

Bayesian Regression Games: A Framework for Adversarial Robustness via Variational Inequality

Chris Junchi Li[◇]

Department of Electrical Engineering and Computer Sciences[◇]
University of California, Berkeley

September 2, 2024

Abstract

This paper introduces a variational inequality (VI) approach to solving Bayesian regression games, which model adversarial robustness in machine learning applications. We extend existing results on Bayesian equilibria by providing sufficient conditions for their existence and uniqueness in a general class of convex and smooth Bayesian games. Specifically, we demonstrate that computing a Bayesian equilibrium can be reduced to solving an infinite-dimensional VI under certain conditions, allowing us to derive efficient algorithms with theoretical guarantees. We also explore two special settings—a finite Bayesian prior and a quadratic adversarial loss—where the VI reduces to a high-dimensional Euclidean VI or a nonconvex stochastic optimization problem, respectively. Our proposed algorithms, including a projected reflected gradient with inertial extrapolation and a randomized block coordinate descent method, achieve strong convergence with rigorous complexity analysis. These results advance the theoretical understanding of Bayesian regression games and provide practical tools for their computation.

Keywords: Bayesian Regression Games; Variational Inequality; Adversarial Robustness; Game Theory; Equilibrium Computation

1 Introduction

Adversarial robustness has become a critical concern in modern machine learning, where models are frequently deployed in environments with malicious actors or noisy data. Over the past few decades, this topic has attracted significant attention across various disciplines, including statistics, optimization, and machine learning. The need for robust models stems from two primary considerations: first, the requirement to maintain performance in the presence of adversarial perturbations, and second, the potential for adversarially robust models to exhibit better generalization and interpretability properties compared to their non-robust counterparts.

Traditionally, adversarial robustness has been studied within the framework of zero-sum games, where a learner and an adversary are pitted against each other with perfectly antagonistic objectives. While zero-sum games provide a solid foundation for analyzing adversarial interactions, they often fail to capture the complexities of real-world scenarios where the goals of the learner and the adversary are not strictly opposed. For example, in applications such as spam detection and network security, the adversary’s objectives may not directly negate those of the learner, leading to overly conservative strategies when using a minimax approach.

To address these limitations, we adopt a more flexible general-sum game framework, specifically focusing on Bayesian regression games. In this setting, the learner models the adversary’s uncertain behavior through a Bayesian prior, allowing for a more nuanced understanding of adversarial interactions. This framework enables the learner to optimize performance while accounting for the adversary’s possible strategies, even when the adversary’s exact cost function is unknown.

Despite its conceptual appeal, the Bayesian regression game framework presents significant challenges in both theoretical analysis and practical computation. Existing results on the existence and uniqueness of Bayesian equilibria are often restricted to specific cases, such as quadratic cost functions, and the algorithms proposed in prior work lack rigorous guarantees of convergence. To advance this line of research, we explore the use of variational inequality (VI) methods to characterize and compute Bayesian equilibria in a broader class of games.

Backgrounds. Adversarially robust models have seen a tremendous surge in research activity by various communities over the past decades, including statistics [Huber(2004)], optimization [Ben-Tal et al.(2009)], and machine learning [Globerson and Roweis(2006), Biggio and Roli(2018)]. In machine learning, there has been renewed interest in the topic driven by work on adversarial methods in deep learning [Bruna et al.(2014), Goodfellow et al.(2015)]. There are two main considerations underlying this line of work: (1) real-world deployments of machine-learning methods require robustness to malicious data and it is an ongoing challenge to provide such robustness [Madry et al.(2018), Wong et al.(2020)]; (2) adversarially robust models may generalize better [Zhu et al.(2020)] and have better interpretability properties than non-robust methods [Tsipras et al.(2019), Santurkar et al.(2019)]. To this end, adversarially robust models are often preferred by practitioners in real-world applications.

From a game-theoretic point of view, adversarially robust models naturally form a two-player game between a learner and an adversary. Both players choose actions simultaneously, and their interacting dynamics constitute a *noncooperative game* [Nash(1951)]. The learner’s action is to select the best set of model parameters while maximizing the test accuracy; the adversary’s action is to impose a perturbation on the input data distribution while paying a perturbation cost. Based on this general framework, existing approaches have mostly been restricted to the *zero-sum* case, in which the learner and adversary have fully conflicting goals. In this case, the learner is best off by choosing a *minimax strategy*; i.e., minimizing the worst-case cost over the action space of the adversary. For classification and regression problems, properties of the minimax solutions have been derived under a variety of assumptions [Lanckriet et al.(2002), EL Ghaoui et al.(2003), Globerson and Roweis(2006), Sayed and Chen(2002), Teo et al.(2008)].

On the positive side, the minimax solutions are computationally tractable in several specific settings [Sayed and Chen(2002)], or can be approximated through a convex relaxation [Teo et al.(2008)]. However, they come with a few limitations. First, zero-sum games are not flexible enough to capture cases in which the learner and the adversary do not have perfectly antagonistic goals [Brückner et al.(2012)]. For example, a credit card defrauder’s goal of maximizing the illicit profit made from exploiting phished account information via spam emails is not the exact inverse of an email service provider’s goal of achieving a close-to-zero false positives rate at spam recognition. In these cases, a minimax strategy can make overly pessimistic assumptions about the adversary’s behavior and lead to an optimal outcome. One approach to filling this gap is by relaxing the zero-sum assumption to *general-sum* games, in which the learner is fully aware of the costs of the adversary [Brückner et al.(2012)].

Moreover, in many practical applications, the cost function of the adversary is simply unknown to the learner. For instance, it is hard for an internet security service provider to know exactly the profit of the attackers. Therefore, the standard minimax solutions, which require *complete information* of the game, become infeasible to compute. This strong assumption of complete information can be lifted via a Bayesian game-theoretic framework [Harsanyi(1967)]. This has been pursued, for example, by [Großhans et al.(2013)], who propose *Bayesian regression games*. In this class of games, the learner’s uncertainty regarding the adversary’s costs is reflected in a

Bayesian prior over the parameters of the cost function. These authors derive sufficient conditions for the existence and uniqueness of a Bayesian equilibrium, as well as a graduated optimization algorithm to compute the equilibrium.

Despite the appealing conceptual framework, the sufficient conditions for the equilibrium uniqueness for Bayesian regression games that are known to date [Großhans et al.(2013), Theorem 2] are designed for quadratic cost functions. This excludes other common choices, e.g., logistic or smooth hinge functions, which are widely used in practice. Moreover, the algorithms studied in this line of work are heuristic, providing no theoretical guarantee of convergence. Thus we are motivated to tackle the following important open questions: *Can we generalize existing sufficient conditions for the existence and uniqueness of a Bayesian equilibrium in Bayesian regression games to more general cost functions? Can we develop efficient algorithms to compute the equilibrium?*

Contributions. We make several key contributions in this paper. First, we establish sufficient conditions for the existence and uniqueness of Bayesian equilibria in a broad class of convex and smooth Bayesian games by utilizing a variational inequality (VI) approach. This approach reformulates equilibrium computation as solving an infinite-dimensional VI in a Hilbert space, leveraging classical optimization results. Second, we consider two specific cases—finite Bayesian priors and quadratic adversarial losses—and demonstrate how the infinite-dimensional VI can be reduced to a high-dimensional Euclidean VI or a nonconvex Euclidean stochastic optimization problem. Third, we propose novel algorithms, including projected reflected gradient techniques with inertial extrapolation and randomized block coordinate descent, providing strong theoretical guarantees for convergence, with detailed analysis of their iteration complexity and convergence properties.

Organization. In Section 2, we present background on Bayesian regression games and prove that the computation of a Bayesian equilibrium amounts to solving an infinite-dimensional VI. In Section 3, we treat the existence and uniqueness of Bayesian equilibria. In Section 4, we consider two special settings of a finite Bayesian prior and a quadratic adversarial loss. We propose specific algorithms and analyze their convergence and iteration complexity in Section 5. We conclude the paper in Section 6.

Notation. We use bold lower-case letters such as \mathbf{x} to denote vectors, upper-case letters such as X to denote matrices, and calligraphic upper-case letters such as \mathcal{X} to denote sets. The notation $[n]$ refers to $\{1, 2, \dots, n\}$ for some integer $n > 0$. The symbols $\mathbf{0}_n$ and $\mathbf{0}_{n \times m}$ refer to the vector and the matrix in \mathbb{R}^n and $\mathbb{R}^{n \times m}$ whose entries are all zeroes. We let $\mathbb{E}[\cdot]$ denote an expectation and use $\mathbb{E}_q[\cdot]$ to denote an expectation over a distribution q . For a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, we let $\nabla f(\mathbf{x})$ denote the gradient of f at \mathbf{x} . For a vector $\mathbf{x} \in \mathbb{R}^d$, we denote $\|\mathbf{x}\|$ as its ℓ_2 -norm. For a matrix $X \in \mathbb{R}^{d \times r}$, we let $\|X\|_F$ denote its Frobenius norm. As an abuse of notation, we denote $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$ as the inner product between two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\langle X, Y \rangle = \text{Trace}(X^\top Y)$ as the inner product between two matrices $X, Y \in \mathbb{R}^{d \times r}$, where $\text{Trace}(\cdot)$ stands for the trace of a matrix. For $\mathcal{X} \subseteq \mathbb{R}^d$, we let $D_{\mathcal{X}}$ denote its diameter: $D_{\mathcal{X}} = \max_{\mathbf{x}, \mathbf{x}' \in \mathcal{X}} \|\mathbf{x} - \mathbf{x}'\|$. Given $\epsilon > 0$, $a = O(b(\epsilon))$ stands for the upper bound $a \leq C \cdot b(\epsilon)$, where $C > 0$ is independent of ϵ . Similarly, $a = \tilde{O}(b(\epsilon))$ indicates that the inequality may depend on a logarithmic function of $1/\epsilon$, where $C > 0$ is independent of ϵ .

1.1 Related Work

We refer to [Kolter and Madry(2018)] as a reference point for the burgeoning literature on adversarially robust models. Despite the attention devoted to this topic, [Großhans et al.(2013)] is one of very few papers that models the general-sum nature and the uncertainty in the adversary via the classical formalism of Bayesian games.

Bayesian games. Bayesian game has been a classical approach in game theory to model the situation of asymmetric or incomplete information in games. In his seminal work, [Harsanyi(1967)] proved the existence of a *Bayesian Nash equilibrium* in finite Bayesian games given a common prior and common knowledge of that prior among all the players. Subsequently, [Mertens and Zamir(1985)] introduced a relaxed notion of a “universal prior space”, which is a sufficiently large space that captures players’ higher-order beliefs. Further, [Aghassi and Bertsimas(2006)] considered the special case where the payoffs are drawn from a bounded uncertainty set but the distribution is fully unknown.

More recently, the appealing formulation of Bayesian games which captures the players’ uncertainty has led to many works in economics settings, in particular for auctions. Specifically, each bidder has a private valuation function that expresses complex preferences over all subsets of items, and bidders have beliefs about the valuation functions of the other bidders, in the form of probability distributions [Myerson(1985)]. In this setting, a Bayesian equilibrium can be viewed as an approximation for the optimal social welfare value. Unfortunately, most existing results on the complexity of finding a Bayesian equilibrium in various auctions are negative [Christodoulou et al.(2008), Bhawalkar and Roughgarden(2011), Feldman et al.(2013)]. Computing a Bayesian equilibrium is in PP and even finding an ϵ -approximate Bayesian equilibrium is NP-hard when $\epsilon > 0$ is small [Cai and Papadimitriou(2014)].

Equilibrium existence and computation. The existence of a mixed strategy Nash equilibrium is well known in finite games with complete information [Nash(1950)]. Such existence results are derived via the Brouwer fixed-point theorem [Kakutani(1941)], which suggests intuitively that a fixed-point iteration might be an efficient approach for computing a Nash equilibrium. Similar existence and uniqueness results are derived for concave games with complete information [Rosen(1965)]. However, computation of the equilibrium has been a challenging problem. [Chen et al.(2009)] recently proved that the problem of finding a Nash equilibrium for even the simplest two-player general-sum games is PPAD-complete [Papadimitriou(1994)]. Further complexity results have been established for equilibrium computation in games with complete information under various assumptions [Gilboa and Zemel(1989), Megiddo and Papadimitriou(1991), Conitzer and Sandholm(2003), Conitzer and Sandholm(2008), Daskalakis et al.(2009), Rubinstein(2018)]. In Bayesian games, the complexity of deciding the existence of a pure Bayesian equilibrium is in general NP-hard [Conitzer and Sandholm(2003), Gottlob et al.(2007)]. Nonetheless, a few Bayesian games are computationally tractable. Two canonical examples are: (i) tree-games [Singh et al.(2004)], where the cost function depends only on the actions/types of neighboring players and the interaction formed by the neighborhood relation is a tree; (ii) two-player zero-sum Bayesian games with finite prior. The counterfactual regret minimization (CFR) algorithm was proposed with solid theoretical guarantee [Zinkevich et al.(2007)].

2 Bayesian Regression Games

In this section, we present the setup and equilibrium concepts for Bayesian regression games with general cost functions.

2.1 Basic setup

We consider a general-sum game between a *learner* of a regression model and a *data generator* who is able to perturb the data distribution. Let $(X, \mathbf{y}) \in \mathbb{R}^{n \times m} \times \mathbb{R}^n$ be a pair of data matrix and target vector, which are generated by the data generator at training time. Denote each row of X as \mathbf{x}_i , each entry of \mathbf{y} as y_i for $i \in [n]$. We assume that all pairs of instances $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ are drawn from an unknown distribution μ over $\mathcal{X} \times \mathcal{Y}$. At testing time, the data generator provides new instances drawn from another distribution $\bar{\mu}$ defined on $\mathcal{X} \times \mathcal{Y}$, but might not be μ . These instances are generally not available at the training time.

For the learner, we denote the action space as $\mathcal{W} \subseteq \mathbb{R}^m$, which is the space of the regression parameters. We denote an instance-specific weight by $c_l(\mathbf{x}, y) \geq 0$. Then, the learner's cost at testing time is the weighted average loss, i.e. $\theta_l(\mathbf{w}, \bar{\mu}, c_l) = \int c_l(\mathbf{x}, y) f_l(\mathbf{w}, \mathbf{x}, y) d\bar{\mu}(\mathbf{x}, y)$, where $\mathbf{w} \in \mathcal{W}$ is a parameter of the prediction model and f_l is the cost function.

For the data generator, intuitively, the goal is to manipulate the input data in order to achieve certain targeted predictions. Therefore, the costs of the data generator contain two parts. The first part is a cost of performing the manipulation, and the second part is a cost which quantifies the loss between the actual predictions and the data generator's targeted predictions. For the cost of manipulation, the manipulation of the input data is reflected in the difference between the distributions μ and $\bar{\mu}$. We denote the cost of such manipulation for the data generator as $\Omega_d(\mu, \bar{\mu})$. For the cost on the predictions, we denote the vector of target values for n data points as $z(\mathbf{x}, y) \in \mathcal{Z} \subseteq \mathbb{R}^n$, and f_d as the cost function. Similarly to the learner, we also allow the data generator to have an instance-specific weight $c_d(\mathbf{x}, y)$. Then, the data generator's costs are defined by $\theta_d(\mathbf{w}, \bar{\mu}, c_d) = \int c_d(\mathbf{x}, y) f_d(\mathbf{w}, \mathbf{x}, z(\mathbf{x}, y)) d\bar{\mu}(\mathbf{x}, y) + \Omega_d(\mu, \bar{\mu})$.

Note that the theoretical costs of both players depend on the unknown distributions μ and $\bar{\mu}$. Thus, we focus on the regularized empirical counterparts of the theoretical costs based on the training samples $(X, \mathbf{y}, \mathbf{z})$, where $\mathbf{y} \in \mathbb{R}^n$ and $\mathbf{z} \in \mathbb{R}^n$ are the empirical versions of y and $z(\mathbf{x}, y)$ respectively. The empirical counterpart of the term $\Omega_d(\mu, \bar{\mu})$ is represented by the difference between the training matrix X and a perturbed matrix \bar{X} that would be the outcome of applying the transformation which translates μ into $\bar{\mu}$ to X . For simplicity, we denote this by $\Omega_d(X, \bar{X})$. We also denote the empirical instance-specific weights for the two players by $\mathbf{c}_l \in \mathbb{R}^n$ and $\mathbf{c}_d \in \mathbb{R}^n$. Then, the empirical costs of the learner and the data generator are given by

$$\begin{aligned}\hat{\theta}_l(\mathbf{w}, \bar{X}, \mathbf{c}_l) &= \sum_{i=1}^n c_{l,i} f_l(\mathbf{w}, \bar{\mathbf{x}}_i, y_i) + \Omega_l(\mathbf{w}) \\ \hat{\theta}_d(\mathbf{w}, \bar{X}, \mathbf{c}_d) &= \sum_{i=1}^n c_{d,i} f_d(\mathbf{w}, \bar{\mathbf{x}}_i, z_i) + \Omega_d(X, \bar{X})\end{aligned}\tag{1}$$

where $\Omega_l(\mathbf{w})$ is a regularization term. With a focus on machine-learning applications, we illustrate the above setup with a few common loss functions, regularization terms and constraint sets. Besides the quadratic case [Großhans et al.(2013)], we provide a motivating example using logistic function and the unit ball constraint sets. This is encouraged by the fact that the logistic loss functions are more suitable than quadratic ones for the binary classification/regression problems, such as the email spam filtering and network security detection.

Example 1 (quadratic). $f_l(\mathbf{w}, \mathbf{x}, y) = (\mathbf{x}^\top \mathbf{w} - y)^2$, $\Omega_l(\mathbf{w}) = \|\mathbf{w}\|^2$, $f_d(\mathbf{w}, \mathbf{x}, z) = (\mathbf{x}^\top \mathbf{w} - z)^2$ and $\Omega_d(X, \bar{X}) = \|X - \bar{X}\|_F^2$; $\mathcal{W} = \mathbb{R}^m$, $\mathcal{X} = \mathbb{R}^{n \times m}$ and $\mathcal{Y} = \mathcal{Z} = \mathbb{R}^n$.

Example 2 (logistic). $f_l(\mathbf{w}, \mathbf{x}, y) = \log(1 + \exp(-y\mathbf{x}^\top \mathbf{w}))$, $\Omega_l(\mathbf{w}) = \|\mathbf{w}\|^2$, $f_d(\mathbf{w}, \mathbf{x}, z) = \log(1 + \exp(-z\mathbf{x}^\top \mathbf{w}))$ and $\Omega_d(X, \bar{X}) = \|X - \bar{X}\|_F^2$; $\mathcal{W} = \{\mathbf{w} \in \mathbb{R}^m \mid \|\mathbf{w}\| \leq 1\}$, $\mathcal{X} = \{X \in \mathbb{R}^{n \times m} \mid \|X\|_F \leq 1\}$ and $\mathcal{Y} = \mathcal{Z} = \{-1, 1\}^n$.

2.2 Bayesian regression game

The previous basic setup describes a general-sum two-player regression game. Now we present the formulation of the Bayesian regression game based on it. First, note that the cost functions (see Eq (1)) depend on the actions of both players—the parameters \mathbf{w} and the transformation manifested in \bar{X} . Further, the cost function of the data generator depends on the instance-specific weight \mathbf{c}_d . We are interested in a setting where these weights are private information of the data generator, and are unknown to the learner. Instead, the learner is only informed about the instance-specific weight \mathbf{c}_d through a Bayesian prior $q(\mathbf{c}_d)$. As argued by [Großhans et al.(2013)], this asymmetry of uncertainty is crucial. By modeling the learner’s lack of information about the data generator, we intend to make the learner more robust to new adversarial examples. Thus, this setting is naturally formulated as a *two-player general-sum Bayesian game*.

We denote this Bayesian regression game by the tuple $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$. From the learner’s viewpoint, \mathbf{c}_d is a random variable that is drawn from a Bayesian prior q at testing time. During training, the data generator commits to a parametric strategy, $\sigma : \mathbb{R}^n \rightarrow \mathcal{X}$, which maps a value of \mathbf{c}_d (unknown to the learner) to a transformation reflected in \bar{X} . In other words, the action space of the data generator Σ contains all the functions from \mathbb{R}^n to \mathcal{X} . We also define Bayesian Equilibrium and its approximation. We denote the best responses of the learner and the data generator as: $\mathbf{w}^*[\sigma] = \operatorname{argmin}_{\mathbf{w} \in \mathcal{W}} \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_l)]$ and $\sigma^*[\mathbf{w}] = \operatorname{argmin}_{\bar{X} \in \mathcal{X}} \hat{\theta}_d(\mathbf{w}, \bar{X}, \mathbf{c}_d)$.

Definition 1 (Bayesian Equilibrium). *The strategy profile $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ is a Bayesian equilibrium for the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ if, for a.e., $\omega \in \Omega$, the following statement holds true: $(\mathbf{w}_*, \sigma_*(\mathbf{c}_d)) = (\mathbf{w}^*[\sigma_*], \sigma^*[\mathbf{w}_*](\mathbf{c}_d))$, or equivalently, $\mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] \leq \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)]$ for all $\mathbf{w} \in \mathcal{W}$ and $\hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d) \leq \hat{\theta}_d(\mathbf{w}_*, \bar{X}, \mathbf{c}_d)$ for all $\bar{X} \in \mathcal{X}$.*

Definition 2 (ϵ -Bayesian Equilibrium). *The strategy profile $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ is an ϵ -Bayesian equilibrium for the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ if, for a.e., $\omega \in \Omega$, the following statement holds true: $\|\mathbf{w} - \mathbf{w}_*\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d)\|_F^2] \leq \epsilon$, where the strategy profile (\mathbf{w}_*, σ_*) is a Bayesian equilibrium.*

When the distribution q is a point mass, it is clear that no uncertainty exists and a Bayesian equilibrium is a Nash equilibrium. This corresponds to a two-player general-sum game which is still more general than the minimax strategy.

2.3 Infinite-dimensional variational inequality model

We show that we can regard a Bayesian equilibrium as a solution of an infinite-dimensional VI [Kinderlehrer and Stampacchia(2000)] (Eq. (3)). This characterization is not only necessary but sufficient under certain assumptions. We adapt a proof technique in [Ui(2016)], and specialized the VI model to Bayesian regression games. We make the following assumption throughout this paper.

Assumption 1. *The Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ is convex and smooth: (1) The cost function $\hat{\theta}_l(\cdot, \bar{X}, \mathbf{c}_l) : \mathcal{W} \rightarrow \mathbb{R}$ is convex and continuously differentiable for each $\bar{X} \in \mathcal{X}$, and $\|\nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \bar{X}, \mathbf{c}_l)\|^2 < +\infty$ for each $(\mathbf{w}, \bar{X}) \in \mathcal{W} \times \mathcal{X}$; (2) The cost function $\hat{\theta}_d(\mathbf{w}, \cdot, \mathbf{c}_d) : \mathcal{X} \rightarrow \mathbb{R}$*

is convex and continuously differentiable for each $\mathbf{w} \in \mathcal{W}$, and $\mathbb{E}_q[\|\nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \bar{X}, \mathbf{c}_d)\|_F^2] < +\infty$ for each $(\mathbf{w}, \bar{X}) \in \mathcal{W} \times \mathcal{X}$; **(3)** The action spaces \mathcal{W} and \mathcal{X} are both closed and convex.

We provide some intuitions for these assumptions. First, note that a Bayesian equilibrium is characterized by the coordinate-wise minimization of the cost functions $\hat{\theta}_l$ and $\hat{\theta}_d$ over the action spaces \mathcal{W} and \mathcal{X} . Thus, it is necessary to impose convexity conditions on the cost functions and a few moment conditions on the Bayesian prior q . Furthermore, our moment conditions are implied by the assumptions $\mathbb{E}_q[c_{d,i}] < +\infty$ for all $i \in [n]$ made in [Großhans et al.(2013), Theorem 1 and 2] and are thus slightly weaker. Finally, even if the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ is not smooth, we can derive a similar first-order condition using the subgradients of $\hat{\theta}_l(\cdot, \bar{X}, \mathbf{c}_l)$ and $\hat{\theta}_d(\mathbf{w}, \cdot, \mathbf{c}_d)$ and regard a Bayesian equilibrium as a solution of a multi-valued VI [Kinderlehrer and Stampacchia(2000)]; see Eq. (3) for the details.

Proposition 1. *Under Assumption 1, the strategy profile $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ is a Bayesian equilibrium for the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ if and only if, for a.e. $\omega \in \Omega$, the following statement holds true:*

$$\begin{aligned} \mathbb{E}_q[\langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l) \rangle] &\geq 0, \quad \forall \mathbf{w} \in \mathcal{W} \\ \langle \bar{X} - \sigma_*(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d) \rangle &\geq 0, \quad \forall \bar{X} \in \mathcal{X} \end{aligned} \quad (2)$$

Proposition 2. *Under Assumption 1, the strategy profile $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ is a Bayesian equilibrium for the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ iff for $\forall (\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$, we have*

$$\mathbb{E}_q[\langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l) \rangle] + \langle \sigma(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d) \rangle \geq 0 \quad (3)$$

For a game with complete information, it is common to study the existence, uniqueness and computation of a Nash equilibrium by regarding it as a solution of a VI in a finite-dimensional space. This approach dates back to [Lions and Stampacchia(1967)] and has been thoroughly studied in the optimization literature [Facchinei and Pang(2007)]. While the VI approach can be formally extended to an infinite-dimensional space [Kinderlehrer and Stampacchia(2000)], it has not been recognized as a useful analytical tool to study games with incomplete information. The only work that we are aware of in this vein is [Ui(2016)], who gives a sufficient condition for the existence and uniqueness of a Bayesian equilibrium by regarding it as a solution of an infinite-dimensional VI. The focus in that work is, however, a general setting without any consideration of the computation of an equilibrium.

2.4 Proof of Proposition 1

Proof of Proposition 1. Note that $\hat{\theta}_l(\cdot, \bar{X}, \mathbf{c}_l) : \mathcal{W} \rightarrow \mathbb{R}$ is convex and continuously differentiable for each $\bar{X} \in \mathcal{X}$. By the Lebesgue monotone convergence theorem, we have

$$\nabla_{\mathbf{w}} \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] = \mathbb{E}_q[\nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] \quad (4)$$

For $\mathbf{w} \in \mathcal{W}$, we let $g(t) = \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_* + t(\mathbf{w} - \mathbf{w}_*), \sigma_*(\mathbf{c}_d), \mathbf{c}_l)]$. Since $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ is a Bayesian equilibrium, we have $g(t) \geq g(0)$ for all $t \in \mathbb{R}$. This implies that $g'(0) \geq 0$. By definition,

$$\begin{aligned} g'(0) &= \left\langle \mathbf{w} - \mathbf{w}_*, \left\{ \nabla_{\mathbf{w}} \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_* + t(\mathbf{w} - \mathbf{w}_*), \sigma_*(\mathbf{c}_d), \mathbf{c}_l)]|_{t=0} \right\} \right\rangle \\ &= \langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] \rangle \\ &\stackrel{\text{Eq. (4)}}{=} \langle \mathbf{w} - \mathbf{w}_*, \mathbb{E}_q[\nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] \rangle = \mathbb{E}_q[\langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l) \rangle] \end{aligned}$$

Putting these pieces together yields the first inequality in Eq. (2). In addition, notice that the cost function $\hat{\theta}_d(\mathbf{w}, \cdot, \mathbf{c}_d) : \mathcal{X} \rightarrow \mathbb{R}$ is convex and continuously differentiable for each $\mathbf{w} \in \mathcal{W}$ and no expectation is involved now: By the similar argument, we obtain the second inequality in Eq. (2).

Conversely, we show that Eq. (2) guarantees that the strategy profile $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ is a Bayesian equilibrium. Note that the function $g(t) = \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_* + t(\mathbf{w} - \mathbf{w}_*), \sigma_*(\mathbf{c}_d), \mathbf{c}_l)]$ is convex since $\hat{\theta}_l(\cdot, \bar{X}, \mathbf{c}_l) : \mathcal{W} \rightarrow \mathbb{R}$ is convex for each $\bar{X} \in \mathcal{X}$. Thus, $g(t) \geq g(0) + tg'(0)$ for each $t \in \mathbb{R}$. By definition, we have

$$\mathbb{E}_q[\hat{\theta}_l(\mathbf{w}, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] = g(1) \geq g(0) + tg'(0) = \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] + \mathbb{E}_q[\langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l) \rangle]$$

Combining the above inequality with the first inequality in Eq. (2), we have

$$\mathbb{E}_q[\hat{\theta}_l(\mathbf{w}, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)] \geq \mathbb{E}_q[\hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l)], \quad \text{for all } \mathbf{w} \in \mathcal{W}$$

Using a similar argument and the second inequality in Eq. (2), we have

$$\hat{\theta}_d(\mathbf{w}_*, \bar{X}, \mathbf{c}_d) \geq \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d), \quad \text{for all } \bar{X} \in \mathcal{X}$$

This completes the proof. \square

2.5 Proof of Proposition 2

Proof of Proposition 2. We first show the “only if” direction. Indeed, let $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ be a Bayesian equilibrium for the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$, we derive from Proposition 1 that Eq. (2) holds true. This implies that, for all $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ and for a.e. $\omega \in \Omega$,

$$\begin{aligned} \mathbb{E}_q[\langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l) \rangle] &\geq 0 \\ \langle \sigma(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d) \rangle &\geq 0 \end{aligned}$$

Summing up the above two inequalities and taking the expectation over the distribution q yields the desired inequality in Eq. (3). Then it suffices to show the “if” direction. Specifically, we show that if $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ is not a Bayesian equilibrium for the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$, then Eq. (2) does not hold true. By Proposition 1, if $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ is not a Bayesian equilibrium, we have

$$\mathbb{E}_q[\langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l) \rangle] < 0 \quad \text{for some } \mathbf{w} \in \mathcal{W} \quad (5)$$

or there exists $E \subseteq \Omega$ with $\mathbb{P}(E) > 0$ such that, for each $\omega \in E$,

$$\langle \bar{X} - \sigma_*(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d) \rangle < 0 \quad \text{for some } \bar{X} \in \mathcal{X} \quad (6)$$

Let $(\mathbf{w}', \sigma') \in \mathcal{W} \times \Sigma$ be defined by

$$\begin{aligned} \mathbf{w}' &= \begin{cases} \mathbf{w} & \text{if Eq. (5) holds true} \\ \mathbf{w}_* & \text{otherwise} \end{cases} \\ \sigma'(\mathbf{c}_d(\omega)) &= \begin{cases} \bar{X} & \text{if Eq. (6) holds true and } \omega \in E \\ \sigma_*(\mathbf{c}_d(\omega)) & \text{otherwise} \end{cases} \end{aligned}$$

By simple calculations, we have

$$\mathbb{E}_q \left[\left\langle \mathbf{w}' - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l) \right\rangle + \left\langle \sigma'(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d) \right\rangle \right] < 0$$

This completes the proof. \square

3 Equilibrium Existence and Uniqueness

In this section, we prove existence and uniqueness results by a variational inequality (VI) formulation. Based on the infinite-dimensional VI formulation (see Eq. (3)), we prove a set of sufficient conditions for the existence and uniqueness of a Bayesian equilibrium. Our results generalize the existing work to Bayesian regression games with general convex cost functions.

For the Bayesian regression game G with its equilibrium defined by Eq. (3), we define a Hilbert space \mathcal{H} consisting of (an equivalence class of) functions $\beta : \mathbb{R}^n \mapsto \mathbb{R}^m \times \mathbb{R}^{n \times m}$ with the inner product by $\langle (\mathbf{w}, \sigma), (\mathbf{w}', \sigma') \rangle_{\mathcal{H}} = \mathbb{E}_q[\langle \mathbf{w}(\mathbf{c}_d), \mathbf{w}'(\mathbf{c}_d) \rangle] + \langle \sigma(\mathbf{c}_d), \sigma'(\mathbf{c}_d) \rangle < +\infty$. Note that each element in $\mathcal{W} \subseteq \mathbb{R}^m$ can be regarded as a constant function from \mathbb{R}^n to \mathbb{R}^m whose value is this element. We denote the set of these constant functions by $\Sigma_{\mathcal{W}}$ (an equivalence class of \mathcal{W}) and define a mapping $T : \Sigma_{\mathcal{W}} \times \Sigma \rightarrow \mathcal{H}$ by

$$T \left(\begin{matrix} \mathbf{w} \\ \sigma(\cdot) \end{matrix} \right) = \left(\begin{matrix} \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma(\cdot), \mathbf{c}_l) \\ \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma(\cdot), \cdot) \end{matrix} \right) \in \mathcal{H} \quad (7)$$

Thus, the computation of a Bayesian equilibrium is equivalent to solving a VI in the space \mathcal{H} . This allows us to analyze the existence and uniqueness of a Bayesian equilibrium under the VI framework. For example, the existence of a Bayesian equilibrium is guaranteed by the continuity and monotonicity of T as well as some additional conditions on $\mathcal{W} \times \Sigma$.

Definition 3 (Monotonicity). *Let \mathcal{H} be a Hilbert space with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$, we define $\mathcal{S} \subseteq \mathcal{H}$ as a closed and convex set and $T : \mathcal{S} \rightarrow \mathcal{H}$ as a mapping. Then, T is **monotone** if $\langle T\beta - T\beta', \beta - \beta' \rangle_{\mathcal{H}} \geq 0$ for each $\beta, \beta' \in \mathcal{S}$; T is **strictly monotone** if $\langle T\beta - T\beta', \beta - \beta' \rangle_{\mathcal{H}} > 0$ for each $\beta, \beta' \in \mathcal{S}$ with $\beta \neq \beta'$; T is **λ -strongly monotone** ($\lambda > 0$) if $\langle T\beta - T\beta', \beta - \beta' \rangle_{\mathcal{H}} \geq \lambda \|\beta - \beta'\|_{\mathcal{H}}^2$ for each $\beta, \beta' \in \mathcal{S}$.*

We summarize the existence and uniqueness results in the following two theorems.

Theorem 1 (Existence). *Suppose that the mapping T defined by Eq. (7) is continuous and monotone, and the action space $\mathcal{W} \times \Sigma$ is nonempty, closed and convex. If $\mathcal{W} \times \Sigma$ is compact, or there exists $(\mathbf{w}_0, \sigma_0) \in \mathcal{W} \times \Sigma$ such that, for all $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ satisfying $\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2] \rightarrow +\infty$, the following statement holds true:*

$$\frac{\mathbb{E}_q[\langle \mathbf{w} - \mathbf{w}_0, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_l) \rangle] + \langle \sigma(\mathbf{c}_d) - \sigma_0(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_d) \rangle]}{\sqrt{\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2]}} \rightarrow +\infty \quad (8)$$

Then, the VI in Eq. (3) has at least one solution.

Corollary 1. *Under Assumption 1, if T defined by Eq. (7) is monotone and there exists $(\mathbf{w}_0, \sigma_0) \in \mathcal{W} \times \Sigma$ such that, for all $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ satisfying $\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2] \rightarrow +\infty$, Eq. (8) holds. Then, the Bayesian regression game G has at least one Bayesian equilibrium.*

Corollary 2. *Under Assumption 1, we assume that T defined by Eq. (7) is monotone and the action space $\mathcal{W} \times \Sigma$ is compact. Then, the Bayesian regression game G has at least one Bayesian equilibrium.*

Theorem 2 (Existence and Uniqueness). *Suppose that T defined by Eq. (7) is continuous and strictly monotone, and $\mathcal{W} \times \Sigma$ is nonempty, closed and convex. If $\mathcal{W} \times \Sigma$ is compact, or there exists $(\mathbf{w}_0, \sigma_0) \in \mathcal{W} \times \Sigma$ such that, for all $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ satisfying $\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2] \rightarrow +\infty$, Eq. (8) holds true. The VI in Eq. (3) has a unique solution.*

Corollary 3. *Under Assumption 1, if T defined by Eq. (7) is λ -strongly monotone. Then, the Bayesian regression game G has a unique Bayesian equilibrium.*

Corollary 4. *Under Assumption 1, if the mapping T defined by Eq. (7) is strictly monotone and the action space $\mathcal{W} \times \Sigma$ is compact. Then, the Bayesian regression game G has a unique Bayesian equilibrium.*

The monotonicity condition in Corollary 4 generalizes [Großhans et al.(2013), Theorem 2], which is a special case with quadratic loss functions. Our VI approach also supplies an intuitive yet rigorous justification for Eq. (8), demonstrating that it arises from the strict monotonicity of the mapping T .

3.1 Proof of Theorem 1

We provide a key notion of monotonicity which is pivotal in the classical VI literature and summarize in Proposition 3 the celebrated existence theorem for an infinite-dimensional VI.

Definition 4 (Monotonicity). *Let \mathcal{H} be a Hilbert space with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$, we define $\mathcal{S} \subseteq \mathcal{H}$ as a closed and convex set and $T : \mathcal{S} \rightarrow \mathcal{H}$ as a mapping. Then,*

- (i) *T is monotone if $\langle T\beta - T\beta', \beta - \beta' \rangle_{\mathcal{H}} \geq 0$ for each $\beta, \beta' \in \mathcal{S}$.*
- (ii) *T is strictly monotone if $\langle T\beta - T\beta', \beta - \beta' \rangle_{\mathcal{H}} > 0$ for each $\beta, \beta' \in \mathcal{S}$ with $\beta \neq \beta'$.*
- (iii) *T is λ -strongly monotone ($\lambda > 0$) if $\langle T\beta - T\beta', \beta - \beta' \rangle_{\mathcal{H}} > \lambda \|\beta - \beta'\|_{\mathcal{H}}^2$ for each $\beta, \beta' \in \mathcal{S}$.*

Proposition 3 (Browder-Hartman-Stampachia). *Let \mathcal{H} be a Hilbert space with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$. Define $\mathcal{S} \subseteq \mathcal{H}$ as a nonempty, closed and convex set, and $T : \mathcal{S} \rightarrow \mathcal{H}$ as a monotone mapping. If the following conditions hold true:*

- (i) *The mapping $t \mapsto \langle T((1-t)\beta + t\beta'), \alpha \rangle_{\mathcal{H}}$ from $[0, 1]$ to \mathbb{R} is continuous for all $\beta, \beta' \in \mathcal{S}$ and $\alpha \in \mathcal{H}$.*
- (ii) *The set \mathcal{S} is compact, or there exists $\beta_0 \in \mathcal{S}$ such that $\frac{\langle T\beta, \beta - \beta_0 \rangle_{\mathcal{H}}}{\|\beta\|_{\mathcal{H}}} \rightarrow +\infty$ as $\|\beta\|_{\mathcal{H}} \rightarrow +\infty$.*

Then, there exists $\bar{\beta} \in \mathcal{S}$ such that $\langle T\bar{\beta}, \beta - \bar{\beta} \rangle_{\mathcal{H}} \geq 0$ for all $\beta \in \mathcal{S}$.

Let us consider the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ and its Bayesian equilibrium in terms of Eq. (3). Then, we can define a Hilbert space \mathcal{H} consisting of (an equivalence class of) a function $\beta : \mathbb{R}^n \mapsto \mathbb{R}^m \times \mathbb{R}^{n \times m}$ with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ by

$$\left\langle \begin{pmatrix} \mathbf{w} \\ \sigma \end{pmatrix}, \begin{pmatrix} \mathbf{w}' \\ \sigma' \end{pmatrix} \right\rangle_{\mathcal{H}} = \mathbb{E}_q[\langle \mathbf{w}(\mathbf{c}_d), \mathbf{w}'(\mathbf{c}_d) \rangle + \langle \sigma(\mathbf{c}_d), \sigma'(\mathbf{c}_d) \rangle] \quad (9)$$

Note that each element in $\mathcal{W} \subseteq \mathbb{R}^m$ can be regarded as a constant function from \mathbb{R}^n to \mathbb{R}^m whose value equal to this element. Then, we denote the set of these constant function by $\Sigma_{\mathcal{W}}$ (an equivalence class of \mathcal{W}) and define the mapping $T : \Sigma_{\mathcal{W}} \times \Sigma \rightarrow \mathcal{H}$ as follows,

$$T \begin{pmatrix} \mathbf{w} \\ \sigma(\cdot) \end{pmatrix} = \begin{pmatrix} \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma(\cdot), \mathbf{c}_l) \\ \nabla_{\hat{X}} \hat{\theta}_d(\mathbf{w}, \sigma(\cdot), \cdot) \end{pmatrix} \in \mathcal{H} \quad (10)$$

To this end, the computation of a Bayesian equilibrium is equivalent to solving a VI in the space \mathcal{H} . This allows us to analyze the existence and uniqueness of a Bayesian equilibrium under the Browder-Hartman-Stampacchia's VI framework (cf. Proposition 3). For example, the existence of a Bayesian equilibrium is guaranteed by the continuity and monotonicity of T as well as some additional conditions on $\mathcal{W} \times \Sigma$.

To prove the existence of a solution, we show that Eq. (3) is a special case of the Browder-Hartman-Stampacchia VIs. Indeed, we set a Hilbert space \mathcal{H} consisting of (an equivalence class of) a function $\beta : \mathbb{R}^n \mapsto \mathbb{R}^m \times \mathbb{R}^{n \times m}$ with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ defined by Eq. (9). By abuse of notation, any element $\mathbf{w} \in \mathcal{W}$ define a constant function $\mathbf{w}(\cdot) \in \Sigma_{\mathcal{W}}$ and vice versa. This implies that $\Sigma_{\mathcal{W}}$ is an equivalent class of \mathcal{W} . Thus, $\mathcal{S} = \Sigma_{\mathcal{W}} \times \Sigma$ is a nonempty, closed and convex subset of \mathcal{H} . In addition, a mapping $T : \mathcal{S} \mapsto \mathcal{H}$ defined by Eq. (10) is continuous and monotone. Putting these pieces together with either the compactness of \mathcal{S} or the certain condition with Eq. (8) yields that all the assumptions in Proposition 3 hold true and Eq. (3) is a special case of the Browder-Hartman-Stampacchia VIs. Therefore, we conclude from Proposition 3 that the VI in Eq. (3) has at least one solution. This completes the proof.

3.2 Proof of Corollary 1

Under Assumption 1, Proposition 2 implies that a Bayesian equilibrium must be a solution of the VI in Eq. (3). Thus, it suffices to verify the assumptions in Theorem 1. Indeed, Assumption 1 guarantees that T defined by Eq. (10) is continuous. In addition, T is monotone and there exists $(\mathbf{w}_0, \sigma_0) \in \mathcal{W} \times \Sigma$ such that, for all $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ satisfying $\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2] \rightarrow +\infty$, Eq. (8) holds true. Therefore, we conclude the desired result.

3.3 Proof of Corollary 2

The proof is nearly the same as that of Corollary 1 and the only difference is that, we assume the compactness of the action space $\mathcal{W} \times \Sigma$, which is another condition of Theorem 1. Thus, by using Theorem 1 again, we conclude the desired result.

3.4 Proof of Theorem 2

Theorem 1 guarantees the existence of at least one solution of the VI in Eq. (3). Thus, it suffices to show that at most one solution exists if T defined by Eq. (10) is strictly monotone. Using the proof by contradiction, suppose that the VI in Eq. (3) has two different solutions. Given the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ is defined by Eq. (9), we have

$$\left\langle \begin{pmatrix} \mathbf{w}_1 - \mathbf{w}_2 \\ \sigma_1 - \sigma_2 \end{pmatrix}, T \begin{pmatrix} \mathbf{w}_1 \\ \sigma_1 \end{pmatrix} \right\rangle_{\mathcal{H}} \geq 0, \quad \left\langle \begin{pmatrix} \mathbf{w}_2 - \mathbf{w}_1 \\ \sigma_2 - \sigma_1 \end{pmatrix}, T \begin{pmatrix} \mathbf{w}_2 \\ \sigma_2 \end{pmatrix} \right\rangle_{\mathcal{H}} \geq 0$$

Summing up the above two inequalities yields:

$$\left\langle \begin{pmatrix} \mathbf{w}_1 - \mathbf{w}_2 \\ \sigma_1 - \sigma_2 \end{pmatrix}, T \begin{pmatrix} \mathbf{w}_1 \\ \sigma_1 \end{pmatrix} - T \begin{pmatrix} \mathbf{w}_2 \\ \sigma_2 \end{pmatrix} \right\rangle_{\mathcal{H}} \geq 0$$

Since a mapping $T : \Sigma_{\mathcal{W}} \times \Sigma \rightarrow \mathbb{R}^m \times \Sigma_0$ is strictly monotone, we have $\mathbf{w}_1 = \mathbf{w}_2$ and $\sigma_1 = \sigma_2$ almost everywhere. This leads to a contradiction and completes the proof.

3.5 Proof of Corollary 3

Under Assumption 1, Proposition 2 implies that a Bayesian equilibrium must be a solution of the VI in Eq. (3). Thus, it suffices to verify the assumptions in Theorem 2. Note that T defined by Eq. (10) is λ -strongly monotone and thus strictly monotone. In order to prove the uniqueness of a Bayesian equilibrium using Theorem 2, it remains to show that there exists $(\mathbf{w}_0, \sigma_0) \in \mathcal{W} \times \Sigma$ such that, for all $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ satisfying $\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2] \rightarrow +\infty$, Eq. (8) holds true. Since T is λ -strongly monotone, we have

$$\begin{aligned} & \mathbb{E}_q \left[\left\langle \mathbf{w} - \mathbf{w}_0, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_l) \right\rangle + \left\langle \sigma(\mathbf{c}_d) - \sigma_0(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_d) \right\rangle \right] \\ & \stackrel{\text{Eq. (9)}}{=} \left\langle \begin{pmatrix} \mathbf{w} - \mathbf{w}_0 \\ \sigma - \sigma_0 \end{pmatrix}, T \begin{pmatrix} \mathbf{w} \\ \sigma \end{pmatrix} \right\rangle_{\mathcal{H}} \\ & \geq \left\langle \begin{pmatrix} \mathbf{w} - \mathbf{w}_0 \\ \sigma - \sigma_0 \end{pmatrix}, T \begin{pmatrix} \mathbf{w}_0 \\ \sigma_0 \end{pmatrix} \right\rangle_{\mathcal{H}} + \lambda (\|\mathbf{w} - \mathbf{w}_0\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d) - \sigma_0(\mathbf{c}_d)\|_F^2]) \end{aligned}$$

Let $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$ satisfy $\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2] \rightarrow +\infty$ and $(\mathbf{w}_0, \sigma_0) \in \mathcal{W} \times \Sigma$ be fixed, we have

$$\frac{\left\langle \begin{pmatrix} \mathbf{w} - \mathbf{w}_0 \\ \sigma - \sigma_0 \end{pmatrix}, T \begin{pmatrix} \mathbf{w}_0 \\ \sigma_0 \end{pmatrix} \right\rangle_{\mathcal{H}}}{\sqrt{\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2]}} \geq -C, \quad \text{for some universal constant } C > 0$$

and

$$\frac{\|\mathbf{w} - \mathbf{w}_0\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d) - \sigma_0(\mathbf{c}_d)\|_F^2]}{\sqrt{\|\mathbf{w}\|^2 + \mathbb{E}_q[\|\sigma(\mathbf{c}_d)\|_F^2]}} \rightarrow +\infty$$

This implies that Eq. (8) holds true, which completes the proof.

3.6 Proof of Corollary 4

The proof is nearly the same as that of Corollary 3 and we need to verify the assumptions in Theorem 2. Note that T defined by Eq. (10) is strictly monotone and $\mathcal{W} \times \Sigma$ is compact. Thus, by using Theorem 2, we conclude the desired result.

4 Two Special Settings

In this section, we consider two special settings of a finite Bayesian prior and a quadratic adversarial loss. This is primarily for practical purposes, for which the computation of a Bayesian equilibrium reduces to solving a high-dimensional VI or solving a nonconvex stochastic optimization problem, both in Euclidean space.

Case I: Finite Bayesian prior. Let $K > 0$ be an integer, assume that the Bayesian prior q is a distribution with support $\{\mathbf{v}_1, \dots, \mathbf{v}_K\}$; i.e., $q(\mathbf{c}_d = \mathbf{v}_k) = p_k > 0$ for all $k \in [K]$, $\sum_{i=1}^K p_k = 1$. Then, the VI in Eq. (3) becomes

$$\sum_{k=1}^K p_k [\langle \mathbf{w} - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{v}_k), \mathbf{c}_l) \rangle + \langle \sigma(\mathbf{v}_k) - \sigma_*(\mathbf{v}_k), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{v}_k), \mathbf{v}_k) \rangle] \geq 0 \quad (11)$$

for all $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$. By definition, $\sigma, \sigma_\star \in \Sigma$ are both mappings from \mathbb{R}^n to \mathcal{X} in an infinite-dimensional space. When the Bayesian prior is finite, Eq. (12) implies that σ can be fully represented by $(\sigma(\mathbf{v}_1), \sigma(\mathbf{v}_2), \dots, \sigma(\mathbf{v}_K))$ where the range $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_K\}$ is finite and known. For simplicity, define $\sigma^k = \sigma(\mathbf{v}_k)$ and $\sigma_\star^k = \sigma_\star(\mathbf{v}_k)$ for all $k \in [K]$. Then, the VI in Eq. (11) can be reformulated as follows:

$$\sum_{k=1}^K p_k [\langle \mathbf{w} - \mathbf{w}_\star, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_\star, \sigma_\star^k, \mathbf{c}_l) \rangle + \langle \sigma^k - \sigma_\star^k, \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_\star, \sigma_\star^k, \mathbf{v}_k) \rangle] \geq 0 \quad (12)$$

for all $\mathbf{w} \in \mathcal{W}$ and $\sigma^k \in \mathcal{X}$ for all $k \in [K]$. Note that the VI in Eq. (12) is a high-dimensional VI in Euclidean space. Indeed, $\mathbf{w} \in \mathbb{R}^m$ and $\sigma^k \in \mathbb{R}^{n \times m}$ for all $k \in [K]$, implying that the total number of unknown variables is $m + mnK$.

Case II: Quadratic adversarial loss. Consider $f_d(\mathbf{w}, \mathbf{x}, z) = (\mathbf{x}^\top \mathbf{w} - z)^2$, $\Omega_d(X, \bar{X}) = \|X - \bar{X}\|_F^2$, where $\mathbf{x} \in \mathbb{R}^m$, and $X, \bar{X} \in \mathcal{X} = \mathbb{R}^{n \times m}$, we have $\sigma^\star[\mathbf{w}](\mathbf{c}_d) = \operatorname{argmin}_{\bar{X} \in \mathbb{R}^{n \times m}} \sum_{i=1}^n c_{d,i} (\bar{\mathbf{x}}_i^\top \mathbf{w} - z_i)^2 + \|X - \bar{X}\|_F^2$. By [Großhans et al.(2013), Lemma 1], $\sigma^\star[\mathbf{w}](\mathbf{c}_d) = X - (\operatorname{diag}(\mathbf{c}_d)^{-1} + \|\mathbf{w}\|^2 I_n)^{-1} (X\mathbf{w} - \mathbf{z})\mathbf{w}^\top$. Equivalently, we have

$$[\sigma^\star[\mathbf{w}](\mathbf{c}_d)]_i = \mathbf{x}_i - \frac{c_{d,i}(\mathbf{x}_i^\top \mathbf{w} - z_i)}{1 + \|\mathbf{w}\|^2 c_{d,i}} \mathbf{w} \in \mathbb{R}^m \quad (13)$$

Putting these pieces together with the best response of the learner yields the following stochastic optimization problem:

$$\min_{\mathbf{w} \in \mathcal{W}} \mathbb{E}_q \left[\sum_{i=1}^n c_{l,i} f_l \left(\mathbf{w}, \mathbf{x}_i - \frac{c_{d,i}(\mathbf{x}_i^\top \mathbf{w} - z_i)}{1 + \|\mathbf{w}\|^2 c_{d,i}} \mathbf{w}, y_i \right) \right] + \Omega_l(\mathbf{w}) \quad (14)$$

The computation of a Bayesian equilibrium thus reduces to the solution of a nonconvex stochastic optimization problem. Standard gradient descent approaches can not be applied for solving the optimization problem in Eq. (14) since the integration with respect to a Bayesian prior q does not have a closed-form expression in general. Nonetheless, the Bayesian prior q is known and accessible through drawing samples, a stochastic-gradient-based algorithm can be applied and particular examples include stochastic gradient descent (SGD) [Robbins and Monro(1951), Bottou(1998)] or its adaptive variants AdaGrad and ADAM [Duchi et al.(2011), Kingma and Ba(2015)].

5 Algorithms and Their Convergence

In this section, we propose specific algorithms and analyze their convergence and iteration complexity.

5.1 Projected Reflected Gradient with Inertial Extrapolation

We present the *projected reflected gradient with inertial extrapolation* (PRG-IE) algorithm for solving Eq. (3). We make the following assumption throughout this section and we discuss the intuitions of them afterwards.

Algorithm 1 Projected Reflected Gradient with Inertial Extrapolation (PRG-IE)

- 1: **Input:** smoothness parameter $L > 0$; Bayesian prior q ; weight $\mathbf{c}_l \in \mathbb{R}^n$
- 2: **Initialize:** $\tilde{\mathbf{w}}_0, \tilde{\mathbf{w}}_1, \mathbf{w}_1 \in \mathcal{W}$, $\tilde{\sigma}_0, \tilde{\sigma}_1, \sigma_1 \in \Sigma$ and $0 < \gamma < \min\{1, \frac{1}{100L}\}$
- 3: **for** $t = 1, 2, \dots, T$ **do**
- 4: Compute $\delta_t \leftarrow 1/t$
- 5: Compute $(\tilde{\mathbf{w}}_{t+1}, \tilde{\sigma}_{t+1})$ by

$$\begin{aligned}\tilde{\mathbf{w}}_{t+1} &\leftarrow P_{\mathcal{W}}(\mathbf{w}_t - \gamma \nabla_{\mathbf{w}} \hat{\theta}_l(2\tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_{t-1}, 2\tilde{\sigma}_t - \tilde{\sigma}_{t-1}, \mathbf{c}_l)) \\ \tilde{\sigma}_{t+1} &\leftarrow P_{\Sigma}(\sigma_t - \gamma \nabla_{\bar{X}} \hat{\theta}_d(2\tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_{t-1}, 2\tilde{\sigma}_t - \tilde{\sigma}_{t-1}, \mathbf{c}_d))\end{aligned}$$

- 6: Compute $(\mathbf{w}_{t+1}, \sigma_{t+1})$ by

$$\begin{pmatrix} \mathbf{w}_{t+1} \\ \sigma_{t+1} \end{pmatrix} \leftarrow \frac{\delta_t}{2} \begin{pmatrix} \mathbf{w}_t \\ \sigma_t \end{pmatrix} + (1 - \delta_t) \begin{pmatrix} \tilde{\mathbf{w}}_{t+1} \\ \tilde{\sigma}_{t+1} \end{pmatrix}$$

- 7: **end for**
-

Assumption 2. *The Bayesian regression game G satisfies the following: (i) \mathcal{W} and \mathcal{X} are both compact with $\mathbf{0}_m \in \mathcal{W}$ and $\mathbf{0}_{n \times m} \in \mathcal{X}$; (ii) Given that $\mathbf{c}_l \in \mathbb{R}^n$ is fixed, we have*

$$\begin{aligned}&\mathbb{E}_q[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_l) - \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}', \sigma'(\mathbf{c}_d), \mathbf{c}_l), \mathbf{w} - \mathbf{w}' \rangle \\ &+ \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_d) - \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}', \sigma'(\mathbf{c}_d), \mathbf{c}_d), \sigma(\mathbf{c}_d) - \sigma'(\mathbf{c}_d) \rangle] \\ &> 0, \quad \forall (\mathbf{w}, \sigma) \neq (\mathbf{w}', \sigma')\end{aligned}$$

(iii) *There exists a constant $L > 0$ such that $\nabla_{\mathbf{w}} \hat{\theta}_l(\cdot, \cdot, \mathbf{c}_l) : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}^m$ is L -Lipschitz for each $\mathbf{c}_l \in \mathbb{R}^n$ and $\nabla_{\bar{X}} \hat{\theta}_d(\cdot, \cdot, \mathbf{c}_d) : \mathcal{W} \times \mathcal{X} \rightarrow \mathbb{R}^{n \times m}$ is L -Lipschitz for each $\mathbf{c}_d \in \mathbb{R}^n$, i.e. for each $(\mathbf{w}, \bar{X}), (\mathbf{w}', \bar{X}') \in \mathcal{W} \times \mathcal{X}$, we have*

$$\begin{aligned}\|\nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \bar{X}, \mathbf{c}_l) - \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}', \bar{X}', \mathbf{c}_l)\| &\leq L(\|\mathbf{w} - \mathbf{w}'\| + \|\bar{X} - \bar{X}'\|_F) \\ \|\nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \bar{X}, \mathbf{c}_d) - \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}', \bar{X}', \mathbf{c}_d)\|_F &\leq L(\|\mathbf{w} - \mathbf{w}'\| + \|\bar{X} - \bar{X}'\|_F)\end{aligned}$$

Assumption 2 is standard in optimization and game theory. The second condition can be interpreted as the strict monotonicity of T defined by Eq. (7) in terms of the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$. The third condition imposes a weak Lipschitz condition on $\nabla_{\mathbf{w}} \hat{\theta}_l(\cdot, \cdot, \mathbf{c}_l)$ and $\nabla_{\bar{X}} \hat{\theta}_d(\cdot, \cdot, \mathbf{c}_d)$. The third condition is necessary for the existence of at least one Bayesian equilibrium and the feasibility of the sequence $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ generated by Algorithm 1. Note that $\mathbf{0}_m \in \mathcal{W}$ and $\mathbf{0}_{n \times m} \in \mathcal{X}$ are not restrictive since \mathcal{W} and \mathcal{X} are commonly taken to be the ball of a norm.

This approach combines Malitsky's projected reflected gradient algorithm [Malitsky(2015)] with Halpern-type inertial extrapolation [Halpern(1967)]. Such an integration has the following advantages: (i) it only performs one projection at each iteration; (ii) it achieves the strong convergence to one Bayesian equilibrium thanks to the extrapolation step. At each iteration, the projections $P_{\mathcal{W}}(\cdot)$ and $P_{\Sigma}(\cdot)$ are required and these operations are based on the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$, which necessitates the expectation with respect to a Bayesian prior q . In applications, variational inference or Markov chain Monte Carlo can be used to approximate this expectation.

Theorem 3. *Under Assumption 1 and 2, the sequence $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 0} \in \mathcal{W} \times \Sigma$ generated by Algorithm 1 satisfies $\|\mathbf{w}_t - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2] \rightarrow 0$, where the point $(\mathbf{w}_\star, \sigma_\star) \in \mathcal{W} \times \Sigma$ is a unique Bayesian equilibrium.*

We make some comments in the sequel. First, the sequence $\delta_t = 1/t$ in Algorithm 1 is one specific choice and more general choices can be considered, as long as they satisfy $\delta_t \rightarrow 0$ and $\sum_{t=1}^{+\infty} \delta_t = +\infty$. The formula for updating the sequence $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ can also be generalized as follows:

$$\begin{pmatrix} \mathbf{w}_{t+1} \\ \sigma_{t+1} \end{pmatrix} \leftarrow \delta_t \cdot g \begin{pmatrix} \mathbf{w}_t \\ \sigma_t \end{pmatrix} + (1 - \delta_t) \begin{pmatrix} \tilde{\mathbf{w}}_{t+1} \\ \tilde{\sigma}_{t+1} \end{pmatrix}$$

where $g : \mathcal{H} \mapsto \mathcal{H}$ is a contraction with the parameter $\kappa \in (0, 1)$. We set $g(\mathbf{x}) = \mathbf{x}/2$ for simplicity. Second, the strict monotonicity can be relaxed to monotonicity and it can be shown that the sequence $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ strongly converges to a particular Bayesian equilibrium.

The rest of this section devotes to proving Theorem 3. We start by reviewing two preliminary results in the literature which are established as [Saejung and Yotkaew(2012), Lemma 2.6] and the Minty's lemma in [Bauschke and Combettes(2011)] respectively.

Lemma 1. *Let $\{s_t\}_{t \geq 0}$ be a sequence of nonnegative real numbers, $\{a_t\}_{t \geq 0}$ be a sequence in $(0, 1)$ such that $\sum_{t=0}^{+\infty} a_t = +\infty$ and $\{b_t\}_{t \geq 0}$ be a sequence of real numbers. Suppose that $s_{t+1} \leq (1 - a_t)s_t + a_t b_t$ for all $t \geq 0$. If $\limsup_{j \rightarrow +\infty} b_{t_j} \leq 0$ for every subsequence $\{s_{t_j}\}_{j \geq 0}$ of $\{s_t\}_{t \geq 0}$ satisfying that $\liminf_{j \rightarrow +\infty} (s_{t_j+1} - s_{t_j}) \geq 0$, then $s_t \rightarrow 0$ as $t \rightarrow +\infty$.*

Lemma 2 (Minty). *Let $T : \mathcal{H} \rightarrow \mathcal{H}$ be a continuous and monotone mapping on a closed and convex subset \mathcal{S} . Then $\hat{\mathbf{x}}$ is a solution of the VI in Eq. (21) if and only if $\langle \mathbf{x} - \hat{\mathbf{x}}, T(\mathbf{x}) \rangle_{\mathcal{H}} \geq 0$ for all $\mathbf{x} \in \mathcal{S}$.*

First, we show that the sequences $\{(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t)\}_{t \geq 1}$ and $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ generated by Algorithm 1 are both bounded in the following lemma.

Lemma 3. *Under Assumption 1 and 2, there exists a constant $M > 0$ such that the sequences $\{(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t)\}_{t \geq 0}$ and $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ generated by Algorithm 1 satisfies that*

$$\|\tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_t(\mathbf{c}_d)\|_F^2] \leq M \quad \text{and} \quad \|\mathbf{w}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d)\|_F^2] \leq M \quad \text{for all } t \geq 1$$

We then present an important descent inequality for the sequences $\{(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t)\}_{t \geq 0}$ and $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ generated by Algorithm 1. Note that our result can not be derived from [Malitsky(2015), Lemma 3.1] due to the Halpern-type inertial extrapolation.

Lemma 4. *Under Assumption 1 and 2, the sequences $\{(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t)\}_{t \geq 0}$ and $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ generated by Algorithm 1 satisfies that*

$$\begin{aligned} & \|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2] \\ & \leq \|\mathbf{w}_t - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2] \\ & \quad - \left(\frac{1}{2} - 2\gamma L \right) (\|\tilde{\mathbf{w}}_{t+1} - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2]) \\ & \quad + 6\gamma L (\|\tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_t(\mathbf{c}_d) - \tilde{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2]) \\ & \quad + (4 + 6\gamma L) (\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2]) \\ & \quad + 4 (\|\mathbf{w}_{t-1} - \tilde{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\sigma_{t-1}(\mathbf{c}_d) - \tilde{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2]) \\ & \quad - 4\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_t - \mathbf{w}_\star \rangle + \langle \nabla_{\tilde{\mathbf{x}}} \hat{\theta}_d(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle \right] \\ & \quad + 2\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t-1} - \mathbf{w}_\star \rangle + \langle \nabla_{\tilde{\mathbf{x}}} \hat{\theta}_d(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t-1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle \right] \\ & \quad - (1 - 6\gamma L) (\|\tilde{\mathbf{w}}_t - \mathbf{w}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_t(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2]) \end{aligned}$$

where the point $(\mathbf{w}_\star, \sigma_\star) \in \mathcal{W} \times \Sigma$ is a unique Bayesian equilibrium, and the auxiliary sequence $(\bar{\mathbf{w}}_t, \bar{\sigma}_t)$ is defined by $(\bar{\mathbf{w}}_t, \bar{\sigma}_t) = (2\tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_{t-1}, 2\tilde{\sigma}_t - \tilde{\sigma}_{t-1})$.

Based on the boundedness results in Lemma 3 and the descent inequality in Lemma 4, we provide two additional inequalities for the sequences $\{(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t)\}_{t \geq 0}$ and $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ generated by Algorithm 1 in the following lemma.

Lemma 5. *Under Assumption 1 and 2, the sequences $\{(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t)\}_{t \geq 0}$ and $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ generated by Algorithm 1 satisfies that*

$$r_{t+1} \leq r_t + \delta_{t-1}M_1 - \left(\frac{1}{2} - 8\gamma L\right) (\|\tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d)\|_F^2]) \quad (15)$$

$$- (1 - 6\gamma L) (\|\bar{\mathbf{w}}_t - \mathbf{w}_t\|^2 + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2])$$

$$r_{t+1} \leq \left(1 - \frac{3\delta_t}{4}\right) r_t + \frac{3\delta_t}{4} (4\delta_t M_1 - 2\langle \mathbf{w}_\star, \mathbf{w}_{t+1} - \mathbf{w}_\star \rangle - 2\mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]) \quad (16)$$

for all $t \geq 2$ and some constant $M_1 > 0$. The residue sequence $\{r_t\}_{t \geq 0}$ is defined as:

$$r_t = \|\mathbf{w}_t - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2] + 6\gamma L (\|\tilde{\mathbf{w}}_t - \bar{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_t(\mathbf{c}_d) - \bar{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2])$$

$$+ 2\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{l}(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t-1} - \mathbf{w}_\star \rangle + \langle \nabla_{\tilde{\sigma}} \hat{l}(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t-1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle \right]$$

where $(\mathbf{w}_\star, \sigma_\star) \in \mathcal{W} \times \Sigma$ is a unique Bayesian equilibrium, and the auxiliary sequence $(\bar{\mathbf{w}}_t, \bar{\sigma}_t)$ is defined in Lemma 4.

We are ready for the

Proof of Theorem 3. It suffices to prove that $r_t \rightarrow 0$ as $t \rightarrow +\infty$. Suppose that $\{r_{t_j}\}_{j \geq 0}$ is any of the subsequences of the whole sequence $\{r_t\}_{t \geq 0}$ and satisfies that $\liminf_{j \rightarrow +\infty} (r_{t_{j+1}} - r_{t_j}) \geq 0$. From Eq. (15) in Lemma 5, we have

$$\limsup_{j \rightarrow +\infty} \left[\left(\frac{1}{2} - 8\gamma L\right) (\|\tilde{\mathbf{w}}_{t_j+1} - \bar{\mathbf{w}}_{t_j}\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t_j+1}(\mathbf{c}_d) - \bar{\sigma}_{t_j}(\mathbf{c}_d)\|_F^2]) \right. \\ \left. + (1 - 6\gamma L) (\|\bar{\mathbf{w}}_{t_j} - \mathbf{w}_{t_j}\|^2 + \mathbb{E}_q[\|\bar{\sigma}_{t_j}(\mathbf{c}_d) - \sigma_{t_j}(\mathbf{c}_d)\|_F^2]) \right] \\ \leq \limsup_{j \rightarrow +\infty} (r_{t_j} - r_{t_j+1} + \delta_{t_j}M_1) \leq \limsup_{j \rightarrow +\infty} (r_{t_j} - r_{t_j+1}) + \limsup_{j \rightarrow +\infty} \delta_{t_j}M_1 \\ \leq -\liminf_{j \rightarrow +\infty} (r_{t_j+1} - r_{t_j}) + \limsup_{j \rightarrow +\infty} \delta_{t_j}M_1 \leq \limsup_{j \rightarrow +\infty} \delta_{t_j}M_1$$

Since $\delta_t \rightarrow 0$ and $0 < \gamma < \min\{1, \frac{1}{100L}\}$, we have

$$\limsup_{j \rightarrow +\infty} (\|\tilde{\mathbf{w}}_{t_j+1} - \bar{\mathbf{w}}_{t_j}\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t_j+1}(\mathbf{c}_d) - \bar{\sigma}_{t_j}(\mathbf{c}_d)\|_F^2]) = 0$$

$$\limsup_{j \rightarrow +\infty} (\|\bar{\mathbf{w}}_{t_j} - \mathbf{w}_{t_j}\|^2 + \mathbb{E}_q[\|\bar{\sigma}_{t_j}(\mathbf{c}_d) - \sigma_{t_j}(\mathbf{c}_d)\|_F^2]) = 0$$

By Lemma 3, the sequence $\{(\tilde{\mathbf{w}}_{t_j+1}, \tilde{\sigma}_{t_j+1})\}_{j \geq 0}$ is bounded. This implies that there exists a subsequence $\{(\tilde{\mathbf{w}}_{t_{j_i}+1}, \tilde{\sigma}_{t_{j_i}+1})\}_{i \geq 0}$ such that $(\tilde{\mathbf{w}}_{t_{j_i}+1}, \tilde{\sigma}_{t_{j_i}+1})$ weakly converges to some point $(\tilde{\mathbf{w}}, \tilde{\sigma}) \in \mathcal{W} \times \Sigma$, and

$$\limsup_{j \rightarrow +\infty} (-2\langle \mathbf{w}_\star, \mathbf{w}_{t_j+1} - \mathbf{w}_\star \rangle - 2\mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t_j+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]) \quad (17)$$

$$= \lim_{i \rightarrow +\infty} \left(-2\langle \mathbf{w}_\star, \mathbf{w}_{t_{j_i}+1} - \mathbf{w}_\star \rangle - 2\mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t_{j_i}+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle] \right)$$

$$= -2\langle \mathbf{w}_\star, \tilde{\mathbf{w}} - \mathbf{w}_\star \rangle - 2\mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \tilde{\sigma}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]$$

It is also clear that $(\bar{\mathbf{w}}_{t_{j_i}}, \bar{\sigma}_{t_{j_i}})$ and $(\mathbf{w}_{t_{j_i}}, \sigma_{t_{j_i}})$ both weakly converge to $(\tilde{\mathbf{w}}, \tilde{\sigma})$. Recall that the optimality condition of updating $(\tilde{\mathbf{w}}_{t+1}, \tilde{\sigma}_{t+1})$ in Eq. (22) is:

$$0 \leq \mathbb{E}_q \left[\langle \mathbf{w} - \tilde{\mathbf{w}}_{t+1}, \tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t + \gamma \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l) \rangle \quad \text{for each } (\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma \right. \\ \left. + \langle \sigma(\mathbf{c}_d) - \tilde{\sigma}_{t+1}(\mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d) + \gamma \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d) \rangle \right]$$

Equivalently, we have

$$0 \leq \mathbb{E}_q [\langle \mathbf{w} - \tilde{\mathbf{w}}_{t+1}, \tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t \rangle + \langle \sigma(\mathbf{c}_d) - \tilde{\sigma}_{t+1}(\mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d) \rangle] \quad \text{for each } (\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma \\ + \gamma \mathbb{E}_q \left[\langle \mathbf{w} - \bar{\mathbf{w}}_t, \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l) \rangle + \langle \sigma(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d) \rangle \right] \\ + \gamma \mathbb{E}_q \left[\langle \bar{\mathbf{w}}_t - \tilde{\mathbf{w}}_{t+1}, \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l) \rangle + \langle \bar{\sigma}_t(\mathbf{c}_d) - \tilde{\sigma}_{t+1}(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d) \rangle \right]$$

Using the first condition in Assumption 2 with $(\mathbf{w}', \sigma') = (\bar{\mathbf{w}}_t, \bar{\sigma}_t)$ and the fact that $\gamma > 0$, we have

$$\gamma \mathbb{E}_q \left[\langle \mathbf{w} - \bar{\mathbf{w}}_t, \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l) \rangle + \langle \sigma(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d) \rangle \right] \\ \leq \gamma \mathbb{E}_q \left[\langle \mathbf{w} - \bar{\mathbf{w}}_t, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_t, \sigma_t(\mathbf{c}_d), \mathbf{c}_l) \rangle + \langle \sigma(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_t, \sigma_t(\mathbf{c}_d), \mathbf{c}_d) \rangle \right]$$

Putting these two inequalities together with $t = t_{j_i}$ yields that, for each $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$,

$$0 \leq \mathbb{E}_q \left[\langle \mathbf{w} - \tilde{\mathbf{w}}_{t_{j_i}+1}, \tilde{\mathbf{w}}_{t_{j_i}+1} - \mathbf{w}_{t_{j_i}} \rangle + \langle \sigma(\mathbf{c}_d) - \tilde{\sigma}_{t_{j_i}+1}(\mathbf{c}_d), \tilde{\sigma}_{t_{j_i}+1}(\mathbf{c}_d) - \sigma_{t_{j_i}}(\mathbf{c}_d) \rangle \right] \\ + \gamma \mathbb{E}_q \left[\langle \mathbf{w} - \bar{\mathbf{w}}_{t_{j_i}}, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_l) \rangle + \langle \sigma(\mathbf{c}_d) - \bar{\sigma}_{t_{j_i}}(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_d) \rangle \right] \\ + \gamma \mathbb{E}_q \left[\langle \bar{\mathbf{w}}_{t_{j_i}} - \tilde{\mathbf{w}}_{t_{j_i}+1}, \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_{t_{j_i}}, \bar{\sigma}_{t_{j_i}}(\mathbf{c}_d), \mathbf{c}_l) \rangle + \langle \bar{\sigma}_{t_{j_i}}(\mathbf{c}_d) - \tilde{\sigma}_{t_{j_i}+1}(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_{t_{j_i}}, \bar{\sigma}_{t_{j_i}}(\mathbf{c}_d), \mathbf{c}_d) \rangle \right]$$

Letting $i \rightarrow +\infty$ in the above inequality, for each $(\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma$, we have

$$\mathbb{E}_q \left[\left\langle \mathbf{w} - \tilde{\mathbf{w}}, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_l) \right\rangle + \left\langle \sigma(\mathbf{c}_d) - \tilde{\sigma}(\mathbf{c}_d), \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma(\mathbf{c}_d), \mathbf{c}_d) \right\rangle \right] \geq 0 \quad (18)$$

Using Lemma 2 and Eq. (18), the point $(\tilde{\mathbf{w}}, \tilde{\sigma})$ is the solution of the VI in Eq. (3). Under Assumption 1 and 2, the VI in Eq. (3) has a unique solution. Thus, $(\tilde{\mathbf{w}}, \tilde{\sigma}) = (\mathbf{w}_\star, \sigma_\star)$ is a unique Bayesian equilibrium.

Finally, we consider Lemma 1 with $s_t = r_t$, $a_t = 3\delta_t/4$ and

$$b_t = 4\delta_t M_1 - 2\langle \mathbf{w}_\star, \mathbf{w}_{t+1} - \mathbf{w}_\star \rangle - 2\mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]$$

More specifically, we have (i) the sequence $\{s_t\}_{t \geq 0}$ is nonnegative and $\{a_t\}_{t \geq 0}$ is a sequence in $(0, 1)$ satisfying $\sum_{t=0}^{+\infty} a_t = +\infty$; (ii) Eq. (16) implies that $s_{t+1} \leq (1 - a_t)s_t + a_t b_t$ for all $t \geq 2$; (iii) for every subsequence $\{s_{t_j}\}_{j \geq 0}$ of $\{s_t\}_{t \geq 0}$ satisfying that $\liminf_{j \rightarrow +\infty} (s_{t_{j+1}} - s_{t_j}) \geq 0$, Eq. (17) and the fact that $\delta_t \rightarrow 0$ implies that

$$\limsup_{j \rightarrow +\infty} b_{t_j} = \limsup_{j \rightarrow +\infty} (\delta_{t_j} M_1 - 2\langle \mathbf{w}_\star, \mathbf{w}_{t_j+1} - \mathbf{w}_\star \rangle - 2\mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t_j+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]) = 0$$

Therefore, we conclude that $r_t \rightarrow 0$ as $t \rightarrow +\infty$. This completes the proof. \square

Algorithm 2 Projected Gradient with Randomized Block Coordinate (PG-RBC)

- 1: **Input:** strongly monotone parameter $\lambda > 0$; finite Bayesian prior $\{p_k, \mathbf{v}_k\}_{k=1}^K$; weight $\mathbf{c}_l \in \mathbb{R}^n$
- 2: **Initialize:** $\mathbf{w}_0 \in \mathcal{W}$, $\sigma_0^k \in \mathcal{X}$ for all $k \in [K]$ and $\gamma_0 > \frac{1}{2\lambda}$
- 3: **for** $t = 0, 1, 2, \dots, T - 1$ **do**
- 4: Randomly pick up an index $j_t \in [K]$ according to $\mathbb{P}(j_t = k) = p_k$ for all $k \in [K]$
- 5: Compute $(\mathbf{w}_{t+1}, \sigma_{t+1}^1, \dots, \sigma_{t+1}^K)$ by

$$\begin{aligned} \mathbf{w}_{t+1} &\leftarrow P_{\mathcal{W}}(\mathbf{w}_t - \gamma_t \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_t, \sigma_t^{j_t}, \mathbf{c}_l)) \\ \sigma_{t+1}^k &\leftarrow \begin{cases} P_{\mathcal{X}}(\sigma_t^k - \gamma_t \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_t, \sigma_t^k, \mathbf{v}_k)), & \text{if } k = j_t \\ \sigma_t^k, & \text{otherwise} \end{cases} \end{aligned}$$

- 6: Compute $\gamma_{t+1} \leftarrow \gamma/(t+1)$
 - 7: **end for**
-

5.2 Projected Gradient with Randomized Block Coordinate

We present a *projected gradient with randomized block coordinate* (PG-RBC) algorithm for solving Eq. (12). To ease the analysis, we make the following assumption throughout this section before providing more intuitions for them.

Assumption 3. *The Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ satisfies: (i) \mathcal{W} and \mathcal{X} are both compact such that there exists a positive constant G such that $\|\nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \bar{X}, \mathbf{c}_l)\| \leq G$ for each $(\mathbf{w}, \bar{X}) \in \mathcal{W} \times \mathcal{X}$ and each $\mathbf{c}_l \in \mathbb{R}^n$ fixed, and $\|\nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \bar{X}, \mathbf{v})\|_F \leq G$ for each $(\mathbf{w}, \bar{X}) \in \mathcal{W} \times \mathcal{X}$ and each $\mathbf{v} \in \{\mathbf{v}_1, \dots, \mathbf{v}_K\}$; (ii) Given a fixed \mathbf{c}_l , there exists a positive constant λ such that¹*

$$\begin{aligned} &\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \bar{X}, \mathbf{c}_l) - \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}', \bar{X}', \mathbf{c}_l), \mathbf{w} - \mathbf{w}' \rangle + \sum_{k=1}^K p_k \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \bar{X}, \mathbf{v}_k) - \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}', \bar{X}', \mathbf{v}_k), \bar{X} - \bar{X}' \rangle \\ &\geq \lambda (\|\mathbf{w} - \mathbf{w}'\|^2 + \|\bar{X} - \bar{X}'\|_F^2), \quad \forall (\mathbf{w}, \bar{X}), (\mathbf{w}', \bar{X}') \in \mathcal{W} \times \mathcal{X} \end{aligned}$$

In Assumption 3, the first condition can be interpreted as the strong monotonicity of T defined by Eq. (7) when the Bayesian prior q is finite. The second condition naturally holds true if $\nabla_{\mathbf{w}} \hat{\theta}_l(\cdot, \cdot, \mathbf{c}_l)$ and $\nabla_{\bar{X}} \hat{\theta}_d(\cdot, \cdot, \mathbf{v})$ are continuous for each $\mathbf{c}_l \in \mathbb{R}^n$ fixed and each $\mathbf{v} \in \{\mathbf{v}_1, \dots, \mathbf{v}_K\}$.

The proposed approach combines projected gradient algorithm with randomized block coordinate update [Nesterov(2012), Wright(2015)]. This design is more efficient than deterministic VI algorithms since the per iteration cost is $O(nm)$ which does not depend on K . Thus, our approach is favorable in application problems when the parameter K is large.

Theorem 4. *Under Assumption 1 and 3, the iterates $\{(\mathbf{w}_t, \sigma_t^1, \dots, \sigma_t^K)\}_{t \geq 0}$ generated by Algorithm 2 satisfy $\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_\star\|^2 + \sum_{k=1}^K \|\sigma_t^k - \sigma_\star^k\|_F^2] = O(1/t)$ where $(\mathbf{w}_\star, \sigma_\star^1, \dots, \sigma_\star^K)$ is a unique Bayesian equilibrium.*

We make some further comments. Since we do not assume any smoothness condition in Assumption 3, the iteration complexity of $O(1/t)$ is the best possible we can hope for all the deterministic

¹The constant λ can depend on \mathbf{c}_l but is independent of the choice of (\mathbf{w}, \bar{X}) and (\mathbf{w}', \bar{X}') .

and stochastic algorithms in general²; see [Nemirovsky(1983)] for the reference. To this end, Theorem 4 demonstrates that the complexity bound of Algorithm 2 is tight in terms of the iteration number.

The rest of this section devotes to proving Theorem 4. We start by reviewing one preliminary result in the literature which is a fact of sequences first established in [Chung(1954)] (although it does not appear to be widely known).

Lemma 6 ([Chung(1954)]). *Let $\{a_t\}_{t \geq 0}$ be a non-negative sequence such that*

$$a_{t+1} \leq \left(1 - \frac{P}{t^p}\right) a_t + \frac{Q}{t^{p+q}}$$

where $P > q > 0$, $0 < p \leq 1$ and $Q > 0$. Then:

$$a_t \leq \begin{cases} \frac{Q}{P} \frac{1}{t^q}, & \text{if } 0 < p < 1 \\ \frac{Q}{P-q} \frac{1}{t^q}, & \text{if } p = 1 \end{cases}$$

We are ready for the

Proof of Theorem 4. By Corollary 3, the Bayesian regression game $G = (\mathcal{W}, \Sigma, \hat{\theta}_l, \hat{\theta}_d, \mathbf{c}_l, q)$ has a unique Bayesian equilibrium $(\mathbf{w}_*, \sigma_*^1, \dots, \sigma_*^K)$ under Assumption 1 and 3. Define $E_t = \|\mathbf{w}_{t+1} - \mathbf{w}_*\|^2 + \sum_{k=1}^K \|\sigma_{t+1}^k - \sigma_*^k\|_F^2$, we derive from the update formula in Algorithm 2 and the fact that two orthogonal projection mappings $P_{\mathcal{W}}$ and $P_{\mathcal{X}}$ are nonexpansive that

$$\begin{aligned} & \mathbb{E}[E_{t+1} \mid (\mathbf{w}_t, \sigma_t^1, \dots, \sigma_t^K)] \\ &= \sum_{k=1}^K p_k \left[\|P_{\mathcal{W}}(\mathbf{w}_t - \gamma_t \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_t, \sigma_t^k, \mathbf{c}_l)) - \mathbf{w}_*\|^2 + \|P_{\mathcal{X}}(\sigma_t^k - \gamma_t \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_t, \sigma_t^k, \mathbf{v}_k)) - \sigma_*^k\|_F^2 \right] \\ & \quad + \sum_{k=1}^K p_k \left[\sum_{j \neq k} \|\sigma_t^j - \sigma_*^j\|_F^2 \right] \\ & \leq \sum_{k=1}^K p_k \left[\|\mathbf{w}_t - \gamma_t \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_t, \sigma_t^k, \mathbf{c}_l) - \mathbf{w}_*\|^2 + \|\sigma_t^k - \gamma_t \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_t, \sigma_t^k, \mathbf{v}_k) - \sigma_*^k\|_F^2 + \sum_{j \neq k} \|\sigma_t^j - \sigma_*^j\|_F^2 \right] \end{aligned}$$

Using the second condition in Assumption 3, we have

$$\begin{aligned} \|\mathbf{w}_t - \gamma_t \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_t, \sigma_t^k, \mathbf{c}_l) - \mathbf{w}_*\|^2 & \leq \|\mathbf{w}_t - \mathbf{w}_*\|^2 - 2\gamma_t \langle \mathbf{w}_t - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_t, \sigma_t^k, \mathbf{c}_l) \rangle + \gamma_t^2 G^2 \\ \|\sigma_t^k - \gamma_t \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_t, \sigma_t^k, \mathbf{v}_k) - \sigma_*^k\|_F^2 & \leq \|\sigma_t^k - \sigma_*^k\|_F^2 - 2\gamma_t \langle \sigma_t^k - \sigma_*^k, \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_t, \sigma_t^k, \mathbf{v}_k) \rangle + \gamma_t^2 G^2 \end{aligned}$$

Putting these pieces together with the fact that $\sum_{k=1}^K p_k = 1$ yields that

$$\begin{aligned} & \mathbb{E}[E_{t+1} \mid (\mathbf{w}_t, \sigma_t^1, \dots, \sigma_t^K)] \\ & \leq E_t + 2\gamma_t^2 G^2 - 2\gamma_t \sum_{k=1}^K p_k \left[\langle \mathbf{w}_t - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_t, \sigma_t^k, \mathbf{c}_l) \rangle + \langle \sigma_t^k - \sigma_*^k, \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_t, \sigma_t^k, \mathbf{v}_k) \rangle \right] \end{aligned} \quad (19)$$

²Convex optimization problems are a special class of the monotone VI problems. Thus, the problem complexity of (strongly) convex optimization implies that of (strongly) monotone VIs.

Since the point $(\mathbf{w}_*, \sigma_*^1, \dots, \sigma_*^K)$ is a Bayesian equilibrium, we have

$$\sum_{k=1}^K p_k \left[\left\langle \mathbf{w}_t - \mathbf{w}_*, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*^k, \mathbf{c}_l) \right\rangle + \left\langle \sigma_t^k - \sigma_*^k, \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*^k, \mathbf{v}_k) \right\rangle \right] \geq 0 \quad (20)$$

Summing up Eq. (19) and Eq. (20) and using the first condition in Assumption 3, we have

$$\mathbb{E}[E_{t+1} \mid (\mathbf{w}_t, \sigma_t^1, \dots, \sigma_t^K)] \leq (1 - 2\lambda\gamma_t)E_t + 2\gamma_t^2 G^2$$

Taking the expectation of both sides and using the definition of γ_t , we have

$$\mathbb{E}[E_{t+1}] \leq (1 - 2\lambda\gamma_t)\mathbb{E}[E_t] + 2\gamma_t^2 G^2 = \left(1 - \frac{2\lambda\gamma_0}{t}\right) \mathbb{E}[E_t] + \frac{2\gamma_0^2 G^2}{t^2}, \quad \text{for all } t \geq 1$$

Applying Lemma 6 with $P = 2\lambda\gamma_0 > 1$, $Q = 2\gamma_0^2 G^2$ and $p = q = 1$, we have

$$\mathbb{E}[E_t] \leq \frac{2\gamma_0^2 G^2}{2\lambda\gamma_0 - 1} \frac{1}{t} = O\left(\frac{1}{t}\right)$$

This completes the proof. \square

6 Conclusions

We have developed a comprehensive computational theory for Bayesian regression games, linking them to broader Bayesian game theory and variational inequalities while emphasizing algorithmic approaches. We established sufficient conditions for the existence and uniqueness of Bayesian equilibria by framing the problem as an infinite-dimensional variational inequality (VI) in Hilbert space, extending existing results to a wider class of cost functions. Additionally, we explored two special cases—finite Bayesian priors and quadratic adversarial losses—where the infinite-dimensional VI simplifies to a high-dimensional VI or a nonconvex stochastic optimization problem. We proposed efficient algorithms with strong convergence guarantees for these cases and demonstrated the effectiveness of our framework through numerical experiments on real datasets, highlighting its applicability in classification settings.

References

- [Aghassi and Bertsimas(2006)] M. Aghassi and D. Bertsimas. Robust game theory. *Mathematical Programming*, 107(1-2):231–273, 2006.
- [Bauschke and Combettes(2001)] H. H. Bauschke and P. L. Combettes. A weak-to-strong convergence principle for Fejér-monotone methods in Hilbert spaces. *Mathematics of Operations Research*, 26(2):248–264, 2001.
- [Bauschke and Combettes(2011)] H. H. Bauschke and P. L. Combettes. *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*, volume 408. Springer, 2011.
- [Ben-Tal et al.(2009)] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski. *Robust Optimization*, volume 28. Princeton University Press, 2009.
- [Bertsekas and Tsitsiklis(1989)] D. P. Bertsekas and J. N. Tsitsiklis. *Parallel and Distributed Computation: Numerical Methods*, volume 23. Prentice Hall Englewood Cliffs, NJ, 1989.

- [Bhawalkar and Roughgarden(2011)] K. Bhawalkar and T. Roughgarden. Welfare guarantees for combinatorial auctions with item bidding. In *SODA*, pages 700–709. SIAM, 2011.
- [Biggio and Roli(2018)] B. Biggio and F. Roli. Wild patterns: Ten years after the rise of adversarial machine learning. *Pattern Recognition*, 84:317–331, 2018.
- [Bottou(1998)] Léon Bottou. Online learning and stochastic approximations. *On-line learning in neural networks*, 17(9):142, 1998.
- [Brückner et al.(2012)] M. Brückner, C. Kanzow, and T. Scheffer. Static prediction games for adversarial learning problems. *Journal of Machine Learning Research*, 13(1):2617–2654, 2012.
- [Bruna et al.(2014)] J. Bruna, C. Szegedy, I. Sutskever, I. Goodfellow, W. Zaremba, R. Fergus, and D. Erhan. Intriguing properties of neural networks. In *ICLR*, 2014. URL https://openreview.net/forum?id=kklr_MTHMRQjG.
- [Cai and Papadimitriou(2014)] Y. Cai and C. Papadimitriou. Simultaneous bayesian auctions and computational complexity. In *EC*, pages 895–910, 2014.
- [Chen et al.(2009)] X. Chen, X. Deng, and S-H. Teng. Settling the complexity of computing two-player Nash equilibria. *Journal of the ACM (JACM)*, 56(3):1–57, 2009.
- [Christodoulou et al.(2008)] G. Christodoulou, A. Kovács, and M. Schapira. Bayesian combinatorial auctions. In *ICALP*, pages 820–832. Springer, 2008.
- [Chung(1954)] K. L. Chung. On a stochastic approximation method. *The Annals of Mathematical Statistics*, pages 463–483, 1954.
- [Conitzer and Sandholm(2003)] V. Conitzer and T. Sandholm. Complexity results about Nash equilibria. In *IJCAI*, pages 765–771, 2003.
- [Conitzer and Sandholm(2008)] V. Conitzer and T. Sandholm. New complexity results about Nash equilibria. *Games and Economic Behavior*, 63(2):621–641, 2008.
- [Daskalakis et al.(2009)] C. Daskalakis, P. W. Goldberg, and C. H. Papadimitriou. The complexity of computing a Nash equilibrium. *SIAM Journal on Computing*, 39(1):195–259, 2009.
- [Dua and Graff(2017)] D. Dua and C. Graff. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- [Duchi et al.(2011)] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(7), 2011.
- [EL Ghaoui et al.(2003)] L. EL Ghaoui, G. Lanckriet, and G. Natsoulis. Robust classification with interval data. *Technical Report UCB/CSD-03-1279*, 2003.
- [Facchinei and Pang(2007)] F. Facchinei and J-S. Pang. *Finite-dimensional Variational Inequalities and Complementarity Problems*. Springer Science & Business Media, 2007.
- [Feldman et al.(2013)] M. Feldman, H. Fu, N. Gravin, and B. Lucier. Simultaneous auctions are (almost) efficient. In *SOTC*, pages 201–210, 2013.
- [Fercoq and Richtárik(2015)] O. Fercoq and P. Richtárik. Accelerated, parallel, and proximal coordinate descent. *SIAM Journal on Optimization*, 25(4):1997–2023, 2015.
- [Gilboa and Zemel(1989)] I. Gilboa and E. Zemel. Nash and correlated equilibria: Some complexity considerations. *Games and Economic Behavior*, 1(1):80–93, 1989.
- [Globerson and Roweis(2006)] A. Globerson and S. Roweis. Nightmare at test time: robust learning by feature deletion. In *ICML*, pages 353–360, 2006.
- [Goodfellow et al.(2015)] I. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. In *ICLR*, 2015. URL <http://arxiv.org/abs/1412.6572>.

- [Gottlob et al.(2007)] G. Gottlob, G. Greco, and T. Mancini. Complexity of pure equilibria in Bayesian games. In *IJCAI*, pages 1294–1299, 2007.
- [Großhans et al.(2013)] M. Großhans, C. Sawade, M. Brückner, and T. Scheffer. Bayesian games for adversarial regression problems. In *ICML*, pages 55–63, 2013.
- [Güler(1991)] O. Güler. On the convergence of the proximal point algorithm for convex minimization. *SIAM Journal on Control and Optimization*, 29(2):403–419, 1991.
- [Halpern(1967)] B. Halpern. Fixed points of nonexpanding maps. *Bulletin of the American Mathematical Society*, 73(6):957–961, 1967.
- [Harsanyi(1967)] J. C. Harsanyi. Games with incomplete information played by "Bayesian" players, I-III. *Management Science*, 14(3), 1967.
- [Huber(2004)] P. J. Huber. *Robust Statistics*, volume 523. John Wiley & Sons, 2004.
- [Kakutani(1941)] S. Kakutani. A generalization of Brouwer’s fixed point theorem. *Duke Mathematical Journal*, 8(3):457–459, 1941.
- [Kinderlehrer and Stampacchia(2000)] D. Kinderlehrer and G. Stampacchia. *An Introduction to Variational Inequalities and Their Applications*. SIAM, 2000.
- [Kingma and Ba(2015)] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. URL <http://arxiv.org/abs/1412.6980>.
- [Kolter and Madry(2018)] Z. Kolter and A. Madry. Adversarial robustness: Theory and practice. *Tutorial at NeurIPS*, 2018.
- [Korpelevich(1976)] G. M. Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756, 1976.
- [Lanckriet et al.(2002)] G. R. G. Lanckriet, L. EL Ghaoui, C. Bhattacharyya, and M. I. Jordan. A robust minimax approach to classification. *Journal of Machine Learning Research*, 3(Dec):555–582, 2002.
- [Lions and Stampacchia(1967)] J-L. Lions and G. Stampacchia. Variational inequalities. *Communications on Pure and Applied Mathematics*, 20(3):493–519, 1967.
- [Madry et al.(2018)] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018. URL <https://openreview.net/forum?id=rJzIBfZAb>.
- [Malitsky(2015)] Y. Malitsky. Projected reflected gradient methods for monotone variational inequalities. *SIAM Journal on Optimization*, 25(1):502–520, 2015.
- [Malitsky and Semenov(2014)] Y. V. Malitsky and V. V. Semenov. An extragradient algorithm for monotone variational inequalities. *Cybernetics and Systems Analysis*, 50(2):271–277, 2014.
- [Megiddo and Papadimitriou(1991)] N. Megiddo and C. H. Papadimitriou. On total functions, existence theorems and computational complexity. *Theoretical Computer Science*, 81(2):317–324, 1991.
- [Mertens and Zamir(1985)] J-F. Mertens and S. Zamir. Formulation of bayesian analysis for games with incomplete information. *International Journal of Game Theory*, 14(1):1–29, 1985.
- [Myerson(1985)] R. B. Myerson. Bayesian equilibrium and incentive-compatibility: An introduction. *Social Goals and Social Organization: Essays in Memory of Elisha Pazner*, pages 229–260, 1985.
- [Nadezhkina and Takahashi(2006)] N. Nadezhkina and W. Takahashi. Strong convergence theorem by a hybrid method for nonexpansive mappings and Lipschitz-continuous monotone mappings. *SIAM Journal on Optimization*, 16(4):1230–1241, 2006.
- [Nash(1950)] J. Nash. Equilibrium points in n-person games. *Proc. Nat. Acad. Sci.*, 36:48–49, 1950.
- [Nash(1951)] J. Nash. Non-cooperative games. *Annals of Mathematics*, pages 286–295, 1951.

- [Nemirovsky(1983)] A. S. Nemirovsky. *Problem Complexity and Method Efficiency in Optimization*. Wiley, 1983.
- [Nesterov(2012)] Y. Nesterov. Efficiency of coordinate descent methods on huge-scale optimization problems. *SIAM Journal on Optimization*, 22(2):341–362, 2012.
- [Papadimitriou(1994)] C. H. Papadimitriou. On the complexity of the parity argument and other inefficient proofs of existence. *Journal of Computer and System Sciences*, 48(3):498–532, 1994.
- [Pedregosa et al.(2011)] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [Popov(1980)] L. D. Popov. A modification of the Arrow-Hurwicz method for search of saddle points. *Mathematical notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.
- [Richtárik and Takáč(2014)] P. Richtárik and M. Takáč. Iteration complexity of randomized block-coordinate descent methods for minimizing a composite function. *Mathematical Programming*, 144(1-2):1–38, 2014.
- [Robbins and Monro(1951)] H. Robbins and S. Monro. A stochastic approximation method. *The Annals of Mathematical Statistics*, pages 400–407, 1951.
- [Rosen(1965)] J. B. Rosen. Existence and uniqueness of equilibrium points for concave n-person games. *Econometrica: Journal of the Econometric Society*, pages 520–534, 1965.
- [Rubinstein(2018)] A. Rubinstein. Inapproximability of Nash equilibrium. *SIAM Journal on Computing*, 47(3):917–959, 2018.
- [Saejung and Yotkaew(2012)] S. Saejung and P. Yotkaew. Approximation of zeros of inverse strongly monotone operators in Banach spaces. *Nonlinear Analysis: Theory, Methods & Applications*, 75(2):742–750, 2012.
- [Santurkar et al.(2019)] S. Santurkar, A. Ilyas, D. Tsipras, L. Engstrom, B. Tran, and A. Madry. Image synthesis with a single (robust) classifier. In *NeurIPS*, pages 1262–1273, 2019.
- [Sayed and Chen(2002)] A. H. Sayed and H. Chen. A uniqueness result concerning a robust regularized least-squares solution. *Systems & Control Letters*, 46(5):361–369, 2002.
- [Singh et al.(2004)] S. Singh, V. Soni, and M. Wellman. Computing approximate Bayes-Nash equilibria in tree-games of incomplete information. In *EC*, pages 81–90, 2004.
- [Solodov and Svaiter(1999)] M. V. Solodov and B. F. Svaiter. A new projection method for variational inequality problems. *SIAM Journal on Control and Optimization*, 37(3):765–776, 1999.
- [Solodov and Svaiter(2000)] M. V. Solodov and B. F. Svaiter. Forcing strong convergence of proximal point iterations in a Hilbert space. *Mathematical Programming*, 87(1):189–202, 2000.
- [Teo et al.(2008)] C. H. Teo, A. Globerson, S. T. Roweis, and A. J. Smola. Convex learning with invariances. In *NeurIPS*, pages 1489–1496, 2008.
- [Tseng(2000)] P. Tseng. A modified forward-backward splitting method for maximal monotone mappings. *SIAM Journal on Control and Optimization*, 38(2):431–446, 2000.
- [Tsipras et al.(2019)] D. Tsipras, S. Santurkar, L. Engstrom, A. Turner, and A. Madry. Robustness may be at odds with accuracy. In *ICLR*, 2019. URL <https://openreview.net/forum?id=SyxAb30cY7>.
- [Ui(2016)] T. Ui. Bayesian Nash equilibrium and variational inequalities. *Journal of Mathematical Economics*, 63:139–146, 2016.
- [Wong et al.(2020)] E. Wong, L. Rice, and J. Z. Kolter. Fast is better than free: Revisiting adversarial training. In *ICLR*, 2020. URL <https://openreview.net/forum?id=BJx040EFvH>.

- [Wright(2015)] S. J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.
- [Zhu et al.(2020)] C. Zhu, Y. Cheng, Z. Gan, S. Sun, T. Goldstein, and J. Liu. FreeLB: Enhanced adversarial training for natural language understanding. In *ICLR*, 2020. URL <https://openreview.net/forum?id=BygzbyHFvB>.
- [Zinkevich et al.(2007)] M. Zinkevich, M. Johanson, M. Bowling, and C. Piccione. Regret minimization in games with incomplete information. *NeurIPS*, 20:1729–1736, 2007.

A Further Backgrounds

In this section, we provide the basic ideas and some additional background materials for the development of our PRG-IE and PG-RBC algorithms. Some discussions on the relevant algorithms are also included.

PRG-IE: We start with a brief overview of the projected reflected gradient algorithm for solving the variational inequality (VI) in Hilbert space. Let \mathcal{S} be a nonempty, closed and convex set of a Hilbert space \mathcal{H} with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}}$, and $T : \mathcal{H} \rightarrow \mathcal{H}$ be *strictly monotone* and ℓ -smooth for some constant $\ell > 0$: for $\forall \mathbf{x}, \mathbf{x}' \in \mathcal{H}$, $\|T(\mathbf{x}) - T(\mathbf{x}')\|_{\mathcal{H}} \leq \ell \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{H}}$. Then, we consider the problem of finding a point $\mathbf{x}_\star \in \mathcal{S}$ such that

$$\langle \mathbf{x} - \mathbf{x}_\star, T(\mathbf{x}_\star) \rangle_{\mathcal{H}} \geq 0, \quad \text{for all } \mathbf{x} \in \mathcal{S} \quad (21)$$

A projected reflected gradient algorithm $\mathbf{x}_{t+1} \leftarrow P_{\mathcal{S}}(\mathbf{x}_t - \gamma \cdot T(2\mathbf{x}_t - \mathbf{x}_{t-1}))$ can be applied for solving this problem where the stepsize $\lambda \in (0, (\sqrt{2} - 1)/\ell)$, and $P_{\mathcal{S}}(\cdot)$ is the orthogonal projection onto a closed set \mathcal{S} . From the update formula, we see that this algorithm has a very simple and elegant structure, which only requires evaluating $T(\cdot)$ and $P_{\mathcal{S}}(\cdot)$ once at each iteration. Thus, it is more computationally appealing than the Korpelevich’s extragradient algorithm [Korpelevich(1976)], Popov’s modified Arrow-Hurwicz algorithm [Popov(1980)], Tseng’s forward-backward splitting algorithm [Tseng(2000)] and some other algorithms [Solodov and Svaiter(1999), Malitsky and Semenov(2014)].

Note that the VI in Eq. (3) is in the form of Eq. (21) with a Hilbert space \mathcal{H} consisting of (an equivalence class of) a function $\beta : \mathbb{R}^n \mapsto \mathbb{R}^m \times \mathbb{R}^{n \times m}$ with the inner product $\langle \cdot, \cdot \rangle_{\mathcal{H}} : \mathcal{H} \times \mathcal{H} \rightarrow \mathbb{R}$ defined by Eq. (9) and the mapping T defined by Eq. (10). Under Assumption 1, we assume that T is strictly monotone and $\mathcal{W} \times \Sigma$ is compact. Then, [Malitsky(2015), Theorem 3.2] guarantees that the sequence generated by projected reflected gradient algorithm weakly converges to a unique Bayesian equilibrium. However, in the infinite-dimensional setting, strong convergence (or norm convergence) is often much more desirable than weak convergence, since it guarantees that the physically tangible property, the error $\|\mathbf{x}_t - \mathbf{x}_\star\|_{\mathcal{H}}^2$ eventually become arbitrarily small [Bauschke and Combettes(2001)]. The importance of strong convergence is also demonstrated by [Güler(1991)] for convex optimization that the convergence rate of the sequence of objectives $\{f(\mathbf{x}_t)\}_{t \geq 0}$ is better when $\{\mathbf{x}_t\}_{t \geq 0}$ with strong convergence than weak convergence. This encourages the strong convergence theorems for various algorithms in Hilbert space [Solodov and Svaiter(2000), Nadezhkina and Takahashi(2006)].

PG-RBC: The variational inequality (VI) in Eq. (12) is the problem of finding a point $(\mathbf{w}_\star, \sigma_\star^1, \dots, \sigma_\star^K) \in \mathcal{S} = \mathcal{W} \times \mathcal{X} \times \dots \times \mathcal{X}$ such that

$$\sum_{k=1}^K p_k \left[\left\langle \mathbf{w} - \mathbf{w}_\star, \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_\star, \sigma_\star^k, \mathbf{c}_l) \right\rangle + \left\langle \sigma^k - \sigma_\star^k, \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_\star, \sigma_\star^k, \mathbf{v}_k) \right\rangle \right] \geq 0$$

for all $\mathbf{w} \in \mathcal{W}$ and $\sigma^k \in \mathcal{X}$ for all $k \in [K]$. Then, by defining the variable $\mathbf{x} \in \mathbb{R}^m \times \mathbb{R}^{n \times m} \times \dots \times \mathbb{R}^{n \times m}$ and a mapping $T : \mathbb{R}^m \times \mathbb{R}^{n \times m} \times \dots \times \mathbb{R}^{n \times m} \mapsto \mathbb{R}^m \times \mathbb{R}^{n \times m} \times \dots \times \mathbb{R}^{n \times m}$ as follows,

$$T \begin{pmatrix} \mathbf{w} \\ \sigma^1 \\ \vdots \\ \sigma^K \end{pmatrix} = \begin{pmatrix} \sum_{k=1}^K \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}, \sigma^k, \mathbf{c}_l) \\ \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma^1, \mathbf{v}_1) \\ \vdots \\ \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}, \sigma^K, \mathbf{v}_K) \end{pmatrix}$$

we can reformulate the above problem in the compact form as follows,

$$\langle \mathbf{x} - \mathbf{x}_*, T(\mathbf{x}_*) \rangle_{\mathcal{H}} \geq 0, \quad \text{for all } \mathbf{x} \in \mathcal{S} = \mathcal{W} \times \mathcal{X} \times \dots \times \mathcal{X}$$

A projected gradient algorithm $\mathbf{x}_{t+1} \leftarrow P_{\mathcal{S}}(\mathbf{x}_t - \gamma \cdot T(\mathbf{x}_t))$ can be applied but becomes problematic when the problem dimension m , the number of data samples n and the range of a Bayesian prior K are huge. Indeed, the algorithm require performing arithmetic operations of order nmK per iteration and the projection step is another source of inefficiency for huge-size problem. The coordinate update and more generally block coordinate update, which are commonly used to address this issue and improve the computational efficiency, are rooted in the optimization community [Bertsekas and Tsitsiklis(1989)]. During the past decade, the *randomized coordinate update* has emerged as one of the most popular coordinate update schemes and were extensively studied [Nesterov(2012), Richtárik and Takáč(2014), Fercoq and Richtárik(2015)].

B Auxiliary Proofs of Lemmas

B.1 Proof of Lemma 3

Proof of Lemma 3. By the compactness of actions spaces \mathcal{W} and \mathcal{X} , it suffices to show that $(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t) \in \mathcal{W} \times \Sigma$ for all $t \geq 0$ and $(\mathbf{w}_t, \sigma_t) \in \mathcal{W} \times \Sigma$ for all $t \geq 1$. Indeed, the initialization step implies that $\tilde{\mathbf{w}}_0, \tilde{\mathbf{w}}_1, \mathbf{w}_1 \in \mathcal{W}$ and $\tilde{\sigma}_0, \tilde{\sigma}_1, \sigma_1 \in \Sigma$. Then, by the updating formula for $\{(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t)\}_{t \geq 2}$, it is clear that for $(\tilde{\mathbf{w}}_t, \tilde{\sigma}_t) \in \mathcal{W} \times \Sigma$ for all $t \geq 0$. It remains to show that $(\mathbf{w}_t, \sigma_t) \in \mathcal{W} \times \Sigma$ for all $t \geq 2$. Since $\mathbf{0}_{n \times m} \in \mathcal{X}$, we have $\mathbf{0}_{n \times m}(\cdot) \in \Sigma$ where $\mathbf{0}_{n \times m}(\cdot)$ is a constant function from \mathbb{R}^n to $\mathbb{R}^{n \times m}$ with value $\mathbf{0}_{n \times m}$. By the convexity of \mathcal{W} and \mathcal{X} , we have the convexity of $\mathcal{W} \times \Sigma$. By the updating formula for $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 2}$, we find that $(\mathbf{w}_{t+1}, \sigma_{t+1})$ is a convex combination of $(\tilde{\mathbf{w}}_{t+1}, \tilde{\sigma}_{t+1})$, (\mathbf{w}_t, σ_t) and $(\mathbf{0}_m, \mathbf{0}_{n \times m}(\cdot))$. This together with the convexity of $\mathcal{W} \times \Sigma$ implies the desired result. This completes the proof. \square

B.2 Proof of Lemma 4

Proof of Lemma 4. Under Assumption 1 and 2, a unique Bayesian equilibrium $(\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ exists. Given that the sequence $(\bar{\mathbf{w}}_t, \bar{\sigma}_t) = (2\tilde{\mathbf{w}}_t - \tilde{\mathbf{w}}_{t-1}, 2\tilde{\sigma}_t - \tilde{\sigma}_{t-1})$ is defined, we derive from the optimality condition of updating $(\tilde{\mathbf{w}}_{t+1}, \tilde{\sigma}_{t+1})$ and the definition of $P_{\mathcal{W}}$ and P_{Σ} that,

$$\begin{aligned} 0 \leq \mathbb{E}_q \Big[& \langle \mathbf{w} - \tilde{\mathbf{w}}_{t+1}, \tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t + \gamma \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l) \rangle \quad \text{for each } (\mathbf{w}, \sigma) \in \mathcal{W} \times \Sigma \quad (22) \\ & + \langle \sigma(\mathbf{c}_d) - \tilde{\sigma}_{t+1}(\mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d) + \gamma \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d) \rangle \Big] \end{aligned}$$

By rearranging the terms in the above inequality with $(\mathbf{w}, \sigma) = (\mathbf{w}_*, \sigma_*) \in \mathcal{W} \times \Sigma$ and the equality $\langle a, b \rangle = (1/2)(\|a + b\|^2 - \|a\|^2 - \|b\|^2)$, we have

$$\begin{aligned} \|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_*\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d)\|_F^2] &\leq \|\mathbf{w}_t - \mathbf{w}_*\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d)\|_F^2] \\ &\quad - \|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t\|^2 - \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2] - 2\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \mathbf{w}_* \rangle \right. \\ &\quad \left. + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \end{aligned} \quad (23)$$

Moreover, we have

$$\begin{aligned} &\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d)) \rangle \right] \\ &= \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d) \rangle \right] \\ &\quad + \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l), \bar{\mathbf{w}}_t - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d), \bar{\sigma}_t(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \end{aligned}$$

Using the first condition in Assumption 2 with $(\mathbf{w}, \sigma) = (\bar{\mathbf{w}}_t, \bar{\sigma}_t)$ and $(\mathbf{w}', \sigma') = (\mathbf{w}_*, \sigma_*)$, we have

$$\begin{aligned} &\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l), \bar{\mathbf{w}}_t - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d), \bar{\sigma}_t(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \\ &\geq \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l), \bar{\mathbf{w}}_t - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d), \bar{\sigma}_t(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \\ &= 2\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_t - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_t(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \\ &\quad - \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t-1} - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t-1}(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \end{aligned}$$

Putting these two inequalities together with Eq. (23) yields that

$$\begin{aligned} \|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_*\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d)\|_F^2] &\leq \|\mathbf{w}_t - \mathbf{w}_*\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d)\|_F^2] \\ &\quad - \|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t\|^2 - \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2] \\ &\quad - 2\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l) - \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle \right] \\ &\quad - 2\gamma \mathbb{E}_q \left[\langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d) - \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d) \rangle \right] \\ &\quad - 2\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d) \rangle \right] \\ &\quad - 4\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_t - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_t(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \\ &\quad + 2\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_l), \mathbf{w}_{t-1} - \mathbf{w}_* \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_*, \sigma_*(\mathbf{c}_d), \mathbf{c}_d), \sigma_{t-1}(\mathbf{c}_d) - \sigma_*(\mathbf{c}_d) \rangle \right] \end{aligned} \quad (24)$$

Then, we obtain two inequalities by changing the index t in Eq. (22) to $t - 1$ and letting $(\mathbf{w}, \sigma) = (\tilde{\mathbf{w}}_{t-1}, \tilde{\sigma}_{t-1})$ and $(\mathbf{w}, \sigma) = (\tilde{\mathbf{w}}_{t+1}, \tilde{\sigma}_{t+1})$ and add them to obtain that

$$\begin{aligned} 0 &\leq \mathbb{E}_q \left[\langle \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t, \bar{\mathbf{w}}_t - \mathbf{w}_{t-1} + \gamma \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_l) \rangle \right. \\ &\quad \left. + \langle \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \tilde{\sigma}_t(\mathbf{c}_d) - \sigma_{t-1}(\mathbf{c}_d) + \gamma \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_d) \rangle \right] \end{aligned}$$

By the definition of the sequence $(\bar{\mathbf{w}}_t, \bar{\sigma}_t)$, we further have

$$\begin{aligned}
& -2\gamma\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d) \rangle \right] \\
& \leq 2\langle \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t, \tilde{\mathbf{w}}_t - \mathbf{w}_{t-1} \rangle + 2\mathbb{E}_q[\langle \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \tilde{\sigma}_t(\mathbf{c}_d) - \sigma_{t-1}(\mathbf{c}_d) \rangle] \\
& = 2\langle \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t, \bar{\mathbf{w}}_t - \tilde{\mathbf{w}}_t \rangle + 2\langle \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t, \tilde{\mathbf{w}}_{t-1} - \mathbf{w}_{t-1} \rangle + 2\mathbb{E}_q[\langle \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \bar{\sigma}_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d) \rangle] \\
& \quad + 2\mathbb{E}_q[\langle \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \tilde{\sigma}_{t-1}(\mathbf{c}_d) - \sigma_{t-1}(\mathbf{c}_d) \rangle] \\
& = 2\langle \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t, \bar{\mathbf{w}}_t - \mathbf{w}_t \rangle + 2\mathbb{E}_q[\langle \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \bar{\sigma}_t(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d) \rangle] + 2\langle \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t, \mathbf{w}_t - \tilde{\mathbf{w}}_t \rangle \\
& \quad + 2\mathbb{E}_q[\langle \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d) \rangle] + 2\langle \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t, \tilde{\mathbf{w}}_{t-1} - \mathbf{w}_{t-1} \rangle \\
& \quad + 2\mathbb{E}_q[\langle \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d), \tilde{\sigma}_{t-1}(\mathbf{c}_d) - \sigma_{t-1}(\mathbf{c}_d) \rangle]
\end{aligned}$$

Using the equality $\langle a, b \rangle = (1/2)(\|a + b\|^2 - \|a\|^2 - \|b\|^2)$ for the first two terms, and the Young's inequality $\langle a, b \rangle \leq (1/8)\|a\|^2 + 2\|b\|^2$ for the last four terms, we have

$$\begin{aligned}
& -2\gamma\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d) \rangle \right] \\
& \leq \|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2] - \frac{1}{2} (\|\tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d)\|_F^2]) \\
& \quad - (\|\bar{\mathbf{w}}_t - \mathbf{w}_t\|^2 + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2]) + 4(\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2]) \\
& \quad + 4(\|\mathbf{w}_{t-1} - \tilde{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\sigma_{t-1}(\mathbf{c}_d) - \tilde{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2]) \tag{25}
\end{aligned}$$

By using the second condition in Assumption 2, we have

$$\begin{aligned}
& -2\gamma\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_l) - \nabla_{\mathbf{w}} \hat{\theta}_l(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t \rangle \right. \\
& \quad \left. + \langle \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_t, \bar{\sigma}_t(\mathbf{c}_d), \mathbf{c}_d) - \nabla_{\bar{X}} \hat{\theta}_d(\bar{\mathbf{w}}_{t-1}, \bar{\sigma}_{t-1}(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d) \rangle \right] \tag{26} \\
& \leq 2\gamma L (\|\bar{\mathbf{w}}_t - \bar{\mathbf{w}}_{t-1}\| + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \bar{\sigma}_{t-1}(\mathbf{c}_d)\|_F]) (\|\tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\| + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d)\|_F]) \\
& \leq 2\gamma L (\|\bar{\mathbf{w}}_t - \bar{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \bar{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2]) + \|\tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d)\|_F^2]) \\
& \leq 2\gamma L (\|\tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d)\|_F^2]) + 6\gamma L (\|\bar{\mathbf{w}}_t - \mathbf{w}_t\|^2 + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2]) \\
& \quad + 6\gamma L (\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2]) + \|\bar{\mathbf{w}}_t - \bar{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \bar{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2])
\end{aligned}$$

Combining Eq. (24), Eq. (25) and Eq. (26) yields the desired inequality. \square

B.3 Proof of Lemma 5

Proof of Lemma 5. By the updating formula for the sequence $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ and the Jensen's inequality, we have

$$\begin{aligned}
& \|\mathbf{w}_{t+1} - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2] \tag{27} \\
& \leq \frac{\delta_t}{2} (\|\mathbf{w}_t - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2]) + \frac{\delta_t}{2} (\|\mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_\star(\mathbf{c}_d)\|_F^2]) \\
& \quad + (1 - \delta_t) (\|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2])
\end{aligned}$$

Recall that the residue sequence $r_t \geq 0$ is defined by

$$\begin{aligned}
r_t &= \|\mathbf{w}_t - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2] + 6\gamma L (\|\tilde{\mathbf{w}}_t - \bar{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_t(\mathbf{c}_d) - \bar{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2]) \\
& \quad + 2\gamma\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_{t-1} - \mathbf{w}_\star \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_{t-1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle \right]
\end{aligned}$$

Combining this with Eq. (27) and Lemma 4, we have

$$\begin{aligned}
r_{t+1} \leq & r_t - \left(\frac{1}{2} - 8\gamma L \right) (\|\tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d)\|_F^2]) + \frac{\delta_t}{2} (\|\mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_\star(\mathbf{c}_d)\|_F^2]) \\
& + (4 + 6\gamma L) (\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2]) + 4 (\|\mathbf{w}_{t-1} - \tilde{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\sigma_{t-1}(\mathbf{c}_d) - \tilde{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2]) \\
& - 2\gamma \mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_t - \mathbf{w}_\star \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle \right] \\
& - (1 - 6\gamma L) (\|\bar{\mathbf{w}}_t - \mathbf{w}_t\|^2 + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2])
\end{aligned} \tag{28}$$

Since $(\mathbf{w}_\star, \sigma_\star) \in \mathcal{W} \times \Sigma$ is a unique solution of the VI in Eq. (3) under Assumption 1 and 2, we have

$$\mathbb{E}_q \left[\langle \nabla_{\mathbf{w}} \hat{\theta}_l(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_l), \tilde{\mathbf{w}}_t - \mathbf{w}_\star \rangle + \langle \nabla_{\bar{X}} \hat{\theta}_d(\mathbf{w}_\star, \sigma_\star(\mathbf{c}_d), \mathbf{c}_d), \tilde{\sigma}_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle \right] \geq 0$$

Using Lemma 3, the updating formula for the sequence $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ and the uniqueness of $(\mathbf{w}_\star, \sigma_\star) \in \mathcal{W} \times \Sigma$, there exists a constant $M_1 > 0$ such that

$$\begin{aligned}
\|\mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_\star(\mathbf{c}_d)\|_F^2] & \leq M_1 \\
\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2] & \leq \frac{\delta_{t-1}^2 M_1}{18} \leq \frac{\delta_{t-1} M_1}{18}
\end{aligned}$$

Putting these pieces together with Eq. (28) and the facts that $0 < \gamma < \min\{1, \frac{1}{100L}\}$ and the sequence $\{\delta_t\}_{t \geq 1}$ is non-increasing yields Eq. (15). Then, we proceed to prove Eq. (16). Using the inequality $\|a + b\|^2 \leq \|a\|^2 + 2\langle a, b \rangle$, together with the Jensen's inequality and the updating formula for the sequence $\{(\mathbf{w}_t, \sigma_t)\}_{t \geq 1}$ yields that

$$\begin{aligned}
& \|\mathbf{w}_{t+1} - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2] \\
= & \|\delta_t(\mathbf{w}_t/2 - \mathbf{w}_\star/2) + (1 - \delta_t)(\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_\star) - \delta_t(\mathbf{w}_\star/2)\|^2 \\
& + \mathbb{E}_q[\|\delta_t(\sigma_t(\mathbf{c}_d)/2 - \sigma_\star(\mathbf{c}_d)/2) + (1 - \delta_t)(\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)) - \delta_t(\sigma_\star(\mathbf{c}_d)/2)\|_F^2] \\
\leq & \frac{\delta_t}{4} (\|\mathbf{w}_t - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2]) + (1 - \delta_t) (\|\tilde{\mathbf{w}}_{t+1} - \mathbf{w}_\star\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d)\|_F^2]) \\
& - \delta_t (\langle \mathbf{w}_\star, \mathbf{w}_{t+1} - \mathbf{w}_\star \rangle + \mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle])
\end{aligned}$$

Note that $1 - \delta_t \geq 1/2$ for all $t \geq 2$ and $0 < \gamma < \min\{1, \frac{1}{100L}\}$. By combining the above inequality and Lemma 4, we have

$$\begin{aligned}
r_{t+1} \leq & \left(1 - \frac{3}{4}\delta_t \right) r_t - \delta_t (\langle \mathbf{w}_\star, \mathbf{w}_{t+1} - \mathbf{w}_\star \rangle + \mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]) \\
& - \left(\frac{1}{4} - 7\gamma L \right) (\|\tilde{\mathbf{w}}_{t+1} - \bar{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\tilde{\sigma}_{t+1}(\mathbf{c}_d) - \bar{\sigma}_t(\mathbf{c}_d)\|_F^2]) \\
& + (4 + 6\gamma L) (\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2]) \\
& + 4 (\|\mathbf{w}_{t-1} - \tilde{\mathbf{w}}_{t-1}\|^2 + \mathbb{E}_q[\|\sigma_{t-1}(\mathbf{c}_d) - \tilde{\sigma}_{t-1}(\mathbf{c}_d)\|_F^2]) \\
& - \left(\frac{1}{2} - 3\gamma L \right) (\|\bar{\mathbf{w}}_t - \mathbf{w}_t\|^2 + \mathbb{E}_q[\|\bar{\sigma}_t(\mathbf{c}_d) - \sigma_t(\mathbf{c}_d)\|_F^2]), \quad \text{for all } t \geq 2
\end{aligned} \tag{29}$$

Note that $\delta_t \leq \delta_{t-1} \leq 2\delta_t$ for all $t \geq 2$ and there exists a constant $M_1 > 0$ such that

$$\|\mathbf{w}_t - \tilde{\mathbf{w}}_t\|^2 + \mathbb{E}_q[\|\sigma_t(\mathbf{c}_d) - \tilde{\sigma}_t(\mathbf{c}_d)\|_F^2] \leq \frac{\delta_{t-1}^2 M_1}{18} \leq \frac{2\delta_t^2 M_1}{9}$$

Putting the above inequality and the fact that $0 < \gamma < \min\{1, \frac{1}{100L}\}$ together with Eq. (29) yields that

$$\begin{aligned} r_{t+1} &\leq \left(1 - \frac{3}{4}\delta_t\right) r_t + 2\delta_t^2 M_1 - \delta_t (\langle \mathbf{w}_\star, \mathbf{w}_{t+1} - \mathbf{w}_\star \rangle + \mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]) \\ &= \left(1 - \frac{3\delta_t}{4}\right) r_t + \frac{3\delta_t}{4} (4\delta_t M_1 - 2\langle \mathbf{w}_\star, \mathbf{w}_{t+1} - \mathbf{w}_\star \rangle - 2\mathbb{E}_q[\langle \sigma_\star(\mathbf{c}_d), \sigma_{t+1}(\mathbf{c}_d) - \sigma_\star(\mathbf{c}_d) \rangle]) \end{aligned}$$

This completes the proof. □