

# Paper Review: Solving Min-Max Optimization with Hidden Structure via Gradient Descent Ascent

Zixuan Pan

zxp@umich.edu

## Abstract

*Many of the recent advances in machine learning have been inspired by min-max games, however many of these methods lack a solid math background telling about their dynamics. Whether these methods could always converge, and which point they will converge to remains much unexplained. This paper proposes a formulation for a kind of games called hidden convex concave games, which is related to more general min-max games. The paper proposed two main results: if the game is equipped with a strictly convex concave objective function, then the convergence to a von Neumann solution can be guaranteed by designing a Lyapunov function. Meanwhile, stronger regularization can make convergence faster, but could also make the convergence point shift. Also some implementation of this paper's result on real datasets is provided in this review.*

## 1. Introduction

Many of the recent advances in machine learning have been inspired by min-max games, such as generative adversarial networks (GAN) in image generation, and actor-critic algorithm in reinforcement learning [1, 4, 8]. These problems are often solved by applying gradient descent ascent (GDA) to some neural network based models. As a main feature of neural networks, the optimization problems could always be non-convex non-concave, but they are bounded by a convex-concave loss function. This kind of min-max game is defined as Hidden Convex-Concave (HCC) games [3]. A formal mathematical definition of it could be: given  $F : \mathbb{R}^N \rightarrow X \subset \mathbb{R}^n$  and  $G : \mathbb{R}^M \rightarrow Y \subset \mathbb{R}^m$  and a continuous convex-concave function  $L : X \times Y \rightarrow \mathbb{R}$ , such that the min-max game is  $\min_{\theta \in \mathbb{R}^N} \max_{\phi \in \mathbb{R}^M} L(F(\theta), G(\phi))$ .

Studying the dynamics of GDA algorithms in HCC games is critical, as there's a bunch of underlying issues that are related to unstable training of HCC games. The vanilla GAN requires heavily parameter tuning to make the training process stable [1]. Meanwhile, the GDA algorithm may often fall into some meaningless saddle points, which generates undesired solutions. WGAN solves part of the problems by making modifications to the loss function, but it still needs some tricks like gradient clipping or adding extra regularization terms to stabilize the training process [5], which still needs proper parameter tuning. Thus, it's very useful to study the dynamics of HCC games and find how to make GDA converge to a meaningful minima. This could shed light on the advances of such min-max game based models by reducing the parameter tuning work with explainability.

This paper focused on finding appropriate initialization to guarantee stability of the dynamic system and the relation between convergence and regularization.

## 2. Preliminary

### 2.1. Hidden Convex Concave Games

We will start with giving a formal definition of HCC games.

**Definition 1 (HCC games)**  $L : \mathbb{R}^N \times \mathbb{R}^M \rightarrow \mathbb{R}$  is convex concave if for every  $y \in \mathbb{R}^M$   $L(\cdot, y)$  is convex and for every  $x \in \mathbb{R}^N$   $L(x, \cdot)$  is concave. Function  $L$  will be called strictly convex concave if it is convex concave and for every  $x \times y \in \mathbb{R}^N \times \mathbb{R}^M$  either  $L(\cdot, y)$  is strictly convex or  $L(x, \cdot)$  is strictly concave.

As the property of min-max games, we assume two players are involved in the game (solving the optimization problem). The goal of one player is to minimize the objective function  $L$ , while the goal of the other player is to maximize  $L$ . We assume the minimization player is equipped with  $n$  functions  $f_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}$  while the maximization player is equipped with  $m$  functions  $g_j : \mathbb{R}^{m_j} \rightarrow \mathbb{R}$ . In order to study the dynamic, we will assume that  $f_i, g_j, L$  are all  $C^2$  functions, which means twice differentiable. The inputs  $\theta_i \in \mathbb{R}^{n_i}$  and  $\phi_j \in \mathbb{R}^{m_j}$  are grouped in two vectors

$$\begin{aligned} \boldsymbol{\theta} &= [\theta_1 \quad \theta_2 \quad \cdots \quad \theta_n]^\top & \mathbf{F}(\boldsymbol{\theta}) &= [f_1(\theta_1) \quad f_2(\theta_2) \quad \cdots \quad f_n(\theta_n)]^\top \\ \boldsymbol{\phi} &= [\phi_1 \quad \phi_2 \quad \cdots \quad \phi_m]^\top & \mathbf{G}(\boldsymbol{\phi}) &= [g_1(\phi_1) \quad g_2(\phi_2) \quad \cdots \quad g_m(\phi_m)]^\top \end{aligned}$$

In the view of neural network, the  $\phi$  and  $\theta$  can be viewed as some high dimensional inputs of neural networks, while  $F$  and  $G$  are two different networks.

Therefore, the hidden convex concave game is:

$$(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*) = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^N} \arg \max_{\boldsymbol{\phi} \in \mathbb{R}^M} L(\mathbf{F}(\boldsymbol{\theta}), \mathbf{G}(\boldsymbol{\phi})).$$

where  $N = \sum_{i=1}^n n_i$  and  $M = \sum_{j=1}^m m_j$

## 2.2. Dynamics of HCC games

**Definition 2 (Von Neumann Solutions)** Given a convex concave function  $L$ , we denote the stationary points of  $L$  as (global) Nash equilibria of the min-max game. The set of all equilibria is defined as von Neumann solutions of  $L$ , denote as  $Solution(L)$ .

When the game is defined on a convex compact set, existence of at least a solution is guaranteed by von Neumann's minimax theorem [2]. By choosing appropriate  $f_i$  and  $g_i$ , the convex compact setting can be obtained, and thus in the following section we assume there's  $Solution(L)$  is not empty. Then the gradient descent ascent dynamic of HCC games can be defined as follow:

$$\begin{aligned} \dot{\theta}_i &= -\nabla_{\theta_i} L(\mathbf{F}(\boldsymbol{\theta}), \mathbf{G}(\boldsymbol{\phi})) = -\nabla_{\theta_i} f_i(\theta_i) \frac{\partial L}{\partial f_i}(\mathbf{F}(\boldsymbol{\theta}), \mathbf{G}(\boldsymbol{\phi})) \\ \dot{\phi}_j &= \nabla_{\phi_j} L(\mathbf{F}(\boldsymbol{\theta}), \mathbf{G}(\boldsymbol{\phi})) = \nabla_{\phi_j} g_j(\phi_j) \frac{\partial L}{\partial g_j}(\mathbf{F}(\boldsymbol{\theta}), \mathbf{G}(\boldsymbol{\phi})) \end{aligned} \tag{1}$$

## 2.3. Reparametrization

The stability of HCC games is guaranteed under appropriate initialization, which will be discussed later. Here we'd like to make some definitions of the initialization and solutions of dynamic system first.

**Lemma 3 (Unique Solution of Dynamic System)** Let  $k : \mathbb{R}^d \rightarrow \mathbb{R}$  be a  $C^2$  function. Let  $h : \mathbb{R} \rightarrow \mathbb{R}$  be a  $C^1$  function and  $\mathbf{x}(t)$  denote the unique solution of the dynamical system  $\Sigma_1$ . Then the unique solution for dynamical system  $\Sigma_2$  is  $\mathbf{z}(t) = \mathbf{x} \left( \int_0^t h(s) ds \right)$

$$\left\{ \begin{array}{l} \dot{\mathbf{x}} = \nabla k(\mathbf{x}) \\ \mathbf{x}(0) = \mathbf{x}_{init} \end{array} \right\} : \Sigma_1 \quad \left\{ \begin{array}{l} \dot{\mathbf{z}} = h(t) \nabla k(\mathbf{z}) \\ \mathbf{z}(0) = \mathbf{x}_{init} \end{array} \right\} : \Sigma_2$$

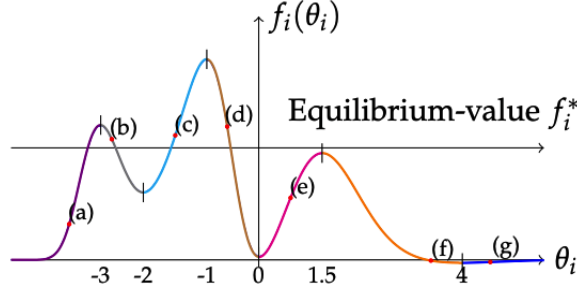


Figure 1. Relation between initialization and GDA trajectory.

See Appendix A for proof.

By choosing  $h(t) = -\partial L(\mathbf{F}(t), \mathbf{G}(t))/\partial f_i$  and  $h(t) = \partial L(\mathbf{F}(t), \mathbf{G}(t))/\partial g_j$  respectively, we can connect the dynamics of each  $\theta_i$  and  $\phi_j$  under Equation (1) to gradient ascent on  $f_i$  and  $g_j$ . Applying Lemma 1, we get that trajectories of  $\theta_i$  and  $\phi_j$  under Equation (1) are restricted to be subsets of the corresponding gradient ascent trajectories with the same initializations. In fact, Lemma 1 guarantees the monotonicity of GDA trajectories. As shown in Figure 1  $\theta_i(t)$  can not escape the purple section if it is initialized at (a) neither the orange section if it is initialized at (f). This limits the attainable values that  $f_i(t)$  and  $g_j(t)$  can take for a specific initialization.

**Definition 4 (Image of Dynamic System)** For each initialization  $\mathbf{x}(0)$  of  $\Sigma_1$ ,  $\text{Im}_k(\mathbf{x}(0))$  is the image of  $k \circ \mathbf{x} : \mathbb{R} \rightarrow \mathbb{R}$ .

Applying Definition 2 in Figure 1,  $\text{Im}_{f_i}(\theta_i(0)) = (f_i(-2), f_i(-1))$  if  $\theta_i$  is initialized at (c). Additionally, observe that in each colored section  $f_i(\theta_i(t))$  uniquely identifies  $\theta_i(t)$ . Generally, even in the case that  $\theta_i$  are vectors, Lemma 1 implies that for a given  $\theta_i(0)$ ,  $f_i(\theta_i(t))$  uniquely identifies  $\theta_i(t)$ . As a result we get that a new dynamical system involving only  $f_i, g_j$  and initializations.

**Theorem 5 (GDA without Intermediate States)** For each initialization  $(\theta(0), \phi(0))$  of Equation (1), there are  $C^1$  functions  $X_{\theta_i(0)}, X_{\phi_j(0)}$  such that  $\theta_i(t) = X_{\theta_i(0)}(f_i(t))$  and  $\phi_j(t) = X_{\phi_j(0)}(g_j(t))$ . If  $(\theta(t), \phi(t))$  satisfy Equation (1) then  $f_i(t) = f_i(\theta_i(t))$  and  $g_j(t) = g_j(\phi_j(t))$  satisfy

$$\begin{aligned} \dot{f}_i &= - \left\| \nabla_{\theta_i} f_i \left( X_{\theta_i(0)}(f_i) \right) \right\|^2 \frac{\partial L}{\partial f_i}(\mathbf{F}, \mathbf{G}) \\ \dot{g}_j &= \left\| \nabla_{\phi_j} g_j \left( X_{\phi_j(0)}(g_j) \right) \right\|^2 \frac{\partial L}{\partial g_j}(\mathbf{F}, \mathbf{G}) \end{aligned} \quad (2)$$

The proof is given in Appendix B.

By determining the ranges of  $f_i$  and  $g_j$ , an initialization clearly dictates if a von Neumann solution is attainable. In Figure 1 for example, any point of the pink, orange or blue colored section like (e), (f) or (g) can not converge to a von Neumann solution with  $f_i(\theta_i) = f_i^*$ . The notion of *safety* captures which initializations can converge to a given element of  $\text{Solution}(L)$ .

**Definition 6 (Safe Initialization)** The initialization  $(\theta(0), \phi(0))$  is called safe for a  $(\mathbf{p}, \mathbf{q}) \in \text{Solution}(L)$  if  $\phi_i(0)$  and  $\theta_j(0)$  are not stationary points of  $f_i$  and  $g_j$  respectively and  $p_i \in \text{Im}_{f_i}(\theta_i(0))$  and  $q_j \in \text{Im}_{g_j}(\phi_j(0))$ .

## 2.4. Stability of Dynamic Systems

Let  $f : D \rightarrow \mathbb{R}^n$  be a locally Lipschitz map from a domain  $D \subset \mathbb{R}^n$  to  $\mathbb{R}^n$ . We consider dynamical systems of the form

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) \quad (\star)$$

A point  $\bar{x}$  for which  $f(\bar{x}) = 0$  is called a fixed point. We will be interested in the following notions of stability for the fixed point points of Equation  $(\star)$ .

**Definition 7 (Stability properties [7])** *The fixed point  $\mathbf{x} = \mathbf{0}$  of Equation  $(\star)$  is stable if, for each  $\epsilon > 0$ , there is a  $\delta = \delta(\epsilon) > 0$  such that*

$$\|\mathbf{x}(0)\| < \delta \implies \|\mathbf{x}(t)\| < \epsilon \quad \forall t \geq 0$$

*unstable if it is not stable*

*asymptotically stable if it is stable and  $\delta$  can be chosen such that*

$$\|\mathbf{x}(0)\| < \delta \implies \lim_{t \rightarrow \infty} \mathbf{x}(t) = \mathbf{0}$$

Lyapunov Function plays an important role in determining the stability of fixed points of a dynamic system.

**Theorem 8 (Lyapunov Theorem [7])** *Let  $\mathbf{x} = \mathbf{0}$  be a fixed point point for Equation  $(\star)$  and  $D \subset \mathbb{R}^n$  be a domain containing  $\mathbf{x} = \mathbf{0}$ . Let  $V : D \rightarrow \mathbb{R}$  be a continuously differentiable function such that*

$$\begin{aligned} V(\mathbf{0}) = 0 \text{ and } V(\mathbf{x}) > 0 \text{ in } D - \{\mathbf{0}\} \\ \dot{V}(\mathbf{x}) \leq 0 \text{ in } D \end{aligned}$$

*then  $\mathbf{x} = \mathbf{0}$  is stable. Moreover if*

$$\dot{V}(\mathbf{x}) < 0 \text{ in } D - \{\mathbf{0}\}$$

*then  $\mathbf{x} = \mathbf{0}$  is asymptotically stable.*

Besides, some notions of stability will be used in the following sections. We call an equilibrium  $\mathbf{x}^*$  of an autonomous dynamical system  $\dot{\mathbf{x}} = \mathcal{D}(\mathbf{x}(t))$  stable if for every neighborhood  $U$  of  $\mathbf{x}^*$  there is a neighborhood  $V$  of  $\mathbf{x}^*$  such that if  $\mathbf{x}(0) \in V$  then  $\mathbf{x}(t) \in U$  for all  $t \geq 0$ . We call a set  $S$  asymptotically stable if there exists a neighborhood  $\mathcal{R}$  such that for any initialization  $\mathbf{x}(0) \in \mathcal{R}$ ,  $\mathbf{x}(t)$  approaches  $S$  as  $t \rightarrow +\infty$ . If  $\mathcal{R}$  is the whole space the set globally asymptotically stable.

## 3. Lyapunov Function for GDA of HCC Games

### 3.1. General Case

The section will be about designing a Lyapunov Function for Equation (2).

**Lemma 9** *If  $L$  is convex concave and  $(\phi(0), \theta(0))$  is a safe for  $(\mathbf{p}, \mathbf{q}) \in \text{Solution}(L)$ , then the following quantity is non-increasing under the dynamics of Equation (3):*

$$H(\mathbf{F}, \mathbf{G}) = \sum_{i=1}^N \int_{p_i}^{f_i} \frac{z - p_i}{\|\nabla f_i(X_{\theta_i(0)}(z))\|^2} dz + \sum_{j=1}^M \int_{q_j}^{g_j} \frac{z - q_j}{\|\nabla g_j(X_{\phi_j(0)}(z))\|^2} dz$$

See Appendix C for proof.

$$\begin{aligned}\dot{H} &\leq L(\mathbf{p}, \mathbf{G}) - L(\mathbf{F}, \mathbf{G}) + L(\mathbf{F}, \mathbf{G}) - L(\mathbf{F}, \mathbf{q}) \\ &\leq L(\mathbf{p}, \mathbf{G}) - L(\mathbf{p}, \mathbf{q}) + L(\mathbf{p}, \mathbf{q}) - L(\mathbf{F}, \mathbf{q}) \leq 0\end{aligned}$$

The last inequality holds since  $(\mathbf{p}, \mathbf{q}) \in \text{Solution}(L)$ .

If  $(\mathbf{p}, \mathbf{q})$  is a saddle point of  $L$  then  $L(\mathbf{p}, \mathbf{G}) \leq L(\mathbf{p}, \mathbf{q}) \leq L(\mathbf{F}, \mathbf{q})$

**Theorem 10** *If  $L$  is convex concave and  $(\phi(0), \theta(0))$  is a safe for  $(\mathbf{p}, \mathbf{q}) \in \text{Solution}(L)$ , then  $(\mathbf{p}, \mathbf{q})$  is stable for Equation (2).*

See Appendix D for proof.

**Theorem 11 (Safety properties of Sigmoid Functions)** *If  $f_i$  and  $g_j$  are sigmoid functions and  $L$  is convex concave and there is a  $\alpha(\phi(0), \theta(0))$  that is safe for  $(\mathbf{p}, \mathbf{q}) \in \text{Solution}(L)$ , then  $(\mathbf{F}^{-1}(\mathbf{p}), \mathbf{G}^{-1}(\mathbf{q}))$  is stable for Equation (1).*

See Appendix 11 for Proof.

With these theorems, the stability of such a dynamic system is guaranteed with a pair  $(p, q)$  in the solution space.

**Theorem 12 (Bound of HCC games)** *Let  $(\mathbf{p}, \mathbf{q}) \in \text{Solution}(L)$ . Let  $R_{f_i}$  and  $R_{g_j}$  be the set of regular values of  $f_i$  and  $g_j$  respectively. Assume that there is a  $\xi > 0$  such that  $[p_i - \xi, p_i + \xi] \subseteq R_{f_i}$  and  $[q_j - \xi, q_j + \xi] \subseteq R_{g_j}$ . Define*

$$r(t) = \|\mathbf{F}(\theta(t)) - \mathbf{p}\|^2 + \|\mathbf{G}(\phi(t)) - \mathbf{q}\|^2.$$

*If  $f_i$  and  $g_j$  are proper functions, then for every  $\epsilon > 0$ , there is an  $\delta > 0$  such that*

$$r(0) < \delta \implies \forall t \geq 0 : r(t) < \epsilon.$$

See Appendix F for the proof.

This theorem provides that the dynamic system over time is bounded. However, it does not guarantee convergence to a certain point. Thus, we will further discuss  $L$  as a strictly convex concave function.

### 3.2. Hidden Strictly Convex Concave Games

Here we continue to study properties of hidden strictly convex concave games with  $L$  as a strictly convex concave function.

**Lemma 13 (Locally Asymptotically Stability of Strictly Convex Concave Functions)** *Let  $L$  be strictly convex concave and  $Z \subset \text{Solution}(L)$  is the non empty set of equilibria of  $L$  for which  $(\theta(0), \phi(0))$  is safe. Then  $Z$  is locally asymptotically stable for Equation (2).*

See Appendix G for proof.

**Theorem 14 (Convergence of Hidden Strictly Convex Concave Games)** *Let  $L$  be strictly convex concave and  $Z \subset \text{Solution}(L)$  is the non empty set of equilibria of  $L$  for which  $(\theta(0), \phi(0))$  is safe. Under the dynamics of Equation (1)  $(\mathbf{F}(\theta(t)), \mathbf{G}(\phi(t)))$  converges to a point in  $Z$  as  $t \rightarrow \infty$ .*

The theorem above guarantees convergence to a von Neumann solution for all initializations that are safe for at least one element of  $\text{Solution}(L)$ . However, this is not the same as global asymptotic stability. To get even stronger guarantees, we can assume that all initializations are safe. In this case it is straightforward to get a global asymptotic stability result:

**Corollary 15** *Let  $L$  be strictly convex concave and assume that all initializations are safe for at least one element of Solution ( $L$ ). The following set is globally asymptotically stable for continuous GDA dynamics.*

$$\{(\boldsymbol{\theta}^*, \boldsymbol{\phi}^*) \in \mathbb{R}^n \times \mathbb{R}^m : (F(\boldsymbol{\theta}^*), G(\boldsymbol{\phi}^*)) \in \text{Solution}(L)\}$$

### 3.3. Convergence via Regularization

Regularization has become a widely used technique in neural networks and also in GANs [1, 4]. It could also have impact on HCC games. Assume a utility  $L(\mathbf{x}, \mathbf{y})$  that is convex concave but not strictly. Here we will propose a modified utility  $L'$  that is strictly convex strictly concave. Specifically we will choose

$$L'(\mathbf{x}, \mathbf{y}) = L(\mathbf{x}, \mathbf{y}) + \frac{\lambda}{2}\|\mathbf{x}\|^2 - \frac{\lambda}{2}\|\mathbf{y}\|^2$$

The choice of the parameter  $\lambda$  captures the trade-off between convergence to the original equilibrium of  $L$  and convergence speed. On the one hand, invoking the implicit function theorem, we get that for small  $\lambda$  the equilibria of  $L$  are not significantly perturbed.

**Theorem 16 (Equilibria of Invertible Hessian)** *If  $L$  is a convex concave function with invertible Hessians at all its equilibria, then for each  $\epsilon > 0$  there is a  $\lambda > 0$  such that  $L'$  has equilibria that are  $\epsilon$ -close to the ones of  $L$ .*

See Appendix I for proof.

Note that invertibility of the Hessian means that  $L$  must have a unique equilibrium. On the other hand increasing  $\lambda$  increases the rate of convergence of safe initializations to the perturbed equilibrium

**Theorem 17 (Bounded by Regularization)** *Let  $(\boldsymbol{\theta}(0), \boldsymbol{\phi}(0))$  be a safe initialization for the unique equilibrium of  $L'(\mathbf{p}, \mathbf{q})$ . If*

$$r(t) = \|\mathbf{F}(\boldsymbol{\theta}(t)) - \mathbf{p}\|^2 + \|\mathbf{G}(\boldsymbol{\phi}(t)) - \mathbf{q}\|^2$$

*then there are initialization dependent constants  $c_0, c_1 > 0$  such that  $r(t) \leq c_0 \exp(-\lambda c_1 t)$ .*

See Appendix J for proof.

This theorem guarantees that the dynamic is bounded by regularization

## 4. Applications

### 4.1. Hidden strictly convex-concave games

A famous example of HCC game is GAN. The conclusion of the paper can be applied to all variants of GANs [1, 4, 8]. In the vanilla GAN architecture, as it is commonly referred, our goal is to obtain a generator, such that it can generate a distribution  $p_G$  that is close to an input data distribution  $p_{\text{data}}$ . To find such a generator function, we can use a discriminator  $D$  that judge if an input is real or fake. For the case of a discrete  $p_{\text{data}}$  over a set  $\mathcal{N}$ , the minimax problem of GAN is the following:

$$\min_{p_G(x) \geq 0, \sum_{x \in \mathcal{N}} p_G(x) = 1} \max_{D \in (0,1)^{|\mathcal{N}|}} V(G, D) = \sum_{x \in \mathcal{N}} p_{\text{data}}(x) \log(D(x)) + \sum_{x \in \mathcal{N}} p_G(x) \log(1 - D(x))$$

The problem above can be formulated as a constrained strictly convex-concave hidden game. On the one hand, for a fixed discriminator  $D^*$ , the  $V(G, D^*)$  is linear over the  $p_G(x)$ . On the other hand, for a fixed generator  $G^*$ ,  $V(G^*, D)$  is strongly-concave. We can implement the inequality constraints on both the generator probabilities and discriminator using sigmoid activations. For the equality constraint  $\sum_{x \in \mathcal{N}} p_G(x) = 1$  we can introduce a Lagrange multiplier. Theretically, as given in the paper, when having effectively removed the constraints, we

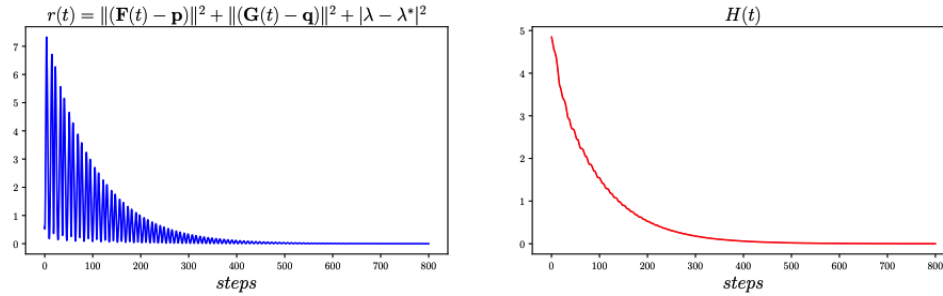


Figure 2. Theoretical l2 distance and Lyapunov function to the equilibria.

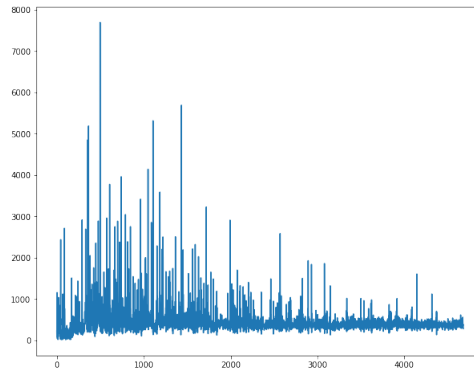


Figure 3. l2 distance on MNIST.

can see in Figure 2, the dynamics of Equation (1) converge to the unique equilibrium of the game. Meanwhile, the dynamic of GAN on MNIST dataset is implemented<sup>1</sup>. An approximation is done to obtain  $r(t)$ , as  $p_{data}$  is hard to obtain, where  $\|G(\phi(t)) - q\|$  term is approximated with  $\|D(G(x))\|$ . See Figure 3 for the l2 distance on MNIST. The result is similar to the theoretical one that l2 distance does not converge monotonically to 0 (left on Figure 2).

#### 4.2. Hidden Convex Concave Games with Regularization

Convergence with respect to regularization strength for vanilla GAN is given in Figure 4. Generally a larger regularization strength gives a faster convergence speed. However, a faster convergence cannot always guarantee converging to a von Neumann solution. Figure 5,6 shows that the generation with larger regularization could fail, indicating not converging to a meaningful solution.

Another example of GAN is WGAN [1]. One of the contributions of is to show that WGANs trained with Stochastic GDA can learn the parameters of Gaussian distributions whose samples are transformed by non-linear activation functions. The original WGAN formulation has a Lipschitz constraint in the discriminator function which is done by weight clipping. However, model performance is very sensitive to the weight clipping ratio. Thus, [5] replaced this constraint with a quadratic regularizer. The min-max problem for the case of one-dimensional Gaussian  $\mathcal{N}(0, \alpha_*^2)$  and linear discriminator  $D_v(x) = v^\top x$  with  $x^2$  activation is:

$$\begin{aligned}
 \min_{\alpha \in \mathbb{R}} \max_{v \in \mathbb{R}} V_{\text{WGAN}}(G_\alpha, D_v) &= \mathbb{E}_{\mathbf{X} \sim p_{\text{data}}}[D(\mathbf{X})] - \mathbb{E}_{\mathbf{X} \sim p_G}[D(\mathbf{X})] - v^2/2 \\
 &= \mathbb{E}_{x \sim \mathcal{N}(0, \alpha_*^2)}[vx] - \mathbb{E}_{x \sim \mathcal{N}(0, \alpha^2)}[vx] - v^2/2 \\
 &= (\alpha_*^2 - \alpha^2)v - v^2/2
 \end{aligned}$$

<sup>1</sup>See <https://github.com/zxp46/EECS559-Final-Project> for the code.

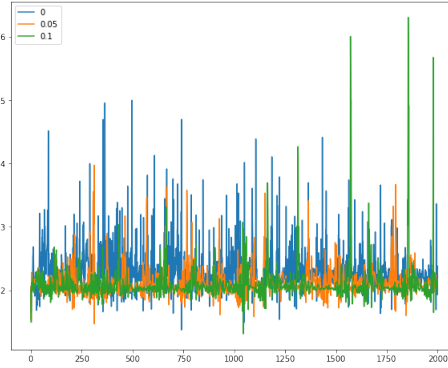


Figure 4. Convergence and regularization strength.

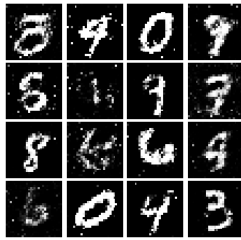


Figure 5. Generation samples for  $\lambda = 0$ .

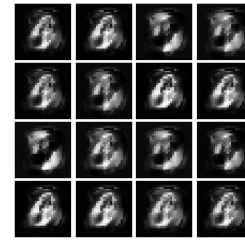


Figure 6. Generation samples for  $\lambda = 0.1$ .

Observe that  $V_{\text{WGAN}}$  is not convex-concave but it can be posed as a hidden strictly convex-concave game with  $G(\alpha) = (\alpha_*^2 - \alpha^2)$  and  $F(v) = v$ . When computing expectations analytically without sampling, Theorem 14 still guarantees convergence.

## 5. Discussion and Conclusion

This paper discusses about the dynamic of HCC games, where the main contribution is proposing that if it is a hidden strictly convex concave game, then the convergence to a meaningful von Neumann solution can be guaranteed by designing a Lyapunov function. Meanwhile, stronger regularization can make convergence faster, but could also make the convergence point shift. For more general cases, as also discussed in [6], the convergence to a meaningful point is not definite. Another limitation of this paper is that the Lyapunov function is hard to implement with code, as the distribution of generated samples and original samples are difficult to obtain. But overall, this paper shed light on formulating a better theoretical background for some zero-sum games.

## References

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan, 2017. 1, 6, 7
- [2] Hichem Ben-El-Mechaiekh and Robert Dimand. A simpler proof of the von neumann minimax theorem. *American Mathematical Monthly*, 118, 08 2011. 2
- [3] Lampros Flokas, Emmanouil-Vasileios Vlatakis-Gkaragkounis, and Georgios Piliouras. Solving min-max optimization with hidden structure via gradient descent ascent. *arXiv preprint arXiv:2101.05248*, 2021. 1
- [4] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. 1, 6
- [5] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. *Advances in neural information processing systems*, 30, 2017. 1, 7



- [6] Ya-Ping Hsieh, Panayotis Mertikopoulos, and Volkan Cevher. The limits of min-max optimization algorithms: Convergence to spurious non-critical sets. In *International Conference on Machine Learning*, pages 4337–4348. PMLR, 2021.  
8
- [7] H.K. Khalil. *Nonlinear Systems*. Pearson Education. Prentice Hall, 2002. 4
- [8] Vijay Konda and John Tsitsiklis. Actor-critic algorithms. *Advances in neural information processing systems*, 12, 1999.  
1, 6

### A. Proof of Lemma 3

*Proof.*

$$\begin{aligned} \mathbf{z}(0) &= \mathbf{x} \left( \int_0^0 h(s) ds \right) = \mathbf{x}(0) = \mathbf{x}_{\text{init}} \\ \dot{\mathbf{z}} &= \dot{\mathbf{x}} \left( \int_0^t h(s) ds \right) \times \frac{d \left[ \int_0^t h(s) ds \right]}{dt} \\ &= \nabla k \left( \mathbf{x} \left( \int_0^t h(s) ds \right) \right) h(t) = \nabla k(\mathbf{z}) h(t) \end{aligned}$$

gives a unique solution for  $\Sigma_2$

### B. Proof of Theorem 5

*Proof.* Let us first study a simpler dynamical system ( $\Sigma^*$ ) with unique solution of  $\gamma_{\theta_i(0)}(t)$ .

$$(\Sigma^*) \equiv \left\{ \begin{array}{l} \dot{\mathbf{z}} = \nabla f_i(\mathbf{z}) \\ \mathbf{z}(0) = \boldsymbol{\theta}_i(0) \end{array} \right\}$$

It is easy to observe that:

$$\dot{f}_i = \nabla f(\mathbf{z}) \dot{\mathbf{z}} = \|\nabla f(\mathbf{z})\|^2$$

If  $\boldsymbol{\theta}_i(0)$  is a stationary point of  $f_i$  then the trajectory of  $z$  is a single point. But the trajectory of  $\boldsymbol{\theta}_i$  under the dynamics of Equation (1) is also a single point so we can pick the following function

$$X_{\boldsymbol{\theta}_i(0)}(f_i) = \boldsymbol{\theta}_i(0).$$

On the other hand if  $\boldsymbol{\theta}_i(0)$  is not a stationary point of  $f_i$ ,  $f_i$  continuously increases along the trajectory of ( $\Sigma^*$ ). Therefore  $A_{\boldsymbol{\theta}_i(0)}(t) = f_i(\gamma_{\boldsymbol{\theta}_i(0)}(t))$  is an increasing function and therefore invertible. Let us call  $A_{\boldsymbol{\theta}_i(0)}^{-1}(f_i)$  the inverse. Let's recall now the  $\boldsymbol{\theta}_i$  part of the dynamical system of interest Equation (1)

$$\dot{\boldsymbol{\theta}}_i = -\nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i) \frac{\partial L}{\partial f_i}(\mathbf{F}(\boldsymbol{\theta}), \mathbf{G}(\phi))$$

initialized at  $\boldsymbol{\theta}_i(0)$ . Applying Lemma 1 for the first equation with

$$h(t) = -\frac{\partial L}{\partial f_i}(\mathbf{F}(\boldsymbol{\theta}(t)), \mathbf{G}(\phi(t)))$$

we have that under the GDA dynamics:

$$\boldsymbol{\theta}_i(t) = \gamma_{\boldsymbol{\theta}_i(0)} \left( \int_0^t h(s) ds \right) \quad (P)$$

Thus it holds

$$f_i(\boldsymbol{\theta}_i(t)) = f_i \left( \gamma_{\boldsymbol{\theta}_i(0)} \left( \int_0^t h(s) ds \right) \right) = A_{\boldsymbol{\theta}_i(0)} \left( \int_0^t h(s) ds \right)$$

or equivalently

$$\int_0^t h(s) ds = A_{\boldsymbol{\theta}_i(0)}^{-1}(f_i(\boldsymbol{\theta}_i(t)))$$

Plugging in back to Equation (P)

$$\boldsymbol{\theta}_i(t) = \gamma_{\boldsymbol{\theta}_i(0)} \left( A_{\boldsymbol{\theta}_i(0)}^{-1} (f_i(\boldsymbol{\theta}_i(t))) \right)$$

Therefore we can pick

$$X_{\boldsymbol{\theta}_i(0)}(f_i) = \gamma_{\boldsymbol{\theta}_i(0)} \left( A_{\boldsymbol{\theta}_i(0)}^{-1} (f_i) \right)$$

which is  $C^1$  as composition of  $C^1$  functions. We can perform an equivalent analysis for  $\phi_j(0)$  and  $g_j$  to pick  $C^1$  function  $X_{\phi_j(0)}$ . Let us now track the time derivative of  $f_i(\boldsymbol{\theta}_i)$  and  $g_j(\phi_j)$

$$\begin{aligned} \dot{f}_i &= \nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i) \dot{\boldsymbol{\theta}}_i = \|\nabla_{\boldsymbol{\theta}_i} f_i(\boldsymbol{\theta}_i)\|^2 \frac{\partial L}{\partial f_i}(\mathbf{F}, \mathbf{G}) \\ \dot{g}_j &= \nabla_{\phi_j} g_j(\phi_j) \dot{\phi}_j = \|\nabla_{\phi_j} g_j(\phi_j)\|^2 \frac{\partial L}{\partial g_j}(\mathbf{F}, \mathbf{G}) \end{aligned}$$

We can now replace  $\boldsymbol{\theta}_i = X_{\boldsymbol{\theta}_i(0)}(f_i)$  and  $\phi_j = X_{\phi_j(0)}(g_j)$  to get the equations required.

### C. Proof for Lemma 9

*Proof.* Simple substitution gets us the following

$$\begin{aligned} \dot{H} &= - \sum_{i=1}^N (f_i - p_i) \frac{\partial L}{\partial f_i}(\mathbf{F}, \mathbf{G}) + \sum_{j=1}^M (g_j - q_j) \frac{\partial L}{\partial g_j}(\mathbf{F}, \mathbf{G}) \\ &= - \langle \mathbf{F} - \mathbf{p}, \nabla_{\mathbf{F}} L(\mathbf{F}, \mathbf{G}) \rangle + \langle \mathbf{G} - \mathbf{q}, \nabla_{\mathbf{G}} L(\mathbf{F}, \mathbf{G}) \rangle \end{aligned}$$

For convex  $L(\cdot, \mathbf{G})$  and concave  $L(\mathbf{F}, \cdot)$ , we have according to the definition:

$$\begin{aligned} - \langle \mathbf{F} - \mathbf{p}, \nabla_{\mathbf{F}} L(\mathbf{F}, \mathbf{G}) \rangle &\leq L(\mathbf{p}, \mathbf{G}) - L(\mathbf{F}, \mathbf{G}) \\ \langle \mathbf{G} - \mathbf{q}, \nabla_{\mathbf{G}} L(\mathbf{F}, \mathbf{G}) \rangle &\leq L(\mathbf{F}, \mathbf{G}) - L(\mathbf{F}, \mathbf{q}) \end{aligned}$$

### D. Proof for Theorem 10

*Proof.* Leveraging Lemma 2, there is a function  $H$  which is well defined in  $D = \{\text{Im}_{f_i}(\boldsymbol{\theta}_i(0))\}_{i=1}^N \times \{\text{Im}_{g_j}(\phi_j(0))\}_{j=1}^M$  and in this domain  $\dot{H} \leq 0$ . Given the safety conditions we know that  $(\mathbf{p}, \mathbf{q}) \in D$ . Observe that for the proposed function, it holds that  $H(\mathbf{p}, \mathbf{q}) = 0$ . Also for each  $f_i$  and  $g_j$  term in  $H$  we know that it has its minimum of value 0 at the corresponding  $p_i$  and  $q_j$ . We can deduce this by taking the derivative of each term to study its monotonicity. For example, the  $f_i$  terms are strictly increasing in  $f_i > p_i$  and strictly decreasing in  $f_i < p_i$ . Thus for all  $D - \{(\mathbf{p}, \mathbf{q})\}$ ,  $H > 0$ . Applying Theorem 8 for the continuously differentiable  $H$  we have that  $(\mathbf{p}, \mathbf{q})$  is stable for Equation (2).

### E. Proof of Theorem 11

*Proof.* Firstly, we recall the property of sigmoid's gradient:

$$\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x)).$$

Thus the transformed dynamical system in the operator space can be written as:

$$(T) := \left\{ \begin{array}{l} \dot{f}_i = -f_i^2(1-f_i)^2 \frac{\partial L}{\partial f_i}(\mathbf{F}, \mathbf{G}) \\ \dot{g}_j = g_j^2(1-g_j)^2 \frac{\partial L}{\partial g_j}(\mathbf{F}, \mathbf{G}) \end{array} \right\}$$

Notice that

1. The dynamical system  $(T)$  in the operator space is independent of the initial conditions. In fact, the dynamical system of  $(T)$  and the one of Equation (1), called  $(\Sigma)$  for short, are diffeomorphic for all initializations, not just a specific trajectory.

2. Since  $(\boldsymbol{\theta}(0), \boldsymbol{\phi}(0))$  is safe, using Theorem 2 we get that  $(\boldsymbol{p}, \boldsymbol{q})$  is stable for  $(T)$ . We would like to prove that for every open neighborhood  $V$  of  $(\boldsymbol{F}^{-1}(\boldsymbol{p}), \boldsymbol{G}^{-1}(\boldsymbol{q}))$  there exists an open neighborhood  $U$  of  $(\boldsymbol{F}^{-1}(\boldsymbol{p}), \boldsymbol{G}^{-1}(\boldsymbol{q}))$  such that

$$(\boldsymbol{\theta}_{\text{init}}, \boldsymbol{\phi}_{\text{init}}) \in U \implies \forall t \geq 0 : (\boldsymbol{\theta}(t), \boldsymbol{\phi}(t)) \in V.$$

Applying the diffeomorphism  $\gamma = \gamma_{\Sigma \rightarrow T}$  between GDA dynamics of  $(\Sigma)$  and  $(T)$ ,  $\gamma(V)$  is an open neighborhood of  $(\boldsymbol{p}, \boldsymbol{q})$  since  $V$  is open and  $\gamma((\boldsymbol{F}^{-1}(\boldsymbol{p}), \boldsymbol{G}^{-1}(\boldsymbol{q}))) \equiv (\boldsymbol{p}, \boldsymbol{q}) \in \gamma(V)$ . By Item 2, since  $(\boldsymbol{p}, \boldsymbol{q})$  is stable for  $(T)$  there is an open neighborhood  $\tilde{U}$  of  $(\boldsymbol{p}, \boldsymbol{q})$  such that:

$$(\boldsymbol{F}_{\text{init}}, \boldsymbol{G}_{\text{init}}) \in \tilde{U} \implies \forall t \geq 0 : (\boldsymbol{F}(t), \boldsymbol{G}(t)) \in \gamma(V)$$

or equivalently

$$\gamma(\boldsymbol{\theta}_{\text{init}}, \boldsymbol{\phi}_{\text{init}}) \in \tilde{U} \implies \forall t \geq 0 : \gamma(\boldsymbol{\theta}(t), \boldsymbol{\phi}(t)) \in \gamma(V)$$

Indeed, using the inverse diffeomorphism  $\gamma^{-1}$ , we can establish that for  $U = \gamma^{-1}(\tilde{U})$  it holds that

$$(\boldsymbol{\theta}_{\text{init}}, \boldsymbol{\phi}_{\text{init}}) \in U \implies \forall t \geq 0 : (\boldsymbol{\theta}(t), \boldsymbol{\phi}(t)) \in V$$

## F. Proof of Theorem 12

*Proof.* Let us define the following sets

$$\begin{aligned} \forall i \in [n] : A_i &= \{ \boldsymbol{\theta}_i \in \mathbb{R}^{n_i} \\ \forall j \in [m] : B_j &= \left\{ \begin{array}{l} f_i(\boldsymbol{\theta}_i) \in [p_i - \xi, p_i + \xi] \\ \boldsymbol{\phi}_j \in \mathbb{R}^{m_j} \end{array} \right. \quad g_j(\boldsymbol{\phi}_j) \in [q_j - \xi, q_j + \xi] \end{aligned}$$

Since  $f_i$  and  $g_j$  are proper  $A_i$  and  $B_j$  are compact sets. Thus, the continuous functions  $\|\nabla f_i(\boldsymbol{\theta}_i)\|^2$  and  $\|\nabla g_j(\boldsymbol{\phi}_j)\|^2$  have a minimum and maximum value on  $A_i$  and  $B_j$  respectively. Let us call  $K_{f_i}$  and  $K_{g_j}$  the maxima and  $\kappa_{f_i}$  and  $\kappa_{g_j}$  the minima. Observe that the minima and maxima must be all greater than zero since  $[p_i - \xi, p_i + \xi]$  and  $[q_j - \xi, q_j + \xi]$  are regular values. Let us define

$$\begin{aligned} \kappa &= \min \left\{ \min_{1 \leq i \leq n} \kappa_{f_i}, \min_{1 \leq j \leq m} \kappa_{g_j} \right\} \\ K &= \max \left\{ \max_{1 \leq i \leq n} K_{f_i}, \max_{1 \leq j \leq m} K_{g_j} \right\} \end{aligned}$$

where  $K \geq \kappa > 0$  as we discussed. Let us create the following set

$$S = \{ (\boldsymbol{\theta}, \boldsymbol{\phi}) \in \mathbb{R}^N \times \mathbb{R}^M \mid \forall i \in [n] : \boldsymbol{\theta}_i \in A_i, \quad \forall j \in [m] : \boldsymbol{\phi}_j \in B_j \}$$

We can prove that every  $(\boldsymbol{\theta}, \boldsymbol{\phi}) \in S$  is a safe initialization for  $(\boldsymbol{p}, \boldsymbol{q})$ . Of course, every  $\boldsymbol{\theta}_i$  and  $\boldsymbol{\phi}_j$  are not stationary points of  $f_i$  and  $g_j$  respectively. We also need to prove that the equilibrium  $(\boldsymbol{p}, \boldsymbol{q})$  is feasible. We will prove this by contradiction. Let there be a  $(\boldsymbol{\theta}, \boldsymbol{\phi}) \in S$  such that  $(\boldsymbol{p}, \boldsymbol{q})$  is not feasible. Without loss of generality we can assume that there is an  $i \in [n]$  such that  $p_i \notin \text{Im}_{f_i}(\boldsymbol{\theta}_i)$ . The case for the  $g_j$  is symmetrical. Along the gradient ascent trajectory of  $f_i$  with initialization at  $\boldsymbol{\theta}_i$ , observe that  $f_i(t)$  cannot attain an infimum or a supremum in  $[p_i - \xi, p_i + \xi]$  because there are no stationary points of  $f_i$  in  $A_i$ . Observe also that at initialization

$f_i(\boldsymbol{\theta}_i) \in [p_i - \xi, p_i + \xi]$ . Thus Let us pick an initialization  $(\boldsymbol{\theta}(0), \phi(0))$  such that  $r(0) \leq \xi^2$ . It is clear that  $(\boldsymbol{\theta}(0), \phi(0)) \in S$  and so it is safe for  $(\mathbf{p}, \mathbf{q})$ . We can do the same steps as in Theorem 2 to prove that the function  $H(\mathbf{F}, \mathbf{G})$  below does not increase under the dynamics of Equation (1):

$$H(\mathbf{F}, \mathbf{G}) = \sum_{i=1}^N \int_{p_i}^{f_i} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz + \sum_{j=1}^M \int_{q_j}^{g_j} \frac{z - q_j}{\|\nabla g_j(X_{\boldsymbol{\phi}_j(0)}(z))\|^2} dz$$

Observe that since  $(\boldsymbol{\theta}(0), \phi(0)) \in S$  we have that the interval between  $p_i$  and  $f_i(\boldsymbol{\theta}_i(0))$  belongs in  $[p_i - \xi, p_i + \xi]$  and  $\|\nabla f_i(\cdot)\|^2 \geq \kappa$  in this interval. Thus we can write

$$\frac{(f_i(\boldsymbol{\theta}_i(0)) - p_i)^2}{2\kappa} \geq \int_{p_i}^{f_i(\boldsymbol{\theta}_i(0))} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz$$

Repeating the same argument for all  $f_i$  and  $g_j$  we have that

$$\frac{r(0)}{2\kappa} \geq H(\mathbf{F}(\boldsymbol{\theta}(0)), \mathbf{G}(\phi(0))) \geq H(\mathbf{F}(\boldsymbol{\theta}(t)), \mathbf{G}(\phi(t)))$$

Let us pick  $r(0) < \min\{\xi^2, \xi^2 \frac{\kappa}{K}\} = \xi^2 \frac{\kappa}{K}$ . We already know that trajectories start in  $S$ . We will prove that they also remain in  $S$ . We will do this by contradiction. If a trajectory escaped  $S$ , then without loss of generality this means that there is at least an  $i \in [n]$  such that at some  $t > 0$ ,  $f_i(\boldsymbol{\theta}_i(t)) \notin [p_i - \xi, p_i + \xi]$ . The case of  $g_j$  is similar. Clearly we have that

$$\int_{p_i}^{f_i(\boldsymbol{\theta}_i(t))} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz \geq \min \left\{ \int_{p_i}^{p_i - \xi} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz, \int_{p_i}^{p_i + \xi} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz \right\}$$

As above, we have that the gradients in the integrals of the right hand side are less or equal than  $K$  so

$$\int_{p_i}^{f_i(\boldsymbol{\theta}_i(t))} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz \geq \frac{\xi^2}{2K}.$$

The terms of  $H$  are all non-negative so we have that

$$\frac{r(0)}{2\kappa} \geq H(\mathbf{F}(\boldsymbol{\theta}(t)), \mathbf{G}(\phi(t))) \geq \int_{p_i}^{f_i(\boldsymbol{\theta}_i(t))} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz \geq \frac{\xi^2}{2K}.$$

But  $r(0) < \xi^2 \frac{\kappa}{K}$ , a contradiction. So the trajectories will stay in  $S$ . We can then write

$$\int_{p_i}^{f_i(\boldsymbol{\theta}_i(t))} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz \geq \frac{(f_i(\boldsymbol{\theta}_i(t)) - p_i)^2}{2K}.$$

Repeating the same argument for all  $f_i$  and  $g_j$  we have that

$$\frac{r(0)}{2\kappa} \geq H(\mathbf{F}(\boldsymbol{\theta}(t)), \mathbf{G}(\phi(t))) \geq \frac{r(t)}{2K}.$$

For every  $\epsilon > 0$ , there is a positive  $\delta = \frac{\min\{\xi^2, \epsilon\}\kappa}{K}$  such that

$$r(0) < \delta \implies r(t) < \epsilon.$$

## G. Proof of Lemma 13

*Proof.* Pick a point  $(\mathbf{p}, \mathbf{q}) \in Z$ . Since our initialization is safe for this saddle point, we can construct the  $H$  function as in Theorem 5 and prove that it has the following property

$$\dot{H} \leq 0 \text{ in } D = \{\text{Im}_{f_i}(\boldsymbol{\theta}_i(0))\}_{i=1}^N \times \{\text{Im}_{g_j}(\boldsymbol{\phi}_j(0))\}_{j=1}^M$$

If  $(\mathbf{F}(\boldsymbol{\theta}(0)), \mathbf{G}(\boldsymbol{\phi}(0))) = (\mathbf{p}, \mathbf{q})$  then the theorem holds trivially. Otherwise, take a ball  $B$  centered at the equilibrium with a small enough radius such that it is contained in the interior of  $D$ .

$$H_0 = \min_{(\mathbf{F}, \mathbf{G}) \in \partial B} H(\mathbf{F}, \mathbf{G})$$

$$\Omega = \{(\mathbf{F}, \mathbf{G}) \in B \mid H(\mathbf{F}, \mathbf{G}) \leq H_0/2\}$$

We know that in both of the cases  $H_0 > 0$  from Theorem 2 .

Since  $\dot{H} \leq 0$ , starting in  $\Omega$ , it implies that  $H(\mathbf{F}(t), \mathbf{G}(t)) \leq H_0$  for  $t \geq 0$ , so  $\Omega$  is forward invariant. Since  $\Omega \subset D$  we know that it is bounded.  $\Omega$  is closed since it is a sublevel set of a continuous function. Notice that the restriction of  $\Omega$  on  $B$  does not affect the above properties since  $\Omega$  is in the interior of  $B$ . Thus  $\Omega$  is a compact forward invariant set, satisfying the requirement of Theorem 9

Let  $E = \{(\mathbf{F}, \mathbf{G}) \in B \mid \dot{H}(\mathbf{F}, \mathbf{G}) = 0\}$ . Without loss of generality we can assume that  $L(\cdot, \mathbf{q})$  is strictly convex as the case of  $L(\mathbf{p}, \cdot)$  being strictly concave is similar. In the following inequality

$$\dot{H} \leq L(\mathbf{p}, \mathbf{G}) - L(\mathbf{p}, \mathbf{q}) + L(\mathbf{p}, \mathbf{q}) - L(\mathbf{F}, \mathbf{q}) \leq 0$$

we know that  $L(\mathbf{p}, \mathbf{G}) - L(\mathbf{p}, \mathbf{q}) \leq 0$  and  $L(\mathbf{p}, \mathbf{q}) - L(\mathbf{F}, \mathbf{q}) \leq 0$ . So  $\dot{H} = 0$  implies  $L(\mathbf{p}, \mathbf{G}) = L(\mathbf{p}, \mathbf{q}) = L(\mathbf{F}, \mathbf{q})$ . By the strict convexity of  $L(\cdot, \mathbf{q})$  we know that this means that  $\mathbf{F} = \mathbf{p}$ . Let  $\mathcal{M}$  be the largest invariant set inside  $E$ . By the properties of  $\mathcal{M}$  being invariant subset of  $E$  we have

$$(\mathbf{F}(0), \mathbf{G}(0)) \in \mathcal{M} \implies \forall t : \mathbf{F}(t) = \mathbf{p} \text{ and } L(\mathbf{p}, \mathbf{G}(t)) = L(\mathbf{p}, \mathbf{q})$$

Taking the time derivatives on each of the constant quantities, they should be zero.

$$\dot{f}_i = 0 \implies \forall i \in [N] : \left\| \nabla_{\boldsymbol{\theta}_i} f_i (X_{\boldsymbol{\theta}_i(0)}(p_i)) \right\|^2 \frac{\partial L}{\partial f_i}(\mathbf{p}, \mathbf{G}) = 0$$

$$L(\mathbf{p}, \dot{\mathbf{G}}(t)) = 0 \implies \sum_{j=1}^M \left\| \nabla_{\boldsymbol{\phi}_j} g_j (X_{\boldsymbol{\phi}_j(0)}(g_j)) \right\|^2 \left[ \frac{\partial L}{\partial g_j}(\mathbf{p}, \mathbf{G}) \right]^2 = 0$$

We know that  $\left\| \nabla_{\boldsymbol{\theta}_i} f_i (X_{\boldsymbol{\theta}_i(0)}(p_i)) \right\| \neq 0$  by the safety conditions and that  $\left\| \nabla_{\boldsymbol{\phi}_j} g_j (X_{\boldsymbol{\phi}_j(0)}(g_j)) \right\|^2 \neq 0$  inside  $D$  again by safety conditions. This implies

$$\forall i \in [N] : \frac{\partial L}{\partial f_i}(\mathbf{p}, \mathbf{G}) = 0$$

$$\forall j \in [M] : \frac{\partial L}{\partial g_j}(\mathbf{p}, \mathbf{G}) = 0$$

Thus  $\mathcal{M}$  contains only stationary points of  $L$  so  $\mathcal{M} \subseteq \text{Solution}(L)$ . In addition  $\mathcal{M} \subseteq D$  so only stationary points of  $L$  for which the initialization is safe are allowed so  $\mathcal{M} \subseteq Z$ . Applying Theorem 9 we have that for any initialization of Equation (3) inside  $\Omega$ , as  $t \rightarrow \infty$   $(\mathbf{F}(t), \mathbf{G}(t))$  approaches  $\mathcal{M}$  and thus  $Z$  is locally asymptotically stable for Equation (3).

## H. Proof of Theorem 14

*Proof.* Again let's pick a point  $(\mathbf{p}, \mathbf{q}) \in Z$ . Since our initialization is safe for this saddle point, we can construct the  $H$  function as in Theorem 5 and prove that it has the following property

$$\dot{H} \leq 0 \text{ in } D = \{\text{Im}_{f_i}(\boldsymbol{\theta}_i(0))\}_{i=1}^N \times \{\text{Im}_{g_j}(\boldsymbol{\phi}_j(0))\}_{j=1}^M$$

If  $(\mathbf{F}(\boldsymbol{\theta}(0)), \mathbf{G}(\boldsymbol{\phi}(0))) = (\mathbf{p}, \mathbf{q})$  then the theorem holds trivially. Otherwise define

$$\begin{aligned} H_0 &= H(\mathbf{F}(\boldsymbol{\theta}(0)), \mathbf{G}(\boldsymbol{\phi}(0))) \\ \Omega &= \{(\mathbf{F}, \mathbf{G}) \in D \mid H(\mathbf{F}, \mathbf{G}) \leq H_0\} \end{aligned}$$

where we know that  $H_0 > 0$  from Theorem 2. Let us assume that indeed  $\Omega$  is in the interior of  $D$ . Then, applying the same argumentation as in Lemma 3 combined with Theorem 2, all fixed points in  $Z$  are stable. So applying Theorem 10 we get that the trajectory initialized at  $(\mathbf{F}(\boldsymbol{\theta}(0)), \mathbf{G}(\boldsymbol{\phi}(0))) \in \Omega$  converges to a point in  $Z$ . It remains to prove our assertion about the set  $\Omega$ : *Claim 1.*  $\Omega$  is in the interior of  $D$ . *Proof.* We will argue that as  $(\mathbf{F}, \mathbf{G})$  approaches the boundary of  $D$ , the value of  $H$  should become unbounded. If this is true then for the finite upper bound of  $H_0$ ,  $\Omega$  should have no points close to the boundary of  $H$  and thus it should be in the interior.

As  $(\mathbf{F}, \mathbf{G})$  approach the boundary of  $D$ , at least one of the variables  $f_i$  or  $g_j$  approaches the endpoints points of  $\text{Im}_{f_i}(\boldsymbol{\theta}_i(0))$  or  $\text{Im}_{g_j}(\boldsymbol{\phi}_j(0))$  respectively. We will study the case of  $f_i$  since the case of  $g_j$  is symmetrical. The endpoint  $f_{is}$  can be either the supremum or the infimum of the gradient ascent trajectory on  $f_i$  or  $\pm\infty$  if they do not exist. Let  $f_{is}$  be the supremum or  $\infty$  depending on if the former exists. We can take the gradient ascent dynamics and apply Lemma 3 to get

$$\dot{f}_i = \|\nabla_{\boldsymbol{\theta}_i} f_i(X_{\boldsymbol{\theta}_i(0)}(f_i))\|^2$$

We know that  $f_i(\boldsymbol{\theta}_i(t))$  goes to  $f_{is}$  when initialized at  $f_i(\boldsymbol{\theta}_i(0))$ . Let us define the following function

$$a(f_i) = \int_{p_i}^{f_i} \frac{1}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz$$

Observe that  $\dot{a} = 1$ , thus  $\lim_{t \rightarrow \infty} a(f_i(t)) = \infty$ . In other words

$$\lim_{t \rightarrow \infty} \int_{p_i}^{f_i(t)} \frac{1}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz = \int_{p_i}^{f_{is}} \frac{1}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz = \infty$$

Symmetrically if  $f_{is}$  is the infimum or  $-\infty$ , then the limit above would be  $-\infty$ . In either case

$$f_i \rightarrow f_{is} \implies \int_{p_i}^{f_i} \frac{z - p_i}{\|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(z))\|^2} dz \rightarrow \infty$$

For the last step it is important to note that  $p_i$  is not at the boundary of  $D$  based on the safety conditions. Therefore as  $(\mathbf{F}, \mathbf{G})$  approach the boundary of  $D$  in the dynamics of Equation (3), at least one of the terms of  $H$  goes to infinity. Also note that all the terms of  $H$  are individually nonnegative so no matter what the other variables in  $(\mathbf{F}, \mathbf{G})$  are doing they cannot stop  $H \rightarrow \infty$ .

## I. Proof of Theorem 16

*Proof.* For any choice of  $\lambda > 0$  we have that  $L'$  is strictly convex strictly concave so the KKT conditions are sufficient to determine its equilibria.

$$\begin{aligned} \frac{\partial L(\mathbf{x}, \mathbf{y})}{\partial x_i} + \lambda x_i &= 0 \\ \frac{\partial L(\mathbf{x}, \mathbf{y})}{\partial y_j} - \lambda y_j &= 0 \end{aligned}$$

We can view the above set of constraints as a single vector constraint  $r(\lambda, \mathbf{x}, \mathbf{y}) = \mathbf{0}$ . Note that by assumption of the Hessians being invertible at all equilibria,  $L$  has a unique equilibrium  $(\mathbf{x}^*, \mathbf{y}^*)$ . Clearly we have that  $r(0, \mathbf{x}^*, \mathbf{y}^*) = \mathbf{0}$ . Observe that for the Jacobian of  $r$  at  $(0, \mathbf{x}^*, \mathbf{y}^*)$  with respect to  $(\mathbf{x}, \mathbf{y})$  we have that

$$D_{(\mathbf{x}, \mathbf{y})} r(0, \mathbf{x}^*, \mathbf{y}^*) = \nabla^2 L(\mathbf{x}^*, \mathbf{y}^*)$$

and thus it is invertible. Invoking the Implicit function Theorem, there is a differentiable function  $g$ , defined in a small enough neighborhood of 0, that takes a  $\lambda$  and returns  $g(\lambda) = (\mathbf{x}(\lambda), \mathbf{y}(\lambda))$  such that  $r(\lambda, g(\lambda)) = \mathbf{0}$ . Thus for a small enough  $\lambda$ , we have that  $g$  returns the corresponding equilibria of  $L'$ . By continuity of  $g$ , for all  $\epsilon$  there is a  $\delta > 0$

$$\forall 0 < \lambda < \delta : \|\mathbf{x}(\lambda) - \mathbf{x}(0)\|^2 + \|\mathbf{y}(\lambda) - \mathbf{y}(0)\|^2 \leq \epsilon^2$$

But  $(\mathbf{x}(0), \mathbf{y}(0)) = (\mathbf{x}^*, \mathbf{y}^*)$  so the equilibrium of  $L'$  has an  $\epsilon$ -close equilibrium of  $L$  for  $\lambda < \delta$ . By strict convexity strict concavity of  $L'$ , it has a unique equilibrium as well. So the equilibria of  $L'$  and  $L$  are  $\epsilon$ -close to each other.

## J. Proof of Theorem 17

*Proof.* Following the same analysis with the strict convex concave analysis of the previous section, if  $(\mathbf{F}(\boldsymbol{\theta}(0)), \mathbf{G}(\boldsymbol{\phi}(0))) = (\mathbf{p}, \mathbf{q})$  then the theorem holds trivially. Otherwise, since our initialization is safe for  $(\mathbf{p}, \mathbf{q})$ , we can construct the  $H$  function as in Theorem 2 and prove that it has the following property in  $D = \{\text{Im}_{f_i}(\boldsymbol{\theta}_i(0))\}_{i=1}^N \times \{\text{Im}_{g_j}(\boldsymbol{\phi}_j(0))\}_{j=1}^M$

$$\begin{aligned} \dot{H} &\leq L'(\mathbf{p}, \mathbf{G}) - L'(\mathbf{p}, \mathbf{q}) + L'(\mathbf{p}, \mathbf{q}) - L'(\mathbf{F}, \mathbf{q}) \\ &\leq -\frac{\lambda}{2} (\|\mathbf{F}(\boldsymbol{\theta}(t)) - \mathbf{p}\|^2 + \|\mathbf{G}(\boldsymbol{\phi}(t)) - \mathbf{q}\|^2) \\ &\leq -\frac{\lambda}{2} r(t) \end{aligned}$$

Where the second step follows from  $L'(\mathbf{p}, \cdot)$  being  $\lambda$  strongly concave and  $L'(\cdot, \mathbf{q})$  being  $\lambda$  strongly convex and  $\mathbf{q}, \mathbf{p}$  being the corresponding optima of these functions since  $(\mathbf{p}, \mathbf{q})$  is an equilibrium. Let us define

$$\begin{aligned} H_0 &= H(\mathbf{F}(\boldsymbol{\theta}(0)), \mathbf{G}(\boldsymbol{\phi}(0))) \\ \Omega &= \{(\mathbf{F}, \mathbf{G}) \in D \mid H(\mathbf{F}, \mathbf{G}) \leq H_0\} \end{aligned}$$

where we know that  $H_0 > 0$  from Theorem 2. Additionally, we can apply Claim 1 even in the new dynamics, so  $\Omega$  is in the interior of  $D$ . Since  $\dot{H} \leq 0$ , starting in  $\Omega$ , it implies that  $H(\mathbf{F}(\boldsymbol{\theta}(t)), \mathbf{G}(\boldsymbol{\phi}(t))) \leq H_0$  for  $t \geq 0$ , so  $(\mathbf{F}(t), \mathbf{G}(t))$  stays in  $\Omega$ . Additionally,  $\Omega$  is closed since it is a sublevel set of a continuous function. Notice that the restriction of  $\Omega$  on  $D$  does not affect the above properties since  $\Omega$  is in the interior of  $D$ . Thus  $\Omega$  is a compact forward invariant set. For a safe initialization  $(\boldsymbol{\theta}(0), \boldsymbol{\phi}(0))$ , the following continuous functions must have a minimum and maximum value on  $\Omega$  respectively.

$$\begin{aligned} K_{f_i} &\geq \|\nabla f_i(X_{\boldsymbol{\theta}_i(0)}(\cdot))\|^2 \geq \kappa_{f_i} \\ K_{g_j} &\geq \|\nabla g_j(X_{\boldsymbol{\phi}_j(0)}(\cdot))\|^2 \geq \kappa_{g_j} \end{aligned}$$

Observe that the minima and maxima must be all greater than zero, since both  $\|\nabla_{\boldsymbol{\phi}_j} g_j(X_{\boldsymbol{\phi}_j(0)}(g(t)))\|, \|\nabla_{\boldsymbol{\theta}_i} f_i(X_{\boldsymbol{\theta}_i(0)}(f(t)))\|$  go to 0 as this happens only at the boundaries of  $D$  which are outside  $\Omega$ . Let us define

$$\begin{aligned} \kappa &= \min \left\{ \min_{1 \leq i \leq n} \kappa_{f_i}, \min_{1 \leq j \leq m} \kappa_{g_j} \right\} \\ K &= \max \left\{ \max_{1 \leq i \leq n} K_{f_i}, \max_{1 \leq j \leq m} K_{g_j} \right\} \end{aligned}$$



Observe that  $K \geq \|\nabla f_i (X_{\theta_i(0)}(\cdot))\|^2 \geq \kappa$  in this interval. Thus we can write

$$\frac{(f_i(\theta_i(t)) - p_i)^2}{2\kappa} \geq \int_{p_i}^{f_i(\theta_i(t))} \frac{z - p_i}{\|\nabla f_i (X_{\theta_i(0)}(z))\|^2} dz \geq \frac{(f_i(\theta_i(t)) - p_i)^2}{2K}$$

Repeating the same argument for all  $f_i$  and  $g_j$  we have that

$$\frac{r(t)}{2\kappa} \geq H(\mathbf{F}(\theta(t)), \mathbf{G}(\phi(t))) \geq \frac{r(t)}{2K}$$

Thus we can extend our analysis

$$\dot{H} \leq -\lambda r(t) \leq -\frac{2\kappa\lambda}{2} H(t) \Rightarrow H(t) \leq H_0 e^{-\lambda\kappa t} \Rightarrow r(t) \leq 2 \times K \times H_0 e^{-\lambda\kappa t}$$