
Double-Step Alternating Extragradient with Increasing Timescale Separation for Finding Local Minimax Points: Provable Improvements

Kyuwon Kim¹ Donghwan Kim¹

Abstract

In nonconvex-nonconcave minimax optimization, *two-timescale* gradient methods have shown their potential to find local minimax (optimal) points, provided that the timescale separation between the min and the max player is sufficiently large. However, existing two-timescale variants of gradient descent ascent and extragradient methods face two shortcomings, especially when we search for *non-strict* local minimax points that are prevalent in modern overparameterized setting. In specific, (i) these methods can be unstable at some *non-strict* local minimax points even with sufficiently large timescale separation, and even (ii) computing a proper amount of timescale separation is infeasible in practice. To remedy these two issues, we propose to incorporate two simple but provably effective schemes, *double-step alternating* update and *increasing* timescale separation, into the two-timescale extragradient method, respectively. Under mild conditions, we show that the proposed methods converge to *non-strict* local minimax points that all existing two-timescale methods fail to converge.

1. Introduction

The significance of minimax problems in the machine learning community has grown considerably, since generative adversarial network (GAN) (Goodfellow et al., 2014), adversarial training (Madry et al., 2018), multi-agent reinforcement learning (Wai et al., 2018), fair classification (Martinez et al., 2020) and sharpness-aware minimization (Foret et al., 2021), are formulated as

$$\min_x \max_y f(x, y).$$

¹Department of Mathematical Sciences, KAIST, Daejeon, Republic of Korea. Correspondence to: Kyuwon Kim <kkw4053@kaist.ac.kr>, Donghwan Kim <donghwan@kaist.ac.kr>.

However, solving a nonconvex-nonconcave minimax problem is known to be problematic, even when using a gradient descent ascent (GDA) method that is a natural extension of a gradient descent method in minimization. This stands in stark contrast to the remarkable success of the gradient descent method in machine learning, a success underpinned by theoretical results; under mild assumptions, the gradient descent finds local minimum and escapes strict saddle points with probability one (Lee et al., 2016; 2019). Therefore, the goal of this paper is to establish a comparable theory in nonconvex-nonconcave minimax optimization.

Nonconvex-nonconcave minimax problems in most machine learning applications are *sequential* games and thus have an intrinsic order between the min-player x and max-player y (Fiez et al., 2020; Jin et al., 2020). While such order is negligible in a convex-concave setting, if we disregard this order in a nonconvex-nonconcave setting, for example in GAN, the undesirable mode collapse phenomenon can arise (Goodfellow, 2016). Nevertheless, a proper definition of local optimal point in minimax problems that takes account of the intrinsic order between the players was not widely recognized until (Fiez et al., 2020; Jin et al., 2020). Accordingly, Jin et al. (2020) proposed the first proper notion of local optimal points, named *local minimax points*, for the *sequential* games, which differs from and includes the commonly employed notion of Nash equilibrium in *simultaneous* games.

To find such local minimax (optimal) points, Jin et al. (2020) considered a *two-timescale* GDA (Heusel et al., 2017) that updates with different step sizes (timescales) for each variable x and y . In specific, Jin et al. (2020) analyzed that the GDA with a sufficiently large timescale separation can find local minimax points, using dynamical system theory. However, it faces two challenges. The first limitation arises from the *non-degeneracy* assumption on $\nabla_{yy}^2 f$ required in (Jin et al., 2020), which neglects the *non-strict* local minimax points that are ubiquitous in the modern overparameterized setting (Cooper, 2021; Liu et al., 2022). The second is that Jin et al. (2020) have not specified how large one should choose an appropriate timescale separation to guarantee the convergence to the local minimax points.

The partial answers to resolve the aforementioned two lim-

itations were provided by Chae et al. (2024b) and Fiez & Ratliff (2021); Li et al. (2022), respectively. First, Chae et al. (2024b) removed the non-degeneracy assumption of $\nabla_{\mathbf{y}\mathbf{y}}^2 f$, and demonstrated that, under mild assumptions, the two-timescale extragradient (EG) method can converge to a set of local minimax points, especially including the *non-strict* local minimax points, that is larger than that of the two-timescale GDA method. Nevertheless, Chae et al. (2024b) state that there still exist some *non-strict* local minimax points that two-timescale EG cannot converge to, due to its insufficient stability. Second, Fiez & Ratliff (2021) and Li et al. (2022) specified an appropriate timescale separation needed to guarantee a local convergence to *strict* local minimax points. However, computing it is infeasible in general as it requires the second-order derivative information of the function f at the local minimax point, and even how one can generalize their result to *non-strict* local minimax points remains open.

This paper thus focuses on addressing these two issues of existing two-timescale gradient methods, which are further detailed in Section 4. This paper then proposes to integrate two simple but provably effective schemes namely *double-step alternating* update and *increasing* timescale separation, in Sections 5 and 6, respectively. Sections 5 and 6 provide improved local convergence analyses in terms of the local minimax points, and these are generalized to a global statement in Section 7. Our contributions via dynamical system theory can be summarized as below.

- **Local Convergence to Local Minimax Points**

- (a) **Double-Step Alternating Update:** In Section 5, we present that updating the min and the max player alternately, which enhances stability from a dynamical system perspective, proves particularly helpful in addressing the first issue of the (simultaneous-update) two-timescale EG. We would like to highlight that the resulting *double-step alternating* EG (Alt2-EG) method, in Algorithm 1, entails a slight yet essential deviation from the usual alternating scheme, which will be detailed later. Built upon its spectral analysis, we show that, under mild conditions, the Alt2-EG with sufficiently large timescale separation is stable at non-strict local minimax points, for those that are unstable for other existing two-timescale methods.
- (b) **Increasing Timescale Separation:** In Section 6, to remedy the second issue, we suggest to simply increase the timescale separation indefinitely as iteration goes. In particular, we demonstrate that the Alt2-EG with *increasing* timescale separation (Alt2-EG-ITS) is stable at non-strict local minimax points that are stable for Alt2-EG with (sufficiently large) fixed timescale separation (Alt2-EG-FTS) in Section 5. We would

like to emphasize here that, unlike the technique being simple and straightforward, analyzing its stability via dynamical system theory is rather complicated, since the resulting system is *non-autonomous*.

- **Global Convergence to Local Minimax Points**

In Section 7, we show that both Alt2-EG-FTS and Alt2-EG-ITS globally find first-order stationary points under a star-convex-star-concave setting. Combined with the aforementioned local convergence analyses, we claim that, under mild conditions, both Alt2-EG-FTS and Alt2-EG-ITS can globally converge to local minimax points under the aforementioned nonconvex-nonconcave setting, while our current analysis for the latter has some limitation due to its *non-autonomous* property.

2. Related Work

Two-Timescale Gradient Methods The vanilla gradient descent ascent (GDA) method may not converge to local minimax points (Daskalakis & Panageas, 2018; Jin et al., 2020). To resolve this non-convergence, the *two-timescale* GDA (Heusel et al., 2017) has been widely studied. For example, under a nonconvex-strongly-concave setting, Lin et al. (2020) established that the GDA with a timescale separation of the order $\Theta(\kappa_{\mathbf{y}}^2)$, where $\kappa_{\mathbf{y}}$ is a global condition number for \mathbf{y} , globally finds a first-order stationary point. However, since the set of first-order stationary points is considerably larger than that of the local minimax (optimal) points, it is crucial for a method to only converge to local minimax points.

After introducing the definition of the local minimax point in their paper, Jin et al. (2020) showed that, under the non-degeneracy assumption on $\nabla_{\mathbf{y}\mathbf{y}}^2 f$, the two-timescale GDA locally converges to *strict* local minimax points. Fiez & Ratliff (2021) then specified a proper value of timescale separation needed for such guarantee, which was missing in (Jin et al., 2020). This was further refined in (Li et al., 2022), which demonstrated that GDA with a timescale separation of the order $\Theta(\tilde{\kappa}_{\mathbf{y}})$ locally converges to *strict* local minimax points. Here, $\tilde{\kappa}_{\mathbf{y}}$ is a local condition number for \mathbf{y} at the local minimax point, which is not available in practice.

Recently, Chae et al. (2024b) generalized the local convergence result of (Jin et al., 2020) by removing the crucial non-degeneracy assumption needed in (Jin et al., 2020; Fiez & Ratliff, 2021; Li et al., 2022). This forward step is essential, given that *non-strict* optimal points are everywhere in the modern overparameterized setting (Cooper, 2021; Liu et al., 2022). In specific, Chae et al. (2024b) proved that there exist non-strict local minimax points that the two-timescale EG converges to, while the two-timescale GDA cannot. However, as pointed out in (Chae et al., 2024b) there

still exist non-strict local minimax points that two-timescale EG does not converge to, due to its insufficient stability from dynamical system perspective.

Alternating GDA The *alternating* gradient descent ascent (Alt-GDA), which updates each variables \mathbf{x} and \mathbf{y} sequentially, has been found to be more stable and sometimes faster than the plain *simultaneous* update GDA (Sim-GDA) under various settings (Mescheder et al., 2018; Zhang et al., 2022; Lee et al., 2024). For instance, in a bilinear setting, Sim-GDA moves far away from the optimal point indefinitely, whereas Alt-GDA also does not converge but oscillates in a stable cycle around the optimal point (Mescheder et al., 2018). Moreover, under a strongly-convex-strongly-concave setting, the Sim-GDA locally converges to the optimal point with an iteration complexity of $O(\kappa^2)$, where κ is a condition number, while the Alt-GDA achieves a better near-optimal bound $O(\kappa)$ (Zhang et al., 2022). Very recently and concurrently, it is further demonstrated that the Alt-GDA and its variant (called Alex-GDA) achieve faster rates even globally, compared to that of Sim-GDA (Lee et al., 2024). This paper utilizes these stabilizing effect of the alternating update, from a dynamical system perspective, to improve the stability guarantee of the two-timescale EG.

GDA with Increasing Timescale Separation Under a nonconvex-strongly-concave setting, Li et al. (2023) studied an AdaGrad-like *parameter-free* GDA method that adapts step sizes for each min and max player, which has a practical importance in deep neural network training. Here, the step sizes are determined by past gradient information, *without requiring problem parameters* such as the condition number needed to determine the appropriate amount of timescale separation in (Lin et al., 2020). In regard of (Lin et al., 2020), the resulting adaptive GDA will likely suffer from the non-convergence, since the adaptive step sizes may not have sufficiently large timescale separation. Therefore, Li et al. (2023) suggested to modify the adaptive step size rule so that the timescale separation increases indefinitely as iteration goes. Similar to (Lin et al., 2020), Li et al. (2023) analyzed the global convergence to a first-order stationary point under a nonconvex-strongly-concave setting. Under a similar *parameter-free* but a more general nonconvex-nonconcave setting, this paper establishes a local convergence to the local minimax points, especially for the Alt2-EG with increasing timescale separation, from a *non-autonomous* dynamical system perspective.

3. Preliminaries

3.1. Notations and Problem Setting

We use the notation $\mathbf{z} := (\mathbf{x}, \mathbf{y})$ to represent a concatenation of the minimization variable $\mathbf{x} \in \mathbb{R}^{d_1}$ and the maximization

variable $\mathbf{y} \in \mathbb{R}^{d_2}$. The saddle-gradient operator of the objective function f will be denoted by $\mathbf{F} := (\nabla_{\mathbf{x}}f, -\nabla_{\mathbf{y}}f)$. For convenience, we denote the second derivatives of f by $\mathbf{A} = \nabla_{\mathbf{x}\mathbf{x}}^2f$, $\mathbf{B} = \nabla_{\mathbf{y}\mathbf{y}}^2f$, and $\mathbf{C} = \nabla_{\mathbf{x}\mathbf{y}}^2f$. Then, the Jacobian of the saddle-gradient \mathbf{F} can be expressed as

$$\mathbf{H} := D\mathbf{F} = \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} \end{bmatrix} \in \mathbb{R}^{(d_1+d_2) \times (d_1+d_2)}. \quad (1)$$

Note that the second derivatives \mathbf{A} , \mathbf{B} , and \mathbf{C} are matrix valued functions of a $(d_1 + d_2)$ -dimensional input point. However, since the input vector will be clear from the context, we simply use \mathbf{A} , \mathbf{B} , and \mathbf{C} to denote the function values. We denote the set of all eigenvalues of a square matrix \mathbf{A} by $\text{spec}(\mathbf{A})$, the smallest eigenvalue of \mathbf{A} by $\lambda_{\min}(\mathbf{A})$, and the spectral radius of \mathbf{A} by $\rho(\mathbf{A})$.

Most of the time, we will impose the following standard smoothness assumption on f .

Assumption 1 (Smoothness of f). *Let $f \in C^1$, and there exist positive constants $L_{\mathbf{x}}$ and $L_{\mathbf{y}}$ such that, for all $(\mathbf{u}, \mathbf{v}), (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_1+d_2}$,*

$$\begin{aligned} \|\nabla_{\mathbf{x}}f(\mathbf{u}, \mathbf{v}) - \nabla_{\mathbf{x}}f(\mathbf{x}, \mathbf{y})\| &\leq L_{\mathbf{x}}\|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|, \\ \|\nabla_{\mathbf{y}}f(\mathbf{u}, \mathbf{v}) - \nabla_{\mathbf{y}}f(\mathbf{x}, \mathbf{y})\| &\leq L_{\mathbf{y}}\|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|. \end{aligned}$$

This implies that $\|\mathbf{F}(\mathbf{z}) - \mathbf{F}(\mathbf{z}')\| \leq L\|\mathbf{z} - \mathbf{z}'\|$ for $L := \sqrt{L_{\mathbf{x}}^2 + L_{\mathbf{y}}^2}$ and for all $\mathbf{z}, \mathbf{z}' \in \mathbb{R}^{d_1+d_2}$.

For the global analysis in Section 7, we will further impose the following nonconvex-nonconcave condition.

Assumption 2 (Star-convex-star-concave property of f). *Let $f \in C^1$, and f is star-convex-star-concave, i.e., there exists a stationary point $\mathbf{z}^* := (\mathbf{x}^*, \mathbf{y}^*)$ that satisfies*

$$\begin{aligned} f(\mathbf{x}^*, \mathbf{y}) &\geq f(\mathbf{x}, \mathbf{y}) + \langle \nabla_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}), \mathbf{x}^* - \mathbf{x} \rangle, \\ f(\mathbf{x}, \mathbf{y}) &\geq f(\mathbf{x}^*, \mathbf{y}) + \langle \nabla_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{y}), \mathbf{x} - \mathbf{x}^* \rangle, \\ f(\mathbf{x}, \mathbf{y}^*) &\leq f(\mathbf{x}, \mathbf{y}) + \langle \nabla_{\mathbf{y}}f(\mathbf{x}, \mathbf{y}), \mathbf{y}^* - \mathbf{y} \rangle, \\ f(\mathbf{x}, \mathbf{y}) &\leq f(\mathbf{x}, \mathbf{y}^*) + \langle \nabla_{\mathbf{y}}f(\mathbf{x}, \mathbf{y}^*), \mathbf{y} - \mathbf{y}^* \rangle. \end{aligned}$$

for all $\mathbf{x} \in \mathbb{R}^{d_1}$ and $\mathbf{y} \in \mathbb{R}^{d_2}$.

This implies that the Minty variational inequality (MVI) condition (Minty, 1967), i.e., $\langle \mathbf{F}(\mathbf{z}), \mathbf{z} - \mathbf{z}^* \rangle \geq 0$ for all $\mathbf{z} \in \mathbb{R}^{d_1+d_2}$, holds.

3.2. Restricted Schur Complement

Chae et al. (2024b) defined the following matrix that reduces to the standard Schur complement $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top$ if \mathbf{B} is non-degenerate. This is needed in next subsection for expressing the property of the local minimax point.

Definition 1 (Chae et al. (2024b, Definition 4)). *For $f \in C^2$, the restricted Schur complement is defined as $\mathbf{S}_{\text{res}}(\mathbf{H}) := \mathbf{U}^\top(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{U}$, with the matrix \mathbf{U} defined below.*

Let $r := \text{rank}(B)$. As B is symmetric, it is orthogonally diagonalizable into $B = P\Delta P^\top$, where $\Delta = \text{diag}\{\delta_1, \dots, \delta_r, 0, \dots, 0\}$ and P is an orthogonal matrix. Define C_1 and C_2 to be submatrices of CP , corresponding to the r leftmost columns and $d_2 - r$ rightmost columns of CP , respectively. Then, U is defined to be a matrix whose columns form an orthonormal basis for $\mathcal{R}(C_2)^\perp$. The matrix U is not unique in general, but since the spectrum of $S_{\text{res}}(H)$ only matters later, U is selected to be any one of possible choices.

3.3. Local Minimax Point

Jin et al. (2020) introduced the following new notion of local optimality for minimax problems.

Definition 2 (Jin et al. (2020, Definition 14)). *A point $(\mathbf{x}^*, \mathbf{y}^*)$ is said to be a local minimax point if there exists $\delta_0 > 0$ and a function h satisfying $h(\delta) \rightarrow 0$ as $\delta \rightarrow 0$ such that, for any $\delta \in (0, \delta_0]$ and any (\mathbf{x}, \mathbf{y}) satisfying $\|\mathbf{x} - \mathbf{x}^*\| \leq \delta$ and $\|\mathbf{y} - \mathbf{y}^*\| \leq \delta$, we have*

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq \max_{\mathbf{y}' : \|\mathbf{y}' - \mathbf{y}^*\| \leq h(\delta)} f(\mathbf{x}, \mathbf{y}').$$

Nevertheless, most of the existing literature, such as (Fiez & Ratliff, 2021; Jin et al., 2020; Wang et al., 2020), only focused on finding a *strict* local minimax point that is a stationary point satisfying the second-order *sufficient* condition of local minimax point (Jin et al., 2020):

$$\begin{aligned} [\nabla_{\mathbf{x}\mathbf{x}}^2 f - \nabla_{\mathbf{x}\mathbf{y}}^2 f (\nabla_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \nabla_{\mathbf{y}\mathbf{x}}^2 f](\mathbf{x}^*, \mathbf{y}^*) &> \mathbf{0}, \\ \text{and } \nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*) &< \mathbf{0}. \end{aligned}$$

Considering that *non-strict* solutions are prevalent in modern overparameterized model training (Cooper, 2021; Liu et al., 2022), Chae et al. (2024a) recently studied constructing a gradient method that can find a *non-strict* (local) minimax point. More precisely, Chae et al. (2024a) (and also this paper) modestly¹ aim to find a stationary point that satisfies the following second-order necessary condition given in (Chae et al., 2024b, Proposition 3.2):

- (Second-order necessary) For $f \in C^2$, any local minimax point \mathbf{z}^* satisfies $\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \preceq \mathbf{0}$. If the function $h(\delta)$ in Definition 2 satisfies $\limsup_{\delta \rightarrow 0^+} h(\delta)/\delta < \infty$, then $S_{\text{res}}(DF(\mathbf{x}^*, \mathbf{y}^*)) \succeq \mathbf{0}$.

Note that this condition does not require a restrictive non-degeneracy condition on $\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$, unlike that in (Jin

¹Already in minimization, finding a local minimizer is NP-Hard in the worst-case (Murty & Kabadi, 1987). So instead, a gradient descent method is similarly shown to find a stationary point that satisfies the second-order necessary condition of a local minimizer (Lee et al., 2016; 2019).

et al., 2020), although it adds a mild condition on the function $h(\delta)$; see (Chae et al., 2024b, Remark 3.3) on this matter.

Built upon their necessary condition, Chae et al. (2024b) presented the concept of *strict non-minimax* point that one hopes to escape, which is analogous to the *strict saddle* point in minimization (Lee et al., 2016).

Definition 3 (Chae et al. (2024b, Definition 5)). *For $f \in C^2$, a stationary point \mathbf{z}^* is said to be a strict non-minimax point of f if $\lambda_{\min}(S_{\text{res}}(DF(\mathbf{z}^*))) < 0$ or $\lambda_{\min}(-\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}^*)) < 0$. We denote the set of strict non-minimax points by \mathcal{T}^* .*

When necessary, we will impose the following Assumption 3 (or its stronger version below) at stationary points \mathbf{z}^* as (Chae et al., 2024b). This is weaker than the non-degeneracy condition on $\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}^*)$, which was crucial in (Jin et al., 2020; Fiez & Ratliff, 2021; Li et al., 2022).

Assumption 3. *Let $f \in C^2$, and for a stationary point \mathbf{z}^* in consideration, at least one of the matrices $S_{\text{res}}(DF(\mathbf{z}^*))$ and $\nabla_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{z}^*)$ is non-degenerate.*

Assumption 3'. *Assumption 3 holds, and $DF(\mathbf{z}^*)$ is non-degenerate.*

To aid understanding, a simple example with a non-strict local minimax point is provided in Section 8, accompanied with numerical experiment.

3.4. Stability of Dynamical Systems

From a dynamical system perspective, this paper analyzes the *asymptotic (or exponential) stability* and the instability of the method in a form

$$\mathbf{z}_{k+1} = \mathbf{w}_k(\mathbf{z}_k), \quad (2)$$

at an equilibrium \mathbf{z}^* , where k denotes the iteration number.

Definition 4. *The equilibrium point \mathbf{z}^* of (2) is*

- **(Lyapunov) stable** if, for each $\epsilon > 0$, there is $\delta = \delta(\epsilon, k_0) > 0$ such that $\|\mathbf{z}_{k_0} - \mathbf{z}^*\| < \delta$ implies $\|\mathbf{z}_k - \mathbf{z}^*\| < \epsilon$, $\forall k \geq k_0 \geq 0$,
- **asymptotically stable** if it is stable and there is a positive constant $c = c(k_0)$ such that $\mathbf{z}_k \rightarrow \mathbf{z}^*$ as $k \rightarrow \infty$, for all $\|\mathbf{z}_{k_0} - \mathbf{z}^*\| < c$,
- **exponentially stable** if it is asymptotically stable and there are constants $\beta, c > 0$ such that $\|\mathbf{z}_k - \mathbf{z}^*\| \leq e^{-\beta k} \|\mathbf{z}_0 - \mathbf{z}^*\|$, for all $\|\mathbf{z}_0 - \mathbf{z}^*\| < c$, $\forall k \geq 1$,
- **unstable**, if it is not stable.

If a dynamical system is *autonomous*, i.e., \mathbf{w}_k is invariant with respect to the iteration number k , we can characterize

its exponential stability and instability at an equilibrium point \mathbf{z}^* by just examining the spectrum of the Jacobian of \mathbf{w}_k at \mathbf{z}^* as below.

Proposition 3.1 (Galor (2007), Polyak (1987)). *Let $\mathbf{w} \in C^1$, and \mathbf{z}^* be an equilibrium of $\mathbf{z}_{k+1} = \mathbf{w}(\mathbf{z}_k)$. Then,*

- \mathbf{z}^* is exponentially stable iff $\rho(D\mathbf{w}(\mathbf{z}^*)) < 1$
- \mathbf{z}^* is unstable if $\rho(D\mathbf{w}(\mathbf{z}^*)) > 1$.

This paper also analyzes a *non-autonomous* system, i.e., a system \mathbf{w}_k that is not invariant with respect to k , in Section 6. Proposition 3.1 cannot be applied here, so we directly analyze the asymptotic stability and the convergence rate of the considered method.

4. Stability and Limitations of the Two-Timescale EG

In this section, we review the stability of the two-timescale EG and its two limitations in (Chae et al., 2024b).

4.1. Stability of the Two-Timescale EG

Chae et al. (2024b) studied the two-timescale EG:

$$\mathbf{z}_{k+1} = \mathbf{w}_\tau(\mathbf{z}_k) := \mathbf{z}_k - \eta \mathbf{\Lambda}_\tau \mathbf{F}(\mathbf{z}_k - \eta \mathbf{\Lambda}_\tau \mathbf{F}(\mathbf{z}_k)), \quad (3)$$

where $\eta > 0$ is a step size, $\tau \geq 1$ is a fixed timescale separation parameter and $\mathbf{\Lambda}_\tau := \text{diag}\{(1/\tau)\mathbf{I}, \mathbf{I}\}$. Based on Proposition 3.1, the stability of the two-timescale EG (3) at an equilibrium \mathbf{z}^* depends on whether or not the spectral radius of its Jacobian matrix $D\mathbf{w}_\tau$ at \mathbf{z}^* is smaller than 1. However, since directly computing $\rho(D\mathbf{w}_\tau)$ is complicated, Chae et al. (2024b) demonstrated that it is equivalent to examine the relationship between the spectrum of

$$\mathbf{H}_\tau := \mathbf{\Lambda}_\tau D\mathbf{F} = \begin{bmatrix} \frac{1}{\tau}\mathbf{A} & \frac{1}{\tau}\mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} \end{bmatrix}. \quad (4)$$

and the set $\mathcal{P}_\eta := \{(x, y) \in \mathbb{C} \mid (\eta x - \frac{1}{2})^2 + \eta^2 y^2 + \frac{3}{4} < \sqrt{1 + 3\eta^2 y^2}\}$ (see Figure 1) as below.

Proposition 4.1 (Chae et al. (2024b), Proposition 5.5). *Let $f \in C^2$. An equilibrium point \mathbf{z}^* is exponentially stable, i.e., $\rho(D\mathbf{w}_\tau(\mathbf{z}^*)) < 1$, if and only if $\text{spec}(\mathbf{H}_\tau(\mathbf{z}^*)) \subset \mathcal{P}_\eta$.*

This does not necessarily imply that the two-timescale EG is stable at the desired local minimax point, so Chae et al. (2024b) next relate the stability condition and the second-order necessary condition when τ is sufficiently large.

4.2. Eigenvalue Characteristics of \mathbf{H}_τ and its Relation to Second-Order Necessary Condition

Chae et al. (2024b) analyzed the behavior of the eigenvalues of \mathbf{H}_τ in terms of τ as below, under Assumption 3. This

reduces to (Jin et al., 2020, Lemma 40) when we strictly impose the non-degeneracy condition on $\nabla_{\mathbf{y}\mathbf{y}}^2 f$.

Theorem 4.2 (Chae et al. (2024b), Theorem 4.3). *Under Assumption 3, for $\tau \geq 1$ and $\epsilon := 1/\tau$, it is possible to construct continuous functions $\lambda_j(\epsilon)$, $j = 1, \dots, d_1 + d_2$ so that they are the $d_1 + d_2$ complex eigenvalues λ_j of \mathbf{H}_τ in (4) with the following asymptotics as $\epsilon \rightarrow 0^+$:*

- (i) $|\lambda_j - i\sigma_j\sqrt{\epsilon}| = o(\sqrt{\epsilon})$, $j = 1, \dots, q$
 $|\lambda_{j+d_1} + i\sigma_j\sqrt{\epsilon}| = o(\sqrt{\epsilon})$,
- (ii) $|\lambda_{j+q} - \epsilon\mu_j| = o(\epsilon)$, $j = 1, \dots, d_1 - q$,
- (iii) $|\lambda_{j+d_1+q} - \nu_j| = o(1)$, $j = 1, \dots, r$,

where $q := \text{rank}(\mathbf{C}_2)$, which are nonzero for all $\epsilon > 0$, while the $(d_2 - r - q)$ remaining λ_j 's are 0. Here, μ_j are the eigenvalues of $\mathbf{S}_{\text{res}}(\mathbf{H})^2$, ν_j are the nonzero eigenvalues of $-\mathbf{B}$, and σ_j are the singular values of \mathbf{C}_2 .

The key distinction between the spectral analysis of \mathbf{H}_τ in (Jin et al., 2020, Lemma 40) and Theorem 4.2 is the existence of additional type (i) eigenvalues (and the $d_2 - r - q$ number of zero eigenvalues), which are irrelevant to the second-order necessary condition. On the other hand, the type (ii) and type (iii) eigenvalues are associated with $\mathbf{S}_{\text{res}}(\mathbf{H})$ and $-\mathbf{B}$, respectively, so they are directly connected to the second-order necessary condition.

4.3. Stability of the Two-Timescale EG at Local Minimax Points

When the point \mathbf{z}^* satisfies the second-order necessary condition of local minimax points, by Theorem 4.2, the type (ii) and type (iii) eigenvalues lie in the set \mathcal{P}_η , for a sufficiently large τ (or equivalently, for a sufficiently small ϵ), as desired. On the other hand, consider the *strict non-minimax* point \mathbf{z}^* that, by definition, does not satisfy the second-order necessary condition. Then, at least one of the type (ii) and type (iii) eigenvalues lie outside the closure of \mathcal{P}_η , even with a large τ , as one wishes.

Let us now focus particularly on the eigenvalues of type (i) in Theorem 4.2. As they are irrelevant to the second-order necessary condition, ideally we want these to be contained in the set \mathcal{P}_η , for a sufficiently large τ , and do not cause instability of the two-timescale EG.

What we know from Theorem 4.2 is that the type (i) eigenvalue $\lambda_j(\epsilon)$ converges to 0 as $\epsilon \rightarrow 0^+$, and asymptotically approaches the imaginary axis. Since the target set \mathcal{P}_η , illustrated in Figure 1, has the imaginary axis as its tangential line at the origin, Theorem 4.2 alone is not sufficient to determine the stability of the two-timescale EG.

²Although $\mathbf{S}_{\text{res}}(\mathbf{H})$ is not unique due to its non-unique choice of matrix \mathbf{U} , its eigenvalues remain the same regardless of the choice of \mathbf{U} due to matrix similarity.

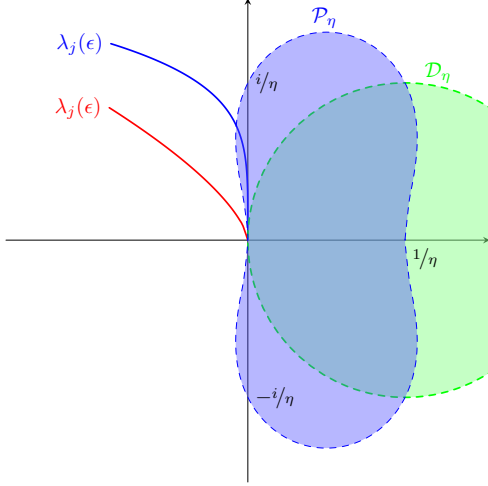


Figure 1. A target set \mathcal{P}_η of the two-timescale EG, and two representative scenarios of type (i) eigenvalues, approaching 0 from the left-half plane, in Theorem 4.2. (As a comparison, we added a similarly derived target set $\mathcal{D}_\eta := \{(x, y) \in \mathbb{C} \mid (\eta x - 1)^2 + \eta^2 y^2 < 1\}$ of the two-timescale GDA $\mathbf{z}_{k+1} = \mathbf{z}_k - \eta \mathbf{\Lambda}_\tau \mathbf{F}(\mathbf{z}_k)$.)

Figure 1 illustrates two representative examples of the type (i) eigenvalue asymptotics, where the blue-colored one converges to 0 from inside the set \mathcal{P}_η , while the other (colored red) converges to 0 from outside the set.³ This observation necessitated Chae et al. (2024b) to further investigate the curvature of both the type (i) eigenvalue $\lambda_j(\epsilon)$ and the target set \mathcal{P}_η around the origin, for completing the stability analysis of the two-timescale EG. This, however, only revealed the fact that the two-timescale EG cannot avoid the red-colored case, resulting in instability for such corresponding non-strict local minimax points.

4.4. Two Limitations of the Two-Timescale EG

We summarize the two limitations of the two-timescale EG, which we address in Sections 5 and 6, respectively.

- 1) There are type (i) eigenvalues $\lambda_j(\epsilon)$ that are not related to the second-order necessary condition but not contained in the set \mathcal{P}_η , even for a large τ , which consequently leads to instability of the two-timescale EG at certain non-strict local minimax points.
- 2) Type (ii) and (iii) eigenvalues become associated with the second-order necessary condition for sufficiently large τ , but its specific value is not available in practice.

³The GDA is unstable even for the blue-colored eigenvalue asymptotic, as the target set \mathcal{D}_η of GDA does not cover a region nearby the imaginary axis as much as the set \mathcal{P}_η of EG. This is illustrated in Figure 1, which shows the superiority of EG over GDA in this context (Chae et al., 2024b).

5. Addressing the First Limitation: Double-Step Alternating Update

In this section, we investigate two stabilizing approaches, namely explicit regularization and implicit regularization (by alternating update), that shift the type (i) eigenvalues towards the set \mathcal{P}_η without affecting the eigenvalue asymptotics of the type (ii) and (iii). In specific, we look for ways to shift the type (i) eigenvalues to the right in a complex plane by the order of $\Omega(\sqrt{\epsilon})$, needed by Theorem 4.2.

5.1. Explicit Regularization for Two-Timescale EG

Since the type (i) eigenvalues are critically related to $\nabla_{\mathbf{y}\mathbf{y}}^2 f$, we consider explicitly adding a regularization term $-\frac{c}{\gamma} \|\mathbf{y}\|^2$ to an objective function $f(\mathbf{x}, \mathbf{y})$, where c and γ are positive constants, while the latter depends on τ . The saddle gradient operator of the resulting regularized function is

$$\mathbf{F}_{\text{reg},\gamma}(\mathbf{x}, \mathbf{y}) := \left(\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) + \frac{2c}{\gamma} \mathbf{y} \right).$$

Although this explicit regularization has a drawback that the stationary points of the regularized function differ from those of the original function f , we proceed as it is straightforward and provides an insight on how we should choose γ in terms of τ .

We first consider the case “ $\gamma = \tau$ ”, where the Jacobian of the timescaled $\mathbf{\Lambda}_\tau \mathbf{F}_{\text{reg},\tau}$ at its stationary point is

$$\mathbf{H}_{\text{reg}1,\tau} := \mathbf{\Lambda}_\tau D\mathbf{F}_{\text{reg},\tau} = \begin{bmatrix} \epsilon \mathbf{A} & \epsilon \mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} + 2c\epsilon \mathbf{I} \end{bmatrix} \quad (5)$$

where $\epsilon = 1/\tau$. The eigenvalues of $\mathbf{H}_{\text{reg}1,\tau}$ have the following asymptotic behaviors.

Proposition 5.1. *Under the same setting as Theorem 4.2, the eigenvalues λ_j of $\mathbf{H}_{\text{reg}1,\tau}$ behave the same as those of \mathbf{H}_τ in Theorem 4.2, which are nonzero for all but finitely many $\epsilon > 0$, except for the eigenvalues λ_{j+d_1+r+q} being $2c\epsilon$ for $j = 1, \dots, d_2 - r - q$ (instead of 0 in Theorem 4.2).*

The result is what we anticipated as we perturbed only in the order of $o(\sqrt{\epsilon})$. So, we next consider the choice “ $\gamma = \sqrt{\tau}$ ”, where the eigenvalues of the corresponding Jacobian

$$\mathbf{H}_{\text{reg}2,\tau} := \mathbf{\Lambda}_\tau D\mathbf{F}_{\text{reg},\sqrt{\tau}} = \begin{bmatrix} \epsilon \mathbf{A} & \epsilon \mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} + 2c\sqrt{\epsilon} \mathbf{I} \end{bmatrix}$$

have the following asymptotic behaviors.

Theorem 5.2. *Under the same setting as Theorem 4.2, the eigenvalues λ_j of $\mathbf{H}_{\text{reg}2,\tau}$ behave the same as those of \mathbf{H}_τ in Theorem 4.2, except for the type (i) eigenvalues being*

$$(i) \quad \left| \lambda_j - \left(c + \sqrt{c^2 - \sigma_j^2} \right) \sqrt{\epsilon} \right| = o(\sqrt{\epsilon}).$$

$$\left| \lambda_{j+d_1} - \left(c - \sqrt{c^2 - \sigma_j^2} \right) \sqrt{\epsilon} \right| = o(\sqrt{\epsilon}),$$

for $j = 1, \dots, q$, which are nonzero for all but finitely many $\epsilon > 0$, and the eigenvalues λ_{j+d_1+r+q} being $2c\sqrt{\epsilon}$ for $j = 1, \dots, d_2 - r - q$.

Here, regardless of the radicand of $\sqrt{c^2 - \sigma_j^2}$ being positive or negative, it is straightforward that the type (i) eigenvalues converge to 0 as $\epsilon \rightarrow 0^+$ from the right-half plane (but not necessarily approaches the positive real axis asymptotically). Therefore, they now lie in the set \mathcal{P}_η , unlike the previous results in Theorem 4.2 and Proposition 5.1. This implies that when an appropriate explicit regularization of the order $\Theta(\sqrt{\epsilon})$ is provided, we will not encounter the eigenvalue examples in Figure 1, approaching 0 from the left-half plane. Nevertheless, as mentioned before, stationary points of $\mathbf{F}_{\text{reg}, \sqrt{\epsilon}}$ differ from those of \mathbf{F} , which is not usually desirable, so we next investigate an approach that provides a similar regularization effect without changing stationary points.

5.2. Double-Step Alternating Update in Two-Timescale EG

Consider an operator for an alternating update:

$$\mathbf{F}_{\text{alt}, \gamma}(\mathbf{x}, \mathbf{y}) := \left(\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}} f \left(\mathbf{x} - \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), \mathbf{y} \right) \right),$$

where $\nabla_{\mathbf{y}} f$ is computed after the update of \mathbf{x} . Here, $\eta > 0$ is a step size and γ is a positive constant that depends on τ . Note that $\mathbf{F}_{\text{alt}, \gamma}$ reduces to \mathbf{F} as $\gamma \rightarrow \infty$. In addition, it is obvious that \mathbf{F} and $\mathbf{F}_{\text{alt}, \gamma}$ share same stationary points, unlike $\mathbf{F}_{\text{reg}, \gamma}$.

Before we proceed to analyze the eigenvalue asymptotics of $\mathbf{F}_{\text{alt}, \gamma}$, we illustrate the corresponding alternating method:

$$\begin{aligned} \mathbf{z}_{k+1} &= \mathbf{w}_{\text{alt}, \tau_k, \gamma_k}(\mathbf{z}_k) \\ &:= \mathbf{z}_k - \eta \mathbf{\Lambda}_{\tau_k} \mathbf{F}_{\text{alt}, \gamma_k}(\mathbf{z}_k - \eta \mathbf{\Lambda}_{\tau_k} \mathbf{F}_{\text{alt}, \gamma_k}(\mathbf{z}_k)), \end{aligned} \quad (6)$$

which is constructed by replacing \mathbf{F} by $\mathbf{F}_{\text{alt}, \gamma}$ in (3). Here, we use τ_k and γ_k , instead of τ and γ , so that they can change over iteration, which will be useful in Section 6. The proposed update (6) takes a stationary point of \mathbf{F} as its equilibrium, see Appendix B.4. We further write down this method in terms of f , in Algorithm 1. Note that this method reduces to the (simultaneous-update) two-timescale EG in (Chae et al., 2024b) if we choose $(\tau_k, \gamma_k) = (\tau, \infty)$ for all k .

We would like to emphasize here that Algorithm 1 deviates from the usual alternating update scheme, as we allow τ_k and γ_k to take different values. In alignment with the spectral analysis in Section 5.1, we will soon show that allowing $\tau_k \neq \gamma_k$ is essential for our purpose. To highlight this deviation, we name this method in Algorithm 1 as a *double-step*

Algorithm 1 Double-step alternating extragradient method with timescale separation (Alt2-EG-TS)

Input: $\mathbf{x}_0 \in \mathbb{R}^{d_1}, \mathbf{y}_0 \in \mathbb{R}^{d_2}, \tau_k \in (1, \infty), \gamma_k \in (1, \infty)$
for $k = 0, 1, \dots$ **do**

$$\mathbf{u}_k = \mathbf{x}_k - \frac{\eta}{\tau_k} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)$$

$$\mathbf{v}_k = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f \left(\mathbf{x}_k - \frac{\eta}{\gamma_k} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k), \mathbf{y}_k \right)$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \frac{\eta}{\tau_k} \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)$$

$$\mathbf{y}_{k+1} = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f \left(\mathbf{u}_k - \frac{\eta}{\gamma_k} \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k), \mathbf{v}_k \right)$$

end for

alternating extragradient method with timescale separation (Alt2-EG-TS). In the rest of this section, we consider fixed constants $(\tau_k, \gamma_k) = (\tau, \gamma)$ for all k , and we name the corresponding method as a double-step alternating extragradient method with *fixed* timescale separation (Alt2-EG-FTS).

Remark 5.1. A concurrent work by Lee et al. (2024) also studied and improved the alternating GDA, especially by adding extrapolation steps, which provided an accelerated rate of convergence in a strongly-convex-strongly-concave problem. Interestingly, our double-step alternating scheme can be also viewed as taking an extrapolation step, like (Lee et al., 2024), since we can rewrite our step as

$$\mathbf{x}_k - \frac{\eta}{\gamma_k} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k) = \left(1 - \frac{\tau_k}{\gamma_k} \right) \mathbf{x}_k + \frac{\tau_k}{\gamma_k} \mathbf{u}_k,$$

where the choice of $\frac{\tau_k}{\gamma_k}$ being larger than 1, in the next subsection, yields extrapolation.

5.3. Implicit Regularization via Double-Step Alternating Update

This section investigates an implicit regularization provided by the double-step alternating update. Similar to Proposition 5.1, the choice “ $\gamma = \tau$ ” in $\mathbf{F}_{\text{alt}, \gamma}$ does not lead to any change in the eigenvalue asymptotics; see Appendix B.5.

We now focus on the choice “ $\gamma = \sqrt{\tau}$ ”, hoping for an appropriate shift of the type (i) eigenvalues. The Jacobian of the timescaled $\mathbf{H}_{\text{alt}, \tau} := \mathbf{\Lambda}_{\tau} D \mathbf{F}_{\text{alt}, \sqrt{\tau}}$, at its stationary point is, where $\epsilon = 1/\tau$,

$$\mathbf{H}_{\text{alt}, \tau} = \begin{bmatrix} \epsilon \mathbf{A} & \epsilon \mathbf{C} \\ -\mathbf{C}^\top + \eta \sqrt{\epsilon} \mathbf{C}^\top \mathbf{A} & -\mathbf{B} + \eta \sqrt{\epsilon} \mathbf{C}^\top \mathbf{C} \end{bmatrix},$$

whose eigenvalues have the following asymptotic behaviors.

Theorem 5.3. *Under the same setting as Theorem 4.2, the eigenvalues λ_j of $\mathbf{H}_{\text{alt}, \tau}$ behave the same as those of \mathbf{H}_{τ} in Theorem 4.2, except for the type (i) eigenvalues being*

$$(i) \left| \lambda_j - \left(\frac{\eta\sigma_j^2}{2} + \frac{\sqrt{\eta^2\sigma_j^4 - 4\sigma_j^2}}{2} \right) \sqrt{\epsilon} \right| = o(\sqrt{\epsilon}),$$

$$\left| \lambda_{j+d_1} - \left(\frac{\eta\sigma_j^2}{2} - \frac{\sqrt{\eta^2\sigma_j^4 - 4\sigma_j^2}}{2} \right) \sqrt{\epsilon} \right| = o(\sqrt{\epsilon}),$$

for $j = 1, \dots, q$.

Similar to Theorem 5.2 for the explicit regularization with $\gamma = \sqrt{\tau}$, the type (i) eigenvalues here approach 0 from the right-half plane and thus from inside the set \mathcal{P}_η . This implies that the double-step alternating update with $\gamma = \sqrt{\tau}$ implicitly induces a regularizing behavior similar to the explicit regularization with $\gamma = \sqrt{\tau}$. We would like to emphasize again that a clear advantage of the double-step alternating, over the explicit regularization, is that the resulting method takes a stationary point of F as its equilibrium.

5.4. Stability Analysis of Alt2-EG-FTS

Built upon the previous spectral analysis, we show that, under Assumption 3', a stationary point that satisfies the second-order necessary condition is *exponentially stable* for the Alt2-EG-FTS with $(\tau_k, \gamma_k) = (\tau, \sqrt{\tau})$, for a sufficiently large τ , under mild conditions.

Theorem 5.4. *Suppose Assumptions 1 and 3' hold and $f \in C^2$. Then, an equilibrium point z^* satisfies $S_{\text{res}} \succeq 0$ and $B \preceq 0$ if and only if there exists some τ^* such that, for any $\tau > \tau^*$, z^* is an exponentially stable point of Alt2-EG-FTS with $(\tau_k, \gamma_k) = (\tau, \sqrt{\tau})$, for any step size $0 < \eta < 1/L$.*

In addition, we show that, under Assumption 3, the Alt2-EG-FTS with $(\tau_k, \gamma_k) = (\tau, \sqrt{\tau})$ almost surely escapes the strict non-minimax points, which we hope to avoid.

Theorem 5.5. *Let z^* be a strict non-minimax point, i.e., $z^* \in \mathcal{T}^*$. Under Assumptions 1, 3, and $0 < \eta < (\sqrt{5}-1)/2\sqrt{2}L$, there exists $\tau^* > 0$ such that for any $\tau > \tau^*$, the set of initial points that converge to z^* by Alt2-EG-FTS with $(\tau_k, \gamma_k) = (\tau, \sqrt{\tau})$ has measure zero. Moreover, if \mathcal{T}^* is finite, then there exists $\tau^* > 0$ such that for any $\tau > \tau^*$, $\mu(\{z_0 : \lim_{k \rightarrow \infty} w_{\text{alt}, \tau, \sqrt{\tau}}^k(z_0) \in \mathcal{T}^*\}) = 0$.*

We ultimately want that there are no points, other than strict non-minimax points, that the method escapes almost surely. Fortunately, the stability results above imply that such can be obtained if we choose step size $0 < \eta < (\sqrt{5}-1)/2\sqrt{2}L$ for the Alt2-EG-FTS with $(\tau_k, \gamma_k) = (\tau, \sqrt{\tau})$ and a sufficiently large τ . This statement will be further generalized to a global statement in Section 7.

6. Addressing the Second Limitation: Increasing Timescale Separation

In this section, we consider increasing the timescale separation, and specifically, we analyze the stability of the double-

step alternating extragradient with *increasing* timescale separation (Alt2-EG-ITS) from a *non-autonomous* dynamical system perspective.

6.1. Increasing Timescale Separation for Alt2-EG

Choosing a properly large value of τ is practically infeasible, so we consider increasing the coefficients (τ_k, γ_k) of Alt2-EG-ITS indefinitely as k increases. We first let $\gamma_k = \sqrt{\tau_k}$, based on the arguments in Section 5, and the only thing left to determine is the rate at which we should increase τ_k .

Our global convergence analysis in Section 7 requires the sequence τ_k to not increase faster than \sqrt{k} to warrant global convergence; see Theorem 7.1. Therefore, we focused on the choice $\tau_k = k^{1/(2+2c)} \approx \sqrt{k}$ for any positive constant c , and we leave further investigation on τ_k as future work.

6.2. Stability Analysis of Alt2-EG-ITS

Since Proposition 3.1 cannot be applied to the *non-autonomous* Alt2-EG-ITS, we derive an analogous *asymptotic* stability result that is specifically designed for the Alt2-EG-ITS ($w_{\text{alt}, \tau_k, \gamma_k}$), under Assumption 3'.

Theorem 6.1. *Suppose Assumptions 1 and 3' hold and $f \in C^3$. Let z^* be an equilibrium point satisfies $S_{\text{res}} \succeq 0$ and $B \preceq 0$. Then, z^* is an asymptotically stable point of Alt2-EG-ITS with $(k^{1/(2+2c)}, k^{1/(4+4c)})$ for any $c > 0$ and $0 < \eta < 1/L$, with a rate $O\left(\frac{1}{\sqrt{k}}e^{-2\sqrt{k}}\right)$.*

Note that, due to decreasing step sizes (as we iteratively increase both τ_k and γ_k), we have a rate of convergence slower than the *exponential* rate of Alt2-EG-FTS in Theorem 5.4, while not requiring one to choose a properly large value of τ as in Theorem 5.4. Moreover, we needed a slightly stronger condition $w_{\text{alt}, \tau_k, \gamma_k} \in C^2$ (and thus $f \in C^3$), compared to $w \in C^1$ in Proposition 3.1. We leave relaxing such condition as a future work.

Our next step would be to investigate the instability of Alt2-EG-ITS at strict non-minimax points that we would like to avoid, as for the Alt2-EG-FTS in the previous section. We, however, leave this as a future work. This is because Theorem 5.5 for the Alt2-EG-FTS uses the stable manifold theorem (Shub, 1987, Theorem III.7) for an *autonomous* system, and to the best of our knowledge, how one can generalize it to a *non-autonomous* system is not known yet.

7. Global Convergence of Alt2-EG-TS

So far, our convergence and avoidance results remained local. To generalize these statements globally, we first show that both Alt2-EG-FTS and Alt2-EG-ITS globally find first-order stationary points, under Assumption 2.

Theorem 7.1. *Under Assumptions 1 and 2, consider the*

Alt2-EG-TS methods with $\tau_k \geq 1$ and $\gamma_k \geq 1$. Then,

- Alt2-EG-FTS with (τ, γ) and $0 < \eta < \frac{\sqrt{\gamma}}{\sqrt{2+\gamma}L}$ satisfies $\lim_{k \rightarrow \infty} \|\mathbf{F}_{\text{alt},\gamma}(\mathbf{x}_k, \mathbf{y}_k)\| = 0$,
- Alt2-EG-ITS with (τ_k, γ_k) and $0 < \eta < \frac{1}{L}$, for any sequence τ_k satisfying $\sum_{k=0}^{\infty} \frac{1}{\tau_k^2} = \infty$, satisfies $\liminf_{k \rightarrow \infty} \|\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{x}_k, \mathbf{y}_k)\| = 0$.

Although $\lim_{k \rightarrow \infty} \|\mathbf{F}_{\text{alt},\gamma}(\mathbf{x}_k, \mathbf{y}_k)\| = 0$ only implies that any accumulation point of the iterates $\{\mathbf{x}_k, \mathbf{y}_k\}$ is a stationary point (see Appendix D.2), it is possible that the iterates of Alt2-EG-FTS converge to a stationary point. For such convergent case, based on Theorems 5.5 and 7.1 and under their settings (such as Assumptions 1, 2 and 3), we have the following global statement.

- The iterates of the Alt2-EG-FTS with sufficiently large τ globally and almost surely converges to a stationary point that satisfies the second-order necessary condition of local minimax point.

(Note that the invertibility of $D\mathbf{F}$, in Assumption 3', is not needed here, which was needed in analyzing the exponential stability in Theorem 5.4. This implies that Alt2-EG-FTS can even converge to *non-strict* local minimax points even with *degenerate* $D\mathbf{F}$.)

The two-timescale EG (Chae et al., 2024b) does not achieve this statement for some *non-strict* local minimax points, as discussed in Section 4.2. Therefore, our result is a lot closer to being analogous to the well-known result of (Lee et al., 2016; 2019), in minimization, that the iterates of the gradient descent method (if they converge) globally and almost surely converges to a stationary point that satisfies the second-order necessary condition of local minimizer.

On the other hand, since we do not have a result of avoiding a strict non-minimax point for the Alt2-EG-ITS, unlike Theorem 5.5 for Alt2-EG-FTS, we have a relatively weaker global statement, based on Theorems 6.1 and 7.1 and under their settings (such as Assumptions 1, 2, 3' and $f \in C^3$).

- Alt2-EG-ITS globally finds a stationary point, and it also locally converges to a stationary point that satisfies the second-order necessary condition of the local minimax point with non-degenerate $D\mathbf{F}$.

8. Example and Experiment

We consider a simple example of a non-strict local minimax point to verify our theory.

Example 1. Consider the function $f(x, y) = -x^2 + 2xy$. This has a unique stationary point $(0, 0)$, which is a non-strict local minimax point. Since this optimal point further

satisfies Assumption 3', by Theorems 5.4 and 6.1, it is asymptotically stable for both Alt2-EG-FTS with $(\tau, \sqrt{\tau})$ and sufficiently large τ , and Alt2-EG-ITS with $(k^{1/(2+2c)}, k^{1/(4+4c)})$ for any $c > 0$, whereas it is unstable for the (vanilla) two-timescale EG; see Section E.1 for the proof.

We ran our proposed Alt2-EG-FTS with $(\tau, \sqrt{\tau})$ and Alt2-EG-ITS with $(k^{1/(2+2c)}, k^{1/(4+4c)})$, and compared them with existing two-timescale gradient methods (GDA-FTS and EG-FTS). We also performed Alt2-EG-TS with (τ, τ) , named Alt1-EG-FTS here, which uses a standard alternating scheme that differs from our double-step alternating scheme. We consider two choices of $\tau = 2$ and 20, where $\tau = 20$ corresponds to a sufficiently large τ , whereas $\tau = 2$ leads to an insufficient timescale separation. As expected, our Alt2-EG-FTS with $\tau = 20$ and Alt2-EG-ITS converge to the local minimax point, while the others do not.

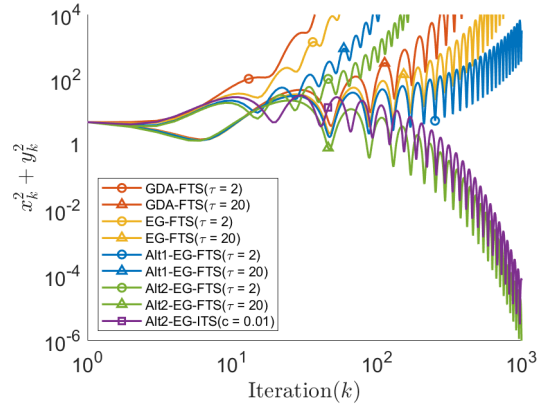


Figure 2. Numerical results with $f(x, y) = -x^2 + 2xy$.

9. Conclusion

We proposed to incorporate double-step alternating update and increasing timescale separation schemes into two-timescale extragradient method—the first method developed for finding *non-strict* local minimax points—to address its two limitations (summarized in Section 4.4). We have then demonstrated that the proposed Alt2-EG-TS method is capable of finding non-strict local minimax points, which cannot be found by existing methods. Therefore, our work, built upon the initial step of (Chae et al., 2024b), is a step closer to establishing a convergence theory in minimax optimization, comparable to the well established theory of gradient descent in minimization problems (Lee et al., 2016; 2019).

Yet, our analysis requires mild but somewhat restrictive conditions, such as Assumption 3, and the escaping behavior around strict non-minimax points for the increasing timescale separation scheme remains unknown. We leave investigating these issues as an interesting future work.

Acknowledgments

This work was supported in part by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No.2019R1A5A1028324, 2022R1C1C1003940), and the Samsung Science & Technology Foundation grant (No. SSTF-BA2101-02).

Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

References

- Apostol, T. M. *Calculus, Volume 1*. John Wiley & Sons, 1991.
- Chae, J., Kim, K., and Kim, D. Two-timescale extragradient for finding local minimax points. In *Proc. Intl. Conf. on Learning Representations*, 2024a.
- Chae, J., Kim, K., and Kim, D. Two-timescale extragradient for finding local minimax points. In *Proc. Intl. Conf. on Learning Representations*, 2024b.
- Cooper, Y. Global minima of overparameterized neural networks. *SIAM Journal on Mathematics of Data Science*, 3(2):676–691, 2021.
- Daskalakis, C. and Panageas, I. The limit points of (optimistic) gradient descent in min-max optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- Fiez, T. and Ratliff, L. Local convergence analysis of gradient descent ascent with finite timescale separation. In *Proc. Intl. Conf. on Learning Representations*, 2021.
- Fiez, T., Chasnov, B., and Ratliff, L. Implicit learning dynamics in Stackelberg games: Equilibria characterization, convergence analysis, and empirical study. In *International Conference on Machine Learning*, pp. 3133–3144. PMLR, 2020.
- Foret, P., Kleiner, A., Mobahi, H., and Neyshabur, B. Sharpness-aware minimization for efficiently improving generalization. In *Proc. Intl. Conf. on Learning Representations*, 2021.
- Galor, O. *Discrete dynamical systems*. Springer Science & Business Media, 2007.
- Goodfellow, I. NIPS 2016 Tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial networks. In *Neural Info. Proc. Sys.*, 2014.
- Heusel, M., Ramsauer, H., Nessler, T. U. B., and Hochreiter, S. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *Neural Info. Proc. Sys.*, 2017.
- Jin, C., Netrapalli, P., and Jordan, M. I. What is local optimality in nonconvex-nonconcave minimax optimization? In *Proc. Intl. Conf. Mach. Learn.*, 2020.
- Lee, J., Cho, H., and Yun, C. Fundamental benefit of alternating updates in minimax optimization. In *International Conference on Machine Learning*, 2024.
- Lee, J. D., Simchowitz, M., Jordan, M. I., and Recht, B. Gradient descent only converges to minimizers. In *Conference on learning theory*, pp. 1246–1257. PMLR, 2016.
- Lee, J. D., Panageas, I., Piliouras, G., Simchowitz, M., Jordan, M. I., and Recht, B. First-order methods almost always avoid strict saddle points. *Mathematical Programming*, 176:311–337, 2019.
- Li, H., Farnia, F., Das, S., and Jadbabaie, A. On convergence of gradient descent ascent: A tight local analysis. In *International Conference on Machine Learning*, pp. 12717–12740. PMLR, 2022.
- Li, X., Yang, J., and He, N. Tiada: A time-scale adaptive algorithm for nonconvex minimax optimization. In *The Eleventh International Conference on Learning Representations*, 2023.
- Lin, T., Jin, C., and Jordan, M. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020.
- Liu, C., Zhu, L., and Belkin, M. Loss landscapes and optimization in over-parameterized non-linear systems and neural networks. *Applied and Computational Harmonic Analysis*, 59:85–116, 2022.
- Madry, A., Makelov, A., Schmidt, L., Tsipras, D., and Vladu, A. Towards deep learning models resistant to adversarial attacks. In *Proc. Intl. Conf. on Learning Representations*, 2018.
- Martinez, N., Bertran, M., and Sapiro, G. Minimax pareto fairness: A multi objective perspective. In *International Conference on Machine Learning*, pp. 6755–6764. PMLR, 2020.

- Mescheder, L., Geiger, A., and Nowozin, S. Which training methods for GANs do actually converge. In *Proc. Intl. Conf. Mach. Learn.*, 2018.
- Minty, G. J. On the generalization of a direct method of the calculus of variations. *Bull. Amer. Math. Soc.*, 73: 314–321, 1967.
- Murty, K. G. and Kabadi, S. N. Some NP-complete problems in quadratic and nonlinear programming. *Mathematical Programming*, 39(2):117–129, 1987.
- Nesterov, Y. *Lectures on convex optimization*, volume 137 of *Springer Optimization and Its Applications*. Springer, second edition, 2018.
- Polyak, B. T. *Introduction to optimization*. New York, Optimization Software,, 1987.
- Shub, M. *Global Stability of Dynamical Systems*. Springer, 1987.
- Tao, T. *Analysis II*. Springer, third edition, 2016.
- Wai, H.-T., Yang, Z., Wang, Z., and Hong, M. Multi-agent reinforcement learning via double averaging primal-dual optimization. *Advances in Neural Information Processing Systems*, 31, 2018.
- Wang, Y., Zhang, G., and Ba, J. On solving minimax optimization locally: A follow-the-ridge approach. In *Proc. Intl. Conf. on Learning Representations*, 2020.
- Zedek, M. Continuity and location of zeros of linear combinations of polynomials. *Proceedings of the American Mathematical Society*, 16(1):78–84, 1965.
- Zhang, G., Wang, Y., Lessard, L., and Grosse, R. B. Near-optimal local convergence of alternating gradient descent-ascent for minimax optimization. In *International Conference on Artificial Intelligence and Statistics*, pp. 7659–7679. PMLR, 2022.

A. Proofs for Section 3

A.1. Proofs for Proposition 3.1

In proving Proposition 3.1, we need the following stable manifold theorem.

Lemma A.1 (Shub (1987, Theorem III.7)). *Let z^* be a fixed point for the C^r local diffeomorphism $\phi : U \rightarrow E$ where U is a neighborhood of z^* in the Banach space E . Suppose that $E = E_{cs} \oplus E_u$, where E_{cs} is the invariant subspace corresponding to the eigenvalues of $D\phi(z^*)$ whose magnitude is less than or equal to 1, and E_u is the invariant subspace corresponding to eigenvalues of $D\phi(z^*)$ whose magnitude is greater than 1. Then there exists a C^r embedded disc W_{loc}^u tangent to the E_u at z^* called the local unstable center manifold. Additionally, there exists a neighborhood B of z^* such that $W_{loc}^u = \{z \in B \mid \phi^k(z) \in B \text{ for all } k \leq 0 \text{ and } d(\phi^k(z), z^*) \text{ tends to zero exponentially}\}$ where $\phi^{-1} : W_{loc}^u \rightarrow W_{loc}^u$ is a contraction mapping.*

Proof of Proposition 3.1. The proof of the first statement can be found in both Galor (2007, Theorem 4.8) and Polyak (1987, Theorem 2.1.2.1), and we are only left to prove the second statement.

Suppose that $\rho(D\mathbf{w}(z^*)) > 1$. Then, by Lemma A.1, there exists a disk W_{loc}^u . Clearly, z^* is contained in W_{loc}^u , since $\mathbf{w}^k(z^*) = z^*$ for any $k \leq 0$. Take $\epsilon = \frac{1}{2} \max_{z \in W_{loc}^u} d(z, z^*)$. Then, there exists $z \in W_{loc}^u$ such that $z \notin D_\epsilon(z^*)$, where $D_\epsilon(z^*)$ is a disc centered at z^* with radius ϵ .

For the sake of contradiction, suppose that z^* of dynamics $\mathbf{w}(\cdot)$ is (Lyapunov) stable. Then, for ϵ discussed above and any given $k_0 \in \mathbb{N}$, there exists $\delta > 0$ such that $\|z_{k_0} - z^*\| < \delta$ implies $\|z_k - z^*\| < \epsilon$ for all $k \geq k_0$. Then, by the definition of W_{loc}^u , there exists $n_0 \in \mathbb{N}$ such that $\mathbf{w}^{-n}(z) \in D_\delta(z^*)$ for all $n \geq n_0$. Then, for $m := \max\{k_0, n_0\} + 1$, we have $z = \mathbf{w}^m(\mathbf{w}^{-m}(z)) \in \mathbf{w}^m(D_\delta(z^*)) \subset D_\epsilon(z^*)$, which is absurd. Therefore, we conclude that $\rho(D\mathbf{w}(z^*)) > 1$ implies that z^* is not stable. □

B. Proofs for Section 5

In proving theorems in Section 5, we need the following lemmas.

Lemma B.1. *Under Assumption 1, the operator $\mathbf{F}_{alt,\gamma}$ with $\gamma > 0$ satisfies, for all $(\mathbf{u}, \mathbf{v}), (\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{d_1+d_2}$,*

$$\|\mathbf{F}_{alt,\gamma}(\mathbf{u}, \mathbf{v}) - \mathbf{F}_{alt,\gamma}(\mathbf{x}, \mathbf{y})\| \leq \sqrt{L_x^2 + L_y^2 \left(1 + \frac{1}{\gamma}\right) \left(1 + L_x^2 \frac{\eta^2}{\gamma}\right)} \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|.$$

Proof. For simplicity, let us denote $\bar{\mathbf{u}} := \mathbf{u} - \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{u}, \mathbf{v})$ and $\bar{\mathbf{x}} := \mathbf{x} - \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$. Then, we have

$$\begin{aligned} \|\mathbf{F}_{alt,\gamma}(\mathbf{u}, \mathbf{v}) - \mathbf{F}_{alt,\gamma}(\mathbf{x}, \mathbf{y})\|^2 &= \|\nabla_{\mathbf{x}} f(\mathbf{u}, \mathbf{v}) - \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})\|^2 + \|\nabla_{\mathbf{y}} f(\bar{\mathbf{u}}, \mathbf{v}) - \nabla_{\mathbf{y}} f(\bar{\mathbf{x}}, \mathbf{y})\|^2 \\ &\leq L_x^2 \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2 + L_y^2 \|(\bar{\mathbf{u}}, \mathbf{v}) - (\bar{\mathbf{x}}, \mathbf{y})\|^2 \\ &\leq L_x^2 \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2 + L_y^2 \|\bar{\mathbf{u}} - \bar{\mathbf{x}}\|^2 + L_y^2 \|\mathbf{v} - \mathbf{y}\|^2 \end{aligned}$$

Since $\bar{\mathbf{u}} - \bar{\mathbf{x}} = \mathbf{u} - \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{u}, \mathbf{v}) - \mathbf{x} + \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$ holds, and using the Young's inequality

$$\|\mathbf{a} + \mathbf{b}\|^2 \leq \left(1 + \frac{1}{\gamma}\right) \|\mathbf{a}\|^2 + (1 + \gamma) \|\mathbf{b}\|^2,$$

we have

$$\begin{aligned} \|\bar{\mathbf{u}} - \bar{\mathbf{x}}\|^2 &\leq \left(1 + \frac{1}{\gamma}\right) \|\mathbf{u} - \mathbf{x}\|^2 + (1 + \gamma) \frac{\eta^2}{\gamma^2} \|\nabla_{\mathbf{x}} f(\mathbf{u}, \mathbf{v}) - \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})\|^2 \\ &\leq \left(1 + \frac{1}{\gamma}\right) \|\mathbf{u} - \mathbf{x}\|^2 + L_x^2 \frac{\eta^2}{\gamma} \left(1 + \frac{1}{\gamma}\right) \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2. \end{aligned}$$

Therefore, we have

$$\begin{aligned}
 & \| \mathbf{F}_{\text{alt},\gamma}(\mathbf{u}, \mathbf{v}) - \mathbf{F}_{\text{alt},\gamma}(\mathbf{x}, \mathbf{y}) \|^2 \\
 & \leq L_x^2 \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2 + L_y^2 \left(1 + \frac{1}{\gamma}\right) \|\mathbf{u} - \mathbf{x}\|^2 + L_x^2 L_y^2 \frac{\eta^2}{\gamma} \left(1 + \frac{1}{\gamma}\right) \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2 + L_y^2 \|\mathbf{v} - \mathbf{y}\|^2 \\
 & = \left(L_x^2 + L_y^2 + L_x^2 L_y^2 \frac{\eta^2}{\gamma} \left(1 + \frac{1}{\gamma}\right) \right) \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2 + \frac{L_y^2}{\gamma} \|\mathbf{u} - \mathbf{x}\|^2 \\
 & \leq \left(L_x^2 + L_y^2 \left(1 + \frac{1}{\gamma}\right) \left(1 + L_x^2 \frac{\eta^2}{\gamma}\right) \right) \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2,
 \end{aligned}$$

and this completes the proof. \square

Lemma B.2 (Zedek (1965, Theorem 1)). *Given a polynomial $p_n(z) := \sum_{k=0}^n a_k z^k$, $a_n \neq 0$, an integer $m \geq n$ and a number $\epsilon > 0$, there exists a number $\delta > 0$ such that whenever the $m + 1$ complex numbers b_k , $0 \leq k \leq m$, satisfy the inequalities*

$$|b_k - a_k| < \delta \quad \text{for } 0 \leq k \leq n, \quad \text{and} \quad |b_k| < \delta \quad \text{for } n + 1 \leq k \leq m,$$

then the roots β_k , $1 \leq k \leq m$, of the polynomial $q_m(z) := \sum_{k=0}^m b_k z^k$ can be labeled in such a way as to satisfy, with respect to the zeros α_k , $1 \leq k \leq n$, of $p_n(z)$, the inequalities

$$|\beta_k - \alpha_k| < \epsilon \quad \text{for } 1 \leq k \leq n, \quad \text{and} \quad |\beta_k| > \frac{1}{\epsilon} \quad \text{for } n + 1 \leq k \leq m.$$

Lemma B.3 (Chae et al. (2024b, Lemma E.1 and Corollary E.2)). *A (possibly complex) number μ is an eigenvalue of the restricted Schur complement \mathbf{S}_{res} if and only if it is a root of the equation*

$$\det \begin{bmatrix} \mu \mathbf{I} - \mathbf{A} & -\mathbf{C}_1 & -\mathbf{L}_q \\ \mathbf{C}_1^\top & -\mathbf{D} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \mathbf{0} \end{bmatrix} = 0, \quad (7)$$

where the definitions of \mathbf{D} and \mathbf{L}_q are given in the proof of Proposition 5.1 below. Moreover, if \mathbf{S}_{res} is invertible, the equation (7) does not have $\mu = 0$ as a solution.

Lemma B.4 (Jin et al. (2020, Lemma 40)). *Let \mathbf{A} , \mathbf{B} and \mathbf{C} respectively be $d_1 \times d_1$ symmetric, $d_2 \times d_2$ non-degenerate symmetric, and $d_1 \times d_2$ matrices. Then, the $d_1 + d_2$ complex eigenvalues of \mathbf{H}_τ have the following asymptotics as $\epsilon = \frac{1}{\tau} \rightarrow 0+$:*

$$|\lambda_j - \epsilon \mu_j| = o(\epsilon), \quad j = 1, \dots, d_1, \quad |\lambda_{j+d_1} - \nu_j| = o(1), \quad j = 1, \dots, d_2,$$

where $\{\mu_j\}_{j=1, \dots, d_1}$ and $\{\nu_j\}_{j=1, \dots, d_2}$ are the eigenvalues of $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top$ and $-\mathbf{B}$, respectively.

B.1. Proof of Proposition 5.1

Proof of Proposition 5.1. We first consider the case where \mathbf{S}_{res} is non-degenerate in Assumption 3, and we begin with the following observation. Consider a block matrix $\mathbf{Q} = \text{diag}(\mathbf{I}, \mathbf{P}^\top)$ where \mathbf{P} is orthogonal matrix such that $\mathbf{B} = \mathbf{P}\mathbf{\Delta}\mathbf{P}^\top$ for some diagonal matrix $\mathbf{\Delta} = \text{diag}\{\delta_1, \dots, \delta_r, 0, \dots, 0\}$. Then, one can show that $D\mathbf{F}_{\text{reg},\tau}$ is similar to

$$\mathbf{G}_{\text{reg},\tau} := \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{C}_2 \\ -\mathbf{C}_1^\top & \mathbf{D} + 2c\epsilon\mathbf{I} & \mathbf{0} \\ -\mathbf{C}_2^\top & \mathbf{0} & 2c\epsilon\mathbf{I} \end{bmatrix}$$

where $\mathbf{D} = \text{diag}(-\delta_1, \dots, -\delta_r)$ and $\epsilon = \frac{1}{\tau}$, since $\mathbf{G}_{\text{reg},\tau} = \mathbf{Q}D\mathbf{F}_{\text{reg},\tau}\mathbf{Q}^\top$. Here, the matrix \mathbf{C}_2 may not be of full column rank matrix, and we further refine the similarity statement below.

Let $q := \text{rank}(\mathbf{C}_2)$, and let $\mathbf{C}_2 = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ be the (full) singular value decomposition of \mathbf{C}_2 , where for some invertible diagonal matrix $\mathbf{\Sigma}_q$, it holds that

$$\mathbf{\Sigma} = \begin{bmatrix} \mathbf{\Sigma}_q & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

Then, by defining $L_q := U \begin{bmatrix} \Sigma_q \\ \mathbf{0} \end{bmatrix}$, we have $U\Sigma = [L_q \ \mathbf{0}]$. Thus, for $\tilde{Q} = \text{diag}\{I, I, V^\top\}$ we have

$$\tilde{Q}G_{\text{reg},\tau}\tilde{Q}^\top = \begin{bmatrix} A & C_1 & C_2V \\ -C_1^\top & D + 2c\epsilon I & \mathbf{0} \\ -V^\top C_2^\top & \mathbf{0} & 2c\epsilon I \end{bmatrix} = \begin{bmatrix} A & C_1 & L_q & \mathbf{0} \\ -C_1^\top & D + 2c\epsilon I & \mathbf{0} & \mathbf{0} \\ -L_q^\top & \mathbf{0} & 2c\epsilon I & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & 2c\epsilon I \end{bmatrix}$$

and therefore

$$\det(\lambda I - DF_{\text{reg},\tau}) = \det(\lambda I - \tilde{Q}G_{\text{reg},\tau}\tilde{Q}^\top) = (\lambda - 2c\epsilon)^{d_2-r-q} \det \left(\lambda I - \begin{bmatrix} A & C_1 & L_q \\ -C_1^\top & D + 2c\epsilon I & \mathbf{0} \\ -L_q^\top & \mathbf{0} & 2c\epsilon I \end{bmatrix} \right).$$

Hence, we notice that the eigenvalues of $DF_{\text{reg},\tau}$ are either $2c\epsilon$ or the eigenvalues of

$$\Phi_{\text{reg},\tau} := \begin{bmatrix} A & C_1 & L_q \\ -C_1^\top & D + 2c\epsilon I & \mathbf{0} \\ -L_q^\top & \mathbf{0} & 2c\epsilon I \end{bmatrix},$$

and analogously the eigenvalues of $H_{\text{reg},\tau} = \Lambda_\tau DF_{\text{reg},\tau}$ (5) are either $2c\epsilon$ or the eigenvalues of $\Phi_{\text{reg},\tau} := \Lambda_\tau \Phi_{\text{reg},\tau}$. Therefore, characterizing the eigenvalues of $\Phi_{\text{reg},\tau}$ is equivalent to characterizing the nonzero eigenvalues of $H_{\text{reg},\tau}$.

Since the eigenvalues of $\Phi_{\text{reg},\tau}$ are the solutions of the equation

$$0 = p_\epsilon(\lambda) := \det \begin{bmatrix} \lambda I - \epsilon A & -\epsilon C_1 & -\epsilon L_q \\ C_1^\top & \lambda I - D - 2c\epsilon I & \mathbf{0} \\ L_q^\top & \mathbf{0} & (\lambda - 2c\epsilon)I \end{bmatrix} = \det(\lambda I - \Phi_{\text{reg},\tau}), \quad (8)$$

we need to investigate the solutions of the equation. By Lemma B.2, constructing the functions $\lambda_j(\epsilon)$ so that they are continuous is possible, and the eigenvalues converge to the solutions of the equation

$$p_0(\lambda) = \det \begin{bmatrix} \lambda I & \mathbf{0} & \mathbf{0} \\ C_1^\top & \lambda I - D & \mathbf{0} \\ L_q^\top & \mathbf{0} & \lambda I \end{bmatrix} = 0$$

as $\epsilon \rightarrow 0$. Hence, the r eigenvalues of $\Phi_{\text{reg},\tau}$ converge to the r nonzero eigenvalues of $-B$, and the other $d_1 + q$ eigenvalues converge to zero, as $\epsilon \rightarrow 0$.

To investigate the order of eigenvalues that converges to zero further, we begin by observing that, whenever $|\lambda|$ and ϵ are small enough so that $\lambda I - D - 2c\epsilon I$ is invertible, it holds that

$$\begin{aligned} \det(\lambda I - \Phi_{\text{reg},\tau}) &= \det \begin{bmatrix} \lambda I - \epsilon A & -\epsilon C_1 & -\epsilon L_q \\ C_1^\top & \lambda I - D - 2c\epsilon I & \mathbf{0} \\ L_q^\top & \mathbf{0} & \lambda I - 2c\epsilon I \end{bmatrix} \\ &= \det \begin{bmatrix} \lambda I - \epsilon A + \epsilon C_1(\lambda I - D - 2c\epsilon I)^{-1}C_1^\top & \mathbf{0} & -\epsilon L_q \\ \mathbf{0} & \lambda I - D - 2c\epsilon I & \mathbf{0} \\ L_q^\top & \mathbf{0} & \lambda I - 2c\epsilon I \end{bmatrix} \\ &= \det(\lambda I - D - 2c\epsilon I) \det \begin{bmatrix} \lambda I - \epsilon A + \epsilon C_1(\lambda I - D - 2c\epsilon I)^{-1}C_1^\top & -\epsilon L_q \\ L_q^\top & \lambda I - 2c\epsilon I \end{bmatrix}. \end{aligned}$$

This implies that the λ_j , which converges to zero as $\epsilon \rightarrow 0$, is a solution of the following equation

$$0 = \det \begin{bmatrix} \lambda I - \epsilon A + \epsilon C_1(\lambda I - D - 2c\epsilon I)^{-1}C_1^\top & -\epsilon L_q \\ L_q^\top & \lambda I - 2c\epsilon I \end{bmatrix}. \quad (9)$$

Now let us reparametrize (9) by $\lambda = \kappa\sqrt{\epsilon}$ to get

$$\begin{aligned} 0 &= \det \begin{bmatrix} \kappa\sqrt{\epsilon}\mathbf{I} - \epsilon\mathbf{A} + \epsilon\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\epsilon\mathbf{I})^{-1}\mathbf{C}_1^\top & -\epsilon\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\sqrt{\epsilon}\mathbf{I} - 2c\epsilon\mathbf{I} \end{bmatrix} \\ &= \sqrt{\epsilon}^{d_1} \det \begin{bmatrix} \kappa\mathbf{I} - \sqrt{\epsilon}\mathbf{A} + \sqrt{\epsilon}\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\epsilon\mathbf{I})^{-1}\mathbf{C}_1^\top & -\sqrt{\epsilon}\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\sqrt{\epsilon}\mathbf{I} - 2c\epsilon\mathbf{I} \end{bmatrix} \\ &= \sqrt{\epsilon}^{d_1+q} \det \begin{bmatrix} \kappa\mathbf{I} - \sqrt{\epsilon}\mathbf{A} + \sqrt{\epsilon}\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & -\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix}. \end{aligned}$$

Since $(\lambda\mathbf{I} - \mathbf{D} - 2c\epsilon\mathbf{I})^{-1}$ converges to \mathbf{D}^{-1} as $\epsilon \rightarrow 0$, we have that if $\lambda_j \rightarrow 0$ as $\epsilon \rightarrow 0$ then λ_j should be a solution of the equation

$$0 = \det \begin{bmatrix} \kappa\mathbf{I} - \sqrt{\epsilon}\mathbf{A} + \sqrt{\epsilon}\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\epsilon\mathbf{I})^{-1}\mathbf{C}_1^\top & -\mathbf{L}_q \\ \mathbf{L}_q^\top & (\kappa - 2c\sqrt{\epsilon})\mathbf{I} \end{bmatrix} \quad (10)$$

By Lemma B.2, notice that eigenvalues divided by $\sqrt{\epsilon}$ converge to the solutions of

$$0 = \det \begin{bmatrix} \kappa\mathbf{I} & -\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\mathbf{I} \end{bmatrix}. \quad (11)$$

From the fact that \mathbf{L}_q is of full column rank matrix, \mathbf{L}_q has exactly q singular values. Therefore, solutions of (11) are nonzero, and those are exactly $i\sigma_k$, $k = 1, \dots, q$ where σ_k are the nonzero singular values of \mathbf{C}_2 , or equivalently singular values of \mathbf{L}_q , and therefore there are $2q$ instances among λ_j such that each λ_j has order of $\sqrt{\epsilon}$, and has asymptotic $+i\sigma_k\sqrt{\epsilon}$ or $-i\sigma_k\sqrt{\epsilon}$.

So far, we have shown that r eigenvalues have magnitude $\Theta(1)$, and $2q$ eigenvalues have magnitude $\Theta(\sqrt{\epsilon})$. On the other hand, we have

$$\begin{aligned} \det(\Phi_{\text{reg1},\tau}) &= \det \begin{bmatrix} \epsilon\mathbf{A} & \epsilon\mathbf{C}_1 & \epsilon\mathbf{L}_q \\ -\mathbf{C}_1^\top & \mathbf{D} + 2c\epsilon\mathbf{I} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & 2c\epsilon\mathbf{I} \end{bmatrix} \\ &= \epsilon^{d_1} \det \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top & \mathbf{D} + 2c\epsilon\mathbf{I} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & 2c\epsilon\mathbf{I} \end{bmatrix}. \end{aligned} \quad (12)$$

Here, since we assumed that \mathbf{S}_{res} is non-degenerate, by Lemma B.3, the RHS of (12) for $\epsilon = 0$ is not zero. Moreover, RHS of (12) is nonzero for sufficiently small ϵ by Lemma B.2. Therefore, the product of all λ_j of $\Phi_{\text{reg1},\tau}$ should be of order $\Theta(\epsilon^{d_1})$. From these two observations, we know that the product of the remaining $d_1 - q$ eigenvalues should be of order $\Theta(\epsilon^{d_1 - q})$. And we claim that each of these $d_1 - q$ eigenvalues is exactly of order $\Theta(\epsilon)$. To this end, let us examine what properties would the eigenvalues of order $O(\epsilon)$ have. By reparametrizing $\lambda = \mu\epsilon$ in (8), we have

$$\begin{aligned} 0 &= \det \begin{bmatrix} \mu\epsilon\mathbf{I} - \epsilon\mathbf{A} & -\epsilon\mathbf{C}_1 & -\epsilon\mathbf{L}_q \\ \mathbf{C}_1^\top & \mu\epsilon\mathbf{I} - \mathbf{D} - 2c\epsilon\mathbf{I} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \mu\epsilon\mathbf{I} - 2c\epsilon\mathbf{I} \end{bmatrix} \\ &= \epsilon^{d_1} \det \begin{bmatrix} \mu\mathbf{I} - \mathbf{A} & -\mathbf{C}_1 & -\mathbf{L}_q \\ \mathbf{C}_1^\top & \mu\mathbf{I} - \mathbf{D} - 2c\epsilon\mathbf{I} & \mathbf{0} \\ \mathbf{L}_q & \mathbf{0} & \mu\mathbf{I} - 2c\epsilon\mathbf{I} \end{bmatrix}. \end{aligned}$$

Then, by Lemma B.2, μ converges to a root of the equation

$$0 = \det \begin{bmatrix} \mu\mathbf{I} - \mathbf{A} & -\mathbf{C}_1 & -\mathbf{L}_q \\ \mathbf{C}_1^\top & -\mathbf{D} & \mathbf{0} \\ \mathbf{L}_q & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (13)$$

as $\epsilon \rightarrow 0$.

Then, $\mu = 0$ cannot be a root of (13) by Lemma B.3. This implies that there is no λ_j of order $o(\epsilon)$, or equivalently, all eigenvalues of $\Phi_{\text{reg1},\tau}$ are of order $\Omega(\epsilon)$. Therefore, from the fact that a product of $d_1 - q$ eigenvalues of $\Phi_{\text{reg1},\tau}$ is of order $\Theta(\epsilon^{d_1-q})$, each of those eigenvalues is exactly order of $\Theta(\epsilon)$, then the claim follows.

Here, we summarize the following two facts implied by the previous discussions:

- If λ is an eigenvalue of order $O(\epsilon)$, then λ/ϵ converges to a solution of (13) as $\epsilon \rightarrow 0$.
- The right-hand side of (13) is a polynomial of degree $d_1 - q$ in μ , whose solutions are nonzero.

By Lemma B.3, it is now immediate that the $d_1 - q$ eigenvalues of $\Phi_{\text{reg1},\tau}$ that are of order $\Theta(\epsilon)$ is of the form $\lambda(\epsilon) = \mu\epsilon + o(\epsilon)$ for μ that is a solution of (13).

The final claim, asserting that $\lambda_j(\epsilon) \neq 0$ for any j can be deduced from the invertibility of S_{res} by Lemmas B.2 and B.3. More precisely, we have $\det(\Phi_{\text{reg1},\tau}) = \det(\mathbf{A}_\tau) \det(D\Phi_{\text{reg},\tau})$. By Lemma B.2,

$$\det(D\Phi_{\text{reg},\tau}) \rightarrow \det \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top & \mathbf{D} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{as } \epsilon \rightarrow 0.$$

Therefore, the invertibility of S_{res} and Lemma B.3 imply that the right-hand side is not zero. Then for sufficiently large τ , the assertion $\lambda_j(\epsilon) \neq 0$ follows.

For the case where $\nabla_{\mathbf{y}\mathbf{y}}^2 f$ is non-degenerate in Assumption 3, the following Lemma completes the proof.

Lemma B.5. *Let \mathbf{A} , \mathbf{B} and \mathbf{C} respectively be $d_1 \times d_1$ symmetric, $d_2 \times d_2$ non-degenerate symmetric, and $d_1 \times d_2$ matrices.*

Then, the $d_1 + d_2$ complex eigenvalues of $\begin{bmatrix} \epsilon\mathbf{A} & \epsilon\mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} + 2c\epsilon\mathbf{I} \end{bmatrix}$ have the following asymptotics as $\epsilon = \frac{1}{\tau} \rightarrow 0+$:

$$|\lambda_j - \epsilon\mu_j| = o(\epsilon), \quad j = 1, \dots, d_1, \quad |\lambda_{j+d_1} - \nu_j| = o(1), \quad j = 1, \dots, d_2,$$

where $\{\mu_j\}_{j=1, \dots, d_1}$ and $\{\nu_j\}_{j=1, \dots, d_2}$ are the eigenvalues of $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top$ and $-\mathbf{B}$, respectively

Proof. Mimicking the proof of (Jin et al., 2020, Lemma 40), we can demonstrate the statement as follows. By definition of eigenvalues, any eigenvalue λ of $\begin{bmatrix} \epsilon\mathbf{A} & \epsilon\mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} + 2c\epsilon\mathbf{I} \end{bmatrix}$ is the roots of the characteristic equation

$$p_\epsilon(\lambda) := \det \begin{bmatrix} \lambda\mathbf{I} - \epsilon\mathbf{A} & -\epsilon\mathbf{C} \\ \mathbf{C}^\top & \lambda\mathbf{I} + \mathbf{B} - 2c\epsilon\mathbf{I} \end{bmatrix}$$

We can express the equation $p_\epsilon(\lambda)$ as follows.

$$p_\epsilon(\lambda) = p_0(\lambda) + \sum_{i=1}^{d_1+d_2} \epsilon^i p_i(\lambda)$$

where $p_0(\lambda) = \det \begin{bmatrix} \lambda\mathbf{I} & \mathbf{0} \\ \mathbf{C}^\top & \lambda\mathbf{I} + \mathbf{B} \end{bmatrix} = \lambda^{d_1} \det(\lambda\mathbf{I} + \mathbf{B})$ and $p_i(\lambda)$ for $i \geq 1$ are polynomials of order equal to or smaller than $d_1 + d_2$. Then, by Lemma B.2, the roots of $p_\epsilon(\lambda)$ are

$$\begin{aligned} |\lambda_j| &= o(1), \quad 1 \leq j \leq d_1, \\ |\lambda_{j+d_1} - \nu_j| &= o(1), \quad 1 \leq j \leq d_2, \end{aligned}$$

Since \mathbf{B} is non-degenerate, λ_{j+d_1} for $1 \leq j \leq d_2$ are of $\Omega(1)$, and therefore only λ_j for $1 \leq j \leq d_1$ converge to zero as $\epsilon \rightarrow 0$. To investigate the eigenvalues that converge to zero further, we reparametrize $\lambda = \kappa\epsilon$, then we have

$$\begin{aligned} p_\epsilon(\kappa\epsilon) &= \det \begin{bmatrix} \kappa\epsilon\mathbf{I} - \epsilon\mathbf{A} & -\epsilon\mathbf{C} \\ \mathbf{C}^\top & \kappa\epsilon\mathbf{I} + \mathbf{B} - 2c\epsilon\mathbf{I} \end{bmatrix} \\ &= \epsilon^{d_1} \det \begin{bmatrix} \kappa\mathbf{I} - \mathbf{A} & -\mathbf{C} \\ \mathbf{C}^\top & \kappa\mathbf{I} + \mathbf{B} - 2c\mathbf{I} \end{bmatrix}. \end{aligned}$$

This implies that the λ_j , for $1 \leq j \leq d_1$, divided by ϵ is a solution of the following equation

$$0 = \det \begin{bmatrix} \kappa \mathbf{I} - \mathbf{A} & -\mathbf{C} \\ \mathbf{C}^\top & \kappa \mathbf{I} + \mathbf{B} - 2c\epsilon \mathbf{I} \end{bmatrix}.$$

and by Lemma B.2, it converge to a root of

$$0 = \det \begin{bmatrix} \kappa \mathbf{I} - \mathbf{A} & -\mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \quad \text{as } \epsilon \rightarrow 0,$$

which is an eigenvalue of $\mathbf{A} - \mathbf{C}\mathbf{B}^{-1}\mathbf{C}^\top$ by Lemma B.3. These arguments complete the proof. \square

Therefore, there is no difference between the eigenvalues asymptotics of both the (vanilla) two-timescale EG (in Theorem 4.2) and the two-timescale EG with the “ $\gamma = \tau$ ” explicit regularization. \square

B.2. Proof of Theorem 5.2

Proof of Theorem 5.2. We first consider the case where \mathbf{S}_{res} is non-degenerate in Assumption 3. Analogous to the observations in proof of Proposition 5.1, by replacing ϵ with $\sqrt{\epsilon}$, one can deduce that the eigenvalues of $\mathbf{H}_{\text{reg}2,\tau}$ are either $2c\sqrt{\epsilon}$ or the nonzero eigenvalues of $\Phi_{\text{reg}2,\tau} := \Lambda_\tau \Phi_{\text{reg},\sqrt{\tau}}$ where

$$\Phi_{\text{reg},\sqrt{\tau}} := \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top & \mathbf{D} + 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix}.$$

Therefore, our next step is to characterize the eigenvalues of $\Phi_{\text{reg}2,\tau}$, and such eigenvalues are the solutions of the equation

$$0 = p_\epsilon(\lambda) := \det \begin{bmatrix} \lambda \mathbf{I} - \epsilon \mathbf{A} & -\epsilon \mathbf{C}_1 & -\epsilon \mathbf{L}_q \\ \mathbf{C}_1^\top & \lambda \mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & (\lambda - 2c\sqrt{\epsilon})\mathbf{I} \end{bmatrix} = \det(\lambda \mathbf{I} - \Phi_{\text{reg}2,\tau}). \quad (14)$$

By Lemma B.2, constructing the functions $\lambda_j(\epsilon)$ so that they are continuous is possible, and the eigenvalues converge to the solutions of the equation

$$p_0(\lambda) = \det \begin{bmatrix} \lambda \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \mathbf{C}_1^\top & \lambda \mathbf{I} - \mathbf{D} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \lambda \mathbf{I} \end{bmatrix} = 0$$

as $\epsilon \rightarrow 0$. Hence, the r eigenvalues of $\Phi_{\text{reg}2,\tau}$ converge to the r nonzero eigenvalues of $-\mathbf{B}$, and the other $d_1 + q$ eigenvalues converge to zero, as $\epsilon \rightarrow 0$.

To investigate the order of eigenvalues that converges to zero further, we begin by observing that, whenever $|\lambda|$ and ϵ are small enough so that $\lambda \mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I}$ is invertible, it holds that

$$\begin{aligned} \det(\lambda \mathbf{I} - \Phi_{\text{reg}2,\tau}) &= \det \begin{bmatrix} \lambda \mathbf{I} - \epsilon \mathbf{A} & -\epsilon \mathbf{C}_1 & -\epsilon \mathbf{L}_q \\ \mathbf{C}_1^\top & \lambda \mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \lambda \mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix} \\ &= \det \begin{bmatrix} \lambda \mathbf{I} - \epsilon \mathbf{A} + \epsilon \mathbf{C}_1(\lambda \mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & \mathbf{0} & -\epsilon \mathbf{L}_q \\ \mathbf{0} & \lambda \mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \lambda \mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix} \\ &= \det(\lambda \mathbf{I} - \mathbf{D} + 2c\sqrt{\epsilon}\mathbf{I}) \det \begin{bmatrix} \lambda \mathbf{I} - \epsilon \mathbf{A} + \epsilon \mathbf{C}_1(\lambda \mathbf{I} - \mathbf{D} + 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & -\epsilon \mathbf{L}_q \\ \mathbf{L}_q^\top & \lambda \mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix}. \end{aligned}$$

This implies that the λ_j , which converges to zero as $\epsilon \rightarrow 0$, is a solution of the following equation

$$0 = \det \begin{bmatrix} \lambda \mathbf{I} - \epsilon \mathbf{A} + \epsilon \mathbf{C}_1(\lambda \mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & -\epsilon \mathbf{L}_q \\ \mathbf{L}_q^\top & \lambda \mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix}. \quad (15)$$

Now let us reparametrize (15) by $\lambda = \kappa\sqrt{\epsilon}$ to get

$$\begin{aligned} 0 &= \begin{bmatrix} \kappa\sqrt{\epsilon}\mathbf{I} - \epsilon\mathbf{A} + \epsilon\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & -\epsilon\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\sqrt{\epsilon}\mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix} \\ &= \sqrt{\epsilon}^{d_1} \det \begin{bmatrix} \kappa\mathbf{I} - \sqrt{\epsilon}\mathbf{A} + \sqrt{\epsilon}\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & -\sqrt{\epsilon}\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\sqrt{\epsilon}\mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix} \\ &= \sqrt{\epsilon}^{d_1+q} \det \begin{bmatrix} \kappa\mathbf{I} - \sqrt{\epsilon}\mathbf{A} + \sqrt{\epsilon}\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & -\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\mathbf{I} - 2c\mathbf{I} \end{bmatrix}. \end{aligned}$$

Since $(\lambda\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}$ converges to \mathbf{D}^{-1} as $\epsilon \rightarrow 0$, we have that if $\lambda_j \rightarrow 0$ as $\epsilon \rightarrow 0$ then λ_j should be a solution of the equation

$$0 = \det \begin{bmatrix} \kappa\mathbf{I} - \sqrt{\epsilon}\mathbf{A} + \sqrt{\epsilon}\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I})^{-1}\mathbf{C}_1^\top & -\mathbf{L}_q \\ \mathbf{L}_q^\top & (\kappa - 2c)\mathbf{I} \end{bmatrix} \quad (16)$$

By Lemma B.2, notice that eigenvalues divided by $\sqrt{\epsilon}$ converge to the solutions of

$$0 = \det \begin{bmatrix} \kappa\mathbf{I} & -\mathbf{L}_q \\ \mathbf{L}_q^\top & (\kappa - 2c)\mathbf{I} \end{bmatrix}. \quad (17)$$

From the fact that \mathbf{L}_q is of full column rank matrix, \mathbf{L}_q has exactly q singular values. Therefore, solutions of (17) are nonzero, and those are exactly $c \pm \sqrt{c^2 - \sigma_k^2}$, $k = 1, \dots, q$ where σ_k are the nonzero singular values of \mathbf{C}_2 , or equivalently singular values of \mathbf{L}_q . Note that the $c \pm \sqrt{c^2 - \sigma_k^2}$ can be written as $c \pm i\sqrt{\sigma_k^2 - c^2}$ when $\sigma_k > c$, however, we will denote it as $c \pm \sqrt{c^2 - \sigma_k^2}$, to cover both cases. Therefore, there are $2q$ instances among λ_j such that each λ_j has order of $\sqrt{\epsilon}$, and has asymptotic $\sqrt{\epsilon}(c + \sqrt{c^2 - \sigma_k^2})$ or $\sqrt{\epsilon}(c - \sqrt{c^2 - \sigma_k^2})$.

So far, we have shown that r eigenvalues have magnitude $\Theta(1)$, and $2q$ eigenvalues have magnitude $\Theta(\sqrt{\epsilon}(c \pm \sqrt{c^2 - \sigma_k^2}))$. On the other hand, we have

$$\begin{aligned} \det(\Phi_{\text{reg2},\tau}) &= \det \begin{bmatrix} \epsilon\mathbf{A} & \epsilon\mathbf{C}_1 & \epsilon\mathbf{L}_q \\ -\mathbf{C}_1^\top & -\mathbf{D} + 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix} \\ &= \epsilon^{d_1} \det \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top & -\mathbf{D} + 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix}. \end{aligned} \quad (18)$$

Here, since we assumed that \mathbf{S}_{res} is non-degenerate, by Lemma B.3, the RHS of (18) for $\epsilon = 0$ is not zero. Moreover, RHS of (18) is nonzero for sufficiently small ϵ by Lemma B.2. Therefore, the product of all λ_j of $\Phi_{\text{reg2},\tau}$ should be of order $\Theta(\epsilon^{d_1})$. From these two observations, we know that the product of the remaining $d_1 - q$ eigenvalues should be of order $\Theta(\epsilon^{d_1 - q})$. And we claim that each of these $d_1 - q$ eigenvalues is exactly of order $\Theta(\epsilon)$. To this end, let us examine what properties would the eigenvalues of order $O(\epsilon)$ have. By reparametrizing $\lambda = \mu\epsilon$ in (14), we have

$$\begin{aligned} 0 &= \det \begin{bmatrix} \mu\epsilon\mathbf{I} - \epsilon\mathbf{A} & -\epsilon\mathbf{C}_1 & -\epsilon\mathbf{L}_q \\ \mathbf{C}_1^\top & \mu\epsilon\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \mu\epsilon\mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix} \\ &= \epsilon^{d_1} \det \begin{bmatrix} \mu\mathbf{I} - \mathbf{A} & -\mathbf{C}_1 & -\mathbf{L}_q \\ \mathbf{C}_1^\top & \mu\mathbf{I} - \mathbf{D} - 2c\sqrt{\epsilon}\mathbf{I} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \mu\mathbf{I} - 2c\sqrt{\epsilon}\mathbf{I} \end{bmatrix}. \end{aligned}$$

Then, by Lemma B.2, μ converges to a root of the equation

$$0 = \det \begin{bmatrix} \mu\mathbf{I} - \mathbf{A} & -\mathbf{C}_1 & -\mathbf{L}_q \\ \mathbf{C}_1^\top & -\mathbf{D} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad (19)$$

as $\epsilon \rightarrow 0$.

Since the \mathbf{S}_{res} is non-degenerate, the $\mu = 0$ cannot be a root of (19) by Lemma B.3. This particularly showing that there is no λ_j of order $o(\epsilon)$, or equivalently, all eigenvalues of $\Phi_{\text{reg}2,\tau}$ are of order $\Omega(\epsilon)$. Therefore, from the fact that a product of $d_1 - q$ eigenvalues of $\Phi_{\text{reg}2,\tau}$ is of order $\Theta(\epsilon^{d_1-q})$, each of those eigenvalues is exactly order of $\Theta(\epsilon)$, then the claim follows.

Here, we want to summarize the following two facts implied by the previous discussions:

- If λ is an eigenvalue of order $O(\epsilon)$, then λ divided by ϵ converges to a solution of (19) as $\epsilon \rightarrow 0$.
- The right-hand side of (19) is a polynomial of degree $d_1 - q$ in μ , whose solutions are nonzero.

By Lemma B.3, it is now immediate that the $d_1 - q$ eigenvalues of $\Phi_{\text{reg}2,\tau}$ that are of order $\Theta(\epsilon)$ is of the form $\lambda(\epsilon) = \mu\epsilon + o(\epsilon)$ for μ that is a solution of (19).

The final claim, asserting that $\lambda_j(\epsilon) \neq 0$ for any j , can be deduced from the invertibility of \mathbf{S}_{res} with Lemmas B.2 and B.3. More precisely, we have $\det(\Phi_{\text{reg}2,\tau}) = \det(\mathbf{A}_\tau) \det(D\Phi_{\text{reg},\sqrt{\tau}})$. By Lemma B.2,

$$\det(D\Phi_{\text{reg},\sqrt{\tau}}) \rightarrow \det \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top & \mathbf{D} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & \mathbf{0} \end{bmatrix} \quad \text{as } \epsilon \rightarrow 0.$$

Therefore, the invertibility of \mathbf{S}_{res} and Lemma B.3 imply that the RHS is not zero. Then, the assertion $\lambda_j(\epsilon) \neq 0$ follows. For the case where $\nabla_{\mathbf{y}\mathbf{y}}^2 f$ is non-degenerate in Assumption 3, mimicking the proof of Lemma B.5 completes the proof. \square

B.3. Proof of Theorem 5.3

Proof of Theorem 5.3. We first consider the case where \mathbf{S}_{res} is non-degenerate. Recall that, alternating saddle gradient operator is as follow

$$\mathbf{F}_{\text{alt},\gamma}(\mathbf{x}, \mathbf{y}) := (\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}} f(\mathbf{x} - \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), \mathbf{y}))$$

where the η is given step size.

Then the proposed alternating extragradient can be formulated as follows

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \eta \mathbf{A}_{\tau_k} \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}_k - \eta \mathbf{A}_{\tau_k} \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}_k)).$$

By the matrix version of chain rule, the Jacobian of the $\mathbf{A}_{\tau_k} \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z})$ at \mathbf{z}^* is

$$\begin{aligned} \mathbf{H}_{\text{alt},\tau_k,\gamma_k} &= \begin{bmatrix} \epsilon_1 \mathbf{A} & \epsilon_1 \mathbf{C} \\ [-\mathbf{C}^\top & -\mathbf{B}] \begin{bmatrix} \mathbf{I} - \eta \epsilon_2 \mathbf{A} \\ 0 \end{bmatrix} & [-\mathbf{C}^\top & -\mathbf{B}] \begin{bmatrix} -\eta \epsilon_2 \mathbf{C} \\ \mathbf{I} \end{bmatrix} \end{bmatrix} \\ &= \begin{bmatrix} \epsilon_1 \mathbf{A} & \epsilon_1 \mathbf{C} \\ -\mathbf{C}^\top + \eta \epsilon_2 \mathbf{C}^\top \mathbf{A} & -\mathbf{B} + \eta \epsilon_2 \mathbf{C}^\top \mathbf{C} \end{bmatrix}. \end{aligned}$$

where $\epsilon_1 = \frac{1}{\tau_k}$, $\epsilon_2 = \frac{1}{\gamma_k}$ and $D\mathbf{F} = \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} \end{bmatrix}$ for saddle gradient $\mathbf{F} = (\nabla_{\mathbf{x}} f, -\nabla_{\mathbf{y}} f)$.

Analogous to the proof of Proposition 5.1, one can show that $\mathbf{H}_{\text{alt},\tau_k,\gamma_k}$ is similar to

$$\mathbf{G}_{\text{alt},\tau_k,\gamma_k} = \begin{bmatrix} \epsilon_1 \mathbf{A} & \epsilon_1 \mathbf{C}_1 & \epsilon_1 \mathbf{C}_2 \\ -\mathbf{C}_1^\top + \eta \epsilon_2 \mathbf{C}_1^\top \mathbf{A} & \mathbf{D} + \eta \epsilon_2 \mathbf{C}_1^\top \mathbf{C}_1 & \eta \epsilon_2 \mathbf{C}_1^\top \mathbf{C}_2 \\ -\mathbf{C}_2^\top + \eta \epsilon_2 \mathbf{C}_2^\top \mathbf{A} & \eta \epsilon_2 \mathbf{C}_2^\top \mathbf{C}_1 & \eta \epsilon_2 \mathbf{C}_2^\top \mathbf{C}_2 \end{bmatrix},$$

under the same settings and notations.

Then, for $\tilde{Q} = \text{diag}\{I, I, V^\top\}$ we have

$$\begin{aligned} \tilde{Q}G_{\text{alt},\tau_k,\gamma_k}\tilde{Q}^\top &= \begin{bmatrix} \epsilon_1 A & \epsilon_1 C_1 & \epsilon_1 C_2 V \\ -C_1^\top + \eta\epsilon_2 C_1^\top A & D + \eta\epsilon_2 C_1^\top C_1 & \eta\epsilon_2 C_1^\top C_2 V \\ -V^\top C_2^\top + \eta\epsilon_2 V^\top C_2^\top A & \eta\epsilon_2 V^\top C_2^\top C_1 & \eta\epsilon_2 V^\top C_2^\top C_2 V \end{bmatrix} \\ &= \begin{bmatrix} \epsilon_1 A & \epsilon_1 C_1 & \epsilon_1 L_q & \mathbf{0} \\ -C_1^\top + \eta\epsilon_2 C_1^\top A & D + \eta\epsilon_2 C_1^\top C_1 & \eta\epsilon_2 C_1^\top L_q & \mathbf{0} \\ -L_q^\top + \eta\epsilon_2 L_q^\top A & \eta\epsilon_2 L_q^\top C_1 & \eta\epsilon_2 L_q^\top L_q & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} \end{aligned}$$

and therefore

$$\begin{aligned} \det(\lambda I - H_{\text{alt},\tau_k,\gamma_k}) &= \det(\lambda I - \tilde{Q}G_{\text{alt},\tau_k,\gamma_k}\tilde{Q}^\top) \\ &= \lambda^{d_2-r-q} \det \left(\lambda I - \begin{bmatrix} \epsilon_1 A & \epsilon_1 C_1 & \epsilon_1 L_q \\ -C_1^\top + \eta\epsilon_2 C_1^\top A & D + \eta\epsilon_2 C_1^\top C_1 & \eta\epsilon_2 C_1^\top L_q \\ -L_q^\top + \eta\epsilon_2 L_q^\top A & \eta\epsilon_2 L_q^\top C_1 & \eta\epsilon_2 L_q^\top L_q \end{bmatrix} \right). \end{aligned}$$

So, the eigenvalues of $H_{\text{alt},\tau_k,\gamma_k}$ are either zero or the eigenvalues of

$$\Phi_{\text{alt},\tau_k,\gamma_k} := \begin{bmatrix} \epsilon_1 A & \epsilon_1 C_1 & \epsilon_1 L_q \\ -C_1^\top + \eta\epsilon_2 C_1^\top A & D + \eta\epsilon_2 C_1^\top C_1 & \eta\epsilon_2 C_1^\top L_q \\ -L_q^\top + \eta\epsilon_2 L_q^\top A & \eta\epsilon_2 L_q^\top C_1 & \eta\epsilon_2 L_q^\top L_q \end{bmatrix}.$$

Therefore, characterizing the eigenvalues of $\Phi_{\text{alt},\tau_k,\gamma_k}$ is equivalent to characterizing the nonzero eigenvalues of $H_{\text{alt},\tau_k,\gamma_k}$.

From now on, let $\tau_k = \tau$ and $\gamma_k = \sqrt{\tau}$. Then, the eigenvalues of $\Phi_{\text{alt},\tau,\sqrt{\tau}}$ are the solutions of the equation

$$0 = p_\epsilon(\lambda) := \det \begin{bmatrix} \lambda I - \epsilon A & -\epsilon C_1 & -\epsilon L_q \\ C_1^\top - \eta\sqrt{\epsilon} C_1^\top A & \lambda I - D - \eta\sqrt{\epsilon} C_1^\top C_1 & -\eta\sqrt{\epsilon} C_1^\top L_q \\ L_q^\top - \eta\sqrt{\epsilon} L_q^\top A & -\eta\sqrt{\epsilon} L_q^\top C_1 & \lambda I - \eta\sqrt{\epsilon} L_q^\top L_q \end{bmatrix} = \det(\lambda I - \Phi_{\text{alt},\tau,\sqrt{\tau}}). \quad (20)$$

By Lemma B.2, constructing the functions $\lambda_j(\epsilon)$ so that they are continuous is possible, and the eigenvalues converge to the solutions of the equation

$$p_0(\lambda) = \det \begin{bmatrix} \lambda I & \mathbf{0} & \mathbf{0} \\ C_1^\top & \lambda I - D & \mathbf{0} \\ L_q^\top & \mathbf{0} & \lambda I \end{bmatrix} = 0.$$

as $\epsilon \rightarrow 0$. Hence, the r eigenvalues of $\Phi_{\text{alt},\tau,\sqrt{\tau}}$ converge to the r nonzero eigenvalues of $-B$, and the other $d_1 + q$ eigenvalues converge to zero, as $\epsilon \rightarrow 0$.

To investigate the order of eigenvalues that converges to zero further, we begin by observing that, whenever $|\lambda|$ and ϵ are small enough so that $\lambda I - D - \eta\sqrt{\epsilon} C_1^\top C_1$ is invertible, it holds that

$$\begin{aligned} \det(\lambda I - \Phi_{\text{alt},\tau,\sqrt{\tau}}) &= \det \begin{bmatrix} \lambda I - \epsilon A & -\epsilon C_1 & -\epsilon L_q \\ C_1^\top - \eta\sqrt{\epsilon} C_1^\top A & \lambda I - D - \eta\sqrt{\epsilon} C_1^\top C_1 & -\eta\sqrt{\epsilon} C_1^\top L_q \\ L_q^\top - \eta\sqrt{\epsilon} L_q^\top A & -\eta\sqrt{\epsilon} L_q^\top C_1 & \lambda I - \eta\sqrt{\epsilon} L_q^\top L_q \end{bmatrix} \\ &= \det \begin{bmatrix} \lambda I - \epsilon A & -\epsilon C_1 & -\epsilon L_q \\ (1 - \frac{\lambda\eta}{\sqrt{\epsilon}}) C_1^\top & \lambda I - D & \mathbf{0} \\ (1 - \frac{\lambda\eta}{\sqrt{\epsilon}}) L_q^\top & \mathbf{0} & \lambda I \end{bmatrix} \\ &= \det \begin{bmatrix} \lambda I - \epsilon A + \epsilon(1 - \frac{\lambda\eta}{\sqrt{\epsilon}}) C_1 (\lambda I - D)^{-1} C_1^\top & \mathbf{0} & -\epsilon L_q \\ \mathbf{0} & \lambda I - D & \mathbf{0} \\ (1 - \frac{\lambda\eta}{\sqrt{\epsilon}}) L_q^\top & \mathbf{0} & \lambda I \end{bmatrix} \\ &= \det(\lambda I - D) \det \begin{bmatrix} \lambda I - \epsilon A + \epsilon(1 - \frac{\lambda\eta}{\sqrt{\epsilon}}) C_1 (\lambda I - D)^{-1} C_1^\top & -\epsilon L_q \\ (1 - \frac{\lambda\eta}{\sqrt{\epsilon}}) L_q^\top & \lambda I \end{bmatrix}. \end{aligned}$$

This implies that the λ_j , which converges to zero as $\epsilon \rightarrow 0$, is a solution of the following equation

$$0 = \det \begin{bmatrix} \lambda \mathbf{I} - \epsilon \mathbf{A} + \epsilon(1 - \frac{\lambda \eta}{\sqrt{\epsilon}}) \mathbf{C}_1 (\lambda \mathbf{I} - \mathbf{D})^{-1} \mathbf{C}_1^\top & -\epsilon \mathbf{L}_q \\ (1 - \frac{\lambda \eta}{\sqrt{\epsilon}}) \mathbf{L}_q^\top & \lambda \mathbf{I} \end{bmatrix}. \quad (21)$$

Now let us reparametrize (21) by $\lambda = \kappa \sqrt{\epsilon}$ to get

$$\begin{aligned} 0 &= \begin{bmatrix} \kappa \sqrt{\epsilon} \mathbf{I} - \epsilon \mathbf{A} + \epsilon(1 - \kappa \eta) \mathbf{C}_1 (\lambda \mathbf{I} - \mathbf{D})^{-1} \mathbf{C}_1^\top & -\epsilon \mathbf{L}_q \\ (1 - \kappa \eta) \mathbf{L}_q^\top & \kappa \sqrt{\epsilon} \mathbf{I} \end{bmatrix} \\ &= \sqrt{\epsilon}^{d_1} \det \begin{bmatrix} \kappa \mathbf{I} - \sqrt{\epsilon} \mathbf{A} + \sqrt{\epsilon}(1 - \kappa \eta) \mathbf{C}_1 (\lambda \mathbf{I} - \mathbf{D})^{-1} \mathbf{C}_1^\top & -\sqrt{\epsilon} \mathbf{L}_q \\ (1 - \kappa \eta) \mathbf{L}_q^\top & \kappa \sqrt{\epsilon} \mathbf{I} \end{bmatrix} \\ &= \sqrt{\epsilon}^{d_1+q} \det \begin{bmatrix} \kappa \mathbf{I} - \sqrt{\epsilon} \mathbf{A} + \sqrt{\epsilon}(1 - \kappa \eta) \mathbf{C}_1 (\lambda \mathbf{I} - \mathbf{D})^{-1} \mathbf{C}_1^\top & -\mathbf{L}_q \\ (1 - \kappa \eta) \mathbf{L}_q^\top & \kappa \mathbf{I} \end{bmatrix}. \end{aligned}$$

Since $(\lambda \mathbf{I} - \mathbf{D})^{-1}$ converges as $\epsilon \rightarrow 0$, we have that if $\lambda_j \rightarrow 0$ as $\epsilon \rightarrow 0$ then λ_j should be a solution of the equation

$$0 = \det \begin{bmatrix} \kappa \mathbf{I} - \sqrt{\epsilon} \mathbf{A} + \sqrt{\epsilon}(1 - \kappa \eta) \mathbf{C}_1 (\lambda \mathbf{I} - \mathbf{D})^{-1} \mathbf{C}_1^\top & -\mathbf{L}_q \\ (1 - \kappa \eta) \mathbf{L}_q^\top & \kappa \mathbf{I} \end{bmatrix} \quad (22)$$

By Lemma B.2, notice that eigenvalues divided by $\sqrt{\epsilon}$ converge to the solutions of

$$0 = \det \begin{bmatrix} \kappa \mathbf{I} & -\mathbf{L}_q \\ (1 - \kappa \eta) \mathbf{L}_q^\top & \kappa \mathbf{I} \end{bmatrix}. \quad (23)$$

From the fact that \mathbf{L}_q is of full column rank matrix, \mathbf{L}_q has exactly q singular values. Therefore, solutions of (23) are nonzero, and those are exactly $\frac{\eta \sigma_k^2}{2} \pm \frac{\sigma_k \sqrt{\eta^2 \sigma_k^2 - 4}}{2}$, $k = 1, \dots, q$ where σ_k are the nonzero singular values of \mathbf{C}_2 , or equivalently singular values of \mathbf{L}_q . Note that the $\frac{\eta \sigma_k^2}{2} \pm \frac{\sigma_k \sqrt{\eta^2 \sigma_k^2 - 4}}{2}$ can be written as $\frac{\eta \sigma_k^2}{2} \pm i \frac{\sigma_k \sqrt{4 - \eta^2 \sigma_k^2}}{2}$ when $\eta \sigma_k < 2$, however, we will denote it as $\frac{\eta \sigma_k^2}{2} \pm \frac{\sigma_k \sqrt{\eta^2 \sigma_k^2 - 4}}{2}$, to cover both cases. Therefore, there are $2q$ instances among λ_j such that each λ_j has order of $\sqrt{\epsilon}$, and has asymptotic, and has asymptotic $\frac{\eta \sigma_k^2}{2} + \frac{\sigma_k \sqrt{\eta^2 \sigma_k^2 - 4}}{2}$ or $\frac{\eta \sigma_k^2}{2} - \frac{\sigma_k \sqrt{\eta^2 \sigma_k^2 - 4}}{2}$.

So far, we have shown that r eigenvalues have magnitude $\Theta(1)$, and $2q$ eigenvalues have magnitude $\Theta(\sqrt{\epsilon}(\frac{\eta \sigma_k^2}{2} \pm \frac{\sigma_k \sqrt{\eta^2 \sigma_k^2 - 4}}{2}))$. On the other hand, we have

$$\begin{aligned} \det(\Phi_{\text{alt}, \tau, \sqrt{\tau}}) &= \epsilon^{d_1} \det \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top + \eta \sqrt{\epsilon} \mathbf{C}_1^\top \mathbf{A} & \mathbf{D} + \eta \sqrt{\epsilon} \mathbf{C}_1^\top \mathbf{C}_1 & \eta \sqrt{\epsilon} \mathbf{C}_1^\top \mathbf{L}_q \\ -\mathbf{L}_q^\top + \eta \sqrt{\epsilon} \mathbf{L}_q^\top \mathbf{A} & \eta \sqrt{\epsilon} \mathbf{L}_q^\top \mathbf{C}_1 & \eta \sqrt{\epsilon} \mathbf{L}_q^\top \mathbf{L}_q \end{bmatrix} \\ &= \epsilon^{d_1} \det \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top & \mathbf{D} & \mathbf{0} \\ -\mathbf{L}_q & \mathbf{0} & \mathbf{0} \end{bmatrix}. \end{aligned} \quad (24)$$

Here, since we assumed that \mathcal{S}_{res} is non-degenerate, by Lemma B.3, the RHS of (24) is not zero, therefore, the product of all λ_j of $\Phi_{\text{alt}, \tau, \sqrt{\tau}}$ should be of order $\Theta(\epsilon^{d_1})$. From these two observations, we know that product of the remaining $d_1 - q$ eigenvalues should be of order $\Theta(\epsilon^{d_1 - q})$. And we claim that each of these $d_1 - q$ eigenvalues is exactly of order $\Theta(\epsilon)$. To this end, let us examine what properties would the eigenvalues of order $O(\epsilon)$ have. By reparametrizing $\lambda = \mu \epsilon$ in (20), we

have

$$\begin{aligned}
 0 &= \det \begin{bmatrix} \mu\epsilon\mathbf{I} - \epsilon\mathbf{A} & -\epsilon\mathbf{C}_1 & -\epsilon\mathbf{L}_q \\ \mathbf{C}_1^\top - \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{A} & \mu\epsilon\mathbf{I} - \mathbf{D} - \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{C}_1 & -\eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{L}_q \\ \mathbf{L}_q^\top - \eta\sqrt{\epsilon}\mathbf{L}_q^\top\mathbf{A} & -\eta\sqrt{\epsilon}\mathbf{L}_q^\top\mathbf{C}_1 & \mu\epsilon\mathbf{I} - \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{L}_q \end{bmatrix} \\
 &= \epsilon^{d_1} \det \begin{bmatrix} \mu\mathbf{I} - \mathbf{A} & -\mathbf{C}_1 & -\mathbf{L}_q \\ \mathbf{C}_1^\top - \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{A} & \mu\mathbf{I} - \mathbf{D} - \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{C}_1 & -\eta\sqrt{\epsilon}\mathbf{L}_q \\ \mathbf{L}_q^\top - \eta\sqrt{\epsilon}\mathbf{L}_q^\top\mathbf{A} & -\eta\sqrt{\epsilon}\mathbf{L}_q^\top\mathbf{C}_1 & \mu\mathbf{I} - \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{L}_q \end{bmatrix}.
 \end{aligned}$$

Then, by Lemma B.2, μ converges to a root of the equation

$$0 = \det \begin{bmatrix} \mu\mathbf{I} - \mathbf{A} & -\mathbf{C}_1 & -\mathbf{L}_q \\ \mathbf{C}_1^\top & -\mathbf{D} & \mathbf{0} \\ \mathbf{L}_q^\top & \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (25)$$

as $\epsilon \rightarrow 0$.

Then, again $\mu = 0$ cannot be a root of (25) by Lemma B.3. This particularly showing that there is no λ_j of order $o(\epsilon)$, or equivalently, all eigenvalues of $\Phi_{\text{alt},\tau,\sqrt{\tau}}$ are of order $\Omega(\epsilon)$. Therefore, from the fact that a product of $d_1 - q$ eigenvalues of $\Phi_{\text{alt},\tau,\sqrt{\tau}}$ is of order $\Theta(\epsilon^{d_1-q})$, each of those eigenvalues is exactly order of $\Theta(\epsilon)$, then the claim follows.

Here, we want to summarize the previous discussions imply the following two facts:

- If λ is an eigenvalue of order $O(\epsilon)$, then λ divided by ϵ converges to a solution of (25) as $\epsilon \rightarrow 0$.
- The right-hand side of (25) is a polynomial of degree $d_1 - q$ in μ , whose solutions are nonzero.

By Lemma B.3, it is now immediate that the $d_1 - q$ eigenvalues of $\Phi_{\text{alt},\tau,\sqrt{\tau}}$ that are of order $\Theta(\epsilon)$, and is of the form $\lambda(\epsilon) = \mu\epsilon + o(\epsilon)$ for a μ which is solution of (25).

The final claim, asserting that $\lambda_j(\epsilon) \neq 0$ for any j , can be deduced from the invertibility of \mathbf{S}_{res} with Lemmas B.2 and B.3. More precisely, we have $\det(\Phi_{\text{alt},\tau,\sqrt{\tau}}) = \det(\mathbf{\Lambda}_\tau) \det(D\Phi_{\text{alt},\sqrt{\tau}})$ where

$$D\Phi_{\text{alt},\sqrt{\tau}} := \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top + \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{A} & \mathbf{D} + \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{C}_1 & \eta\sqrt{\epsilon}\mathbf{C}_1^\top\mathbf{L}_q \\ -\mathbf{L}_q^\top + \eta\sqrt{\epsilon}\mathbf{L}_q^\top\mathbf{A} & \eta\sqrt{\epsilon}\mathbf{L}_q^\top\mathbf{C}_1 & \eta\sqrt{\epsilon}\mathbf{L}_q^\top\mathbf{L}_q \end{bmatrix}.$$

Since

$$\det(D\Phi_{\text{alt},\sqrt{\tau}}) = \det \begin{bmatrix} \mathbf{A} & \mathbf{C}_1 & \mathbf{L}_q \\ -\mathbf{C}_1^\top & \mathbf{D} & \mathbf{0} \\ -\mathbf{L}_q^\top & \mathbf{0} & \mathbf{0} \end{bmatrix}, \quad (26)$$

the assumption and lemma B.3 implies that the RHS of (26) is not zero. Then, for any $\tau \geq 1$, the assertion $\lambda_j(\epsilon) \neq 0$ follows.

For the case where $\nabla_{\mathbf{y}\mathbf{y}}^2 f$ is non-degenerate in Assumption 3, following the proof of Lemma B.4 with the fact that $\det(\mathbf{H}_{\text{alt},\tau,\sqrt{\tau}}) = \det \begin{bmatrix} \epsilon\mathbf{A} & \epsilon\mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} \end{bmatrix} = \det(\mathbf{H}_\tau)$ completes the proof.

□

B.4. Relationship between the equilibrium of Alt2-EG-TS and the stationary point of F

Proposition B.6. *Under Assumption 1, a point \mathbf{z}^* is an equilibrium point of Alt2-EG-TS (6)*

$$\mathbf{z}_{k+1} = \mathbf{z}_k - \eta\mathbf{\Lambda}_{\tau_k} \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}_k - \eta\mathbf{\Lambda}_{\tau_k} \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}_k))$$

if and only if $\mathbf{F}(\mathbf{z}^*) = \mathbf{0}$, for $0 < \eta < \frac{1}{\sqrt{2}L}$, $\tau_k \geq 1$ and $\gamma_k \geq 2$.

Proof. It is obvious that $\mathbf{F}(\mathbf{z}^*) = \mathbf{0}$ if and only if $\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) = \mathbf{0}$. It is also straightforward that $\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) = \mathbf{0}$ implies $\mathbf{z}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*)) = \mathbf{z}^*$. We are now left to prove the ‘‘only if’’ statement.

Suppose that \mathbf{z}^* is an equilibrium point of Alt2-EG-TS. Then, we have

$$\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*)) = \mathbf{0}.$$

For the sake of contradiction, suppose that $\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) \neq \mathbf{0}$ and let $\mathbf{w}^* := \mathbf{z}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*)$. Note that, by the assumption, $\mathbf{w}^* \neq \mathbf{z}^*$. Then, we have

$$\begin{aligned} \mathbf{z}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) &= \mathbf{z}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*)) - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) \\ &= \mathbf{w}^* - \eta\mathbf{\Lambda}_{\tau_k}\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{w}^*), \end{aligned}$$

hence we have $\mathbf{z}^* - \mathbf{w}^* = \eta\mathbf{\Lambda}_{\tau_k}(\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) - \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{w}^*))$.

Meanwhile, under Assumption 1, one can deduce that $\mathbf{F}_{\text{alt},\gamma_k}$ is $\sqrt{2}L$ -Lipschitz for $0 < \eta < \frac{1}{\sqrt{2}L}$ and $\gamma_k \geq 2$, since

$$\begin{aligned} \|\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{u}, \mathbf{v}) - \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{x}, \mathbf{y})\|^2 &\leq \left(L_x^2 + L_y^2 \left(1 + \frac{1}{\gamma_k} \right) \left(1 + L_x^2 \frac{\eta^2}{\gamma_k} \right) \right) \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2 \quad (\text{by Lemma B.1}) \\ &\leq \left(L_x^2 + L_y^2 \left(1 + \frac{1}{2} \right) \left(1 + \frac{1}{4} \right) \right) \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2 \\ &\leq 2(L_x^2 + L_y^2) \|(\mathbf{u}, \mathbf{v}) - (\mathbf{x}, \mathbf{y})\|^2. \end{aligned} \quad (27)$$

Therefore, for $0 < \eta < \frac{1}{\sqrt{2}L}$, $\tau_k \geq 1$, and $\gamma_k \geq 2$, we have

$$\begin{aligned} \|\mathbf{z}^* - \mathbf{w}^*\| &= \eta \|\mathbf{\Lambda}_{\tau_k}(\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) - \mathbf{F}_{\text{alt},\gamma_k}(\mathbf{w}^*))\| \\ &\leq \sqrt{2}\eta L \|\mathbf{\Lambda}_{\tau_k}\| \|\mathbf{z}^* - \mathbf{w}^*\| \\ &< \|\mathbf{z}^* - \mathbf{w}^*\| \end{aligned}$$

which is absurd. Therefore, we can deduce that $\mathbf{F}_{\text{alt},\gamma_k}(\mathbf{z}^*) = \mathbf{0}$. \square

B.5. Eigenvalue Asymptotic of $\mathbf{\Lambda}_{\tau}\mathbf{F}_{\text{alt},\gamma}$

We present the eigenvalues asymptotic of $\mathbf{\Lambda}_{\tau}\mathbf{F}_{\text{alt},\gamma}$. Following few statements in proof of Theorem 5.3, one can deduce that the eigenvalues of order $\Theta(\sqrt{\epsilon})$, divided by $\sqrt{\epsilon}$, are solutions of the equation

$$0 = \det \begin{bmatrix} \kappa\mathbf{I} - \sqrt{\epsilon}\mathbf{A} + \sqrt{\epsilon}(1 - \kappa\sqrt{\epsilon}\eta)\mathbf{C}_1(\lambda\mathbf{I} - \mathbf{D})^{-1}\mathbf{C}_1^\top & -\mathbf{L}_q \\ (1 - \kappa\sqrt{\epsilon}\eta)\mathbf{L}_q^\top & \kappa\mathbf{I} \end{bmatrix}$$

By Lemma B.2, notice that eigenvalues of order $\Theta(\sqrt{\epsilon})$, divided by $\sqrt{\epsilon}$, converge to the solutions of

$$0 = \det \begin{bmatrix} \kappa\mathbf{I} & -\mathbf{L}_q \\ \mathbf{L}_q^\top & \kappa\mathbf{I} \end{bmatrix},$$

as $\epsilon \rightarrow 0$. Therefore, these eigenvalues have asymptotics that are equivalent to those of the type (i) eigenvalues in Theorem 4.2.

B.6. Proof of Theorem 5.4

Proof of Proposition 5.4. We first investigate the necessary and/or sufficient condition for each type of eigenvalues λ_j of $\mathbf{H}_{\text{alt}2,\tau} = \mathbf{H}_{\text{alt},\tau,\sqrt{\tau}}$, which is categorized in Theorem 5.3, to lie in \mathcal{P}_η .

- (i) First, consider the type (i) eigenvalues $\lambda_j = \left(\frac{\eta\sigma_k^2}{2} \pm \frac{\sqrt{\eta^2\sigma_k^4 - 4\sigma_k^2}}{2} \right) \sqrt{\epsilon} + o(\sqrt{\epsilon})$ for some k . The radicand of $\sqrt{\eta^2\sigma_k^4 - 4\sigma_k^2}$ is negative for $0 < \eta < \frac{1}{L}$, since $\eta\sigma_k < \frac{1}{L} \cdot L = 1$. Therefore, the type (i) eigenvalues is in a form $\lambda_j = \left(\frac{\eta\sigma_k^2}{2} \pm i \frac{\sqrt{4\sigma_k^2 - \eta^2\sigma_k^4}}{2} \right) \sqrt{\epsilon} + o(\sqrt{\epsilon})$. Here, the leading term of real part of λ_j is $\frac{\eta\sigma_k^2}{2} \sqrt{\epsilon}$ that is positive, so $\lambda_j(\epsilon) \in \mathcal{P}_\eta$ for sufficiently large τ .

(ii) Secondly, consider the type (ii) eigenvalues, $\lambda_j = \epsilon\mu_k + o(\epsilon)$ for some k . Since the coefficient of leading term is μ_k , for sufficiently small ϵ , $\mu_k > 0$ implies that $\lambda_j(\epsilon) \in \mathcal{P}_\eta$, and $\mu_k < 0$ implies that $\lambda_j(\epsilon) \notin \mathcal{P}_\eta$. Recall that the μ_k are the eigenvalues of the restricted Schur complement \mathbf{S}_{res} in Theorem 5.3, which are nonzero due to Assumption 3'. Therefore, $\mathbf{S}_{\text{res}} \succeq \mathbf{0}$ if and only if there exists some τ^* such that $\tau > \tau^*$ implies that every λ_j of order $\Theta(\epsilon)$ satisfies $\lambda_j(\epsilon) \in \mathcal{P}_\eta$.

(iii) Finally, consider the type (iii) eigenvalues, $\lambda_j = \nu_k + o(1)$ for some k . By the inequality

$$\|\mathbf{B}\| = \left\| \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{bmatrix} \right\| = \left\| \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ -\mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \right\| \leq \|\mathbf{DF}(\mathbf{z}^*)\| = L,$$

we have $\|\nu_k\| \leq L$ for the eigenvalues of $-\mathbf{B}$. Now, let $0 < \eta < \frac{1}{L}$, and suppose that $\mathbf{B} \not\leq \mathbf{0}$. Then, there exists some k such that $-L \leq \nu_k < 0$. Since the half-open interval $[-L, 0)$ is contained in the complement of $\bar{\mathcal{P}}_\eta$, for sufficiently large τ , we have $\lambda_j(\epsilon) \notin \mathcal{P}_\eta$ clearly. On the other hand, suppose that $\mathbf{B} \preceq \mathbf{0}$. Then $\nu_k > 0$ for all k , and it implies that for $0 < \eta < \frac{1}{L}$ and sufficiently large τ , we have $\lambda_j(\epsilon) \in \mathcal{P}_\eta$ for all j such that $\lambda_j(\epsilon) = \nu_k + o(1)$. Therefore, $\mathbf{B} \preceq \mathbf{0}$ if and only if there exists some $0 < \eta < \frac{1}{L}$ such that τ being sufficiently large implies that every λ_j of order $\Theta(1)$ satisfies $\lambda_j(\epsilon) \in \mathcal{P}_\eta$.

Combining all the previous discussions, we can conclude that $\mathbf{S}_{\text{res}} \succeq \mathbf{0}$ and $\mathbf{B} \preceq \mathbf{0}$ if and only if, for any η satisfying $0 < \eta < \frac{1}{L}$, there exists sufficiently large τ such that $\lambda_j(\epsilon) \in \mathcal{P}_\eta$ for all j . The conclusion then follows from the Proposition 4.1. \square

B.7. Proof of Theorem 5.5

In proving Theorem 5.5, we need the following results.

Proposition B.7. *Under Assumption 1, $f \in C^2$, $0 < \eta < \frac{\sqrt{5}-1}{2\sqrt{2}L}$ and $\tau \geq 4$, we have $\det(D\mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}(\mathbf{z})) \neq 0$ for all \mathbf{z} .*

Proof. We begin by observing that

$$D\mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}(\mathbf{z}) = \mathbf{I} - \eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z} - \eta\mathbf{\Lambda}_\tau \mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z})) (\mathbf{I} - \eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z})).$$

Under Assumption 1 and by (27) with $\frac{\sqrt{5}-1}{2\sqrt{2}L} < \frac{1}{\sqrt{2}L}$, we have $\|\mathbf{DF}_{\text{alt},\sqrt{\tau}}\| \leq \sqrt{2}L$. Hence, whenever $0 < \eta < \frac{\sqrt{5}-1}{2\sqrt{2}L}$ and $\tau \geq 4$, we obtain the bound

$$\begin{aligned} & \left\| \eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z} - \eta\mathbf{\Lambda}_\tau \mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z})) (\mathbf{I} - \eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z})) \right\| \\ & \leq \left\| \eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z} - \eta\mathbf{\Lambda}_\tau \mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z})) \right\| \left\| \mathbf{I} - \eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z}) \right\| \\ & \leq \sqrt{2}\eta L (1 + \sqrt{2}\eta L) \\ & < 1. \end{aligned}$$

It follows that any eigenvalue of $\eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z} - \eta\mathbf{\Lambda}_\tau \mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z})) (\mathbf{I} - \eta\mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}(\mathbf{z}))$ has its magnitude strictly less than 1, and hence, $D\mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}(\mathbf{z})$ cannot have zero eigenvalue. Therefore, $\det(D\mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}(\mathbf{z})) \neq 0$ holds. \square

For convenience, let us define a subset of the complex plane, for a real negative constant $a < 0$,

$$\mathcal{O}_a^\# := \left\{ z \in \mathbb{C} : |z - a| < \frac{|a|}{2} \right\},$$

which is an open disk centered at a with radius $\frac{|a|}{2}$.

Lemma B.8. $\mathcal{O}_a^\# \cap \mathcal{P}_\eta = \emptyset$ for any real negative constant $a < 0$ and real positive constant $\eta > 0$.

Proof. Noticing that \mathcal{O}_a lies in the left half plane, the only region to care about is \mathbb{C}°_- . Thus, if the disk \mathcal{O}_a and the peanut-shaped \mathcal{P}_η do not intersect on that region, then the assertion follows immediately.

Consider a circle centered at origin with radius R , denoted by \mathcal{O}^* . Then, if the circle \mathcal{O}^* and the boundary of \mathcal{O}_a intersects, then it intersects at a point z_1 with a real part $\operatorname{Re} z_1 = \frac{R^2}{2a} + \frac{3a}{8}$. Similarly, if the circle \mathcal{O}^* and the boundary of \mathcal{P}_η intersects, then it intersects at a point z_2 with a real part $\operatorname{Re} z_2 = \frac{1}{4\eta} + \frac{\eta R^2}{4} - \sqrt{\frac{1}{16\eta^2} - \frac{3\eta^2 R^4}{16} + \frac{3R^2}{8}}$. For \mathcal{O}_a and \mathcal{P}_η to have an overlap, there must exist some R such that $\operatorname{Re} z_1 = \operatorname{Re} z_2$. We show that such R does not exist for any $a < 0$ and $\eta > 0$, by proving the following statement

$$\operatorname{Re} z_2 - \operatorname{Re} z_1 = \frac{1}{4\eta} + \frac{\eta R^2}{4} - \left(\frac{R^2}{2a} + \frac{3a}{8}\right) - \sqrt{\frac{1}{16\eta^2} - \frac{3\eta^2 R^4}{16} + \frac{3R^2}{8}} > 0$$

for any $a < 0$, $\eta > 0$ and $R > 0$. This is done by showing that the following

$$\begin{aligned} & \left(\frac{1}{4\eta} + \frac{\eta R^2}{4} - \frac{R^2}{2a} - \frac{3a}{8}\right)^2 - \left(\frac{1}{16\eta^2} - \frac{3\eta^2 R^4}{16} + \frac{3R^2}{8}\right) \\ &= \frac{9a^2}{64} - \frac{3a}{16\eta} + R^2 \left(\frac{1}{8} - \frac{3a\eta}{16} - \frac{1}{4a\eta}\right) + R^4 \left(\frac{1}{4a^2} - \frac{\eta}{4a} + \frac{\eta^2}{4}\right) > 0 \end{aligned}$$

holds for any negative a with positive η and R , and this concludes the proof. \square

Definition 5. Given a C^1 mapping w , the set $\mathcal{A}^*(w) := \{z^* : z^* = w(z^*), \rho(Dw(z^*)) > 1\}$ is the set of **strictly unstable equilibrium points**.

Theorem B.9 (Lee et al. (2019), Theorem 2)). Let w be a C^1 mapping such that $\det(Dw(z)) \neq 0$ for all z . Then the set of initial points that converge to a unstable equilibrium point has (Lebesgue) measure zero, i.e., $\mu(\{z_0 : \lim_{k \rightarrow \infty} w^k(z_0) \in \mathcal{A}^*(w)\}) = 0$.

Proposition B.10. Let z^* be a strict non-minimax point i.e., $z^* \in \mathcal{T}^*$. Under Assumptions 3, there exists a positive constant $\tau^* > 0$ such that $z^* \in \mathcal{A}^*(w_{\text{alt}, \tau, \sqrt{\tau}})$ for any $\tau > \tau^*$.

Proof. By Proposition 4.1, we have

$$\begin{aligned} \mathcal{A}^*(w_{\text{alt}, \tau, \sqrt{\tau}}) &= \{z^* : z^* = w_{\text{alt}, \tau, \sqrt{\tau}}(z^*), \rho(Dw_{\text{alt}, \tau, \sqrt{\tau}}(z^*)) > 1\} \\ &= \{z^* : z^* = w_{\text{alt}, \tau, \sqrt{\tau}}(z^*), \exists \lambda \in \operatorname{spec}(\mathbf{H}_{\text{alt}2, \tau}(z^*)) \text{ s.t. } \lambda \notin \bar{\mathcal{P}}_\eta\}. \end{aligned}$$

For any strict non-minimax point $z^* \in \mathcal{T}^*$, either $\mathbf{S}_{\text{res}}(z^*)$ or $-\mathbf{B}(z^*)$ has at least one strictly negative eigenvalue. First, suppose that $\mathbf{S}_{\text{res}}(z^*)$ has a strictly negative eigenvalue $\mu < 0$. By Theorem 5.3, there exists a constant τ^* such that at least one eigenvalue of $\mathbf{H}_{\text{alt}2, \tau}(z^*)$ lies in a disk $\mathcal{O}_{\mu\epsilon}^\sharp$ for any $\tau > \tau^*$. So by Lemma B.8, we would have $\mathcal{O}_{\mu\epsilon}^\sharp \cap \mathcal{P}_\eta = \emptyset$. On the other hand, suppose that $-\mathbf{B}(z^*)$ has a strictly negative eigenvalue $\nu < 0$. Similarly, by Theorem 5.3, there exists a constant τ^* such that at least one eigenvalue of $\mathbf{H}_{\text{alt}2, \tau}(z^*)$ lies in a disk \mathcal{O}_ν^\sharp for any $\tau > \tau^*$. So by Lemma B.8, we would have $\mathcal{O}_\nu^\sharp \cap \mathcal{P}_\eta = \emptyset$. Therefore, we can conclude that for any $z^* \in \mathcal{T}^*$, there exists a constant τ^* such that $z^* \in \mathcal{A}^*(w_{\text{alt}, \tau, \sqrt{\tau}})$ for any $\tau > \tau^*$. \square

Proof of Theorem 5.5. Because $z^* \in \mathcal{A}^*(w_{\text{alt}, \tau, \sqrt{\tau}})$ implies

$$\left\{z_0 : \lim_{k \rightarrow \infty} w_{\text{alt}, \tau, \sqrt{\tau}}^k(z_0) = z^*\right\} \subset \left\{z_0 : \lim_{k \rightarrow \infty} w_{\text{alt}, \tau, \sqrt{\tau}}^k(z_0) \in \mathcal{A}^*(w_{\text{alt}, \tau, \sqrt{\tau}})\right\},$$

by Theorem B.9, there exists a positive constant $\tau^* > 0$ such that

$$\mu \left(\left\{z_0 : \lim_{k \rightarrow \infty} w_{\text{alt}, \tau, \sqrt{\tau}}^k(z_0) = z^*\right\} \right) = 0$$

for any $\tau > \tau^*$.

Moreover, if \mathcal{T}^* is finite, then a maximum of τ^* for all $z^* \in \mathcal{T}^*$ is also finite. Let us denote such maximum by τ_{\max}^* . Then, for any $\tau > \tau_{\max}^*$ we have $\mathcal{T}^* \subset \mathcal{A}^*(\mathbf{w}_{\text{alt},\tau,\sqrt{\tau}})$. This implies that

$$\left\{ z_0 : \lim_{k \rightarrow \infty} \mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}^k(z_0) \in \mathcal{T}^* \right\} \subset \left\{ z_0 : \lim_{k \rightarrow \infty} \mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}^k(z_0) \in \mathcal{A}^*(\mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}) \right\},$$

for any $\tau > \tau_{\max}^*$, and by Theorem B.9, we can conclude that

$$\mu \left(\left\{ z_0 : \lim_{k \rightarrow \infty} \mathbf{w}_{\text{alt},\tau,\sqrt{\tau}}^k(z_0) \in \mathcal{T}^* \right\} \right) = 0 \quad \square.$$

C. Proofs for Section 6

C.1. Proof of Theorem 6.1

In proving the theorem, we begin with the following observations. Recall that, the dynamical system of the Alt2-EG-ITS method at time k is as follows

$$\mathbf{w}_{\text{alt},\tau_k,\sqrt{\tau_k}}(z) = z - \eta \Lambda_{\tau_k} \mathbf{F}_{\text{alt},\sqrt{\tau_k}}(z - \eta \Lambda_{\tau_k} \mathbf{F}_{\text{alt},\sqrt{\tau_k}}(z)).$$

Then, by the Taylor expansion around z^* , we have

$$\begin{aligned} \mathbf{w}_{\text{alt},\tau_k,\sqrt{\tau_k}}(z) &= \mathbf{w}_{\text{alt},\tau_k,\sqrt{\tau_k}}(z^*) + D\mathbf{w}_{\text{alt},\tau_k,\sqrt{\tau_k}}(z^*)(z - z^*) + o(z - z^*) \\ &= z^* + (\mathbf{I} - \eta \Lambda_{\tau_k} D\mathbf{F}_{\text{alt},\sqrt{\tau_k}}(z^*) + \eta^2 (\Lambda_{\tau_k} D\mathbf{F}_{\text{alt},\sqrt{\tau_k}}(z^*))^2) (z - z^*) + o(z - z^*) \end{aligned} \quad (28)$$

For convenience, let us denote $\mathbf{x}_k := z_k - z^*$, $\mathbf{y}_k := o(z_k - z^*)$ and

$$\mathbf{A}_k := \mathbf{I} - \eta \Lambda_{\tau_k} D\mathbf{F}_{\text{alt},\sqrt{\tau_k}}(z^*) + \eta^2 (\Lambda_{\tau_k} D\mathbf{F}_{\text{alt},\sqrt{\tau_k}}(z^*))^2. \quad (29)$$

Then, for arbitrary k_0 such that $0 \leq k_0 < k$, we have

$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{w}_{\text{alt},\tau_k,\sqrt{\tau_k}}(z_k) - z^* \\ &= \mathbf{A}_k \mathbf{x}_k + \mathbf{y}_k \\ &= \mathbf{A}_k (\mathbf{A}_{k-1} \mathbf{x}_{k-1} + \mathbf{y}_{k-1}) + \mathbf{y}_k \\ &= \dots \\ &= \left(\prod_{j=k_0}^k \mathbf{A}_j \right) \mathbf{x}_{k_0} + \sum_{i=k_0}^k \left(\prod_{j=i+1}^k \mathbf{A}_j \right) \mathbf{y}_i \end{aligned}$$

which implies that

$$\|\mathbf{x}_{k+1}\| \leq \left(\prod_{j=k_0}^k \|\mathbf{A}_j\| \right) \|\mathbf{x}_{k_0}\| + \sum_{i=k_0}^k \left(\prod_{j=i+1}^k \|\mathbf{A}_j\| \right) \|\mathbf{y}_i\|. \quad (30)$$

Then, the rest of the proof is to show that, for a certain increasing sequence of τ_k , the RHS of (30) (and thus $\|\mathbf{x}_k\|$) decreases at a certain rate. Note that, for a simpler setting, where \mathbf{A}_k is fixed for all k with $\|\mathbf{A}_k\| < 1$, Polyak (1987, Theorem 2.1.2.1) showed in a few lines that the RHS of (30) (and thus $\|\mathbf{x}_k\|$) decreases at an exponential rate. However, our proof is not as straightforward as that of (Polyak, 1987, Theorem 2.1.2.1), since here we not only consider \mathbf{A}_k that varies over time k , but also satisfies $\lim_{k \rightarrow \infty} \|\mathbf{A}_k\| = 1$ as shown next.

Lemma C.1. *Suppose Assumptions 1 and 3' hold, and let z^* be an equilibrium point that satisfies the second-order necessary condition of local minimax points. Then, there exists $K > 0$ such that the matrix \mathbf{A}_k (29) of the Alt2-EG-ITS with $(\tau_k, \sqrt{\tau_k})$, where τ_k is increasing and $\lim_{k \rightarrow \infty} \tau_k = \infty$, for any $0 < \eta < 1/L$, satisfies*

$$\|\mathbf{A}_k\| \leq 1 - \frac{2}{\tau_k^{1+c}}$$

for all $k \geq K$ and for any $c > 0$.

Proof. Recall that, under Assumption 3', Theorem 5.3 implies that there are three types of eigenvalue asymptotics of $\mathbf{H}_{\text{alt}2,\tau} := \mathbf{\Lambda}_\tau D\mathbf{F}_{\text{alt},\sqrt{\tau}}$. This directly implies that the matrix \mathbf{A}_k has the following three types of eigenvalue asymptotics

$$\begin{aligned} \text{(i)} \quad & 1 - \eta \left(\left(\frac{\eta\sigma_j^2}{2} \pm i \frac{\sqrt{4\sigma_j^2 - \eta^2\sigma_j^4}}{2} \right) \frac{1}{\sqrt{\tau_k}} + o\left(\frac{1}{\sqrt{\tau_k}}\right) \right) + \eta^2 \left(\left(\frac{\eta\sigma_j^2}{2} \pm i \frac{\sqrt{4\sigma_j^2 - \eta^2\sigma_j^4}}{2} \right) \frac{1}{\sqrt{\tau_k}} + o\left(\frac{1}{\sqrt{\tau_k}}\right) \right)^2 \\ & = 1 - \frac{\eta^2\sigma_j^2}{2\sqrt{\tau_k}} + \frac{\eta^4\sigma_j^4}{4\tau_k} - \eta^2 \frac{4\sigma_j^2 - \eta^2\sigma_j^4}{4\tau_k} + i \left(\mp \eta \frac{\sqrt{4\sigma_j^2 - \eta^2\sigma_j^4}}{2\sqrt{\tau_k}} \pm \frac{\eta^3\sigma_j^2\sqrt{4\sigma_j^2 - \eta^2\sigma_j^4}}{2\tau_k} \right) \\ & \quad - \eta o\left(\frac{1}{\sqrt{\tau_k}}\right) + \eta^2 o\left(\frac{1}{\tau_k}\right) + \eta^3\sigma_m^2 \frac{1}{\sqrt{\tau_k}} o\left(\frac{1}{\sqrt{\tau_k}}\right) \pm i \frac{\eta^2\sqrt{4\sigma_m^2 - \eta^2\sigma_m^4}}{2\sqrt{\tau_k}} o\left(\frac{1}{\sqrt{\tau_k}}\right), \end{aligned}$$

$$\text{(ii)} \quad 1 - \eta \left(\frac{\mu_j}{\tau_k} + o\left(\frac{1}{\tau_k}\right) \right) + \eta^2 \left(\frac{\mu_j}{\tau_k} + o\left(\frac{1}{\tau_k}\right) \right)^2 = 1 - \frac{1}{\tau_k}(\eta\mu_j + o(1)) + \frac{1}{\tau_k^2}(\eta^2\mu_j^2 + o(1)),$$

$$\text{(iii)} \quad 1 - \eta(\nu_j + o(1)) + \eta^2(\nu_j + o(1))^2,$$

where μ_j are the eigenvalues of the restricted Schur complement $\mathbf{S}_{\text{res}}(\mathbf{H})$, ν_j are the nonzero eigenvalues of $-\mathbf{B}$, and σ_j are the singular values of \mathbf{C}_2 . Note that, under Assumption 3', $\mathbf{S}_{\text{res}}(\mathbf{H})$ is invertible.

Since the sequence τ_k is increasing and $\lim_{k \rightarrow \infty} \tau_k = \infty$, for any $c > 0$, there exists sufficiently large K such that the followings satisfy

$$\left| o\left(\frac{1}{\sqrt{\tau_k}}\right) \right| \leq \min_j \left(\frac{\eta\sigma_j^2}{8\sqrt{\tau_k}}, \frac{\sqrt{4\sigma_j^2 - \eta^2\sigma_j^4}}{8} \right) \quad (31)$$

$$\left| o\left(\frac{1}{\tau_k}\right) \right| \leq \frac{\mu_j}{4\tau_k}, \quad (32)$$

$$6.25\eta\mu_j \leq \tau_k, \quad (33)$$

$$\max_j \left(\frac{4}{\eta^2\sigma_j^2}, \frac{8}{\eta\mu_j} \right) \leq \tau_k^c, \quad (34)$$

$$\max_j \left(6.5\eta^2\sigma_j^2 + \eta^3\sigma_j^2 + 1, 5\eta^2\sigma_j^2 + \frac{\eta^3\sigma_j^2}{2} + \frac{\eta\sqrt{4\sigma_j^2 - \eta^2\sigma_j^4}}{8} \right) \leq \sqrt{\tau_k}, \quad (35)$$

$$1 - \frac{1}{2}\eta\nu_j(1 - \eta\nu_j) \leq 1 - \frac{2}{\tau_k^{1+c}}, \quad (36)$$

$$|1 - \eta(\nu_j + o(1)) + \eta^2(\nu_j + o(1))^2| \leq 1 - \frac{\eta\nu_j(1 - \eta\nu_j)}{2}, \quad (37)$$

for all $k > K$, and for all μ_j and σ_j . Using the above inequalities for any $k > K$, we characterize the type (i), (ii) and (iii)

eigenvalues as follows. First, for the type (i) eigenvalue, we have

$$\begin{aligned}
 |\lambda| &\leq \sqrt{\left(1 - \frac{\eta^2 \sigma_j^2}{4\sqrt{\tau_k}}\right)^2 + \left(\eta \frac{\sqrt{4\sigma_j^2 - \eta^2 \sigma_j^4}}{4\sqrt{\tau_k}}\right)^2} \quad \text{by (31) and (35)} \\
 &= \sqrt{1 - \frac{\eta^2 \sigma_j^2}{2\sqrt{\tau_k}} + \frac{\eta^2 \sigma_j^2}{4\tau_k}} \\
 &\leq \sqrt{1 - \frac{\eta^2 \sigma_j^2}{4\sqrt{\tau_k}}} \\
 &\leq 1 - \frac{\eta^2 \sigma_j^2}{2\sqrt{\tau_k}} \\
 &\leq 1 - \frac{2}{\tau_k^{(1+c)/2}} \quad \text{by (34)}.
 \end{aligned}$$

Next, for the type (ii) eigenvalue, we have we have

$$\begin{aligned}
 |\lambda| &\leq \sqrt{\left(1 - \frac{\eta \mu_j}{2\tau_k}\right)^2 + \left(\frac{\eta \mu_j}{4\tau_k}\right)^2} \quad \text{by (32) and (33)} \\
 &= \sqrt{1 - \frac{\eta \mu_j}{\tau_k} + \frac{5\eta^2 \mu_j^2}{16\tau_k^2}} \\
 &\leq \sqrt{1 - \frac{\eta \mu_j}{2\tau_k}} \\
 &\leq 1 - \frac{\eta \mu_j}{4\tau_k} \\
 &\leq 1 - \frac{2}{\tau_k^{1+c}} \quad \text{by (34)}.
 \end{aligned}$$

Finally, for the type (iii) eigenvalue, we have

$$\begin{aligned}
 |\lambda| &\leq 1 - \frac{1}{2}\eta\nu_j(1 - \eta\nu_j) \quad \text{by (37)} \\
 &\leq 1 - \frac{2}{\tau_k^{1+c}} \quad \text{by (36)}.
 \end{aligned}$$

Combining all the previous arguments, the assertion then follows. \square

From now on, let $\tau_k = k^{1/(2+2c)}$. Then, by Lemma C.1, the upper bound of $\|\mathbf{A}_k\|$ is $\|\mathbf{A}_k\| \leq 1 - \frac{2}{\sqrt{k}}$. Our next step is to find a bound of RHS in (30). To accomplish this, we require the following technical Lemma, which will be used later.

Lemma C.2. *For any k and $n \geq k + 1$, there exists a constant M_1 such that $\frac{1}{\sqrt{k}} \prod_{j=k}^n \left(1 - \frac{1}{\sqrt{j}}\right) \leq M_1 \frac{1}{\sqrt{n}}$*

Proof. We will find upper bound of $\prod_{j=i+1}^n \left(1 - \frac{1}{\sqrt{j}}\right)$ via finding tight upper and lower bound of $\prod_{j=2}^n \left(1 - \frac{1}{\sqrt{j}}\right)$.

From the fact that logarithm function is an increasing function on \mathbb{R}^+ , we have

$$\begin{aligned}
 & \int_2^{n+1} \log\left(1 - \frac{1}{\sqrt{x}}\right) dx \\
 &= -\sqrt{n+1} - \frac{1}{2} \log(n+1) + n \log\left(1 - \frac{1}{\sqrt{n+1}}\right) + \sqrt{2} + \log 2 + \sinh^{-1}(1) \\
 &\geq \sum_{j=2}^n \log\left(1 - \frac{1}{\sqrt{j}}\right) \\
 &\geq \int_1^n \log\left(1 - \frac{1}{\sqrt{x}}\right) dx \\
 &= -\sqrt{n} - \frac{\log(n)}{2} + (n-1) \log\left(1 - \frac{1}{\sqrt{n}}\right) + 1.
 \end{aligned}$$

Therefore, for sufficiently large n , we have upper bound

$$\begin{aligned}
 \prod_{j=2}^n \left(1 - \frac{1}{\sqrt{j}}\right) &\leq C_1 e^{-\sqrt{n+1}} \cdot \frac{1}{\sqrt{n+1}} \cdot \left(1 - \frac{1}{\sqrt{n+1}}\right)^n \\
 &\leq C_2 \frac{1}{\sqrt{n+1}} e^{-(\sqrt{n}+\sqrt{n+1})} \\
 &\quad \text{(From the fact that } \left(1 - \frac{1}{\sqrt{n+1}}\right)^n \times e^{\sqrt{n}} \text{ converges to } 1/\sqrt{e}\text{)} \\
 &\leq C_2 \frac{1}{\sqrt{n}} e^{-(\sqrt{n-1}+\sqrt{n})} \tag{38}
 \end{aligned}$$

for some positive constant C_1 and C_2 . Similarly, we have lower bound

$$\begin{aligned}
 \prod_{j=2}^n \left(1 - \frac{1}{\sqrt{j}}\right) &\geq C_3 e^{-\sqrt{n}} \cdot \frac{1}{\sqrt{n}} \cdot \left(1 - \frac{1}{\sqrt{n}}\right)^n \\
 &\geq C_4 \frac{1}{\sqrt{n}} e^{-(\sqrt{n-1}+\sqrt{n})} \\
 &\quad \text{(From the fact that } \left(1 - \frac{1}{\sqrt{n}}\right)^n \times e^{\sqrt{n-1}} \text{ converges to } 1/\sqrt{e}\text{)} \\
 &\geq C_4 \frac{1}{\sqrt{n+1}} e^{-2\sqrt{n+1}} \tag{39}
 \end{aligned}$$

for some positive constant C_3 and C_4 . Using (38) and (39), we can bound the product $\frac{1}{\sqrt{k}} \prod_{i=k}^n \left(1 - \frac{1}{\sqrt{i}}\right)$ as follows.

$$\begin{aligned}
 \frac{1}{\sqrt{k}} \prod_{i=k}^n \left(1 - \frac{1}{\sqrt{i}}\right) &\leq \frac{1}{\sqrt{k}} \cdot \frac{C_2 e^{-(\sqrt{n-1}+\sqrt{n})}}{\frac{C_4}{\sqrt{k}} e^{-2\sqrt{k}}} \\
 &\leq \frac{1}{\sqrt{n}} \cdot \frac{C_2 e^{-(\sqrt{n-1}+\sqrt{n})}}{C_4 e^{-2\sqrt{k}}}
 \end{aligned}$$

Here, for any $n \geq k+1$, the $e^{-(\sqrt{n-1}+\sqrt{n}-2\sqrt{k})} \leq 1$. Therefore, $\frac{1}{\sqrt{k}} \prod_{i=k}^n \left(1 - \frac{1}{\sqrt{i}}\right) \leq M_1 \frac{1}{\sqrt{n}}$ for $M_1 = \frac{C_2}{C_4}$, and this completes the proof. \square

By utilizing the previous Lemma C.2, we can bound the RHS in (30) via mathematical induction. The precise statement is as follows.

Proposition C.3. *Let $f \in C^3$ and $(\tau_k, \gamma_k) = (k^{1/(2+2c)}, k^{1/(4+4c)})$ for $c > 0$. Let \mathbf{z}^* be an equilibrium point that satisfies the necessary condition of local minimax points. Then, under Assumptions 1 and 3', \mathbf{z}^* is asymptotically stable point of Alt2-EG-ITS.*

Proof. Under the assumption $f \in C^3$, the Lagrange's form of the remainder (Apostol, 1991, §7.7) is $o(\mathbf{z} - \mathbf{z}^*) = \frac{\mathbf{w}_k^{(2)}(\xi)}{2!}(\mathbf{z} - \mathbf{z}^*)^2$ for some ξ lies in the closed interval between \mathbf{z} and \mathbf{z}^* . Moreover, $\mathbf{w}_k^{(2)}/2$ can be bounded by some positive constant M_2 on small neighborhood $B_{\delta_1}(\mathbf{z}^*)$. Let $M := \max(1, M_1)$. Then, we will prove that, if for some k -th step, the iterate lies in $B_{\delta_k}(\mathbf{z}^*)$ where $\delta_k := \min\left(\delta_1, \frac{1}{M_2 \max(1, M_1)\sqrt{k}}\right)$, the future step \mathbf{z}_n converge to \mathbf{z}^* as n goes to infinity.

Under aforementioned settings, suppose that \mathbf{z}_k lies in $B_{\delta_k}(\mathbf{z}^*)$ for some k . Then

$$\begin{aligned} \|o(\mathbf{z}_k - \mathbf{z}^*)\| &\leq M_2 \|\mathbf{z}_k - \mathbf{z}^*\|^2 \\ &\leq M_2 \delta_k \|\mathbf{z}_k - \mathbf{z}^*\| \\ &\leq \frac{1}{\max(1, M_1)\sqrt{k}} \|\mathbf{z}_k - \mathbf{z}^*\| \\ &= \frac{1}{M\sqrt{k}} \|\mathbf{z}_k - \mathbf{z}^*\|. \end{aligned}$$

Hence, the first iterate satisfies the following iterates

$$\begin{aligned} \|\mathbf{x}_{k+1}\| &= \|\mathbf{A}_k \mathbf{x}_k + o(\mathbf{x}_k)\| \\ &\leq \|\mathbf{A}_k\| \|\mathbf{x}_k\| + \|o(\mathbf{x}_k)\| \\ &\leq \left(1 - \frac{2}{\sqrt{k}}\right) \|\mathbf{x}_k\| + \frac{1}{M\sqrt{k}} \|\mathbf{x}_k\| \quad (\text{By Lemma C.1}) \\ &\leq \left(1 - \frac{1}{\sqrt{k}}\right) \|\mathbf{x}_k\|. \end{aligned}$$

We use induction on n to prove that $\|\mathbf{x}_n\| \leq \prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\|$. The $n = k$ case is trivial. Suppose that $\|\mathbf{x}_n\| \leq \prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\|$ for some n such that $n \geq k + 1$. Then, we have

$$\begin{aligned} \|\mathbf{x}_{n+1}\| &\leq \|\mathbf{A}_n\| \|\mathbf{x}_n\| + \|o(\mathbf{x}_n)\| \\ &\leq \left(1 - \frac{2}{\sqrt{n}}\right) \prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\| + M_2 \left(\prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right)\right)^2 \|\mathbf{x}_k\|^2 \quad (\text{By Lemma C.1}) \\ &\leq \left(1 - \frac{2}{\sqrt{n}}\right) \prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\| + \left(\prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right)\right)^2 \frac{1}{M\sqrt{k}} \|\mathbf{x}_k\| \\ &\leq \left(1 - \frac{2}{\sqrt{n}}\right) \prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\| + \left(\prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right)\right)^2 \frac{1}{M_1\sqrt{k}} \|\mathbf{x}_k\| \end{aligned}$$

$$\begin{aligned} &\leq \left(1 - \frac{2}{\sqrt{n}}\right) \prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\| + \frac{1}{\sqrt{n}} \prod_{j=k}^{n-1} \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\| \quad (\text{By Lemma C.2}) \\ &\leq \prod_{j=k}^n \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_k\|. \end{aligned}$$

Then the assertion follows from the fact that $\lim_{n \rightarrow \infty} \prod_{j=k}^n \left(1 - \frac{1}{\sqrt{j}}\right) = 0$. \square

We are now ready to prove the asymptotic stability of the Alt2-EG-ITS.

Proof of Theorem 6.1. Suppose that the stationary point \mathbf{z}^* satisfies the second-order necessary condition of local minimax points. Then, by Proposition C.3, Alt2-EG-ITS can converge to the \mathbf{z}^* . Moreover, following a few statements in proof of Proposition C.3, the convergence rate is upper bounded by $\prod_{j=i}^k \left(1 - \frac{1}{\sqrt{j}}\right) \|\mathbf{x}_i\|$, and the product has (tight) upper bound $O\left(\frac{1}{\sqrt{k}} e^{-2\sqrt{k}}\right)$. These arguments complete the proof. \square

Intuitively, as the value of k increases, the neighborhood ensuring local convergence gradually shrinks. However, once it lies inside of the neighborhood, the future iterates converge to the \mathbf{z}^* .

D. Proof for Section 7

D.1. Proof of Theorem 7.1

We need the following lemma for proving Theorem 7.1.

Lemma D.1. *Suppose Assumptions 1 and 2 hold. Then, there exists a stationary point $(\mathbf{x}^*, \mathbf{y}^*)$ that satisfies*

$$\langle \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), \mathbf{x} - \mathbf{x}^* \rangle - \langle \nabla_{\mathbf{y}} f(\bar{\mathbf{x}}, \mathbf{y}), \mathbf{y} - \mathbf{y}^* \rangle \geq f(\mathbf{x}, \mathbf{y}) - f(\bar{\mathbf{x}}, \mathbf{y}) \geq \frac{\eta}{\gamma} \left(1 - \frac{L_{\mathbf{x}} \eta}{2 \gamma}\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})\|^2, \quad (40)$$

for all $\mathbf{x} \in \mathbb{R}^{d_1}$ and $\mathbf{y} \in \mathbb{R}^{d_2}$, where we let $\bar{\mathbf{x}} := \mathbf{x} - \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$.

Proof. Assumption 2 implies that the following four inequalities hold, for any $\mathbf{x} \in \mathbb{R}^{d_1}$ and $\mathbf{y} \in \mathbb{R}^{d_2}$:

$$\begin{aligned} f(\mathbf{x}^*, \mathbf{y}) &\geq f(\mathbf{x}, \mathbf{y}) + \langle \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), \mathbf{x}^* - \mathbf{x} \rangle, \\ f(\bar{\mathbf{x}}, \mathbf{y}^*) &\geq f(\mathbf{x}^*, \mathbf{y}^*), \\ f(\bar{\mathbf{x}}, \mathbf{y}^*) &\leq f(\bar{\mathbf{x}}, \mathbf{y}) + \langle \nabla_{\mathbf{y}} f(\bar{\mathbf{x}}, \mathbf{y}), \mathbf{y}^* - \mathbf{y} \rangle, \\ f(\mathbf{x}^*, \mathbf{y}) &\leq f(\mathbf{x}^*, \mathbf{y}^*). \end{aligned}$$

By summing over the above four inequalities, we have the first inequality of (40). Moreover, the second inequality of (40) can be shown as

$$\begin{aligned} f(\bar{\mathbf{x}}, \mathbf{y}) &\leq f(\mathbf{x}, \mathbf{y}) + \langle \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), \bar{\mathbf{x}} - \mathbf{x} \rangle + \frac{L_{\mathbf{x}}}{2} \|\bar{\mathbf{x}} - \mathbf{x}\|^2 \\ &= f(\mathbf{x}, \mathbf{y}) + \left\langle \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), -\frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) \right\rangle + \frac{L_{\mathbf{x}}}{2} \left\| \frac{\eta}{\gamma} \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) \right\|^2 \\ &= f(\mathbf{x}, \mathbf{y}) - \frac{\eta}{\gamma} \left(1 - \frac{L_{\mathbf{x}} \eta}{2 \gamma}\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})\|^2, \end{aligned}$$

where the first inequality uses Assumption 1 and (Nesterov, 2018, Theorem 2.1.5). \square

We are now ready to show that the Alt2-EG-TS

$$\begin{cases} \mathbf{u}_k = \mathbf{x}_k - \frac{\eta}{\tau_k} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k), \\ \mathbf{v}_k = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f\left(\mathbf{x}_k - \frac{\eta}{\gamma_k} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k), \mathbf{y}_k\right), \end{cases} \quad \begin{cases} \mathbf{x}_{k+1} = \mathbf{x}_k - \frac{\eta}{\tau_k} \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k), \\ \mathbf{y}_{k+1} = \mathbf{y}_k + \eta \nabla_{\mathbf{y}} f\left(\mathbf{u}_k - \frac{\eta}{\gamma_k} \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k), \mathbf{v}_k\right). \end{cases}$$

finds a stationary point under Assumptions 1 and 2, *i.e.*, the smoothness and star-convex-star-concave assumptions on f .

Proof of Theorem 7.1. To prove Theorem 7.1, we begin with the following observation. For simplicity, let us denote $\bar{\mathbf{x}}_k := \mathbf{x}_k - \frac{\eta}{\gamma_k} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)$ and $\bar{\mathbf{u}}_k := \mathbf{u}_k - \frac{\eta}{\gamma_k} \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)$. Then, we have the inequality

$$\begin{aligned} & \|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_{k+1} - \mathbf{y}^*\|^2 \\ &= \left\| \mathbf{x}_k - \frac{\eta}{\tau_k} \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k) - \mathbf{x}^* \right\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k + \eta \nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k) - \mathbf{y}^*\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 - \frac{2\eta}{\tau_k} (\langle \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k), \mathbf{x}_k - \mathbf{x}^* \rangle - \langle \nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k), \mathbf{y}_k - \mathbf{y}^* \rangle) \\ &\quad + \frac{\eta^2}{\tau_k^2} \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 + \frac{\eta^2}{\tau_k} \|\nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k)\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 - \frac{2\eta}{\tau_k} (\langle \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k), \mathbf{x}_k - \mathbf{u}_k \rangle - \langle \nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k), \mathbf{y}_k - \mathbf{v}_k \rangle) \\ &\quad - \frac{\eta}{\gamma_k} \left(1 - \frac{\eta L_{\mathbf{x}}}{2\gamma_k}\right) \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 + \frac{\eta^2}{\tau_k^2} \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 + \frac{\eta^2}{\tau_k} \|\nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k)\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 - \frac{2\eta}{\tau_k} \left(\left\langle \nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k), \frac{\eta}{\tau_k} \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k) \right\rangle - \langle \nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k), -\eta \nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k) \rangle \right) \\ &\quad - \frac{\eta}{\gamma_k} \left(1 - \frac{\eta L_{\mathbf{x}}}{2\gamma_k}\right) \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 + \frac{\eta^2}{\tau_k^2} \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 + \frac{\eta^2}{\tau_k} \|\nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k)\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 + \frac{\eta^2}{\tau_k^2} (\|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k) - \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|^2 - \|\nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|^2) \\ &\quad + \frac{\eta^2}{\tau_k} (\|\nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k) - \nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|^2 - \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|^2) - \frac{\eta}{\gamma_k} \left(1 - \frac{\eta L_{\mathbf{x}}}{2\gamma_k}\right) \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 + \frac{\eta^2}{\tau_k} (\|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k) - \nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|^2 + \|\nabla_{\mathbf{y}} f(\bar{\mathbf{u}}_k, \mathbf{v}_k) - \nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|^2) \\ &\quad - \frac{\eta^2}{\tau_k^2} \|\nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|^2 - \frac{\eta^2}{\tau_k} \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|^2 - \frac{\eta}{\gamma_k} \left(1 - \frac{\eta L_{\mathbf{x}}}{2\gamma_k}\right) \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 \\ &\quad + \frac{\eta^2}{\tau_k} \left(L_{\mathbf{x}}^2 + L_{\mathbf{y}}^2 + L_{\mathbf{x}}^2 L_{\mathbf{y}}^2 \frac{\eta^2}{\gamma_k} \left(1 + \frac{1}{\gamma_k}\right) \right) \|\mathbf{u}_k, \mathbf{v}_k) - (\mathbf{x}_k, \mathbf{y}_k)\|^2 + \frac{\eta^2 L_{\mathbf{y}}^2}{\tau_k \gamma_k} \|\mathbf{u}_k - \mathbf{x}_k\|^2 \\ &\quad - \frac{\eta^2}{\tau_k^2} \|\nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|^2 - \frac{\eta^2}{\tau_k} \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|^2 - \frac{\eta}{\gamma_k} \left(1 - \frac{\eta L_{\mathbf{x}}}{2\gamma_k}\right) \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 \\ &\quad - \frac{\eta^2}{\tau_k^2} \left(1 - \frac{\eta^2 L^2}{\tau_k} \left(1 + \frac{1}{2\gamma_k} \left(1 + \frac{1}{\gamma_k}\right)\right) - \frac{\eta^2 L^2}{\tau_k \gamma_k}\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|^2 \\ &\quad - \frac{\eta^2}{\tau_k} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_k} \left(1 + \frac{1}{\gamma_k}\right)\right)\right) \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|^2 - \frac{\eta}{\gamma_k} \left(1 - \frac{\eta L}{2\gamma_k}\right) \|\nabla_{\mathbf{x}} f(\mathbf{u}_k, \mathbf{v}_k)\|^2 \\ &\leq \|\mathbf{x}_k - \mathbf{x}^*\|^2 + \frac{1}{\tau_k} \|\mathbf{y}_k - \mathbf{y}^*\|^2 \end{aligned}$$

$$-\frac{\eta^2}{\tau_k} \left(1 - \frac{\eta^2 L^2}{\tau_k} \left(1 + \frac{1}{2\gamma_k} \left(3 + \frac{1}{\gamma_k}\right)\right)\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|^2 - \frac{\eta^2}{\tau_k} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_k} \left(1 + \frac{1}{\gamma_k}\right)\right)\right) \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|^2$$

where the first inequality uses Lemma D.1, and the second inequality uses $\tau_k \geq 1$, and the fourth inequality uses the update rules and Lemma B.1. These inequalities lead to the following lemma which is essential in proving the convergence of $\|\mathbf{F}_{\text{alt}, \gamma_k}(\mathbf{x}_k, \mathbf{y}_k)\|$.

Lemma D.2. *Let $\tau_i \geq 1$. Then, the series*

$$\sum \frac{\eta^2}{\tau_i^2} \left(1 - \frac{\eta^2 L^2}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \mathbf{y}_i)\|^2 + \sum \frac{\eta^2}{\tau_i} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_i} \left(1 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_i, \mathbf{y}_i)\|^2 \quad (41)$$

is bounded.

Proof. By taking a telescoping summation, then we have

$$\begin{aligned} & \sum_{i=1}^k \frac{\eta^2}{\tau_i^2} \left(1 - \frac{\eta^2 L^2}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \mathbf{y}_i)\|^2 + \sum_{i=1}^k \frac{\eta^2}{\tau_i} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_i} \left(1 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_i, \mathbf{y}_i)\|^2 \\ & \leq \|\mathbf{x}_1 - \mathbf{x}^*\|^2 + \frac{1}{\tau_1} \|\mathbf{y}_1 - \mathbf{y}^*\|^2 - \|\mathbf{x}_{i+1} - \mathbf{x}^*\|^2 - \frac{1}{\tau_{i+1}} \|\mathbf{y}_{i+1} - \mathbf{y}^*\|^2 \\ & \leq \|\mathbf{x}_1 - \mathbf{x}^*\|^2 + \frac{1}{\tau_1} \|\mathbf{y}_1 - \mathbf{y}^*\|^2. \end{aligned}$$

Therefore, the series is bounded. \square

We will then show that the series increases monotonically as k increases. Then combining monotonicity and boundedness, one can deduce that both summands of (41) converge to zero as $i \rightarrow \infty$.

Lemma D.3. *Let $\tau_i \geq 1$ and $\gamma_i \geq 1$. Then, the series*

$$\sum_{i=j}^k \frac{\eta^2}{\tau_i^2} \left(1 - \frac{\eta^2 L^2}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \mathbf{y}_i)\|^2 + \sum_{i=j}^k \frac{\eta^2}{\tau_i} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_i} \left(1 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_i, \mathbf{y}_i)\|^2 \quad (42)$$

for some fixed j is monotonically increasing.

Proof. To verify the monotonicity of the series, we need to check the positivity of the summands. The coefficient of the norms are positive, when $\gamma_i > \frac{2\eta^2 L^2}{(1-\eta^2 L^2)}$ satisfied, because

$$\begin{aligned} \gamma_i & > \frac{2\eta^2 L^2}{(1-\eta^2 L^2)} \\ & \Rightarrow \gamma_i > \frac{2\eta^2 L^2}{(\tau_i - \eta^2 L^2)} \quad (\text{Since } \tau_i \geq 1) \\ & \Rightarrow \frac{\tau_i - \eta^2 L^2}{\eta^2 L^2} > \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right) \\ & \Rightarrow \frac{1}{\eta^2 L^2} > \frac{1}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right) \\ & \Rightarrow 1 - \frac{\eta^2 L^2}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right) > 0 \end{aligned}$$

and

$$\begin{aligned}
 \gamma_i &> \frac{\eta^2 L^2}{(1 - \eta^2 L^2)} \\
 &\Rightarrow \frac{1 - \eta^2 L^2}{\eta^2 L^2} > \frac{1}{2\gamma_i} \left(1 + \frac{1}{\gamma_i}\right) \\
 &\Rightarrow 1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_i} \left(1 + \frac{1}{\gamma_i}\right)\right) > 0.
 \end{aligned}$$

For fixed (τ, γ) , since the $0 < \eta < \frac{\sqrt{\gamma}}{\sqrt{2+\gamma}L}$ implies the condition $\gamma > \frac{2\eta^2 L^2}{(1-\eta^2 L^2)}$, the both summands are positive. For increasing (τ_i, γ_i) , since the γ_i increases without any upper bound, for sufficiently large (fixed) i , the condition $\gamma_i > \frac{2\eta^2 L^2}{(1-\eta^2 L^2)}$ holds for any $0 < \eta < \frac{1}{L}$. Therefore, the summation is monotonically increasing. \square

We are now ready to prove the convergence of each norms $\|\nabla_{\mathbf{x}} f(\mathbf{x}_k, \mathbf{y}_k)\|$ and $\|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_k, \mathbf{y}_k)\|$. For the case of Alt2-EG-FTS, by Lemmas D.2 and D.3, for $\tau \geq 1$, $\gamma \geq 1$ and $\eta < \frac{\sqrt{\gamma}}{\sqrt{2+\gamma}L}$, the both summands $\frac{\eta^2}{\tau^2} \left(1 - \frac{\eta^2 L^2}{\tau} \left(1 + \frac{1}{2\gamma} \left(3 + \frac{1}{\gamma}\right)\right)\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \mathbf{y}_i)\|^2$ and $\frac{\eta^2}{\tau} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma} \left(1 + \frac{1}{\gamma}\right)\right)\right) \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_i, \mathbf{y}_i)\|^2$ converge to zero as $i \rightarrow \infty$. Then the assertion follows from the fact that the coefficients of the both terms are invariant as i varies.

For the case of Alt2-EG-ITS, since the conditions $\tau_i \geq 1$, $\gamma > \frac{2\eta^2 L^2}{1-\eta^2 L^2}$ for $\eta < \frac{1}{L}$ are satisfied as $i \rightarrow \infty$, both summands $\frac{\eta^2}{\tau_i^2} \left(1 - \frac{\eta^2 L^2}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \mathbf{y}_i)\|^2$ and $\frac{\eta^2}{\tau_i} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_i} \left(1 + \frac{1}{\gamma_i}\right)\right)\right) \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_i, \mathbf{y}_i)\|^2$ converge to zero as $i \rightarrow \infty$ by Lemmas D.2 and D.3. However, the coefficient of the terms also diminishes as $i \rightarrow \infty$, therefore, the only thing we can say about the terms is the limit inferior of the both terms converge to zero.

Therefore, the rest of the proof is to demonstrate that, the summation of the coefficients in (42) are infinite. Since, $1 - \frac{\eta^2 L^2}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right) \rightarrow 1$ as $i \rightarrow \infty$, there exists i_0 such that any $i \geq i_0$ implies the $1 - \frac{\eta^2 L^2}{\tau_i} \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right) > \frac{1}{2}$. Therefore, we have

$$\begin{aligned}
 \sum_{i=i_0}^{\infty} \frac{\eta^2}{\tau_i^2} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_i} \left(3 + \frac{1}{\gamma_i}\right)\right)\right) &\geq \sum_{i=i_0}^{\infty} \frac{\eta^2}{2\tau_i^2} \\
 &= \frac{\eta^2}{2} \sum_{i=i_0}^{\infty} \frac{1}{\tau_i^2} \\
 &= \infty
 \end{aligned}$$

Therefore, $\liminf_{i \rightarrow \infty} \|\nabla_{\mathbf{x}} f(\mathbf{x}_i, \mathbf{y}_i)\| = 0$ holds.

Similarly, we have

$$\begin{aligned}
 \sum_{i=i_1}^{\infty} \frac{\eta^2}{\tau_i} \left(1 - \eta^2 L^2 \left(1 + \frac{1}{2\gamma_i} \left(1 + \frac{1}{\gamma_i}\right)\right)\right) &\geq \sum_{i=i_1}^{\infty} \frac{\eta^2}{2\tau_i} \\
 &= \frac{\eta^2}{2} \sum_{i=i_0}^{\infty} \frac{1}{\tau_i} \\
 &= \infty
 \end{aligned}$$

Therefore, $\liminf_{i \rightarrow \infty} \|\nabla_{\mathbf{y}} f(\bar{\mathbf{x}}_i, \mathbf{y}_i)\| = 0$ holds, and these arguments complete the proof. \square

D.2. Accumulation Point of the Alt2-EG-FTS is a Stationary Point

Proof. Let \tilde{z} be an accumulation point of the sequence $\{z_k\}_{k \geq 0}$. Then, there exists a subsequence $\{z_{k_j}\}_{j \geq 0}$ of original sequence such that $z_{k_j} \rightarrow \tilde{z}$ as $j \rightarrow \infty$ (Tao, 2016, Proposition 1.4.5). By triangle inequality and Lemma B.1, we have

$$\begin{aligned} 0 &\leq \|\mathbf{F}_{\text{alt},\gamma}(\tilde{z})\| \\ &\leq \|\mathbf{F}_{\text{alt},\gamma}(\tilde{z}) - \mathbf{F}_{\text{alt},\gamma}(z_{k_j})\| + \|\mathbf{F}_{\text{alt},\gamma}(z_{k_j})\| \\ &\leq \sqrt{L_x^2 + L_y^2 \left(1 + \frac{1}{\gamma}\right) \left(1 + L_x^2 \frac{\eta^2}{\gamma}\right)} \|\tilde{z} - z_{k_j}\| + \|\mathbf{F}_{\text{alt},\gamma}(z_{k_j})\|. \end{aligned}$$

From the fact that both $\|\tilde{z} - z_{k_j}\|$ and $\|\mathbf{F}_{\text{alt},\gamma}(z_{k_j})\|$ converge to zero as $j \rightarrow \infty$, we have $\|\mathbf{F}_{\text{alt},\gamma}(\tilde{z})\| = 0$, and this completes the proof. \square

E. Proofs for Section 8

E.1. Proof of Example 1

Proof. Consider the function $f(x, y) = -x^2 + 2xy$. Its saddle-gradient is $\mathbf{F}(x, y) = (-2x + 2y, -2x)$, and it has a unique stationary point $(0, 0)$. Moreover, f satisfies Assumption 1 with $L_x = L_y = 2$, so $L := \sqrt{L_x^2 + L_y^2} = 2\sqrt{2}$.

For any $\delta > 0$ and any (x, y) satisfying $|x - 0| \leq \delta$ and $|y - 0| \leq \delta$, the inequality

$$f(0, y) = 0 \leq f(0, 0) = 0 \leq \max_{y' : |y' - 0| \leq \delta} f(x, y') = -x^2 + 2|x\delta|$$

holds, so the stationary point $(0, 0)$ is a local minimax point. In addition, since $\nabla_{yy}f(x, y)$ is degenerate, it is a non-strict local minimax point.

Since

$$D\mathbf{F} = \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ -\mathbf{C}^\top & -\mathbf{B} \end{bmatrix} = \begin{bmatrix} -2 & 2 \\ -2 & 0 \end{bmatrix},$$

we have $\mathbf{C}_2 = [2] \in \mathbb{R}^{1 \times 1}$ and $q = \text{rank}(\mathbf{C}_2) = 1$. Then, the matrix \mathbf{U} is of size 0×0 , and so is $\mathbf{S}_{\text{res}}(D\mathbf{F}) = \mathbf{U}^\top (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top) \mathbf{U}$, which is vacuously positive definite. Therefore, Assumption 3' is satisfied, and thus, by Theorems 5.4 and 6.1, both Alt2-EG-FTS and Alt2-EG-ITS are asymptotically stable at $(0, 0)$, respectively.

Let us now show that the (vanilla) two-timescale EG is unstable at $(0, 0)$ for any choice of timescale separation τ . In particular, we show that the eigenvalues of its Jacobian

$$\mathbf{H}_\tau = \mathbf{\Lambda}_\tau D\mathbf{F} = \begin{bmatrix} \frac{1}{\tau} & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} -2 & 2 \\ -2 & 0 \end{bmatrix} = \begin{bmatrix} -2\epsilon & 2\epsilon \\ -2 & 0 \end{bmatrix},$$

which are $-\epsilon \pm i\sqrt{4\epsilon - \epsilon^2}$, are outside \mathcal{P}_η for any choice of τ , based on Propositions 3.1 and 4.1. By the definition of \mathcal{P}_η , it is enough to show that all $\epsilon > 0$ satisfy

$$\begin{aligned} &(\eta\epsilon + \frac{1}{2})^2 + \eta^2(4\epsilon - \epsilon^2) + \frac{3}{4} > \sqrt{1 + 3\eta^2(4\epsilon - \epsilon^2)} \\ &\Leftrightarrow 1 + \eta\epsilon + 4\eta^2\epsilon > \sqrt{1 + 3\eta^2(4\epsilon - \epsilon^2)} \\ &\Leftrightarrow 1 + \eta^2\epsilon^2 + 16\eta^4\epsilon^2 + 2\eta\epsilon + 8\eta^2\epsilon + 8\eta^3\epsilon^2 > 1 + 12\eta^2\epsilon - 3\eta^2\epsilon^2 \\ &\Leftrightarrow 4\eta^2\epsilon^2 + 16\eta^4\epsilon^2 + 2\eta\epsilon - 4\eta^2\epsilon + 8\eta^3\epsilon^2 > 0 \\ &\Leftrightarrow 4\eta^2\epsilon + 16\eta^4\epsilon + 2\eta - 4\eta^2 + 8\eta^3\epsilon > 0 \\ &\Leftrightarrow \epsilon(4\eta^2 + 16\eta^4 + 8\eta^3) > 4\eta^2 - 2\eta \end{aligned}$$

where the fifth equivalent comes from the $\epsilon > 0$. Since $4\eta^2 - 2\eta \leq 0$ for any $0 < \eta < \frac{1}{L} = \frac{1}{2\sqrt{2}}$, the above inequality indeed holds for any $\epsilon > 0$, which completes the proof. \square