\log(d/\delta)+\sqrt{\log(d/\delta)[\log(d/\delta)+18n]}\}}{3n}$$
$$\le \sqrt{2d_2}b\gamma_{\max}\psi_{\max}^2\phi_{\max}^2(1+\sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}\psi_{\max}) + M$$
$$\quad + \frac{K_1\{\log(d/\delta)+\sqrt{\log(d/\delta)[\log(d/\delta)+18n]}\}}{3n}.$$

Noting that $c = \sqrt{2}\max\{\theta_{\max}\sqrt{d_1},\gamma_{\max}\sqrt{d_2}\}\psi_{\max}^2\phi_{\max}^2(1+\sqrt{d_1 d_2}\gamma_{\max}\theta_{\max}\phi_{\max}\psi_{\max})$, the proof is complete.

## D.6 Proof of Theorem 1

By Assumption 4, we know $\min\{\lambda_\phi,\lambda_\psi\} - 6c_0 M > 0$. Set

$$b_0 = \min\left\{\frac{\lambda_\phi - 6c_0 M}{8\sqrt{2d_2}\gamma_{\max}\phi_{\max}^2\psi_{\max}^2 + 6c_0 c}, \frac{\lambda_\psi - 6c_0 M}{8\sqrt{2d_1}\theta_{\max}\phi_{\max}^2\psi_{\max}^2 + 6c_0 c}\right\}$$

with $c$ defined in Lemma S8. We begin by bounding $\|\theta_{t+1} - \theta_*\|_2^2$. According to the update in Algorithm 1, we have

$$\|\theta_{t+1} - \theta_*\|_2^2 = \|\theta_t - \eta_1 \nabla_\theta L_n(\theta_t,\gamma_t) - \theta_*\|_2^2$$
$$= \|\theta_t - \theta_*\|_2^2 + \eta_1^2\|\nabla_\theta L_n(\theta_t,\gamma_t)\|_2^2 - 2\eta_1\langle\nabla_\theta L_n(\theta_t,\gamma_t),\theta_t - \theta_*\rangle.$$
(S60)

By Lemmas S4 and S5, we obtain

$$\|\nabla_\theta L_n(\theta_t, \gamma_t)\|_2^2$$

$$= \|\nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\theta_t, \gamma_*) + \nabla_\theta L_n(\theta_t, \gamma_*) - \nabla_\theta L_n(\theta_*, \gamma_*) + \nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2$$

$$\leq 3[\|\nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\theta_t, \gamma_*)\|_2^2 + \|\nabla_\theta L_n(\theta_t, \gamma_*) - \nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2 + \|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2]$$

$$\leq 3\widetilde{M}^2 \|\gamma_t - \gamma_*\|_2^2 + 3d_2^2 K^2 \|\theta_t - \theta_*\|_2^2 + 3\|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2. \tag{S61}$$

Next, we have

$$\langle \nabla_\theta L_n(\theta_t, \gamma_t), \theta_t - \theta_* \rangle = \langle \nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_t), \theta_t - \theta_* \rangle + \langle \nabla_\theta L_n(\theta_*, \gamma_*), \theta_t - \theta_* \rangle$$

$$+ \langle \nabla_\theta L_n(\theta_*, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_*), \theta_t - \theta_* \rangle$$

$$\geq \langle \nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_t), \theta_t - \theta_* \rangle - \frac{1}{2w}\|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2$$

$$- \frac{w}{2}\|\theta_t - \theta_*\|_2^2 + \langle \nabla_\theta L_n(\theta_*, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_*), \theta_t - \theta_* \rangle. \tag{S62}$$

By Lemma S6, with probability at least $1 - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_1 K}}$, we have for all $t \in 1, \cdots, T$

$$\langle \nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_t), \theta_t - \theta_* \rangle \geq w\|\theta_t - \theta_*\|_2^2, \tag{S63}$$

From Lemma S8, with probability at least $1 - \delta$, we have

$$|\langle \nabla_\theta L_n(\theta_*, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_*), \theta_t - \theta_* \rangle| \leq \|\nabla_\theta L_n(\theta_*, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_*)\|_2 \|\theta_t - \theta_*\|_2$$

$$\leq \|\nabla_{\theta\gamma}^2 L_n(\theta_*, \bar\gamma_t)\|_2 \|\gamma_t - \gamma_*\|_2 \|\theta_t - \theta_*\|_2$$

$$\leq (cb + M + f(d_1 + d_2, n, \delta, K_1)) \|\gamma_t - \gamma_*\|_2 \|\theta_t - \theta_*\|_2$$

$$\leq \frac{cb + M + f(d_1 + d_2, n, \delta, K_1)}{2}(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_t - \theta_*\|_2^2).$$

Therefore,

$$\langle \nabla_\theta L_n(\theta_*, \gamma_t) - \nabla_\theta L_n(\theta_*, \gamma_*), \theta_t - \theta_* \rangle \geq -\frac{cb + M}{2}(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_t - \theta_*\|_2^2)$$

$$- \frac{f(d_1 + d_2, n, \delta, K_1)}{2}(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_t - \theta_*\|_2^2). \tag{S64}$$

21

Substituting (S63) and (S64) into (S62), we obtain

$$\langle \nabla_\theta L_n(\theta_t, \gamma_t), \theta_t - \theta_* \rangle \geq \left( \frac{w - cb - M}{2} \right) \|\theta_t - \theta_*\|_2^2 - \frac{cb + M}{2} \|\gamma_t - \gamma_*\|_2^2$$

$$- \frac{f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \gamma_*\|_2^2 + \|\theta_t - \theta_*\|_2^2) - \frac{1}{2w} \|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2.$$
(S65)

Plugging (S61) and (S65) into (S60), with probability at least $1 - \delta - d_1 \left( \frac{e}{2} \right)^{-\frac{\lambda_\phi n}{8d_2 K}}$, we have

$$\|\theta_{t+1} - \theta_*\|_2^2 \leq [1 + 3d_2^2 K^2 \eta_1^2 - (w - cb - M)\eta_1] \|\theta_t - \theta_*\|_2^2 + (3\widetilde{M}^2 \eta_1 + cb + M)\eta_1 \|\gamma_t - \gamma_*\|_2^2$$

$$+ \eta_1 f(d_1 + d_2, n, \delta, K_1)(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_t - \theta_*\|_2^2) + (3\eta_1^2 + \eta_1/w) \|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2.$$
(S66)

Now we turn to bound $\|\gamma_{t+1} - \gamma_*\|_2^2$. According to the update in Algorithm 1, we have

$$\|\gamma_{t+1} - \gamma_*\|_2^2 = \|\gamma_t - \eta_2 \nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \gamma_*\|_2^2$$

$$= \|\gamma_t - \gamma_*\|_2^2 + \eta_2^2 \|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t)\|_2^2 - 2\eta_2 \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t), \gamma_t - \gamma_* \rangle.$$
(S67)

By Lemmas S4 and S5, we obtain

$$\|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t)\|_2^2$$

$$= \|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_*, \gamma_t) + \nabla_\gamma L_n(\theta_*, \gamma_t) - \nabla_\gamma L_n(\theta_*, \gamma_*) + \nabla_\gamma L_n(\theta_*, \gamma_*)\|_2^2$$

$$\leq 3[\|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_*, \gamma_t)\|_2 + \|\nabla_\gamma L_n(\theta_*, \gamma_t) - \nabla_\gamma L_n(\theta_*, \gamma_*)\|_2 + \|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2^2]$$

$$\leq 3\widetilde{M}^2 \|\theta_{t+1} - \theta_*\|_2^2 + 3d_1^2 K^2 \|\gamma_t - \gamma_*\|_2^2 + 3\|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2^2.$$
(S68)

Next,

$$\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t), \gamma_t - \gamma_* \rangle = \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_{t+1}, \gamma_*), \gamma_t - \gamma_* \rangle + \langle \nabla_\gamma L_n(\theta_*, \gamma_*), \gamma_t - \gamma_* \rangle$$

$$+ \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_*) - \nabla_\gamma L_n(\theta_*, \gamma_*), \gamma_t - \gamma_* \rangle$$

$$\geq \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_{t+1}, \gamma_*), \gamma_t - \gamma_* \rangle - \frac{1}{2w} \|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2^2$$

$$- \frac{w}{2} \|\gamma_t - \gamma_*\|_2^2 + \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_*) - \nabla_\gamma L_n(\theta_*, \gamma_*), \gamma_t - \gamma_* \rangle.$$
(S69)

By Lemma S6, the following

$$\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_{t+1}, \gamma_*), \gamma_t - \gamma_* \rangle \geq w \|\gamma_t - \gamma_*\|_2^2,$$
(S70)

holds with probability at least $1 - d_2 \left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8 d_1 K}}$, uniformly over $t$. By Lemma S8, with probability at least $1 - \delta$, we have

$$
\begin{aligned}
&|\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_*) - \nabla_\gamma L_n(\theta_*, \gamma_*), \gamma_t - \gamma_* \rangle| \\
&\leq \|\nabla_\gamma L_n(\theta_{t+1}, \gamma_*) - \nabla_\gamma L_n(\theta_*, \gamma_*)\|_2 \|\theta_{t+1} - \theta_*\|_2 \|\gamma_t - \gamma_*\|_2 \\
&\leq \|\nabla_{\theta\gamma}^2 L_n(\bar{\theta}_t, \gamma_*)\|_2 \|\theta_{t+1} - \theta_*\|_2 \|\gamma_t - \gamma_*\|_2 \\
&\leq (cb + M + f(d_1 + d_2, n, \delta, K_1)) \|\theta_{t+1} - \theta_*\|_2 \|\gamma_t - \gamma_*\|_2 \\
&\leq \frac{cb + M + f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \gamma_*\|_2^2 + \|\theta_{t+1} - \theta_*\|_2^2).
\end{aligned}
$$

Therefore,

$$
\begin{aligned}
\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_*) - \nabla_\gamma L_n(\theta_*, \gamma_*), \gamma_t - \gamma_* \rangle \geq {} & -\frac{f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \gamma_*\|_2^2 + \|\theta_{t+1} - \theta_*\|_2^2) \\
& -\frac{cb + M}{2} (\|\gamma_t - \gamma_*\|_2^2 + \|\theta_{t+1} - \theta_*\|_2^2).
\end{aligned}
$$

$$(\text{S71})$$

By (S69), (S70) and (S71), we conclude

$$
\begin{aligned}
\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t), \gamma_t - \gamma_* \rangle \geq {} & \left(\frac{w - cb - M}{2}\right) \|\gamma_t - \gamma_*\|_2^2 - \frac{cb + M}{2} \|\theta_{t+1} - \theta_*\|_2^2 \\
& -\frac{f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \gamma_*\|_2^2 + \|\theta_{t+1} - \theta_*\|_2^2) - \frac{1}{2w} \|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2.
\end{aligned}
$$

$$(\text{S72})$$

Combining (S67), (S68) and (S72), with probability at least $1 - \delta - d_2 \left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8 d_1 K}}$, we have

$$
\begin{aligned}
\|\gamma_{t+1} - \gamma_*\|_2^2 \leq {} & [1 + 3 d_1^2 K^2 \eta_2^2 - (w - cb - M)\eta_2] \|\gamma_t - \gamma_*\|_2^2 + (3 \widetilde{M}^2 \eta_2 + cb + M)\eta_2 \|\theta_{t+1} - \theta_*\|_2^2 \\
& + \eta f(d_1 + d_2, n, \delta, K_1)(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_{t+1} - \theta_*\|_2^2) + (3\eta_2^2 + \eta_2/w)\|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2^2.
\end{aligned}
$$

$$(\text{S73})$$

Under the condition $0 \leq b < b_0$, we can verify $w - 3(cb + M) > 0$. Denote

$$
\eta = \max\{\eta_1, \eta_2\}. \tag{S74}
$$

We define $A_1 = 1 + \max\{3 d_2^2 K^2 \eta_1^2 - (w - cb - M)\eta_1, 3 d_1^2 K^2 \eta_2^2 - (w - cb - M)\eta_2\}$ and $A_2 = 3\widetilde{M}^2 \eta^2 + (cb + M)\eta$. If $(w - cb - M)^2 - 12 d_2^2 K^2 \geq 0$, set

$$
\eta_1 \in \left(0, \min\left\{\frac{w - cb - M - \sqrt{(w - cb - M)^2 - 12 d_2^2 K^2}}{6 d_2^2 K^2}, \frac{w - 3(cb + M)}{3(d_2^2 K^2 + \widetilde{M}^2)}\right\}\right).
$$

Otherwise, set $\eta_1 \in \left(0, \frac{w-3(cb+M)}{3d_2^2 K^2 + \widetilde{M}^2}\right)$. If $(w - cb - M)^2 - 12d_1^2 K^2 \geq 0$, set

$$\eta_2 \in \left(0, \min\left\{\frac{w - cb - M - \sqrt{(w - cb - M)^2 - 12d_1^2 K^2}}{6d_1^2 K^2}, \frac{w - 3(cb + M)}{3(d_1^2 K^2 + \widetilde{M}^2)}\right\}\right).$$

Otherwise, set $\eta_2 \in \left(0, \frac{w-3(cb+M)}{3d_1^2 K^2 + \widetilde{M}^2}\right)$. Then, it is easy to verify

$$A_1 > 0, A_2 > 0 \text{ and } A_1 + 2A_2 \in (0, 1).$$

We define

$$A_{1n} = A_1 + \eta f(d_1 + d_2, n, \delta, K_1) \text{ and } A_{2n} = A_2 + \eta f(d_1 + d_2, n, \delta, K_1). \quad \text{(S75)}$$

Since $A_1, A_2, \eta$ are constants and $f(d_1 + d_2, n, \delta, K_1)$ decreases as $n$ increases, there exists some constant $n_1$ such that when $n > n_1$, we have

$$A_{1n} > 0, A_{2n} > 0 \text{ and } A_{1n} + 2A_{2n} \in (0, 1). \quad \text{(S76)}$$

It follows that

$$A_{1n} \in (0, 1), A_{2n} \in (0, 1), A_{1n} + A_{2n} \in (0, 1) \quad \text{(S77)}$$

and

$$A_{1n} + A_{2n}^2 + A_{1n}A_{2n} = A_{1n} + A_{2n}(A_{1n} + A_{2n}) < A_{1n} + A_{2n} \in (0, 1). \quad \text{(S78)}$$

By (S66), (S73) and Lemma S7, with probability at least $1 - 3\delta - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}} - d_2 \left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$ for $0 < \delta < 1$, we have

$$\|\theta_{t+1} - \theta_*\|_2^2 \leq A_1 \|\theta_t - \theta_*\|_2^2 + A_2 \|\gamma_t - \gamma_*\|_2^2 + (3\eta^2 + \eta/w)\|\nabla_\theta L_n(\theta_*, \gamma_*)\|_2^2$$

$$+ \eta f_3(n, \delta)(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_t - \theta_*\|_2^2) \quad \text{(S79)}$$

$$\leq A_{1n}\|\theta_t - \theta_*\|_2^2 + A_{2n}\|\gamma_t - \gamma_*\|_2^2 + (3\eta^2 + \eta/w)f^2(d_1 + 1, n, \delta, 2\sqrt{d_2}K)$$

and

$$\|\gamma_{t+1} - \gamma_*\|_2^2$$

$$\leq A_1 \|\gamma_t - \gamma_*\|_2^2 + A_2\|\theta_{t+1} - \theta_*\|_2^2 + (3\eta^2 + \eta/w)\|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2^2$$

$$+ \eta f_3(n, \delta)(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_{t+1} - \theta_*\|_2^2)$$

$$= A_{1n}\|\gamma_t - \gamma_*\|_2^2 + A_{2n}\|\theta_{t+1} - \theta_*\|_2^2 + (3\eta^2 + \eta/w)\|\nabla_\gamma L_n(\theta_*, \gamma_*)\|_2^2$$

$$\leq (A_{1n} + A_{2n}^2)\|\gamma_t - \gamma_*\|_2^2 + A_{1n}A_{2n}\|\theta_t - \theta_*\|_2^2 + (A_{2n} + 1)(3\eta^2 + \eta/w)f^2(d_2 + 1, n, \delta, 2\sqrt{d_1}K).$$

$$\text{(S80)}$$

When calculating the above probability, we should note that the probability $\delta$ in (S66) and (S73) is derived from the same event. Since $f(d_1 + 1, n, \delta, 2\sqrt{d_2 K})$ and $f(d_2 + 1, n, \delta, 2\sqrt{d_1 K})$ decrease in $n$, there exists $n_2$ such that when $n \geq n_2$ we have

$$\frac{(A_{1n} + A_{2n})b^2}{2} + (3\eta^2 + \eta/w)f^2(d_1 + 1, n, \delta, 2\sqrt{d_2 K}) \leq \frac{b^2}{2}$$

and

$$\frac{(A_{1n} + A_{2n}^2 + A_{1n}A_{2n})b^2}{2} + (A_{2n} + 1)(3\eta^2 + \eta/w)f^2(d_2 + 1, n, \delta, 2\sqrt{d_1 K}) \leq \frac{b^2}{2}.$$

Note that we have $\|\theta_0 - \theta_*\|_2^2 \leq b^2/2$ and $\|\gamma_0 - \gamma_*\|_2^2 \leq b^2/2$ by initialization. Using (S79) and (S80), we can inductively prove that $\|\theta_t - \theta_*\|_2^2 \leq b^2/2$ and $\|\gamma_t - \gamma_*\|_2^2 \leq b^2/2$ for all $t \geq 0$ when $n \geq n_2$. By (S79) and (S80), we obtain

$$\|\theta_{t+1} - \theta_*\|_2^2 + \|\gamma_{t+1} - \gamma_*\|_2^2 \leq (A_{1n} + A_{2n} + A_{2n}^2)\|\gamma_t - \gamma_*\|_2^2 + (A_{1n} + A_{1n}A_{2n})\|\theta_t - \theta_*\|_2^2$$

$$+ (A_{2n} + 2)(3\eta^2 + \eta/w)f^2(\max\{d_1, d_2\} + 1, n, \delta, 2\sqrt{(d_1 + d_2)K})$$

$$\leq (A_{1n} + A_{2n} + A_{2n}^2)(\|\gamma_t - \gamma_*\|_2^2 + \|\theta_t - \theta_*\|_2^2)$$

$$(A_{2n} + 2)(3\eta^2 + \eta/w)f^2(\max\{d_1, d_2\} + 1, n, \delta, 2\sqrt{(d_1 + d_2)K}),$$

where the second inequality follows from

$$A_{1n} + A_{2n} + A_{2n}^2 - (A_{1n} + A_{1n}A_{2n}) = A_{2n}(1 + A_{2n} - A_{1n}) > 0. \tag{S81}$$

Let $n_0 \geq \max\{n_1, n_2\}$ be chosen sufficiently large such that $d_2\left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_1 K}} < \frac{1}{2n}$ and $d_1\left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}} < \frac{1}{2n}$ hold for all $n \geq n_0$. We denote $\rho = A_{1n_0} + A_{2n_0} + A_{2n_0}^2$. By (S76), we have

$$0 < \rho < A_{1n_0} + 2A_{2n_0} < 1. \tag{S82}$$

By recursion, with probability at least $1 - 3\delta - \frac{1}{n}$, when $n > n_0$, we obtain

$$\|\theta_T - \theta_*\|_2^2 + \|\gamma_T - \gamma_*\|_2^2 \leq \rho^T b^2 + \frac{(A_{2n} + 2)(3w\eta^2 + \eta)f^2(\max\{d_1, d_2\} + 1, n, \delta, 2\sqrt{(d_1 + d_2)K})}{(1 - \rho)w}$$

$$< \rho^T b^2 + \frac{c_2' f^2(\max\{d_1, d_2\} + 1, n, \delta, 2\sqrt{(d_1 + d_2)K})}{(1 - \rho)},$$

25

where $c_2' = (A_{2n_0} + 2)(3w\eta^2 + \eta)/w$. By (S6), there exists some constant $c_1$ such that $c_2' f^2(\max\{d_1, d_2\} + 1, n, \delta, 2\sqrt{(d_1 + d_2)K}) \leq \frac{c_1 \log(1/\delta)}{n}$ since both $d_1$ and $d_2$ are fixed. Since $\delta$ can be any value in $(0, 1)$, the probability can be equivalently expressed as $1 - \delta - \frac{1}{n}$. This completes the proof.

# E  Proofs for Uncertainty Quantification

This section contains the proofs for the theoretical results in Section 4.2 and Section 5.

## E.1  Proof of Lemma 2

By Assumptions 1 and 2, we have $|\sigma_{\gamma_*}(x)\theta_*^\top z| \leq 2\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max}$. Then

$$\mu(\sigma_{\gamma_*}(x)\theta_*^\top z)[1 - \mu(\sigma_{\gamma_*}(x)\theta_*^\top z)] \geq \frac{e^{2\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max}}}{(1 + e^{2\theta_{\max}\gamma_{\max}\phi_{\max}\psi_{\max}})^2} := \kappa > 0.$$

By (S4) and Assumption 3, we have

$$\mathcal{I}_{\gamma\gamma}(\theta_*, \gamma_*) = \mathbb{E}\{\mu(\sigma_{\gamma_*}(x)\theta_*^\top z)[1 - \mu(\sigma_{\gamma_*}(x)\theta_*^\top z)](\theta_*^\top z)^2 \psi(x)\psi^\top(x)\} \succeq \kappa \mathbb{E}[(\theta_*^\top z)^2 \psi(x)\psi^\top(x)] \succeq \kappa\lambda_\psi I,$$

which is positive definite. Under Assumption 5, $\mathcal{I}(\tau_*)$ is positive definite by following Lemma S11.

## E.2  Proof of Theorem 2

Denote $\widehat{\tau}_n = (\widehat{\theta}_n^\top, \widehat{\gamma}_n^\top)^\top$. The proof of Theorem 2 proceeds by first showing that $\|\tau_T - \widehat{\tau}_n\|_2 \to 0$ as $T \to \infty$. Once this consistency is established, the asymptotic normality of $\tau_T$ follows from that of the MLE $\widehat{\tau}_n$, by appropriately controlling the growth rate of $T$ relative to $n$. We first state the asymptotic normality of the MLE $\widehat{\tau}_n$.

**Theorem 7.** *Let the assumptions of Lemma 2 hold. As $n \to \infty$, we have*

$$\sqrt{n}(\widehat{\tau}_n - \tau_*) \xrightarrow{D} \mathcal{N}(0, \mathcal{I}^{-1}(\tau_*)).$$

Theorem 7 provides the classical normal approximation for the MLE, with the covariance structure determined by the inverse of the Fisher information matrix. The proof follows the

standard arguments for asymptotic normality of MLE (see, e.g., Vaart (1998); Casella and Berger (2002)). For completeness, we provide the detailed proof.

*Proof.* We expand the first derivative of the negative log-likelihood around the true value $\tau_*$,

$$\nabla L_n(\widehat{\tau}_n) = \nabla L_n(\tau_*) + \nabla^2 L_n(\tau_*)(\widehat{\tau}_n - \tau_*) + [\nabla^2 L_n(\bar{\tau}) - \nabla^2 L_n(\tau_*)](\widehat{\tau}_n - \tau_*)$$

where

$$\bar{\tau} = c\widehat{\tau}_n + (1 - c)\tau_* \tag{S83}$$

for some $c \in [0, 1]$. Note that $\nabla L_n(\widehat{\tau}_n) = \mathbf{0}$. Rearranging and multiplying by $\sqrt{n}$ gives us

$$\sqrt{n}\{\nabla^2 L_n(\tau_*) + [\nabla^2 L_n(\bar{\tau}) - \nabla^2 L_n(\tau_*)]\}(\widehat{\tau}_n - \tau_*) = -\sqrt{n}\nabla L_n(\tau_*). \tag{S84}$$

By (S2), Assumptions 1 and 2, there exist some positive constants $C_1$ and $C_1'$ such that

$$
\begin{aligned}
\|\nabla^2_{\theta\theta} L_n(\bar{\theta}, \bar{\gamma}) - \nabla^2_{\theta\theta} L_n(\theta_*, \gamma_*)\|_2 &\leq \frac{C_1}{n} \sum_{i=1}^{n} |\sigma^2_{\bar{\gamma}}(x_i) - \sigma^2_{\gamma_*}(x_i)| \\
&= \frac{C_1}{n} \sum_{i=1}^{n} |[\sigma_{\bar{\gamma}}(x_i) + \sigma_{\gamma_*}(x_i)][\sigma_{\bar{\gamma}}(x_i) - \sigma_{\gamma_*}(x_i)]| \\
&\leq C_1' \|\bar{\gamma} - \gamma_*\|_2.
\end{aligned}
$$

By (S4), Assumptions 1 and 2, there exist some positive constants $C_2$ and $C_2'$ such that

$$
\begin{aligned}
\|\nabla^2_{\gamma\gamma} L_n(\bar{\theta}, \bar{\gamma}) - \nabla^2_{\gamma\gamma} L_n(\theta_*, \gamma_*)\|_2 &\leq \frac{C_2}{n} \sum_{i=1}^{n} |(\bar{\theta}^\top z_i)^2 - (\theta_*^\top z_i)^2| \\
&= \frac{C_2}{n} \sum_{i=1}^{n} |[\bar{\theta}^\top z_i + \theta_*^\top z_i][\bar{\theta}^\top z_i - \theta_*^\top z_i]| \\
&\leq C_2' \|\bar{\theta} - \theta_*\|_2.
\end{aligned}
$$

By (S5), Assumptions 1 and 2, there exist some positive constants $C_3, C_3'$ and $C_3''$ such that

$$\|\nabla_{\gamma\theta}^2 L_n(\bar{\theta}, \bar{\gamma}) - \nabla_{\gamma\theta}^2 L_n(\theta_*, \gamma_*)\|_2$$

$$\leq \frac{C_3}{n} \sum_{i=1}^{n} [|\mu(\sigma_{\gamma_*}(x_i)\theta_*^\top z_i) - \mu(\sigma_{\bar{\gamma}}(x_i)\bar{\theta}^\top z_i)| + |\sigma_{\bar{\gamma}}(x_i)\bar{\theta}^\top z_i - \sigma_{\gamma_*}(x_i)\theta_*^\top z_i|]$$

$$\leq \frac{C_3'}{n} \sum_{i=1}^{n} |\sigma_{\bar{\gamma}}(x_i)\bar{\theta}^\top z_i - \sigma_{\gamma_*}(x_i)\theta_*^\top z_i|$$

$$= \frac{C_3'}{n} \sum_{i=1}^{n} |\sigma_{\bar{\gamma}}(x_i)\bar{\theta}^\top z_i - \sigma_{\bar{\gamma}}(x_i)\theta_*^\top z_i + \sigma_{\bar{\gamma}}(x_i)\theta_*^\top z_i - \sigma_{\gamma_*}(x_i)\theta_*^\top z_i|$$

$$\leq C_3''(\|\bar{\theta} - \theta_*\|_2 + \|\bar{\gamma} - \gamma_*\|_2),$$

where the second inequality follows from the mean value theorem. Therefore, there exists some constant $C$ such that

$$\|\nabla^2 L_n(\tau_*) - \nabla^2 L_n(\bar{\tau})\|_2 \leq C(\|\bar{\theta} - \theta_*\|_2 + \|\bar{\gamma} - \gamma_*\|_2) = cC(\|\widehat{\theta}_n - \theta_*\|_2 + \|\widehat{\gamma}_n - \gamma_*\|_2),$$

where the last equality follows from (S83). Therefore, by Lemma S12, we obtain

$$\|\nabla^2 L_n(\tau_*) - \nabla^2 L_n(\bar{\tau})\|_2 \xrightarrow{p} 0. \tag{S85}$$

By the central limit theorem, we have

$$-\sqrt{n}\nabla L_n(\tau_*) \xrightarrow{D} \mathcal{N}(0, \mathcal{I}(\tau_*)). \tag{S86}$$

Combining (S84), (S85) and (S86), by Slutsky's theorem, we have

$$\sqrt{n}(\widehat{\tau}_n - \tau_*) \xrightarrow{D} \mathcal{N}(0, \mathcal{I}^{-1}(\tau_*)).$$

□

Next, we propose the following Theorem that demonstrates that the estimator from Algorithm 1 converges to the MLE when MLE is close to the true parameters.

**Theorem 8.** *Let the conditions in Theorem 1 hold. Suppose* $\|\widehat{\theta}_n - \theta_*\|_2 \leq \epsilon/\sqrt{2}$ *and* $\|\widehat{\gamma}_n - \gamma_*\|_2 \leq \epsilon/\sqrt{2}$ *when* $n \geq n_\epsilon$ *for some constant* $0 \leq \epsilon < b$. *When* $n > \max\{n_0, n_\epsilon\}$, *with probability at least* $1 - \delta - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}} - d_2 \left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$, *we have*

$$\|\theta_T - \widehat{\theta}_n\|_2^2 + \|\gamma_T - \widehat{\gamma}_n\|_2^2 \leq \rho_1^T (b + \epsilon)^2,$$

*where* $0 < \rho_1 < 1$.

The condition $\|\widehat{\theta}_n - \theta_*\|_2 \le \epsilon/\sqrt{2}$ and $\|\widehat{\gamma}_n - \gamma_*\|_2 \le \epsilon/\sqrt{2}$ for large enough $n$ can be guaranteed by the consistency of MLE, see Lemma S12. The proof of Theorem 8 is postponed and presented after the proof of Theorem 2.

Now we prove Theorem 2. By Lemma S12, for any $\epsilon > 0$, there exists $n_\epsilon$ such that $\|\widehat{\theta}_n - \theta_*\|_2 \le \epsilon/\sqrt{2}$ and $\|\widehat{\gamma}_n - \gamma_*\|_2 \le \epsilon/\sqrt{2}$ for all $n > n_\epsilon$ with probability 1. By Theorem 8, with probability at least $1 - \delta - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8 d_2 K}} - d_2 \left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8 d_1 K}}$, we have

$$\sqrt{n}\|\tau_T - \widehat{\tau}_n\|_2 = \sqrt{n}\sqrt{\|\theta_T - \widehat{\theta}_n\|_2^2 + \|\gamma_T - \widehat{\gamma}_n\|_2^2} \le \sqrt{n}\sqrt{\rho_1^T (b + \epsilon)^2} = (b + \epsilon)\sqrt{n\rho_1^T}.$$

Since $n\rho_1^T \to 0$ as $T \to \infty$ and $n \to \infty$, we have

$$\sqrt{n}\|\tau_T - \widehat{\tau}_n\|_2 \xrightarrow{p} 0.$$

Therefore, by Slutsky's theorem and Theorem 7, we have

$$\sqrt{n}(\tau_T - \tau_*) = \sqrt{n}(\tau_T - \widehat{\tau}_n) + \sqrt{n}(\widehat{\tau}_n - \tau_*) \xrightarrow{D} \mathcal{N}(0, \mathcal{I}^{-1}(\tau_*)).$$

## E.3 Proof of Theorem 8

We first bound $\|\theta_{t+1} - \widehat{\theta}_n\|_2^2$. According to the update rule in Algorithm 1, we have

$$\begin{aligned}
\|\theta_{t+1} - \widehat{\theta}_n\|_2^2 &= \|\theta_t - \eta_1 \nabla_\theta L_n(\theta_t, \gamma_t) - \widehat{\theta}_n\|_2^2 \\
&= \|\theta_t - \widehat{\theta}_n\|_2^2 + \eta_1^2 \|\nabla_\theta L_n(\theta_t, \gamma_t)\|_2^2 - 2\eta_1 \langle \nabla_\theta L_n(\theta_t, \gamma_t), \theta_t - \widehat{\theta}_n \rangle.
\end{aligned} \tag{S87}$$

Since $\nabla_\theta L_n(\widehat{\theta}_n, \widehat{\gamma}_n) = 0$. Applying Lemmas S4 and S5 yields

$$\begin{aligned}
\|\nabla_\theta L_n(\theta_t, \gamma_t)\|_2^2 &= \|\nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\theta_t, \widehat{\gamma}_n) + \nabla_\theta L_n(\theta_t, \widehat{\gamma}_n) - \nabla_\theta L_n(\widehat{\theta}_n, \widehat{\gamma}_n)\|_2^2 \\
&\le 2\|\nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\theta_t, \widehat{\gamma}_n)\|_2^2 + 2\|\nabla_\theta L_n(\theta_t, \widehat{\gamma}_n) - \nabla_\theta L_n(\widehat{\theta}_n, \widehat{\gamma}_n)\|_2^2 \\
&\le 2\widetilde{M}^2 \|\gamma_t - \widehat{\gamma}_n\|_2^2 + 2 d_2^2 K^2 \|\theta_t - \widehat{\theta}_n\|_2^2.
\end{aligned} \tag{S88}$$

Next,

$$\langle \nabla_\theta L_n(\theta_t, \gamma_t), \theta_t - \widehat{\theta}_n \rangle = \langle \nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\widehat{\theta}_n, \gamma_t), \theta_t - \widehat{\theta}_n \rangle + \langle \nabla_\theta L_n(\widehat{\theta}_n, \gamma_t) - \nabla_\theta L_n(\widehat{\theta}_n, \widehat{\gamma}_n), \theta_t - \widehat{\theta}_n \rangle. \tag{S89}$$

29

By Lemma S6, with probability at least $1 - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_1 K}}$, the following

$$\langle \nabla_\theta L_n(\theta_t, \gamma_t) - \nabla_\theta L_n(\widehat{\theta}_n, \gamma_t), \theta_t - \widehat{\theta}_n \rangle \geq w \|\theta_t - \widehat{\theta}_n\|_2^2, \tag{S90}$$

holds for all $t \in \{1, \cdots, T\}$. Denote $\bar{\gamma}_t = c\gamma_t + (1-c)\widehat{\gamma}_n$ for $c \in [0,1]$. We have

$$\begin{aligned}
\|\bar{\gamma}_t - \gamma_*\|_2 &= \|c\gamma_t + (1-c)\widehat{\gamma}_n - \gamma_*\|_2 \\
&\leq c\|\gamma_t - \gamma_*\|_2 + (1-c)\|\widehat{\gamma}_n - \gamma_*\|_2 \\
&\leq cb/\sqrt{2} + (1-c)\epsilon/\sqrt{2} \\
&\leq b/\sqrt{2},
\end{aligned}$$

and

$$\begin{aligned}
\|\nabla_{\theta\gamma}^2 L_n(\widehat{\theta}_n, \bar{\gamma}_t)\|_2 &= \|\nabla_{\theta\gamma}^2 L_n(\theta_*, \bar{\gamma}_t) + \nabla_{\theta\gamma}^2 L_n(\widehat{\theta}_n, \bar{\gamma}_t) - \nabla_{\theta\gamma}^2 L_n(\theta_*, \bar{\gamma}_t)\|_2 \\
&\leq \|\nabla_{\theta\gamma}^2 L_n(\theta_*, \bar{\gamma}_t)\|_2 + \|\nabla_{\theta\gamma}^2 L_n(\widehat{\theta}_n, \bar{\gamma}_t) - \nabla_{\theta\gamma}^2 L_n(\theta_*, \bar{\gamma}_t)\|_2 \\
&\leq \|\nabla_{\theta\gamma}^2 L_n(\theta_*, \bar{\gamma}_t)\|_2 + \sqrt{2}c_3\|\widehat{\theta}_n - \theta_*\|_2 \\
&\leq \|\nabla_{\theta\gamma}^2 L_n(\theta_*, \bar{\gamma}_t)\|_2 + c_3\epsilon
\end{aligned}$$

for some positive constant $c_3$. By Lemma S8, with probability at least $1 - \delta$, we have

$$\begin{aligned}
&|\langle \nabla_\theta L_n(\widehat{\theta}_n, \gamma_t) - \nabla_\theta L_n(\widehat{\theta}_n, \widehat{\gamma}_n), \theta_t - \widehat{\theta}_n \rangle| \\
&\leq \|\nabla_\theta L_n(\widehat{\theta}_n, \gamma_t) - \nabla_\theta L_n(\widehat{\theta}_n, \widehat{\gamma}_n)\|_2 \|\theta_t - \widehat{\theta}_n\|_2 \\
&\leq \|\nabla_{\theta\gamma}^2 L_n(\widehat{\theta}, \bar{\gamma}_t)\|_2 \|\gamma_t - \widehat{\gamma}_n\|_2 \|\theta_t - \widehat{\theta}_n\|_2 \\
&\leq (\|\nabla_{\theta\gamma}^2 L_n(\theta_*, \bar{\gamma}_t)\|_2 + c_3\epsilon) \|\gamma_t - \widehat{\gamma}_n\|_2 \|\theta_t - \widehat{\theta}_n\|_2 \\
&\leq (cb + M + c_3\epsilon + f(d_1 + d_2, n, \delta, K_1)) \|\gamma_t - \widehat{\gamma}_n\|_2 \|\theta_t - \widehat{\theta}_n\|_2 \\
&\leq \frac{cb + M + c_3\epsilon + f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \widehat{\gamma}_n\|_2^2 + \|\theta_t - \widehat{\theta}_n\|_2^2).
\end{aligned}$$

Therefore,

$$\begin{aligned}
&\langle \nabla_\theta L_n(\widehat{\theta}_n, \gamma_t) - \nabla_\theta L_n(\widehat{\theta}_n, \widehat{\gamma}_n), \theta_t - \widehat{\theta}_n \rangle \\
&\geq -\frac{cb + M + c_3\epsilon + f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \widehat{\gamma}_n\|_2^2 + \|\theta_t - \widehat{\theta}_n\|_2^2).
\end{aligned} \tag{S91}$$

By (S89), (S90) and (S91), we obtain

$$\langle \nabla_\theta L_n(\theta_t, \gamma_t), \theta_t - \widehat{\theta}_n \rangle \geq \left( w - \frac{cb + M + c_3\epsilon + f(d_1 + d_2, n, \delta, K_1)}{2} \right) \|\theta_t - \widehat{\theta}_n\|_2^2$$
$$- \frac{cb + M + c_3\epsilon + f(d_1 + d_2, n, \delta, K_1)}{2} \|\gamma_t - \widehat{\gamma}_n\|_2^2. \tag{S92}$$

Combining (S87), (S88) and (S92), with probability at least $1 - d_1 \left( \frac{e}{2} \right)^{-\frac{\lambda_\phi n}{8 d_2 K}}$, we have

$$\|\theta_{t+1} - \widehat{\theta}_n\|_2^2 \leq [1 + 2d_2^2 K^2 \eta_1^2 - (2w - cb - M - c_3\epsilon - f(d_1 + d_2, n, \delta, K_1))\eta_1] \|\theta_t - \widehat{\theta}_n\|_2^2$$
$$+ \{2\widetilde{M}^2 \eta_1^2 + [cb + M + c_3\epsilon + f(d_1 + d_2, n, \delta, K_1)]\eta_1\} \|\gamma_t - \widehat{\gamma}_n\|_2^2. \tag{S93}$$

Now we bound $\|\gamma_{t+1} - \widehat{\gamma}_n\|_2^2$. According to the update in Algorithm 1, we have

$$\|\gamma_{t+1} - \widehat{\gamma}_n\|_2^2 = \|\gamma_t - \eta_2 \nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \widehat{\gamma}_n\|_2^2$$
$$= \|\gamma_t - \widehat{\gamma}_n\|_2^2 + \eta_2^2 \|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t)\|_2^2 - 2\eta_2 \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t), \gamma_t - \widehat{\gamma}_n \rangle. \tag{S94}$$

Note that $\nabla_\gamma L_n(\widehat{\theta}_n, \widehat{\gamma}_n) = 0$. By Lemmas S4 and S5, we obtain

$$\|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t)\|_2^2 = \|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n) + \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n) - \nabla_\gamma L_n(\widehat{\theta}_n, \widehat{\gamma}_n)\|_2^2$$
$$\leq 2\|\nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n)\|_2^2 + 2\|\nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n) - \nabla_\gamma L_n(\widehat{\theta}_n, \widehat{\gamma}_n)\|_2^2$$
$$\leq 2d_2^2 K^2 \|\gamma_t - \widehat{\gamma}_n\|_2^2 + 2\widetilde{M}^2 \|\theta_{t+1} - \widehat{\theta}_n\|_2^2. \tag{S95}$$

Next,

$$\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t), \gamma_t - \widehat{\gamma}_n \rangle$$
$$= \langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n), \gamma_t - \widehat{\gamma}_n \rangle + \langle \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n) - \nabla_\gamma L_n(\widehat{\theta}_n, \widehat{\gamma}_n), \gamma_t - \widehat{\gamma}_n \rangle. \tag{S96}$$

By Lemma S6, with probability at least $1 - d_2 \left( \frac{e}{2} \right)^{-\frac{\lambda_\psi n}{8 d_1 K}}$, the following

$$\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t) - \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n), \gamma_t - \widehat{\gamma}_n \rangle \geq w \|\gamma_t - \widehat{\gamma}_n\|_2^2, \tag{S97}$$

holds for all $t \in \{1, \cdots, T\}$. Denote $\bar{\theta}_t = a\theta_t + (1 - a)\widehat{\theta}_n$ for $a \in [0, 1]$. We have

$$\|\bar{\theta}_t - \theta_*\|_2 = \|a\theta_t + (1 - a)\widehat{\theta}_n - \theta_*\|_2$$
$$\leq a\|\theta_t - \theta_*\|_2 + (1 - a)\|\widehat{\theta}_n - \theta_*\|_2$$
$$\leq ab/\sqrt{2} + (1 - a)\epsilon/\sqrt{2}$$
$$\leq b/\sqrt{2},$$

31

and

$$\|\nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \widehat{\gamma}_n)\|_2 = \|\nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \gamma_*) + \nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \widehat{\gamma}_n) - \nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \gamma_*)\|_2$$

$$\leq \|\nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \gamma_*)\|_2 + \|\nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \widehat{\gamma}_n) - \nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \gamma_*)\|_2$$

$$\leq \|\nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \gamma_*)\|_2 + \sqrt{2} c_3' \|\widehat{\gamma}_n - \gamma_*\|_2$$

$$\leq \|\nabla^2_{\theta\gamma} L_n(\bar{\theta}_t, \gamma_*)\|_2 + c_3' \epsilon$$

for some positive constant $c_3'$. By Lemma S8, with probability at least $1 - \delta$, we have

$$|\langle \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n) - \nabla_\gamma L_n(\widehat{\theta}_n, \widehat{\gamma}_n), \gamma_t - \widehat{\gamma}_n \rangle|$$

$$\leq \|\nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n) - \nabla_\gamma L_n(\widehat{\theta}_n, \widehat{\gamma}_n)\|_2 \|\gamma_t - \widehat{\gamma}_n\|_2$$

$$\leq \nabla^2_{\theta\gamma} L_n(\bar{\theta}_{t+1}, \widehat{\gamma}_n) \|\theta_{t+1} - \widehat{\theta}_n\|_2 \|\gamma_t - \widehat{\gamma}_n\|_2$$

$$\leq (\|\nabla^2_{\theta\gamma} L_n(\bar{\theta}_{t+1}, \gamma_*)\|_2 + c_3' \epsilon) \|\theta_{t+1} - \widehat{\theta}_n\|_2 \|\gamma_t - \widehat{\gamma}_n\|_2$$

$$\leq \frac{cb + c_3' \epsilon + M + f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \widehat{\gamma}_n\|_2^2 + \|\theta_{t+1} - \widehat{\theta}_n\|_2^2).$$

Therefore,

$$\langle \nabla_\gamma L_n(\theta_{t+1}, \widehat{\gamma}_n) - \nabla_\gamma L_n(\widehat{\theta}_n, \widehat{\gamma}_n), \gamma_t - \widehat{\gamma}_n \rangle$$
$$\geq -\frac{cb + c_3' \epsilon + M + f(d_1 + d_2, n, \delta, K_1)}{2} (\|\gamma_t - \widehat{\gamma}_n\|_2^2 + \|\theta_{t+1} - \widehat{\theta}_n\|_2^2). \tag{S98}$$

By (S96), (S97) and (S98), we obtain

$$\langle \nabla_\gamma L_n(\theta_{t+1}, \gamma_t), \gamma_t - \widehat{\gamma}_n \rangle \geq \left( w - \frac{cb + c_3' \epsilon + M + f(d_1 + d_2, n, \delta, K_1)}{2} \right) \|\gamma_t - \widehat{\gamma}_n\|_2^2$$
$$- \frac{cb + c_3' \epsilon + M + f(d_1 + d_2, n, \delta, K_1)}{2} \|\theta_{t+1} - \widehat{\theta}_n\|_2^2. \tag{S99}$$

Combining (S94), (S95) and (S99), with probability at least $1 - d_2 \left( \frac{e}{2} \right)^{-\frac{\lambda_\psi n}{8 d_1 K}}$, we have

$$\|\gamma_{t+1} - \widehat{\gamma}_n\|_2^2 \leq [1 + 2 d_1^2 K^2 \eta_2^2 - (2w - cb - c_3' \epsilon - M - f_3(n, \delta)) \eta_2] \|\gamma_t - \widehat{\gamma}_n\|_2^2$$
$$+ \{2 \widetilde{M}^2 \eta_2^2 + [cb + c_3' \epsilon + M + f(d_1 + d_2, n, \delta, K_1)] \eta_2\} \|\theta_{t+1} - \widehat{\theta}_n\|_2^2. \tag{S100}$$

Let $A_{1n}$ and $A_{2n}$ be defined in (S75), and $\eta$ be defined in (S74). When

$$\epsilon \leq \frac{\min\{d_2^2 K^2 \eta_1^2 + w \eta_1, d_1^2 K^2 \eta_2^2 + w \eta_2, \widetilde{M}^2 \eta\}}{\max\{c_3, c_3'\}},$$

32

we can verify $1 + \max\{2d_2^2 K^2 \eta_1^2 - (2w - cb - c_3\epsilon - M - f(d_1 + d_2, n, \delta, K_1))\eta_2, 2d_1^2 K^2 \eta_2^2 - (2w - cb - c_3'\epsilon - M - f(d_1 + d_2, n, \delta, K_1))\eta_2\} \le A_{1n}$ and $2\widetilde{M}^2 \eta^2 + [cb + \max\{c_3, c_3'\}\epsilon + M + f(d_1 + d_2, n, \delta, K_1)]\eta \le A_{2n}$. By (S93) and (S100), with probability at least $1 - \delta - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}} - d_2 \left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$, we can obtain

$$\|\theta_{t+1} - \widehat{\theta}_n\|_2^2 \le A_{1n}\|\theta_t - \widehat{\theta}_n\|_2^2 + A_{2n}\|\gamma_t - \widehat{\gamma}_n\|_2^2 \tag{S101}$$

and

$$\|\gamma_{t+1} - \widehat{\gamma}_n\|_2^2 \le A_{1n}\|\gamma_t - \widehat{\gamma}_n\|_2^2 + A_{2n}\|\theta_{t+1} - \widehat{\theta}_n\|_2^2$$
$$\le (A_{1n} + A_{2n}^2)\|\gamma_t - \widehat{\gamma}_n\|_2^2 + A_{1n}A_{2n}\|\theta_t - \widehat{\theta}_n\|_2^2. \tag{S102}$$

Denote $n_0' = \max\{n_0, n_\epsilon\}$. By (S81), (S101) and (S102), we have

$$\|\theta_{t+1} - \widehat{\theta}_n\|_2^2 + \|\gamma_{t+1} - \widehat{\gamma}_n\|_2^2 \le (A_{1n} + A_{2n} + A_{2n}^2)\|\gamma_t - \widehat{\gamma}_n\|_2^2 + (A_{1n} + A_{1n}A_{2n})\|\theta_t - \widehat{\theta}_n\|_2^2$$
$$\le (A_{1n} + A_{2n} + A_{2n}^2)(\|\gamma_t - \widehat{\gamma}_n\|_2^2 + \|\theta_t - \widehat{\theta}_n\|_2^2)$$
$$\le \rho_1(\|\gamma_t - \widehat{\gamma}_n\|_2^2 + \|\theta_t - \widehat{\theta}_n\|_2^2),$$

where $\rho_1 = A_{1n_0'} + A_{2n_0'} + A_{2n_0'}^2 < 1$ according to (S82). For the initialization, we have $\|\theta_0 - \widehat{\theta}_n\|_2 \le \|\theta_t - \theta_*\|_2 + \|\theta_* - \widehat{\theta}_n\|_2 \le (b+\epsilon)/\sqrt{2}$ and $\|\gamma_t - \widehat{\gamma}_n\|_2 \le \|\gamma_t - \gamma_*\|_2 + \|\gamma_* - \widehat{\gamma}_n\|_2 \le (b+\epsilon)/\sqrt{2}$. Thus, $\|\theta_0 - \widehat{\theta}_n\|_2^2 + \|\gamma_* - \widehat{\gamma}_n\|_2^2 \le (b+\epsilon)^2$. By recursion, with probability at least $1 - \delta - d_1 \left(\frac{e}{2}\right)^{-\frac{\lambda_\phi n}{8d_2 K}} - d_2 \left(\frac{e}{2}\right)^{-\frac{\lambda_\psi n}{8d_1 K}}$, when $n > n_0'$, we obtain

$$\|\theta_T - \widehat{\theta}_n\|_2^2 + \|\gamma_T - \widehat{\gamma}_n\|_2^2 \le \rho_1^T (b+\epsilon)^2.$$

## E.4    Proof of Theorem 3

We denote $v = (\phi^\top(s, a), 0)^\top \in \mathbb{R}^{d_1 + d_2}$. Then, by Assumption 1, we have

$$v^\top(\tau_T - \tau_*) = (\theta_T - \theta_*)^\top \phi(s, a) = r_{\theta_T}(s, a) - r_{\theta_*}(s, a).$$

By Lemma S10, we have

$$v^\top \widetilde{\mathcal{I}}^{-1}(\tau_*)v = \phi^\top(s, a)[\mathcal{I}_{\theta\theta}(\theta_*, \gamma_*) - \mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*)\mathcal{I}_{\gamma\gamma}^{-1}(\theta_*, \gamma_*)\mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*)]^{-1}\phi(s, a)$$
$$= \phi^\top(s, a)\widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\phi(s, a).$$

Therefore, by Theorem 2, as $n \to \infty$ and $T \to \infty$, we have

$$\sqrt{n}[r_{\theta_T}(s, a) - r_{\theta_*}(s, a)] \xrightarrow{D} \mathcal{N}(0, \phi^\top(s, a)\widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\phi(s, a)).$$

33

## E.5  Proof of Theorem 4

Because the samples $\{(x_i, s_i, a_i^{(0)}, a_i^{(1)})\}_{i=1}^n$ are i.i.d., the law of large numbers implies that, as $n \to \infty$, entry-wisely, we have

$$\widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*) \xrightarrow{p} \mathcal{I}_{\theta\theta}(\theta_*, \gamma_*), \ \widehat{\mathcal{I}}_{\theta\gamma}(\theta_*, \gamma_*) \xrightarrow{p} \mathcal{I}_{\theta\gamma}(\theta_*, \gamma_*), \ \widehat{\mathcal{I}}_{\gamma\gamma}(\theta_*, \gamma_*) \xrightarrow{p} \mathcal{I}_{\gamma\gamma}(\theta_*, \gamma_*).$$

Consequently, for sufficiently large $n$, we have

$$\|I - \mathcal{I}_{\theta\theta}^{-1}(\theta_*, \gamma_*)\widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*)\|_2 = \|\mathcal{I}_{\theta\theta}^{-1}(\theta_*, \gamma_*)[\mathcal{I}_{\theta\theta}(\theta_*, \gamma_*) - \widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*)]\|_2$$

$$\leq \|\mathcal{I}_{\theta\theta}^{-1}(\theta_*, \gamma_*)\|_2 \|\mathcal{I}_{\theta\theta}(\theta_*, \gamma_*) - \widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*)\|_2 < 1$$

with probability approaching 1. By Lemma S14, this implies that $\mathcal{I}_{\theta\theta}^{-1}(\theta_*, \gamma_*)\widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*)$ is invertible, and therefore $\widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*)$ is invertible. An analogous argument applies to $\widehat{\mathcal{I}}_{\gamma\gamma}(\theta_*, \gamma_*)$. By the continuous mapping theorem, the inverses of these matrices converge in probability to the inverses of the population matrices. Since $S_\theta^2(\theta_*, \gamma_*) = [\widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*) - \widehat{\mathcal{I}}_{\theta\gamma}(\theta_*, \gamma_*)\widehat{\mathcal{I}}_{\gamma\gamma}^{-1}(\theta_*, \gamma_*)\widehat{\mathcal{I}}_{\theta\gamma}^{\top}(\theta_*, \gamma_*)]^{-1}$, as $n \to \infty$, we have

$$\|S_\theta^2(\theta_*, \gamma_*) - \widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\|_2 \xrightarrow{p} 0. \tag{S103}$$

Note that $\mu(\cdot)[1 - \mu(\cdot)] \leq 1/4$. Then, by Assumptions 1 and 2, there exist some constants $C_1$ and $C_2$, such that

$$\|\widehat{\mathcal{I}}_{\theta\theta}(\theta_T, \gamma_T) - \widehat{\mathcal{I}}_{\theta\theta}(\theta_*, \gamma_*)\|_2 \leq \frac{1}{4n} \left\| \sum_{i=1}^n \{[\gamma_T^\top \psi(x_i)]^2 - [\gamma_*^\top \psi(x_i)]^2\} z_i z_i^\top \right\|_2$$

$$= \frac{1}{4n} \left\| \sum_{i=1}^n \{[\gamma_T^\top \psi(x_i) - \gamma_*^\top \psi(x_i)][\gamma_T^\top \psi(x_i) + \gamma_*^\top \psi(x_i)]\} z_i z_i^\top \right\|_2$$

$$\leq \frac{1}{4n} \sum_{i=1}^n \|\gamma_T - \gamma_*\|_2 \|\psi(x_i)\|_2 |\gamma_T^\top \psi(x_i) + \gamma_*^\top \psi(x_i)| \|_2 \|z_i z_i^\top\|_2$$

$$\leq C_1 \|\gamma_T - \gamma_*\|_2$$

and

$$\|\widehat{\mathcal{I}}_{\gamma\gamma}(\theta_T, \gamma_T) - \widehat{\mathcal{I}}_{\gamma\gamma}(\theta_*, \gamma_*)\|_2 \le \frac{1}{4n} \left\| \sum_{i=1}^n [(\theta_T^\top z)^2 - (\theta^\top z)^2] \psi(x_i) \psi(x_i)^\top \right\|_2$$

$$= \frac{1}{4n} \left\| \sum_{i=1}^n [(\theta_T^\top z - \theta^\top z)(\theta_T^\top z + \theta^\top z)] \psi(x_i) \psi(x_i)^\top \right\|_2$$

$$\le \frac{1}{4n} \sum_{i=1}^n \|\theta_T - \theta\|_2 \|z\|_2 |\theta_T^\top z + \theta^\top z| \|\psi(x_i) \psi(x_i)^\top\|_2$$

$$\le C_2 \|\theta_T - \theta_*\|_2.$$

By Assumptions 1 and 2, there exists some constant $C_3$ such that

$$\|\widehat{\mathcal{I}}_{\theta\gamma}(\theta_T, \gamma_T) - \widehat{\mathcal{I}}_{\theta\gamma}(\theta_*, \gamma_*)\|_2$$

$$= \|\widehat{\mathcal{I}}_{\theta\gamma}(\theta_T, \gamma_T) - \widehat{\mathcal{I}}_{\theta\gamma}(\theta_*, \gamma_T) + \widehat{\mathcal{I}}_{\theta\gamma}(\theta_*, \gamma_T) - \widehat{\mathcal{I}}_{\theta\gamma}(\theta_*, \gamma_*)\|_2$$

$$= \frac{1}{4n} \left\| \sum_{i=1}^n \{[\sigma_{\gamma_T}(x_i)(\theta_T - \theta_*)^\top z_i] + (\gamma_T - \gamma_*)^\top \psi(x_i)(\theta_*^\top z_i)\} \psi(x_i) z_i^\top \right\|_2$$

$$\le \frac{1}{4n} \sum_{i=1}^n [|\sigma_{\gamma_T}(x_i)| \|\theta_T - \theta_*\|_2 \|z_i\|_2 + \|\gamma_T - \gamma_*\|_2 \|\psi(x_i)\|_2 |\theta_*^\top z_i|] \|\psi(x_i) z_i^\top\|_2$$

$$\le C_3(\|\theta_T - \theta_*\|_2 + \|\gamma_T - \gamma_*\|_2).$$

Since

$$S_\theta^2(\theta_T, \gamma_T) = [\widehat{\mathcal{I}}_{\theta\theta}(\theta_T, \gamma_T) - \widehat{\mathcal{I}}_{\theta\gamma}(\theta_T, \gamma_T) \widehat{\mathcal{I}}_{\gamma\gamma}^{-1}(\theta_T, \gamma_T) \widehat{\mathcal{I}}_{\theta\gamma}^\top(\theta_T, \gamma_T)]^{-1},$$

by Theorem 1, we have $\|S_\theta^2(\theta_T, \gamma_T) - S_\theta^2(\theta_*, \gamma_*)\|_2 \xrightarrow{p} 0$ as $n \to \infty$ and $T \to \infty$. Combining (S103), we have

$$\|S_\theta^2(\theta_T, \gamma_T) - \widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\|_2$$

$$= \|S_\theta^2(\theta_T, \gamma_T) - S_\theta^2(\theta_*, \gamma_*) + S_\theta^2(\theta_*, \gamma_*) - \widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\|_2$$

$$\le \|S_\theta^2(\theta_T, \gamma_T) - S_\theta^2(\theta_*, \gamma_*)\|_2 + \|S_\theta^2(\theta_*, \gamma_*) - \widetilde{\mathcal{I}}^{-1}(\theta_*, \gamma_*)\|_2$$

$$\xrightarrow{p} 0.$$

By Theorem 3 and Slutsky's theorem, we have

$$\frac{\sqrt{n}[r_{\theta_T}(s, a) - r_{\theta_*}(s, a)]}{\sqrt{\phi^\top(s, a) S_\theta^2(\theta_T, \gamma_T) \phi(s, a)}} \xrightarrow{D} \mathcal{N}(0, 1).$$

## E.6  Proof of Theorem 5

The proof of Theorem 5 follows the same strategy as used in the proofs for Theorems 3 and 4, and hence is omitted.

## E.7  Proof of Corollary 1

When $T > \frac{\log[(1-\rho)b^2 n] - \log\{c_2 \log(1/\delta)\}}{-\log \rho}$ (equivalent to $T > c \log n$ for some constant $c$), we can verify $\rho^T b^2 < \frac{c_2 \log(1/\delta)}{(1-\rho)n}$. By Theorem 1, with probability at least $1 - \delta - \frac{1}{n}$, we have

$$\|\theta_T - \theta_*\|_2^2 \leq \frac{2c_2 \log(1/\delta)}{(1-\rho)n}.$$

Therefore,

$$\|\theta_T - \theta_*\|_{S^{-2}(\theta_T,\gamma_T)} \leq \|\theta_T - \theta_*\|_2 \|S^{-2}(\theta_T,\gamma_T)\|_2 \leq \frac{c_3}{\sqrt{n}},$$

where $c_3 = \sqrt{\frac{2c_2 \log(1/\delta)}{1-\rho}} \|S^{-2}(\theta_T,\gamma_T)\|_2$. Since $q_{1-\alpha/2}$ can take any positive value, we can choose $\alpha$ such that $q_{1-\alpha/2} = c_3$. The proof is complete.

## E.8  Proof of Lemma 3

By Assumption 1, we have $r_\theta(s,a) = \theta^\top \phi(s,a)$. Then, for any $\theta$, we can write $r_\theta(s,a) = r_{\theta_T}(s,a) + \phi(s,a)^\top (\theta - \theta_T)$. Then minimizing $r_\theta(s,a)$ over $\widetilde{\Theta}$ is equivalent to

$$\min_{\theta \in \widetilde{\Theta}} r_\theta(s,a) = r_{\theta_T}(s,a) + \min_{\|\theta - \theta_T\|_{S_\theta^{-2}(\theta_T,\gamma_T)} \leq \frac{q_{1-\alpha/2}}{\sqrt{n}}} (\theta - \theta_T)^\top \phi(s,a).$$

By the Cauchy-Schwarz inequality, we have

$$\begin{aligned}
(\theta - \theta_T)^\top \phi(s,a) &= (\theta - \theta_T)^\top S_\theta^{-1}(\theta_T,\gamma_T) S_\theta(\theta_T,\gamma_T) \phi(s,a) \\
&\geq -\|S_\theta^{-1}(\theta_T,\gamma_T)(\theta - \theta_T)\|_2 \|S_\theta(\theta_T,\gamma_T)\phi(s,a)\|_2 \\
&= -\sqrt{(\theta - \theta_T)^\top S_\theta^{-2}(\theta_T,\gamma_T)(\theta - \theta_T)} \sqrt{\phi^\top(s,a) S_\theta^2(\theta_T,\gamma_T) \phi(s,a)} \quad \text{(S104)} \\
&\geq -q_{1-\alpha/2} \sqrt{\frac{\phi^\top(s,a) S_\theta^2(\theta_T,\gamma_T) \phi(s,a)}{n}},
\end{aligned}$$

36

where the last inequality is from the definition of $\widetilde{\Theta}$. Pick $\theta = \theta_T - \frac{q_{1-\alpha/2}}{\sqrt{n}} \frac{S_\theta^2(\theta_T, \gamma_T)\phi(s,a)}{\sqrt{\phi(s,a)^\top S_\theta^2(\theta_T, \gamma_T)\phi(s,a)}}$,

which satisfies $\|\theta - \theta_T\|_{S_\theta^{-2}(\theta_T, \gamma_T)} \leq q_{1-\alpha/2}/\sqrt{n}$. We have

$$
\begin{aligned}
(\theta - \theta_T)^\top \phi(s,a) &= -\frac{q_{1-\alpha/2}}{\sqrt{n}} \frac{\phi^\top(s,a) S_\theta^2(\theta_T, \gamma_T)\phi(s,a)}{\sqrt{\phi(s,a)^\top S_\theta^2(\theta_T, \gamma_T)\phi(s,a)}} \\
&= -q_{1-\alpha/2}\sqrt{\frac{\phi^\top(s,a) S_\theta^2(\theta_T, \gamma_T)\,\phi(s,a)}{n}}.
\end{aligned}
\tag{S105}
$$

By (S104) and (S105), we have

$$
\min_{\|\theta - \theta_T\|_{S_\theta^{-2}(\theta_T, \gamma_T)} \leq \frac{q_{1-\alpha/2}}{\sqrt{n}}} (\theta - \theta_T)^\top \phi(s,a) = -q_{1-\alpha/2}\sqrt{\frac{\phi^\top(s,a) S_\theta^2(\theta_T, \gamma_T)\,\phi(s,a)}{n}}.
$$

Therefore,

$$
\min_{\theta \in \widetilde{\Theta}} r_\theta(s,a) = r_{\theta_T}(s,a) - q_{1-\alpha/2}\sqrt{\frac{\phi^\top(s,a) S_\theta^2(\theta_T, \gamma_T)\,\phi(s,a)}{n}}.
$$

This completes the proof.

## E.9  Proof of Theorem 6

For notational simplicity, write $a_* = a_*(s)$ and $a_{pBoN} = a_{pBoN}(s)$. The reward difference between $a_*$ and $a_{pBoN}$ for $s$ is

$$
\begin{aligned}
&r_{\theta_*}(s, a_*) - r_{\theta_*}(s, a_{pBoN}) \\
&= [r_{\theta_*}(s, a_*) - \widehat{r}(s, a_*)] + [\widehat{r}(s, a_*) - \widehat{r}(s, a_{pBoN})] + [\widehat{r}(s, a_{pBoN}) - r_{\theta_*}(s, a_{pBoN})].
\end{aligned}
\tag{S106}
$$

Since $a_{pBoN}$ is the optimal action under $\widehat{r}(s,a)$, i.e., $a_{pBoN}(s) = \arg\max_{a \in \mathcal{A}_N(s)} \widehat{r}(s,a)$, the second difference satisfies

$$
\widehat{r}(s, a_*) - \widehat{r}(s, a_{pBoN}) \leq 0.
\tag{S107}
$$

For the third difference, by (10), we have

$$
\widehat{r}(s, a_{pBoN}) - r_{\theta_*}(s, a_{pBoN}) = \min_{\theta \in \widetilde{\Theta}} r_\theta(s, a_{pBoN}) - r_{\theta_*}(s, a_{pBoN}),
$$

where $\widetilde{\Theta} = \{\theta : \|\theta - \theta_T\|_{S_\theta^{-2}(\theta_T, \gamma_T)} \leq q_{1-\alpha/2}/\sqrt{n}\}$. From Corollary 1, we have $\|\theta_* - \theta_T\|_{S_\theta^{-2}(\theta_T, \gamma_T)} \leq q_{1-\alpha/2}/\sqrt{n}$ with probability at least $1 - \delta - \frac{1}{n}$, i.e., $\mathbb{P}(\theta_* \in \widetilde{\Theta}) \geq 1 - \delta - \frac{1}{n}$. Thus, with probability at least $1 - \delta - \frac{1}{n}$, we have $\min_{\theta \in \widetilde{\Theta}} r_\theta(s, a_{pBoN}) \leq r_{\theta_*}(s, a_{pBoN})$, and then

$$
\widehat{r}(s, a_{pBoN}) - r_{\theta_*}(s, a_{pBoN}) \leq 0.
\tag{S108}
$$

37

Now, we bound the first difference. By Assumption 1, we obtain

$$
\begin{aligned}
r_{\theta_*}(s, a_*) - \widehat{r}(s, a_*) &= \theta_*^\top \phi(s, a_*) - \min_{\theta \in \widetilde{\Theta}} \theta^\top \phi(s, a_*) \\
&= \max_{\theta \in \widetilde{\Theta}} (\theta_* - \theta)^\top \phi(s, a_*) \\
&= \max_{\theta \in \widetilde{\Theta}} (\theta_* - \theta_T + \theta_T - \theta)^\top \phi(s, a_*) \\
&= (\theta_* - \theta_T)^\top \phi(s, a_*) + \max_{\theta \in \widetilde{\Theta}} (\theta_T - \theta)^\top \phi(s, a_*) \\
&\le |(\theta_T - \theta_*)^\top \phi(s, a_*)| + \max_{\theta \in \widetilde{\Theta}} |(\theta_T - \theta)^\top \phi(s, a_*)|.
\end{aligned}
$$

Since $\theta_* \in \widetilde{\Theta}$ with probability at least $1 - \delta - \frac{1}{n}$ by Corollary 1, we have $|(\theta_T - \theta_*)^\top \phi(s, a_*)| \le \max_{\theta \in \widetilde{\Theta}} |(\theta_T - \theta)^\top \phi(s, a_*)|$ with probability at least $1 - \delta - \frac{1}{n}$. Therefore, with probability at least $1 - \delta - \frac{1}{n}$, we have

$$
\begin{aligned}
r_{\theta_*}(s, a_*) - \widehat{r}(s, a_*) &\le 2 \max_{\theta \in \widetilde{\Theta}} |(\theta_T - \theta)^\top \phi(s, a_*)| \\
&\le 2 \max_{\theta \in \widetilde{\Theta}} \|\theta_T - \theta\|_{S_\theta^{-2}(\theta_T, \gamma_T)} \|S_\theta(\theta_T, \gamma_T)\phi(s, a_*)\|_2 \\
&\le 2 q_{1-\alpha/2} \|S_\theta(\theta_T, \gamma_T)\phi(s, a_*)\|_2 / \sqrt{n} \\
&\le 3 q_{1-\alpha/2} \|\widetilde{\mathcal{I}}^{-\frac{1}{2}}(\theta_*, \gamma_*)\phi(s, a_*)\|_2 / \sqrt{n},
\end{aligned}
\tag{S109}
$$

where the third inequality follows from the definition of $\widetilde{\Theta}$, and the last inequality holds when $T$ and $n$ are large enough. Therefore, by (S106), (S107), (S108), (S109), with probability at least $1 - \delta - \frac{1}{n}$, we have

$$
\mathsf{SubOpt}(a_{pBoN}) = \mathbb{E}_s[r_{\theta_*}(s, a_*) - r_{\theta_*}(s, a_{pBoN})] \le 2 q_{1-\alpha/2} \|\widetilde{\mathcal{I}}^{-\frac{1}{2}}(\theta_*, \gamma_*)\mathbb{E}_s \phi(s, a_*(s))\|_2 / \sqrt{n}.
$$

# F  Supporting Lemmas

**Lemma S9.** *(Corollary 5.2 (Tropp, 2012)) Consider a finite sequence $\{\mathbf{X}_k\}$ of independent, random, self-adjoint matrices with dimension $d$ that satisfy*

$$
\mathbf{X}_k \succeq \mathbf{0} \text{ and } \lambda_{\max}(\mathbf{X}_k) \le L \text{ almost surely.}
$$

*Compute the minimum eigenvalue of the sum of expectations, $\mu_{\min} := \lambda_{\min}\left(\sum_k \mathbb{E}\mathbf{X}_k\right)$. Then for $\zeta \in [0, 1]$,*

$$
\mathbb{P}\left\{\lambda_{\min}\left(\sum_k \mathbf{X}_k\right) \le (1-\zeta)\mu_{\min}\right\} \le d\left[\frac{e^{-\zeta}}{(1-\zeta)^{1-\zeta}}\right]^{\mu_{\min}/L}.
$$

**Lemma S10.** *(Theorem 8.5.11, (Harville, 1997)) Let $\boldsymbol{T}$ represent an $m \times m$ matrix, $\boldsymbol{U}$ an $m \times n$ matrix, $\boldsymbol{V}$ an $n \times m$ matrix, and $\boldsymbol{W}$ an $n \times n$ matrix. Suppose that $\boldsymbol{T}$ is nonsingular. Then, $\begin{pmatrix} \boldsymbol{T} & \boldsymbol{U} \\ \boldsymbol{V} & \boldsymbol{W} \end{pmatrix}$ or equivalently $\begin{pmatrix} \boldsymbol{W} & \boldsymbol{V} \\ \boldsymbol{U} & \boldsymbol{T} \end{pmatrix}$, is nonsingular if and only if the $n \times n$ matrix $\boldsymbol{Q} = \boldsymbol{W} - \boldsymbol{V}\boldsymbol{T}^{-1}\boldsymbol{U}$ is nonsingular, in which case*

$$\begin{pmatrix} \boldsymbol{T} & \boldsymbol{U} \\ \boldsymbol{V} & \boldsymbol{W} \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{T}^{-1} + \boldsymbol{T}^{-1}\boldsymbol{U}\boldsymbol{Q}^{-1}\boldsymbol{V}\boldsymbol{T}^{-1} & -\boldsymbol{T}^{-1}\boldsymbol{U}\boldsymbol{Q}^{-1} \\ -\boldsymbol{Q}^{-1}\boldsymbol{V}\boldsymbol{T}^{-1} & \boldsymbol{Q}^{-1} \end{pmatrix},$$

$$\begin{pmatrix} \boldsymbol{W} & \boldsymbol{V} \\ \boldsymbol{U} & \boldsymbol{T} \end{pmatrix}^{-1} = \begin{pmatrix} \boldsymbol{Q}^{-1} & -\boldsymbol{Q}^{-1}\boldsymbol{V}\boldsymbol{T}^{-1} \\ -\boldsymbol{T}^{-1}\boldsymbol{U}\boldsymbol{Q}^{-1} & \boldsymbol{T}^{-1} + \boldsymbol{T}^{-1}\boldsymbol{U}\boldsymbol{Q}^{-1}\boldsymbol{V}\boldsymbol{T}^{-1} \end{pmatrix}.$$

**Lemma S11.** *(Theorem 14.8.5, (Harville, 1997)) Let $\boldsymbol{A} = \begin{pmatrix} \boldsymbol{T} & \boldsymbol{U} \\ \boldsymbol{U}^\top & \boldsymbol{W} \end{pmatrix}$, where $\boldsymbol{T}$ is of dimensions $m \times m$, $\boldsymbol{U}$ of dimensions $m \times n$, and $\boldsymbol{W}$ of dimensions $n \times n$. If $\boldsymbol{T}$ is symmetric positive definite and the Schur complement $\boldsymbol{W} - \boldsymbol{U}^\top\boldsymbol{T}^{-1}\boldsymbol{U}$ of $\boldsymbol{T}$ is positive definite, or if $\boldsymbol{W}$ is symmetric positive definite and the Schur complement $\boldsymbol{T} - \boldsymbol{U}\boldsymbol{W}^{-1}\boldsymbol{U}^\top$ of $\boldsymbol{W}$ is positive definite, then $\boldsymbol{A}$ is positive definite.*

**Lemma S12.** *(Theorem 17, (Ferguson, 1996)) Let $X_1, X_2, \ldots$ be i.i.d. with density $f(x \mid \theta)$, $\theta \in \Theta$, and let $\theta_0$ denote the true value of $\theta$. If*

*(1) $\Theta$ is compact,*

*(2) $f(x \mid \theta)$ is upper semicontinuous in $\theta$ for all $x$,*

*(3) there exists a function $K(x)$ such that $E_{\theta_0}[K(X)] < \infty$ and*

$$U(x, \theta) = \log f(x \mid \theta) - \log f(x \mid \theta_0) \leq K(x), \quad \text{for all } x \text{ and } \theta,$$

*(4) for all $\theta \in \Theta$ and sufficiently small $\rho > 0$, $\sup_{|\theta' - \theta| < \rho} f(x \mid \theta')$ is measurable in $x$,*

*(5) (identifiability) $f(x \mid \theta) = f(x \mid \theta_0)$ (a.e. dv) $\Rightarrow \theta = \theta_0$,*

*then, for any sequence of maximum-likelihood estimates $\widehat{\theta}_n$ of $\theta$,*

$$\widehat{\theta}_n \xrightarrow{a.s.} \theta_0.$$

**Lemma S13.** *(Corollary 6.1.2, (Tropp et al., 2015)) Consider a finite sequence $\{S_k\}$ of independent random matrices with common dimension $d_1 \times d_2$. Assume that each matrix has*

*uniformly bounded deviation from its mean:*

$$\|S_k - \mathbb{E}S_k\| \leq L \quad \text{for each index } k.$$

*Introduce the sum*

$$Z = \sum_k S_k,$$

*and let $v(Z)$ denote the matrix variance statistic of the sum:*

$$v(Z) = \max \left\{ \|\mathbb{E}\left[(Z - \mathbb{E}Z)(Z - \mathbb{E}Z)^*\right]\|, \|\mathbb{E}\left[(Z - \mathbb{E}Z)^*(Z - \mathbb{E}Z)\right]\| \right\}$$

$$= \max \left\{ \left\| \sum_k \mathbb{E}\left[(S_k - \mathbb{E}S_k)(S_k - \mathbb{E}S_k)^*\right] \right\|, \left\| \sum_k \mathbb{E}\left[(S_k - \mathbb{E}S_k)^*(S_k - \mathbb{E}S_k)\right] \right\| \right\}.$$

*Then*

$$\mathbb{E}\|Z - \mathbb{E}Z\| \leq \sqrt{2v(Z)\log(d_1 + d_2)} + \frac{1}{3}L\log(d_1 + d_2).$$

*Furthermore, for all $t \geq 0$,*

$$\mathbb{P}\left\{\|Z - \mathbb{E}Z\| \geq t\right\} \leq (d_1 + d_2) \cdot \exp\left(\frac{-t^2/2}{v(Z) + Lt/3}\right).$$

**Lemma S14.** *(Corollary 5.6.16, (Horn and Johnson, 2012)) A matrix $A$ is nonsingular if there is a matrix norm $\|\cdot\|$ such that $\|I - A\| < 1$.*