

---

# Global Optimality in Bivariate Gradient-based DAG Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Recently, a new class of non-convex optimization problems motivated by the  
2 statistical problem of learning an acyclic directed graphical model from data  
3 has attracted significant interest. While existing work uses standard first-order  
4 optimization schemes to solve this problem, proving the global optimality of such  
5 approaches has proven elusive. The difficulty lies in the fact that unlike other  
6 non-convex problems in the literature, this problem is not “benign”, and possesses  
7 multiple spurious solutions that standard approaches can easily get trapped in. In  
8 this paper, we prove that a simple path-following optimization scheme globally  
9 converges to the global minimum of the population loss in the bivariate setting.

## 10 1 Introduction

11 Over the past decade, non-convex optimization has become a major topic of research within the  
12 machine learning community, in part due to the successes of training large-scale models with simple  
13 first-order methods such as gradient descent—along with their stochastic and accelerated variants—  
14 in spite of the non-convexity of the loss function. A large part of this research has focused on  
15 characterizing which problems have *benign* loss landscapes that are amenable to the use of gradient-  
16 based methods, i.e., there are no spurious local minima, or they can be easily avoided. By now,  
17 several theoretical results have shown this property for different non-convex problems such as:  
18 learning a two hidden unit ReLU network [48], learning (deep) over-parameterized quadratic neural  
19 networks [43, 27], low-rank matrix recovery [19, 13, 3], learning a two-layer ReLU network with  
20 a single non-overlapping convolutional filter [6], semidefinite matrix completion [4, 20], learning  
21 neural networks for binary classification with the addition of a single special neuron [30], and learning  
22 deep networks with independent ReLU activations [26, 11], to name a few.

23 Recently, a new class of non-convex optimization problems due to Zheng et al. [51] have emerged in  
24 the context of learning the underlying structure of a structural equation model (SEM) or Bayesian  
25 network. This underlying structure is typically represented by a directed acyclic graph (DAG), which  
26 makes the learning task highly complex due to its combinatorial nature. In general, learning DAGs is  
27 well-known to be NP-complete [8, 10]. The key innovation in Zheng et al. [51] was the introduction  
28 of a differentiable function  $h$ , whose level set at zero *exactly* characterizes DAGs. Thus, replacing the  
29 challenges of combinatorial optimization by those of non-convex optimization. Mathematically, this  
30 class of non-convex problems take the following general form:

$$\min_{\Theta} f(\Theta) \text{ subject to } h(W(\Theta)) = 0, \quad (1)$$

31 where  $\Theta \in \mathbb{R}^l$  represents the model parameters,  $f : \mathbb{R}^l \rightarrow \mathbb{R}$  is a (possibly non-convex) smooth loss  
32 function (sometimes called a *score function*) that measures the fitness of  $\Theta$ , and  $h : \mathbb{R}^{d \times d} \rightarrow [0, \infty)$   
33 is a smooth **non-convex** function that takes the value of zero if and only if the induced weighted  
34 adjacency matrix of  $d$  nodes,  $W(\Theta)$ , corresponds to a DAG.

Given the smoothness of  $f$  and  $h$ , problem (1) can be solved using off-the-shelf nonlinear solvers, which has driven a series of remarkable developments in structure learning for DAGs. Multiple empirical studies have demonstrated that global or near-global minimizers for (1) can often be found in a variety of settings, such as linear models with Gaussian and non-Gaussian noises [e.g., 51, 34, 1], and non-linear models, represented by neural networks, with additive Gaussian noises [e.g., 29, 52, 49, 1]. The empirical success for learning DAGs via (1), which started with the NOTEARS method of Zheng et al. [51], bears a resemblance to the success of training deep models, which started with AlexNet for image classification.

Importantly, the reader should note that the majority of applications in ML consist of solving a *single unconstrained* non-convex problem. In contrast, the class of problems (1) contains a non-convex constraint. Thus, researchers have considered some type of penalty method such as the augmented Lagrangian [51, 52], quadratic penalty [35], and a log-barrier [1]. In all cases, the penalty approach consists in solving a *sequence* of unconstrained non-convex problems, where the constraint is enforced progressively [see e.g. 2, for background]. In this work, we will consider the following form of penalty:

$$\min_{\Theta} g_{\mu_k}(\Theta) := \mu_k f(\Theta) + h(W(\Theta)). \quad (2)$$

It was shown by Bello et al. [1] that due to the invexity property of  $h$ ,<sup>1</sup> solutions to (2) will converge to a DAG as  $\mu_k \rightarrow 0$ . However, no guarantees on local/global optimality were given.

With the above considerations in hand, one is inevitably led to ask the following questions:

- (i) *Are the loss landscapes  $g_{\mu_k}(\Theta)$  benign for different  $\mu_k$ ?*
- (ii) *Is there a (tractable) solution path  $\{\Theta_k\}$  that converges to a global minimum of (1)?*

Due to the NP-completeness of learning DAGs, one would expect the answer to (i) to be negative in its most general form. Moreover, it is known from the classical theory of constrained optimization [e.g. 2] that if we can *exactly* and *globally* optimize (1) for each  $\mu_k$ , then the answer to (ii) is affirmative. This is not a practical algorithm, however, since the problem (1) is nonconvex. Thus we seek a solution path that can be tractably computed in practice, e.g. by gradient descent.

In this work, we focus on perhaps the simplest setting where interesting phenomena take place. That is, a linear SEM with two nodes (i.e.,  $d = 2$ ),  $f$  is the population least squared loss (i.e.,  $f$  is convex), and  $\Theta_k$  is defined via gradient flow with warm starts. More specifically, we consider the case where  $\Theta_k$  is obtained by following the gradient flow of  $g_{\mu_k}$  with initial condition  $\Theta_{k-1}$ .

Under this setting, to answer (i), it is easy to see that for a large enough  $\mu_k$ , the convex function  $f$  dominates and we can expect a benign landscape, i.e., a (almost) convex landscape. Similarly, when  $\mu_k$  approaches zero, the invexity of  $h$  kicks in and we could expect that all stationary points are (near) global minimizers.<sup>2</sup> That is, at the extremes  $\mu_k \rightarrow \infty$  and  $\mu_k \rightarrow 0$ , the landscapes seem well-behaved, and the reader might wonder if it follows that for any  $\mu_k \in [0, \infty)$  the landscape is well-behaved. We answer the latter in the *negative* and show that there always exists a  $\tau > 0$  where the landscape of  $g_{\mu_k}$  is non-benign for any  $\mu_k < \tau$ , namely, there exist three stationary points: i) A saddle point, ii) A spurious local minimum, and iii) The global minimum. In addition, each of these stationary points have wide basins of attractions, thus making the initialization of the gradient flow for  $g_{\mu_k}$  crucial. Finally, we answer (ii) in the affirmative and provide an explicit scheduling for  $\mu_k$  that guarantees the asymptotic convergence of  $\Theta_k$  to the global minimum of (1). Moreover, we show that this scheduling cannot be arbitrary as there exists a sequence of  $\{\mu_k\}$  that leads  $\{\Theta_k\}$  to a spurious local minimum.

Overall, we establish the first set of results that study the optimization landscape and global optimality for the class of problems (1). We believe that this comprehensive analysis in the bivariate case provides a valuable starting point for future research in more complex settings.

**Remark 1.** *We emphasize that solving (1) in the bivariate case is not an inherently difficult problem. Indeed, when there are only two nodes, there are only two DAGs to distinguish and one can simply*

<sup>1</sup>An invex function is any function where all its stationary points are global minima. It is worth noting that the composite objective in (2) is not necessarily invex, even when  $f$  is convex.

<sup>2</sup>This transition or path, from an optimizer of a simple function to an optimizer of a function that closely resembles the original constrained formulation, is also known as a *homotopy*.

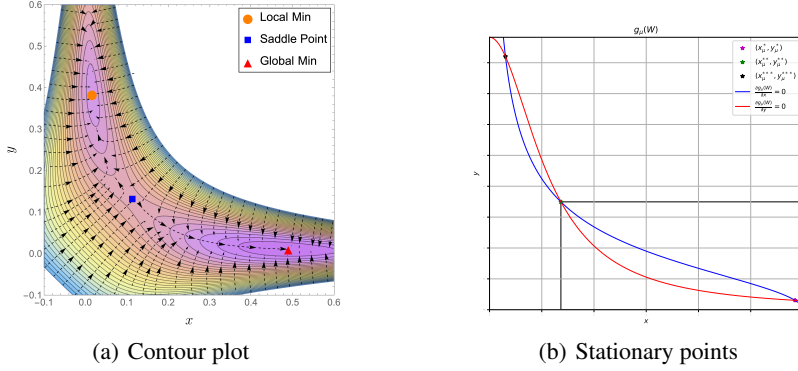


Figure 1: Visualizing the nonconvex landscape. (a) A contour plot of  $g_\mu$  for  $a = 0.5$  and  $\mu = 0.005$  (see Section 2 for definitions). We only show a section of the landscape for better visualization. The solid lines represent the contours, while the dashed lines represent the vector field  $-\nabla g_\mu$ . (b) Stationary points of  $g_\mu$ ,  $r(y; \mu) = 0$  and  $r(x; \mu) = 0$  (see Section 4 for definitions).

fit  $f$  under the only two possible DAGs, and select the model with the lowest value for  $f$ . However, evaluating  $f$  for each possible DAG structure clearly cannot scale beyond 10 or 20 nodes, and is not a standard algorithm for solving (1). Instead, here our focus is on studying how (1) is actually being solved in practice, namely, by solving unconstrained non-convex problems in the form of (2). Previous work suggests that such gradient-based approaches indeed scale well to hundreds and even thousands of nodes [e.g. 51, 34, 1].

## 1.1 Our Contributions

More specifically, we make the following contributions:

1. We present a homotopy-based optimization scheme (Algorithm 2) to find global minimizers of the program (1) by iteratively decreasing the penalty coefficient according to a given schedule. Gradient flow is used to find the stationary points of (2) at each step, starting from the previous solution.
2. We prove that Algorithm 2 converges *globally* (i.e. regardless of initialization for  $W$ ) to the *global* minimum (Theorem 1).
3. We show that the non-convex program (1) is indeed non-benign, and naïve implementation of black-box solvers are likely to get trapped in a bad local minimum. See Figure 1 (a).
4. Experimental results verify our theory, consistently recovering the global minimum of (1), regardless of initialization or initial penalty value. We show that our algorithm converges to the global minimum while naïve approaches can get stuck.

The analysis consists of three main parts: First, we explicitly characterize the trajectory of the stationary points of (2). Second, we classify the number and type of all stationary points (Lemma 1) and use this to isolate the desired global minimum. Finally, we apply Lyapunov analysis to identify the basin of attraction for each stationary point, which suggests a schedule for the penalty coefficient that ensures that the gradient flow is initialized within that basin at the previous solution.

## 1.2 Related Work

The class of problems (1) falls under the umbrella of score-based methods, where given a score function  $f$ , the goal is to identify the DAG structure with the lowest score possible [9, 22]. We shall note that learning DAGs is a very popular structure model in a wide range of domains such as biology [40], genetics [50], and causal inference [44, 39], to name a few.

**Score-based methods that consider the combinatorial constraint.** Given the ample set of score-based methods in the literature, we briefly mention some classical works that attempt to optimize  $f$

by considering the combinatorial DAG constraint. In particular, we have approximate algorithms such as the greedy search method of Chickering et al. [10], order search methods [45, 41, 38], the LP-relaxation method of Jaakkola et al. [24], and the dynamic programming approach of Loh and Bühlmann [31]. There are also exact methods such as GOBNILP [12] and Bene [42], however, these algorithms only scale up to  $\approx 30$  nodes.

**Score-based methods that consider the continuous non-convex constraint  $h$ .** The following works are the closest to ours since they attempt to solve a problem in the form of (1). Most of these developments either consider optimizing different score functions  $f$  such as ordinary least squares [51, 52], the log-likelihood [29, 34], the evidence lower bound [49], a regret function [53]; or consider different differentiable characterizations of acyclicity  $h$  [49, 1]. However, none of the aforementioned works provide any type of optimality guarantee. Few studies have examined the optimization intricacies of problem (1). Wei et al. [47] investigated the optimality issues and provided *local* optimality guarantees under the assumption of convexity in the score  $f$  and linear models. On the other hand, Ng et al. [35] analyzed the convergence to (local) DAGs of generic methods for solving nonlinear constrained problems, such as the augmented Lagrangian and quadratic penalty methods. In contrast to both, our work is the first to study global optimality and the loss landscapes of actual methods used in practice for solving (1).

**Bivariate causal discovery.** Even though in a two-node model the discrete DAG constraint does not pose a major challenge, the bivariate setting has been subject to major research in the area of causal discovery. See for instance [36, 16, 32, 25] and references therein.

**Penalty and homotopy methods.** There exist classical global optimality guarantees for the penalty method if  $f$  and  $h$  were convex functions, see for instance [2, 5, 37]. However, to our knowledge, there are no global optimality guarantees for general classes of non-convex constrained problems, let alone for the specific type of non-convex functions  $h$  considered in this work. On the other hand, homotopy methods (also referred to as continuation or embedding methods) are in many cases capable of finding better solutions than standard first-order methods for non-convex problems, albeit they typically do not come with global optimality guarantees either. When homotopy methods come with global optimality guarantees, they are commonly computationally more intensive as it involves discarding solutions, thus, closely resembling simulated annealing methods, see for instance [15]. Authors in [21] characterize a family of non-convex functions where a homotopy algorithm provably converges to a global optimum. However, the conditions for such family of non-convex functions are difficult to verify and are very restrictive; moreover, their homotopy algorithm involves Gaussian smoothing, making it also computationally more intensive than the procedure we study here. Other examples of homotopy methods in machine learning include [7, 18, 46, 17, 23], in all these cases, no global optimality guarantees are given.

## 2 Preliminaries

The objective  $f$  we consider can be easily written down as follows:

$$f(W) = \frac{1}{2} \mathbb{E}_X [\|X - W^\top X\|_2^2], \quad (3)$$

where  $X \in \mathbb{R}^2$  is a random vector and  $W \in \mathbb{R}^{2 \times 2}$ . Although not strictly necessary for the developments that follow, we begin by introducing the necessary background on linear SEM that leads to this objective and the resulting optimization problem of interest.

**The bivariate model.** Let  $X = (X_1, X_2) \in \mathbb{R}^2$  denote the random variables in the model, and let  $N = (N_1, N_2) \in \mathbb{R}^2$  denote a vector of independent errors. Then a linear SEM over  $X$  is defined as  $X = W_*^\top X + N$ , where  $W_* \in \mathbb{R}^{2 \times 2}$  is a weighted adjacency matrix encoding the coefficients in the linear model. In order to represent a valid Bayesian network for  $X$  [see e.g. 39, 44, for details], the matrix  $W_*$  must be acyclic: More formally, the weighted graph induced by the adjacency matrix  $W_*$  must be a DAG. This (non-convex) acyclicity constraint represents the major computational hurdle that must overcome in practice (cf. Remark 1).

The goal is to recover the matrix  $W_*$  from the random vector  $X$ . Since  $W_*$  is acyclic, we can assume the diagonal of  $W_*$  is zero (i.e. no self-loops). Thus, under the bivariate linear model, it then suffices

162 to consider two parameters  $x$  and  $y$  that define the matrix of parameters<sup>3</sup>

$$W = W(x, y) = \begin{pmatrix} 0 & x \\ y & 0 \end{pmatrix} \quad (4)$$

163 For notational simplicity, we will use  $f(W)$  and  $f(x, y)$  interchangeably, similarly for  $h(W)$  and  
164  $h(x, y)$ . Without loss of generality, we write the underlying parameter as

$$W_* = \begin{pmatrix} 0 & a \\ 0 & 0 \end{pmatrix} \quad (5)$$

165 which implies

$$X = W_*^\top X + N \implies \begin{cases} X_1 = N_1, \\ X_2 = aX_1 + N_2. \end{cases}$$

166 In general, we only require  $N_i$  to have finite mean and variance, hence we *do not* assume Gaussianity.  
167 We assume that  $\text{Var}[N_1] = \text{Var}[N_2]$ , and for simplicity, we consider  $\mathbb{E}[N] = 0$  and  $\text{Cov}[N] = I$ ,  
168 where  $I$  denotes the identity matrix. Finally, in the sequel we assume w.l.o.g. that  $a > 0$ .

169 **The population least squares.** In this work, we consider the population squared loss defined by (3).  
170 If we equivalently write  $f$  in terms of  $x$  and  $y$ , then we have:  $f(W) = ((1-ay)^2 + y^2 + (a-x)^2 + 1)/2$ .  
171 In fact, the population loss can be substituted with empirical loss. In such a case, our algorithm  
172 can still attain the global minimum,  $W_G$ , of problem (6). However, the output  $W_G$  will serve as an  
173 empirical estimation of  $W_*$ . An in-depth discussion on this topic can be found in Appendix B

174 **The non-convex function  $h$ .** We use the continuous acyclicity characterization of Yu et al. [49], i.e.,  
175  $h(W) = \text{Tr}((I + \frac{1}{d}W \circ W)^d) - d$ , where  $\circ$  denotes the Hadamard product. Then, for the bivariate  
176 case, we have  $h(W) = x^2y^2/2$ . We note that the analysis presented in this work is not tailored to  
177 this version of  $h$ , that is, we can use the same techniques used throughout this work for other existing  
178 formulations of  $h$ , such as the trace of the matrix exponential [51], and the log-det formulation [1].  
179 Nonetheless, here we consider that the polynomial formulation of Yu et al. [49] is more amenable for  
180 the analysis.

181 **Remark 2.** *Our restriction to the bivariate case highlights the simplest setting in which this problem*  
182 *exhibits nontrivial behaviour. Extending our analysis to higher dimensions remains a challenging*  
183 *future direction, however, we emphasize that even in two-dimensions this problem is nontrivial. Our*  
184 *approach is similar to that taken in other parts of the literature that started with simple cases (e.g.*  
185 *single-neuron models in deep learning).*

186 **Remark 3.** *It is worth noting that our choice of the population least squares is not arbitrary. Indeed,*  
187 *for linear models with identity error covariance, such as the model considered in this work, it is*  
188 *known that the global minimizer of the population squared loss is unique and corresponds to the*  
189 *underlying matrix  $W_*$ . See Theorem 7 in [31].*

190 Gluing all the pieces together, we arrive to the following version of (1) for the bivariate case:

$$\min_{x,y} f(x, y) := \frac{1}{2}((1-ay)^2 + y^2 + (a-x)^2 + 1) \quad \text{subject to} \quad h(x, y) := \frac{x^2y^2}{2} = 0. \quad (6)$$

191 Moreover, for any  $\mu \geq 0$ , we have the corresponding version of (2) expressed as:

$$\min_{x,y} g_\mu(x, y) := \mu f(x, y) + h(x, y) = \frac{\mu}{2}((1-ay)^2 + y^2 + (a-x)^2 + 1) + \frac{x^2y^2}{2}. \quad (7)$$

192 To conclude this section, we present a visualization of the landscape of  $g_\mu(x, y)$  in Figure 1 (a), for  
193  $a = 0.5$  and  $\mu = 0.005$ . We can clearly observe the non-benign landscape of  $g_\mu$ , i.e., there exists  
194 a spurious local minimum, a saddle point, and the global minimum. In particular, we can see that  
195 the basin of attraction of the spurious local minimum is comparable to that of the global minimum,  
196 which is problematic for a local algorithm such as the gradient flow (or gradient descent) as it can  
197 easily get trapped in a local minimum if initialized in the wrong basin.

---

<sup>3</sup>Following the notation in (1), for the bivariate model we simply have  $\Theta \equiv (x, y)$  and  $W(\Theta) \equiv \begin{pmatrix} 0 & x \\ y & 0 \end{pmatrix}$ .

---

**Algorithm 1:** GradientFlow( $f, z_0$ )

---

1: set  $z(0) = z_0$   
2:  $\frac{d}{dt}z(t) = -\nabla f(z(t))$   
3: **return**  $\lim_{t \rightarrow \infty} z(t)$

---

---

**Algorithm 2:** Homotopy algorithm for solving (1).

---

**Input:** Initial  $W_0 = W(x_0, y_0)$ ,  $\mu_0 \in \left[ \frac{a^2}{4(a^2+1)^3}, \frac{a^2}{4} \right)$

**Output:**  $\{W_{\mu_k}\}_{k=0}^{\infty}$

1  $W_{\mu_0} \leftarrow \text{GradientFlow}(g_{\mu_0}, W_0)$   
2 **for**  $k = 1, 2, \dots$  **do**  
3     Let  $\mu_k = (2/a)^{2/3} \mu_{k-1}^{4/3}$   
4      $W_{\mu_k} \leftarrow \text{GradientFlow}(g_{\mu_k}, W_{\mu_{k-1}})$   
5 **end**

---

### 3 A Homotopy-Based Approach and Its Convergence to the Global Optimum

To fix notation, let us write  $W_k := W_{\mu_k} := \begin{pmatrix} 0 & x^{\mu_k} \\ y^{\mu_k} & 0 \end{pmatrix}$ , and let  $W_G$  denote the global minimizer of (6). In this section, we present our main result, which provides conditions under which solving a series of unconstrained problems (7) with first-order methods will converge to the global optimum  $W_G$  of (6), in spite of facing non-benign landscapes. Recall that from Remark 3, we have that  $W_G = \begin{pmatrix} 0 & a \\ 0 & 0 \end{pmatrix}$ . Since we use gradient flow path to connect  $W_{\mu_k}$  and  $W_{\mu_{k+1}}$ , we specify this path in Procedure 1 for clarity. Although the theory here assumes continuous-time gradient flow with  $t \rightarrow \infty$ , see Section 5 for an iteration complexity analysis for (discrete-time) gradient descent, which is a straightforward consequence of the continuous-time theory.

In Algorithm 2, we provide an explicit regime of initialization for the homotopy parameter  $\mu_0$  and a specific scheduling for  $\mu_k$  such that the solution path found by Algorithm 2 will converge to the global optimum of (6). This is formally stated in Theorem 1, whose proof is given in Section 5.

**Theorem 1.** *For any initialization  $W_0$  and  $a \in \mathbb{R}$ , the solution path provided in Algorithm 2 converges to the global optimum of (6), i.e.,*

$$\lim_{k \rightarrow \infty} W_{\mu_k} = W_G.$$

A few observations regarding Algorithm 2: Observe that when the underlying model parameter  $a \gg 0$ , the regime of initialization for  $\mu_0$  is wider; on the other hand, if  $a$  is closer to zero then the interval for  $\mu_0$  is much narrower. As a concrete example, if  $a = 2$  then it suffices to have  $\mu_0 \in [0.008, 1)$ ; whereas if  $a = 0.1$  then the regime is about  $\mu_0 \in [0.0089, 0.01)$ . This matches the intuition that for a “stronger” value of  $a$  it should be easier to detect the right direction of the underlying model. Second, although in Line 3 we set  $\mu_k$  in a specific manner, it actually suffices to have

$$\mu_k \in \left[ \left( \frac{\mu_{k-1}}{2} \right)^{2/3} (a^{1/3} - \sqrt{a^{2/3} - (4\mu_{k-1})^{1/3}})^2, \mu_{k-1} \right).$$

We simply chose a particular expression from this interval for clarity of presentation; see the proof in Section 5 for details.

As presented, Algorithm 2 is of theoretical nature in the sense that the initialization for  $\mu_0$  and the decay rate for  $\mu_k$  in Line 3 depend on the underlying parameter  $a$ , which in practice is unknown. In Algorithm 3, we present a modification that is independent of  $a$  and  $W_*$ . By assuming instead a lower bound on  $a$ , which is a standard assumption in the literature, we can prove that Algorithm 3 also converges to the global minimum:

**Corollary 1.** *Initialize  $\mu_0 = \frac{1}{27}$ . If  $a > \sqrt{5/27}$  then for any initialization  $W_0$ , Algorithm 3 outputs the global optimal solution to (6), i.e.*

$$\lim_{k \rightarrow \infty} W_{\mu_k} = W_G.$$

For more details on this modification, see Appendix A.

---

**Algorithm 3:** Practical (i.e. independent of  $a$  and  $W_*$ ) homotopy algorithm for solving (1).

---

**Input:** Initial  $W_0 = W(x_0, y_0)$

**Output:**  $\{W_{\mu_k}\}_{k=0}^{\infty}$

```

1  $\mu_0 \leftarrow 1/27$ 
2  $W_{\mu_0} \leftarrow \text{GradientFlow}(g_{\mu_0}, W_0)$ 
3 for  $k = 1, 2, \dots$  do
4   Let  $\mu_k = (2/\sqrt{5\mu_0})^{2/3} \mu_{k-1}^{4/3}$ 
5    $W_{\mu_k} \leftarrow \text{GradientFlow}(g_{\mu_k}, W_{\mu_{k-1}})$ 
6 end

```

---

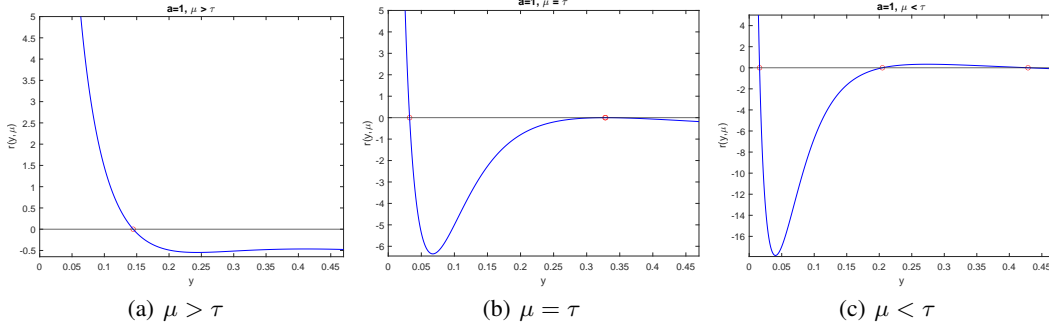


Figure 2: The behavior of  $r(y; \mu)$  for different  $\mu$ .

## 222 4 A Detailed Analysis of the Evolution of the Stationary Points

223 The homotopy approach in Algorithm 2 relies heavily on how the stationary points of (7) behave with  
 224 respect to  $\mu_k$ . In this section, we dive deep into the properties of these critical points.

225 By analyzing the first-order conditions for  $g_\mu$ , we first narrow our attention to the region  $A = \{0 \leq$   
 226  $x \leq a, 0 \leq y \leq \frac{a}{a^2+1}\}$ . By solving the resulting equations, we obtain an equation that only involves  
 227 the variable  $y$ :

$$r(y; \mu) = \frac{a}{y} - \frac{\mu a^2}{(y^2 + \mu)^2} - (a^2 + 1). \quad (8)$$

228 Likewise, we can find an equation only involving the variable  $x$ :

$$t(x; \mu) = \frac{a}{x} - \frac{\mu a^2}{(\mu(a^2 + 1) + x^2)^2} - 1. \quad (9)$$

229 To understand the behavior of the stationary points of  $g_\mu(W)$ , we can examine the characteristics of  
 230  $t(x; \mu)$  in the range  $x \in [0, a]$  and the properties of  $r(y; \mu)$  in the interval  $y \in [0, \frac{a}{a^2+1}]$ .

231 In Figures 2 and 3, we show the behavior of  $r(y; \mu)$  and  $t(x; \mu)$  for  $a = 1$ . Theorems 5 and 6 in the  
 232 appendix establish the existence of a  $\tau > 0$  with the following useful property:

233 **Corollary 2.** *There exists  $\mu < \tau$  such that the equation  $\nabla g_\mu(W) = 0$  has three different solutions,*  
 234 *denoted as  $W_\mu^*, W_\mu^{**}, W_\mu^{***}$ . Then,*

$$\lim_{\mu \rightarrow 0} W_\mu^* = \begin{bmatrix} 0 & a \\ 0 & 0 \end{bmatrix}, \quad \lim_{\mu \rightarrow 0} W_\mu^{**} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad \lim_{\mu \rightarrow 0} W_\mu^{***} = \begin{bmatrix} 0 & 0 \\ \frac{a}{a^2+1} & 0 \end{bmatrix}$$

235 Note that the interesting regime takes place when  $\mu < \tau$ . Then, we characterize the stationary points  
 236 as either local minima or saddle points:

237 **Lemma 1.** *Let  $\mu < \tau$ , then  $g_\mu(W)$  has two local minima at  $W_\mu^*, W_\mu^{***}$ , and a saddle point at  $W_\mu^{**}$ .*

238 With the above results, it has been established that  $W_\mu^*$  converges to the global minimum  $W_G$  as  
 239  $\mu \rightarrow 0$ . In the following section for the proof of Theorem 1, we perform a thorough analysis on how  
 240 to track  $W_\mu^*$  and avoid the local minimum at  $W_\mu^{**}$  by carefully designing the scheduling for  $\mu_k$ .

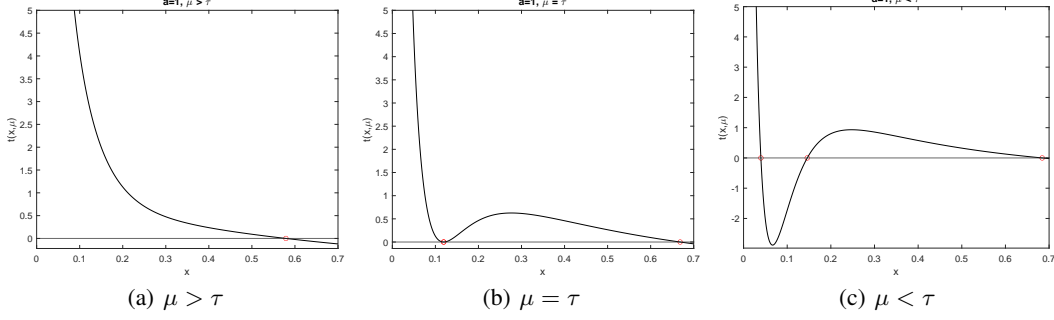


Figure 3: The behavior of  $t(x; \mu)$  for different  $\mu$ .

## 241 5 Convergence Analysis: From continuous to discrete

242 We now discuss the iteration complexity of our method when gradient descent is used in place of  
 243 gradient flow. We begin with some preliminaries regarding the continuous-time analysis.

### 244 5.1 Continuous case: Gradient flow

245 The key to ensuring the convergence of gradient flow to  $W_\mu^*$  is to accurately identify the basin of  
 246 attraction of  $W_\mu^*$ . The following lemma provides a region that lies within such basin of attraction.

247 **Lemma 2.** Define  $B_\mu = \{(x, y) \mid x_\mu^{**} < x \leq a, 0 \leq y < y_\mu^{**}\}$ . Run Algorithm 1 with input  $f =$   
 248  $g_\mu(x, y)$ ,  $z_0 = (x(0), y(0))$  where  $(x(0), y(0)) \in B_\mu$ , then  $\forall t \geq 0$ , we have that  $(x(t), y(t)) \in B_\mu$   
 249 and  $\lim_{t \rightarrow \infty} (x(t), y(t)) = (x_\mu^*, y_\mu^*)$ .

250 In Figure 1 (b), the lower-right rectangle corresponds to  $B_\mu$ . Lemma 2 implies that the gradient flow  
 251 with any initialization inside  $B_{\mu_{k+1}}$  will converge to  $W_{\mu_{k+1}}^*$  at last. Then, by utilizing the previous  
 252 solution  $W_{\mu_k}^*$  as the initial point, as long as it lies within region  $B_{\mu_{k+1}}$ , the gradient flow can converge  
 253 to  $W_{\mu_{k+1}}^*$ , thereby achieving the goal of tracking  $W_{\mu_{k+1}}^*$ . Following the scheduling for  $\mu_k$  prescribed  
 254 in Algorithm 2 provides a sufficient condition to ensure that will happen.

255 The following lemma, with proof in the appendix, is used for the Proof of Theorem 1. It provides a  
 256 lower bound for  $y_\mu^{**}$  and upper bound for  $y_\mu^*$ .

257 **Lemma 3.** If  $\mu < \tau$ , then  $y_\mu^{**} > \sqrt{\mu}$ , and  $\frac{(4\mu)^{1/3}}{2} \left( a^{1/3} - \sqrt{a^{2/3} - (4\mu)^{1/3}} \right) > y_\mu^*$ .

258 **Proof of Theorem 1.** Consider that we are at iteration  $k + 1$  of Algorithm 2, then  $\mu_{k+1} < \mu_k$ . If  
 259  $\mu_k > \tau$  and  $\mu_{k+1} > \tau$ , then there is only one stationary point for  $g_{\mu_k}(x, y)$  and  $g_{\mu_{k+1}}(x, y)$ , thus,  
 260 1 will converge to such stationary point. Hence, let us assume  $\mu_{k+1} \leq \tau$ . From Theorem 6 in the  
 261 appendix, we known that  $x_{\mu_{k+1}}^{**} < x_{\mu_k}^*$ . Then, the following relations hold:

$$y_{\mu_{k+1}}^{**} \stackrel{(1)}{>} \sqrt{\mu_{k+1}} \geq 2 \left( \frac{\mu_k^2}{4a} \right)^{1/3} \stackrel{(2)}{\geq} \frac{(4\mu_k)^{1/3}}{2} \left( a^{1/3} - \sqrt{a^{2/3} - (4\mu_k)^{1/3}} \right) \stackrel{(3)}{>} y_{\mu_k}^*$$

262 Here (1) and (3) are due to Lemma 3, and (2) follows from  $\sqrt{1-x} \geq 1-x$  for  $0 \leq x \leq 1$ . Then it  
 263 implies that  $(x_{\mu_k}^*, y_{\mu_k}^*)$  is in the region  $\{(x, y) \mid x_{\mu_{k+1}}^{**} < x \leq a, 0 \leq y < y_{\mu_{k+1}}^{**}\}$ . By Lemma 2, the  
 264 1 procedure will converge to  $(x_{\mu_{k+1}}^*, y_{\mu_{k+1}}^*)$ . Finally, from Theorems 5 and 6, if  $\lim_{k \rightarrow \infty} \mu_k = 0$ ,  
 265 then  $\lim_{k \rightarrow \infty} x_{\mu_k}^* = a$ ,  $\lim_{k \rightarrow \infty} y_{\mu_k}^* = 0$ , thus, converging to the global optimum, i.e.,

$$\lim_{k \rightarrow \infty} W_{\mu_k} = W_G.$$

### 266 5.2 Discrete case: Gradient Descent

267 In Algorithms 2 and 4, gradient flow is employed to locate the next stationary points, which is not  
 268 practically feasible. A viable alternative is to execute Algorithm 2, replacing the gradient flow with



269 gradient descent. Now, at every iteration  $k$ , Algorithm 6 uses gradient descent to output  $W_{\mu_k, \epsilon_k}$ , a  
 270  $\epsilon_k$  stationary point of  $g_{\mu_k}$ , initialized at  $W_{\mu_{k-1}, \epsilon_{k-1}}$ , and a step size of  $\eta_k = 1/(\mu_k(a^2 + 1) + 3a^2)$ .  
 271 The tolerance parameter  $\epsilon_k$  can significantly influence the behavior of the algorithm and must be  
 272 controlled for different iterations. A convergence guarantee is established via a simplified theorem  
 273 presented here. A more formal version of the theorem and a comprehensive description of the  
 274 algorithm (i.e., Algorithm 6) can be found in Appendix C.

275 **Theorem 2 (Informal).** *For any  $\varepsilon_{\text{dist}} > 0$ , set  $\mu_0$  satisfy a mild condition, and use updating rule  $\epsilon_k =$   
 276  $\min\{\beta a \mu_k, \mu_k^{3/2}\}$ ,  $\mu_{k+1} = (2\mu_k^2)^{2/3} \frac{(a + \epsilon_k/\mu_k)^{2/3}}{(a - \epsilon_k/\mu_k)^{4/3}}$ , and let  $K \equiv K(\mu_0, a, \varepsilon_{\text{dist}}) \in O\left(\ln \frac{\mu_0}{a\varepsilon_{\text{dist}}}\right)$ .  
 277 Then, for any initialization  $W_0$ , following the updated procedure above for  $k = 0, \dots, K$ , we have:*

$$\|W_{\mu_k, \epsilon_k} - W_G\|_2 \leq \varepsilon_{\text{dist}}$$

278 that is,  $W_{\mu_k, \epsilon_k}$  is  $\varepsilon_{\text{dist}}$ -close in Frobenius norm to global optimum  $W_G$ . Moreover, the total number  
 279 of gradient descent steps is upper bounded by  $O\left((\mu_0 a^2 + a^2 + \mu_0) \left(\frac{1}{a^6} + \frac{1}{\varepsilon_{\text{dist}}^6}\right)\right)$ .

## 280 6 Experiments

281 We conducted experiments to verify that Algorithms 2 and 4 both converge to the global minimum of  
 282 (7). Our purpose is to illustrate two main points: First, we compare our updating scheme as given in  
 283 Line 3 of Algorithm 2 against a faster-decreasing updating scheme for  $\mu_k$ . In Figure 4 we illustrate  
 284 how a naive faster decrease of  $\mu$  can lead to spurious a local minimum. Second, in Figure 5, we show  
 285 that regardless of the initialization, Algorithms 2 and 4 always return the global minimum. In the  
 286 supplementary material, we provide additional experiments where the gradient flow is replaced with  
 287 gradient descent. For more details, please refer to Appendix F.

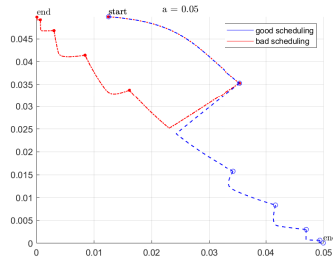


Figure 4: Trajectory of the gradient flow path for two different update rules for  $\mu_k$  with same initialization and  $\mu_0$ . Here, “good scheduling” uses Line 3 of Algorithm 2, while “bad scheduling” uses a faster decreasing scheme for  $\mu_k$  which leads the path to a spurious local minimum.

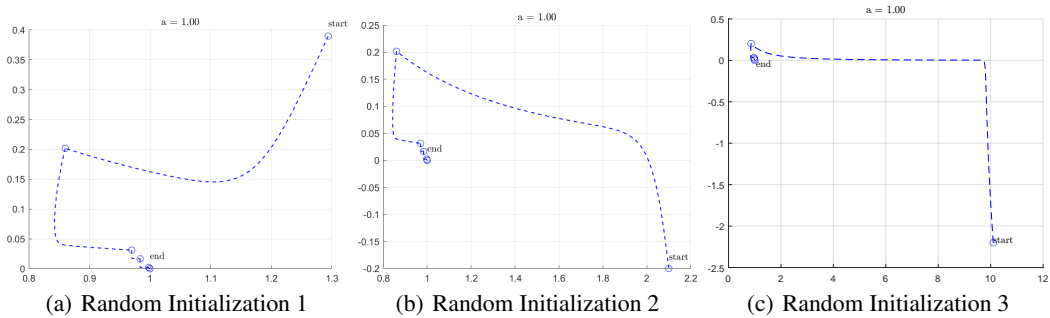


Figure 5: Trajectory of the gradient flow path with the different initializations. We observe that under a proper scheduling for  $\mu_k$ , they all converge to the global minimum.

## References

- [1] Bello, K., Aragam, B. and Ravikumar, P. [2022], DAGMA: Learning dags via M-matrices and a log-determinant acyclicity characterization, in ‘Advances in Neural Information Processing Systems’.
- [2] Bertsekas, D. P. [1997], ‘Nonlinear programming’, *Journal of the Operational Research Society* **48**(3), 334–334.
- [3] Bhojanapalli, S., Neyshabur, B. and Srebro, N. [2016], ‘Global optimality of local search for low rank matrix recovery’, *Advances in Neural Information Processing Systems* **29**.
- [4] Boumal, N., Voroninski, V. and Bandeira, A. [2016], ‘The non-convex burer-monteiro approach works on smooth semidefinite programs’, *Advances in Neural Information Processing Systems* **29**.
- [5] Boyd, S., Boyd, S. P. and Vandenberghe, L. [2004], *Convex optimization*, Cambridge university press.
- [6] Brutzkus, A. and Globerson, A. [2017], Globally optimal gradient descent for a convnet with gaussian inputs, in ‘International conference on machine learning’, PMLR, pp. 605–614.
- [7] Chen, W., Drton, M. and Wang, Y. S. [2019], ‘On causal discovery with an equal-variance assumption’, *Biometrika* **106**(4), 973–980.
- [8] Chickering, D. M. [1996], Learning Bayesian networks is NP-complete, in ‘Learning from data’, Springer, pp. 121–130.
- [9] Chickering, D. M. [2003], ‘Optimal structure identification with greedy search’, *JMLR* **3**, 507–554.
- [10] Chickering, D. M., Heckerman, D. and Meek, C. [2004], ‘Large-sample learning of Bayesian networks is NP-hard’, *Journal of Machine Learning Research* **5**, 1287–1330.
- [11] Choromanska, A., Henaff, M., Mathieu, M., Arous, G. B. and LeCun, Y. [2015], The loss surfaces of multilayer networks, in ‘Artificial intelligence and statistics’, PMLR, pp. 192–204.
- [12] Cussens, J. [2012], ‘Bayesian network learning with cutting planes’, *arXiv preprint arXiv:1202.3713*.
- [13] De Sa, C., Re, C. and Olukotun, K. [2015], Global convergence of stochastic gradient descent for some non-convex matrix problems, in ‘International conference on machine learning’, PMLR, pp. 2332–2341.
- [14] Dontchev, A. L., Rockafellar, R. T. and Rockafellar, R. T. [2009], *Implicit functions and solution mappings: A view from variational analysis*, Vol. 11, Springer.
- [15] Dunlavy, D. M. and O’Leary, D. P. [2005], Homotopy optimization methods for global optimization, Technical report, Sandia National Laboratories (SNL), Albuquerque, NM, and Livermore, CA . . . .
- [16] Duong, B. and Nguyen, T. [2022], Bivariate causal discovery via conditional divergence, in ‘Conference on Causal Learning and Reasoning’.
- [17] Gargiani, M., Zanelli, A., Tran-Dinh, Q., Diehl, M. and Hutter, F. [2020], ‘Convergence analysis of homotopy-sgd for non-convex optimization’, *arXiv preprint arXiv:2011.10298*.
- [18] Garrigues, P. and Ghaoui, L. [2008], ‘An homotopy algorithm for the lasso with online observations’, *Advances in neural information processing systems* **21**.
- [19] Ge, R., Jin, C. and Zheng, Y. [2017], No spurious local minima in nonconvex low rank problems: A unified geometric analysis, in ‘International Conference on Machine Learning’, PMLR, pp. 1233–1242.
- [20] Ge, R., Lee, J. D. and Ma, T. [2016], ‘Matrix completion has no spurious local minimum’, *Advances in neural information processing systems* **29**.

- [21] Hazan, E., Levy, K. Y. and Shalev-Shwartz, S. [2016], On graduated optimization for stochastic non-convex problems, in ‘International conference on machine learning’, PMLR, pp. 1833–1841.
- [22] Heckerman, D., Geiger, D. and Chickering, D. M. [1995], ‘Learning bayesian networks: The combination of knowledge and statistical data’, *Machine learning* **20**(3), 197–243.
- [23] Iwakiri, H., Wang, Y., Ito, S. and Takeda, A. [2022], ‘Single loop gaussian homotopy method for non-convex optimization’, *arXiv preprint arXiv:2203.05717*.
- [24] Jaakkola, T., Sontag, D., Globerson, A. and Meila, M. [2010], Learning bayesian network structure using lp relaxations, in ‘Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics’, Vol. 9 of *Proceedings of Machine Learning Research*, pp. 358–365.
- [25] Jiao, R., Lin, N., Hu, Z., Bennett, D. A., Jin, L. and Xiong, M. [2018], ‘Bivariate causal discovery and its applications to gene expression and imaging data analysis’, *Frontiers Genetics* **9**, 347.
- [26] Kawaguchi, K. [2016], ‘Deep learning without poor local minima’, *Advances in neural information processing systems* **29**.
- [27] Kazemipour, A., Larsen, B. W. and Druckmann, S. [2019], ‘Avoiding spurious local minima in deep quadratic networks’, *arXiv preprint arXiv:2001.00098*.
- [28] Khalil, H. K. [2002], ‘Nonlinear systems third edition’, *Patience Hall* **115**.
- [29] Lachapelle, S., Brouillard, P., Deleu, T. and Lacoste-Julien, S. [2020], Gradient-based neural dag learning, in ‘International Conference on Learning Representations’.
- [30] Liang, S., Sun, R., Lee, J. D. and Srikant, R. [2018], ‘Adding one neuron can eliminate all bad local minima’, *Advances in Neural Information Processing Systems* **31**.
- [31] Loh, P.-L. and Bühlmann, P. [2014], ‘High-dimensional learning of linear causal networks via inverse covariance estimation’, *The Journal of Machine Learning Research* **15**(1), 3065–3105.
- [32] Mooij, J. M., Peters, J., Janzing, D., Zscheischler, J. and Schölkopf, B. [2014], ‘Distinguishing cause from effect using observational data: methods and benchmarks’, *Arxiv*.
- [33] Nesterov, Y. et al. [2018], *Lectures on convex optimization*, Vol. 137, Springer.
- [34] Ng, I., Ghassami, A. and Zhang, K. [2020], On the role of sparsity and dag constraints for learning linear DAGs, in ‘Advances in Neural Information Processing Systems’.
- [35] Ng, I., Lachapelle, S., Ke, N. R., Lacoste-Julien, S. and Zhang, K. [2022], On the convergence of continuous constrained optimization for structure learning, in ‘International Conference on Artificial Intelligence and Statistics’, PMLR, pp. 8176–8198.
- [36] Ni, Y. [2022], ‘Bivariate causal discovery for categorical data via classification with optimal label permutation’, *Arxiv*.
- [37] Nocedal, J. and Wright, S. J. [1999], *Numerical optimization*, Springer.
- [38] Park, Y. W. and Klabjan, D. [2017], ‘Bayesian network learning via topological order’, *The Journal of Machine Learning Research* **18**(1), 3451–3482.
- [39] Pearl, J. [2009], *Causality: Models, Reasoning, and Inference*, 2nd edn, Cambridge University Press.
- [40] Sachs, K., Perez, O., Pe’er, D., Lauffenburger, D. A. and Nolan, G. P. [2005], ‘Causal protein-signaling networks derived from multiparameter single-cell data’, *Science* **308**(5721), 523–529.
- [41] Scanagatta, M., de Campos, C. P., Corani, G. and Zaffalon, M. [2015], Learning bayesian networks with thousands of variables., in ‘NIPS’, pp. 1864–1872.

- 378 [42] Silander, T. and Myllymaki, P. [2006], A simple approach for finding the globally optimal  
379 bayesian network structure, in ‘Proceedings of the 22nd Conference on Uncertainty in Artificial  
380 Intelligence’.
- 381 [43] Soltanolkotabi, M., Javanmard, A. and Lee, J. D. [2018], ‘Theoretical insights into the opti-  
382 mization landscape of over-parameterized shallow neural networks’, *IEEE Transactions on*  
383 *Information Theory* **65**(2), 742–769.
- 384 [44] Spirtes, P., Glymour, C. N., Scheines, R. and Heckerman, D. [2000], *Causation, prediction, and*  
385 *search*, MIT press.
- 386 [45] Teyssier, M. and Koller, D. [2005], Ordering-based search: A simple and effective algorithm for  
387 learning bayesian networks, in ‘Proceedings of the Twenty-First Conference on Uncertainty in  
388 Artificial Intelligence’.
- 389 [46] Wauthier, F. and Donnelly, P. [2015], A greedy homotopy method for regression with nonconvex  
390 constraints, in ‘Artificial Intelligence and Statistics’, PMLR, pp. 1051–1060.
- 391 [47] Wei, D., Gao, T. and Yu, Y. [2020], DAGs with no fears: A closer look at continuous optimization  
392 for learning bayesian networks, in ‘Advances in Neural Information Processing Systems’.
- 393 [48] Wu, C., Luo, J. and Lee, J. D. [2018], No spurious local minima in a two hidden unit relu  
394 network, in ‘International Conference on Learning Representations’.
- 395 [49] Yu, Y., Chen, J., Gao, T. and Yu, M. [2019], Dag-gnn: Dag structure learning with graph neural  
396 networks, in ‘International Conference on Machine Learning’, PMLR, pp. 7154–7163.
- 397 [50] Zhang, B., Gaiteri, C., Bodea, L.-G., Wang, Z., McElwee, J., Podtelezhnikov, A. A., Zhang, C.,  
398 Xie, T., Tran, L., Dobrin, R. et al. [2013], ‘Integrated systems approach identifies genetic nodes  
399 and networks in late-onset alzheimer’s disease’, *Cell* **153**(3), 707–720.
- 400 [51] Zheng, X., Aragam, B., Ravikumar, P. K. and Xing, E. P. [2018], DAGs with NO TEARS:  
401 Continuous optimization for structure learning, in ‘Advances in Neural Information Processing  
402 Systems’.
- 403 [52] Zheng, X., Dan, C., Aragam, B., Ravikumar, P. and Xing, E. [2020], Learning sparse nonpara-  
404 metric DAGs, in ‘International Conference on Artificial Intelligence and Statistics’, PMLR,  
405 pp. 3414–3425.
- 406 [53] Zhu, S., Ng, I. and Chen, Z. [2020], Causal discovery with reinforcement learning, in ‘Interna-  
407 tional Conference on Learning Representations’.

---

**Algorithm 4:** Find a path  $\{W_{\mu_k}\}$  via a particular scheduling for  $\mu_k$  when  $a$  is unknown.

---

**Input:**  $\mu_0 \in \left[\frac{a^2}{4(a^2+1)^3}, \frac{a^2}{4}\right)$ ,  $\varepsilon > 0$

**Output:**  $\{W_{\mu_k}\}_{k=0}^\infty$

```

1  $\hat{a} \leftarrow \sqrt{4(\mu_0 + \varepsilon)}$  //  $\forall \varepsilon \geq 0$  s.t.  $\hat{a} < a$ 
2  $W_{\mu_0} \leftarrow \text{GradientFlow}(g_{\mu_0}, \mathbf{0})$ 
3 for  $k = 1, 2, \dots$  do
4   | Let  $\mu_{k+1} \in \left[(2/\hat{a})^{2/3} \mu_k^{4/3}, \mu_k\right)$ 
5   |  $W_{\mu_{k+1}} \leftarrow \text{GradientFlow}(g_{\mu_{k+1}}, W_{\mu_k})$ 
6 end
7 return  $\{W_{\mu_k}\}_{k=0}^\infty$ 

```

---

## A Practical Implementation of Algorithm 2

We present a practical implementation of our homotopy algorithm in Algorithm 4. The updating scheme for  $\mu_k$  is now independent of the parameter  $a$ , but as presented, the initialization for  $\mu_0$  still depends on  $a$ . This is for the following reason: It is possible to make the updating scheme independent of  $a$  *without imposing any additional assumptions on  $a$* , as evidenced by Lemma 4 below. The initialization for  $\mu_0$ , however, is trickier, and we must consider two separate cases:

1. *No assumptions on  $a$ .* In this case, if  $a$  is too small, then the problem becomes harder and the initial choice of  $\mu_0$  matters.
2. *Lower bound on  $a$ .* If we are willing to accept a lower bound on  $a$ , then there is an initialization for  $\mu_0$  that does not depend on  $a$ .

In Corollary 1, we illustrate this last point with the additional condition that  $a > \sqrt{5/27}$ . This essentially amounts to an assumption on the minimum signal, and is quite standard in the literature on learning SEM.

**Lemma 4.** *Under the assumption  $\frac{a^2}{4(a^2+1)^3} \leq \mu_0 < \frac{a^2}{4}$ , the Algorithm 4 outputs the global optimal solution to (6), i.e.*

$$\lim_{k \rightarrow \infty} W_{\mu_k} = W_G.$$

It turns out that the assumption in Lemma 4 is not overly restrictive, as there exist pre-determined sequences of  $\{\mu_k\}_{k=0}^\infty$  that can ensure the effectiveness of Algorithm 4 for any values of  $a$  greater than a certain threshold.

## B From Population Loss to Empirical Loss

The transformation from population loss to empirical can be thought from two components. First, with a given empirical loss, Algorithms 2 and 3 still achieve the global minimum,  $W_G$ , of problem 6, but now the output from the Algorithm is an empirical estimator  $\hat{a}$ , rather than ground truth  $a$ . Theorem 1 and Corollary 1 would continue to be valid. Second, the global optimum,  $W_G$ , of the empirical loss possess the same DAG structure as the underlying  $W_*$ . The finite-sample findings in Section 5 (specifically, Lemmas 18 and 19) of Loh and Bühlmann [31], which offer sufficient conditions on the sample size to ensure that the DAG structures of  $W_G$  and  $W_*$  are identical.

## C From Continuous to Discrete: Gradient Descent

Previously, gradient flow was employed to address the intermediate problem (7), a method that poses implementation challenges in a computational setting. In this section, we introduce Algorithm 6 that leverages gradient descent to solve (7) in each iteration. This adjustment serves practical considerations. We start with the convergence results of Gradient Descent.

**Definition 1.**  $f$  is  $L$ -smooth, if  $f$  is differentiable and  $\forall x, y \in \text{dom}(f)$  such that  $\|\nabla f(x) - \nabla f(y)\|_2 \leq L\|x - y\|_2$ .

---

**Algorithm 5:** Gradient Descent( $f, \eta, W_0, \epsilon$ )

---

**Input:** function  $f$ , step size  $\eta$ , initial point  $W_0$ , tolerance  $\epsilon$

**Output:**  $W_t$

```
1  $t \leftarrow 0$ 
2 while  $\|\nabla f(W_t)\|_2 > \epsilon$  do
3    $W_{t+1} \leftarrow W_t - \eta \nabla f(W_t)$ 
4    $t \leftarrow t + 1$ 
5 end
```

---

---

**Algorithm 6:** Homotopy algorithm using gradient descent for solving (1).

---

**Input:** Initial  $W_{-1} = W(x_{-1}, y_{-1})$ ,  $\mu_0 \in \left[ \frac{a^2}{4(a^2+1)^3} \frac{(1+\beta)^4}{(1-\beta)^2}, \frac{a^2}{4} \frac{(1-\delta)^3(1-\beta)^4}{(1+\beta)^2} \right)$ ,

$\eta_0 = \frac{1}{\mu_0(a^2+1)+3a^2}$ ,  $\epsilon_0 = \min\{\beta a \mu_0, \mu_0^{3/2}\}$

**Output:**  $\{W_{\mu_k}\}_{k=0}^\infty$

```
1  $W_{\mu_0, \epsilon_0} \leftarrow \text{Gradient Descent}(g_{\mu_0}, \eta_0, W_{-1}, \epsilon_0)$ 
2 for  $k = 1, 2, \dots$  do
3   Let  $\mu_k = (2\mu_{k-1}^2)^{2/3} \frac{(a+\epsilon_{k-1}/\mu_{k-1})^{2/3}}{(a-\epsilon_{k-1}/\mu_{k-1})^{4/3}}$ 
4   Let  $\eta_k = \frac{1}{\mu_k(a^2+1)+3a^2}$ 
5   Let  $\epsilon_k = \min\{\beta a \mu_k, \mu_k^{3/2}\}$ 
6    $W_{\mu_k, \epsilon_k} \leftarrow \text{Gradient Descent}(g_{\mu_k}, \eta_k, W_{\mu_{k-1}}, \epsilon_k)$ 
7 end
```

---

441 **Theorem 3** (Nesterov et al. 33). *If function  $f$  is  $L$ -smooth, then Gradient Descent (Algorithm 5) with*  
442 *step size  $\eta = 1/L$ , finds an  $\epsilon$ -first-order stationary point (i.e.  $\|\nabla f(x)\|_2 \leq \epsilon$ ) in  $2L(f(x^0) - f^*)/\epsilon^2$*   
443 *iterations.*

444 One of the pivotal factors influencing the convergence of gradient descent is the selection of the step  
445 size. Theorem 3 select a step size  $\eta = \frac{1}{L}$ . Therefore, our initial step is to determine the smoothness  
446 of  $g_\mu(W)$  within our region of interest,  $A = \{0 \leq x \leq a, 0 \leq y \leq \frac{a}{a^2+1}\}$ .

447 **Lemma 5.** *Consider the function  $g_\mu(W)$  as defined in Equation 7 within the region  $A = \{0 \leq x \leq$   
448  $a, 0 \leq y \leq \frac{a}{a^2+1}\}$ . It follows that for all  $\mu \geq 0$ , the function  $g_\mu(W)$  is  $\mu(a^2 + 1) + 3a^2$ -smooth.*

449 Since gradient descent is limited to identifying the  $\epsilon$  stationary point of the function. Thus, we study  
450 the gradient of  $g_\mu(W) = \mu f(W) + h(W)$ , i.e.  $\nabla g_\mu(W)$  has the following form

$$\nabla g_\mu(W) = \begin{pmatrix} \mu(x-a) + y^2x \\ \mu(a^2+1)y - a\mu + yx^2 \end{pmatrix}$$

451 As gradient descent is limited to identifying the  $\epsilon$  stationary point of the function, we, therefore, focus  
452 on  $\|\nabla g_\mu(W)\|_2 \leq \epsilon$ . This can be expressed in the subsequent manner:

$$\|\nabla g_\mu(W)\|_2 \leq \epsilon \Rightarrow -\epsilon \leq \mu(x-a) + y^2x < \epsilon \quad \text{and} \quad -\epsilon \leq \mu(a^2+1)y - a\mu + yx^2 \leq \epsilon$$

453 As a result,

$$\{(x, y) \mid \|\nabla g_\mu(W)\|_2 \leq \epsilon\} \subseteq \{(x, y) \mid \frac{\mu a - \epsilon}{\mu + y^2} \leq x \leq \frac{\mu a + \epsilon}{\mu + y^2}, \frac{\mu a - \epsilon}{x^2 + \mu(a^2+1)} \leq y \leq \frac{\mu a + \epsilon}{x^2 + \mu(a^2+1)}\}$$

454 Here we denote such region as  $A_{\mu, \epsilon}$

$$A_{\mu, \epsilon} = \{(x, y) \mid \frac{\mu a - \epsilon}{\mu + y^2} \leq x \leq \frac{\mu a + \epsilon}{\mu + y^2}, \frac{\mu a - \epsilon}{x^2 + \mu(a^2+1)} \leq y \leq \frac{\mu a + \epsilon}{x^2 + \mu(a^2+1)}\} \quad (10)$$

455 Figure 6 and 7 illustrate the region  $A_{\mu, \epsilon}$ .

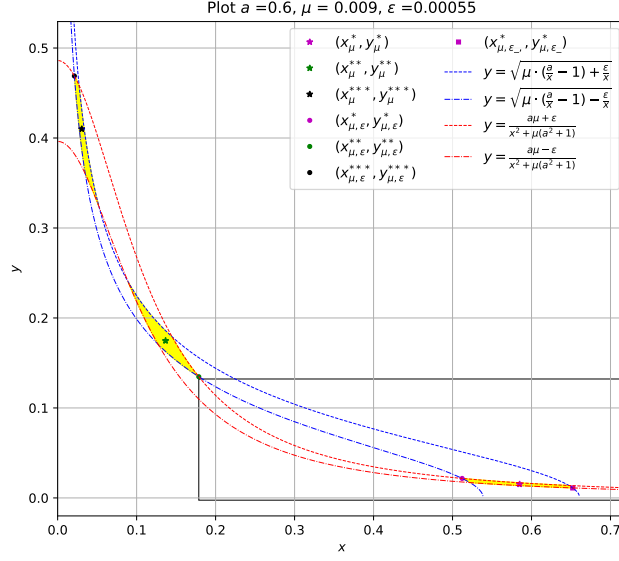


Figure 6: An example of  $A_{\mu, \epsilon}$  is depicted for  $a = 0.6$ ,  $\mu = 0.009$ , and  $\epsilon = 0.00055$ . The yellow region signifies  $\epsilon$  stationary points, denoted as  $A_{\mu, \epsilon}$  and defined by Equation (10).  $A_{\mu, \epsilon}$  is the disjoint union of  $A_{\mu, \epsilon}^1$  and  $A_{\mu, \epsilon}^2$ , which are defined by Equations (21) and (22), respectively.

456

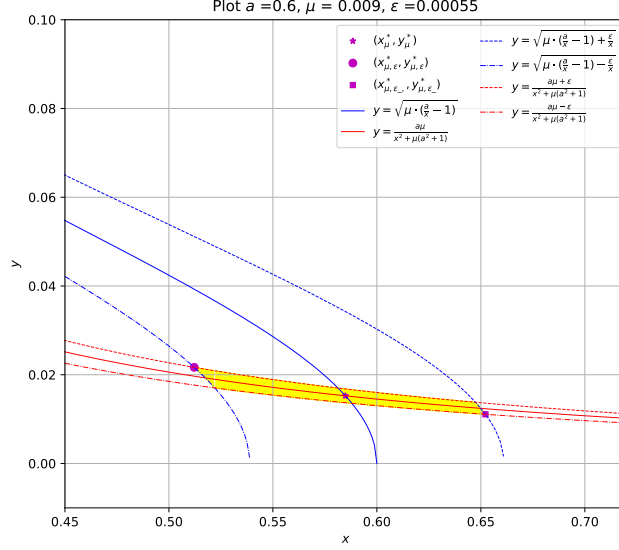


Figure 7: Here is a localized illustration of  $A_{\mu, \epsilon}$  that includes the point  $(x_{\mu}^*, y_{\mu}^*)$ . This region, referred to as  $A_{\mu, \epsilon}^1$ , is defined in Equation (21).

457

458 Given that the gradient descent can only locate  $\epsilon$  stationary points within the region  $A_{\mu, \epsilon}$  during  
 459 each iteration, the boundary of  $A_{\mu, \epsilon}$  becomes a critical component of our analysis. To facilitate clear  
 460 presentation, it is essential to establish some pertinent notations.

$$\begin{cases} x = \frac{\mu a}{\mu + y^2} & (11a) \\ y = \frac{\mu a}{\mu(a^2 + 1) + x^2} & (11b) \end{cases}$$

461 If the system of equations yields only a single solution, we denote this solution as  $(x_\mu^*, y_\mu^*)$ .  
 462 If it yields two solutions, these solutions are denoted as  $(x_\mu^*, y_\mu^*), (x_\mu^{**}, y_\mu^{**})$ , with  $x_\mu^{**} < x_\mu^*$ .  
 463 In the event that there are three distinct solutions to the system of equations, these solutions  
 464 are denoted as  $(x_\mu^*, y_\mu^*), (x_\mu^{**}, y_\mu^{**}), (x_\mu^{***}, y_\mu^{***})$ , where  $x_\mu^{***} < x_\mu^{**} < x_\mu^*$ .

$$\begin{cases} x = \frac{\mu a - \epsilon}{\mu + y^2} & (12a) \\ y = \frac{\mu a + \epsilon}{\mu(a^2 + 1) + x^2} & (12b) \end{cases}$$

465 If the system of equations yields only a single solution, we denote this solution as  $(x_{\mu,\epsilon}^*, y_{\mu,\epsilon}^*)$ .  
 466 If it yields two solutions, these solutions are denoted as  $(x_{\mu,\epsilon}^*, y_{\mu,\epsilon}^*), (x_{\mu,\epsilon}^{**}, y_{\mu,\epsilon}^{**})$ , with  $x_{\mu,\epsilon}^{**} <$   
 467  $x_{\mu,\epsilon}^*$ . In the event that there are three distinct solutions to the system of equations, these  
 468 solutions are denoted as  $(x_{\mu,\epsilon}^*, y_{\mu,\epsilon}^*), (x_{\mu,\epsilon}^{**}, y_{\mu,\epsilon}^{**}), (x_{\mu,\epsilon}^{***}, y_{\mu,\epsilon}^{***})$ , where  $x_{\mu,\epsilon}^{***} < x_{\mu,\epsilon}^{**} < x_{\mu,\epsilon}^*$ .

$$\begin{cases} x = \frac{\mu a + \epsilon}{\mu + y^2} & (13a) \\ y = \frac{\mu a - \epsilon}{\mu(a^2 + 1) + x^2} & (13b) \end{cases}$$

469 If the system of equations yields only a single solution, we denote this solu-  
 470 tion as  $(x_{\mu,\epsilon_-}^*, y_{\mu,\epsilon_-}^*)$ . If it yields two solutions, these solutions are denoted  
 471 as  $(x_{\mu,\epsilon_-}^*, y_{\mu,\epsilon_-}^*), (x_{\mu,\epsilon_-}^{**}, y_{\mu,\epsilon_-}^{**})$ , with  $x_{\mu,\epsilon_-}^{**} < x_{\mu,\epsilon_-}^*$ . In the event that there are  
 472 three distinct solutions to the system of equations, these solutions are denoted as  
 473  $(x_{\mu,\epsilon_-}^*, y_{\mu,\epsilon_-}^*), (x_{\mu,\epsilon_-}^{**}, y_{\mu,\epsilon_-}^{**}), (x_{\mu,\epsilon_-}^{***}, y_{\mu,\epsilon_-}^{***})$ , where  $x_{\mu,\epsilon_-}^{***} < x_{\mu,\epsilon_-}^{**} < x_{\mu,\epsilon_-}^*$ .

474 **Remark 4.** *There always exists at least one solution to the above system of equations. When  $\mu$  is*  
 475 *sufficiently small, the above system of equations always yields three solutions, as demonstrated in*  
 476 *Theorem 5, and Theorem 9.*

477 The parameter  $\epsilon$  can substantially influence the behavior of the systems of equations (12a),(12b) and  
 478 (13a),(13b). A crucial consideration is to ensure that  $\epsilon$  remains adequately small. To facilitate this,  
 479 we introduce a new parameter,  $\beta$ , whose specific value will be determined later. At this stage, we  
 480 merely require that  $\beta$  should lie within the interval  $(0, 1)$ . We further impose a constraint on  $\epsilon$  to  
 481 satisfy the following inequality:

$$\epsilon \leq \beta a \mu \quad (14)$$

482 Following the same procedure when we deal with  $\epsilon = 0$ . Let us substitute (12a) into (12b), then we  
 483 obtain an equation that only involves the variable  $y$

$$r_\epsilon(y; \mu) = \frac{a + \epsilon/\mu}{y} - (a^2 + 1) - \frac{(\mu a - \epsilon)^2/\mu}{(y^2 + \mu)^2} \quad (15)$$

484 Let us substitute (12b) into (12a), then we obtain an equation that only involves the variable  $x$

$$t_\epsilon(x; \mu) = \frac{a - \epsilon/\mu}{x} - 1 - \frac{(\mu a + \epsilon)^2/\mu}{(\mu(a^2 + 1) + x^2)^2} \quad (16)$$

485 Proceed similarly for equations (13a) and (13b).

$$r_{\epsilon_-}(y; \mu) = \frac{a - \epsilon/\mu}{y} - (a^2 + 1) - \frac{(\mu a + \epsilon)^2/\mu}{(y^2 + \mu)^2} \quad (17)$$



$$t_{\epsilon_-}(x; \mu) = \frac{a + \epsilon/\mu}{x} - 1 - \frac{(\mu a - \epsilon)^2/\mu}{(\mu(a^2 + 1) + x^2)^2} \quad (18)$$

Given the substantial role that the system of equations 12a and 12b play in our analysis, the existence of  $\epsilon$  in these equations complicates the analysis, this can be avoided by considering the worst-case scenario, i.e., when  $\epsilon = \beta a \mu$ . With this particular choice of  $\epsilon$ , we can reformulate (15) and (16) as follows, denoting them as  $r_\beta(y; \epsilon)$  and  $r_\beta(x; \epsilon)$  respectively.

$$r_\beta(y; \mu) = \frac{a(1 + \beta)}{y} - (a^2 + 1) - \frac{\mu a^2(1 - \beta)^2}{(y^2 + \mu)^2} \quad (19)$$

$$t_\beta(x; \mu) = \frac{a(1 - \beta)}{x} - 1 - \frac{\mu a^2(1 + \beta)^2}{(\mu(a^2 + 1) + x^2)^2} \quad (20)$$

The functions  $r_\epsilon(y; \mu)$ ,  $r_{\epsilon_-}(y; \mu)$ , and  $r_\beta(y; \mu)$  possess similar properties to  $r(y; \mu)$  as defined in Equation (8), with more details available in Theorem 7 and 8. Additionally, the functions  $t_\epsilon(x; \mu)$ ,  $t_{\epsilon_-}(x; \mu)$ , and  $t_\beta(x; \mu)$  share similar characteristics with  $t(x; \mu)$  as defined in Equation (9), with more details provided in Theorem 9.

As illustrated in Figure 6, the  $\epsilon$ -stationary point region  $A_{\mu, \epsilon}$  can be partitioned into two distinct areas, of which only the lower-right one contains  $(x_{\mu, \epsilon}^*, y_{\mu, \epsilon}^*)$  and it is of interest to our analysis. Moreover,  $(x_{\mu, \epsilon}^*, y_{\mu, \epsilon}^*)$  and  $(x_{\mu, \epsilon}^{**}, y_{\mu, \epsilon}^{**})$  are extremal point of two distinct regions. The upcoming corollary substantiates this intuition.

**Corollary 3.** If  $\mu < \tau$  ( $\tau$  is defined in Theorem 5(v)), assume  $\epsilon$  satisfies (14),  $\beta$  satisfies  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq a^2 + 1$ , systems of equations (12a), (12b) at least have two solutions. Moreover,  $A_{\mu, \epsilon} = A_{\mu, \epsilon}^1 \cup A_{\mu, \epsilon}^2$

$$A_{\mu, \epsilon}^1 = A_{\mu, \epsilon} \cap \{(x, y) \mid x \geq x_{\mu, \epsilon}^*, y \leq y_{\mu, \epsilon}^*\} \quad (21)$$

$$A_{\mu, \epsilon}^2 = A_{\mu, \epsilon} \cap \{(x, y) \mid x \leq x_{\mu, \epsilon}^{**}, y \geq y_{\mu, \epsilon}^{**}\} \quad (22)$$

Corollary 3 suggests that  $A_{\mu, \epsilon}$  can be partitioned into two distinct regions, namely  $A_{\mu, \epsilon}^1$  and  $A_{\mu, \epsilon}^2$ . Furthermore, for every  $(x, y)$  belonging to  $A_{\mu, \epsilon}^1$ , it follows that  $x \geq x_{\mu, \epsilon}^*$  and  $y \leq y_{\mu, \epsilon}^*$ . Similarly, for every  $(x, y)$  that lies within  $A_{\mu, \epsilon}^2$ , the condition  $x \leq x_{\mu, \epsilon}^{**}$  and  $y \geq y_{\mu, \epsilon}^{**}$  holds. The region  $A_{\mu, \epsilon}^1$  represents the “correct” region that gradient descent should identify. In this context, identifying the region equates to pinpointing the extremal points of the region. As a result, our focus should be on the extremal points of  $A_{\mu, \epsilon}^1$  and  $A_{\mu, \epsilon}^2$ , specifically at  $(x_{\mu, \epsilon}^*, y_{\mu, \epsilon}^*)$  and  $(x_{\mu, \epsilon}^{**}, y_{\mu, \epsilon}^{**})$ . Furthermore, the key to ensuring the convergence of the gradient descent to the  $A_{\mu, \epsilon}^1$  is to accurately identify the “basin of attraction” of the region  $A_{\mu, \epsilon}^1$ . The following lemma provides a region within which, regardless of the initialization point of the gradient descent, it converges inside  $A_{\mu, \epsilon}^1$ .

**Lemma 6.** Assume  $\mu < \tau$  ( $\tau$  is defined in Theorem 5(v)),  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq a^2 + 1$ . Define  $B_{\mu, \epsilon} = \{(x, y) \mid x_{\mu, \epsilon}^{**} < x \leq a, 0 \leq y < y_{\mu, \epsilon}^{**}\}$ . Run Algorithm 5 with input  $f = g_\mu(x, y), \eta = \frac{1}{\mu(a^2+1)+3a^2}, W_0 = (x(0), y(0))$ , where  $(x(0), y(0)) \in B_{\mu, \epsilon}$ , then after at most  $\frac{2(\mu(a^2+1)+3a^2)(g_\mu(x(0), y(0)) - g_\mu(x_{\mu, \epsilon}^*, y_{\mu, \epsilon}^*))}{\epsilon^2}$  iterations,  $(x_t, y_t) \in A_{\mu, \epsilon}^1$ .

Lemma 6 can be considered the gradient descent analogue of Lemma 2. It plays a pivotal role in the proof of Theorem 4. In Figure 6, the lower-right rectangle corresponds to  $B_{\mu, \epsilon}$ . Lemma 6 implies that the gradient descent with any initialization inside  $B_{\mu_{k+1}, \epsilon_{k+1}}$  will converge to  $A_{\mu_{k+1}, \epsilon_{k+1}}^1$  at last. Then, by utilizing the previous solution  $W_{\mu_k, \epsilon_k}$  as the initial point, as long as it lies within region  $B_{\mu_{k+1}, \epsilon_{k+1}}$ , the gradient descent can converge to  $A_{\mu_{k+1}, \epsilon_{k+1}}^1$  which is  $\epsilon$  stationary points region that contains  $W_{\mu_{k+1}}^*$ , thereby achieving the goal of tracking  $W_{\mu_{k+1}}^*$ . Following the scheduling for  $\mu_k$  prescribed in Algorithm 6 provides a sufficient condition to ensure that will happen.

We now proceed to present the theorem which guarantees the global convergence of Algorithm 6.

524 **Theorem 4.** If  $\delta \in (0, 1)$ ,  $\beta \in (0, 1)$ ,  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq (1-\delta)(a^2+1)$ , and  $\mu_0$  satisfies

$$\frac{a^2}{4(a^2+1)^3} \leq \frac{a^2}{4(a^2+1)^3} \frac{(1+\beta)^4}{(1-\beta)^2} \leq \mu_0 \leq \frac{a^2}{4} \frac{(1-\delta)^3(1-\beta)^4}{(1+\beta)^2} \leq \frac{a^2}{4}$$

525 Set the updating rule

$$\begin{aligned} \epsilon_k &= \min\{\beta a \mu_k, \mu_k^{3/2}\} \\ \mu_{k+1} &= (2\mu_k^2)^{2/3} \frac{(a + \epsilon_k/\mu_k)^{2/3}}{(a - \epsilon_k/\mu_k)^{4/3}} \end{aligned}$$

526 Then  $\mu_{k+1} \leq (1-\delta)\mu_k$ . Moreover, for any  $\varepsilon_{\text{dist}} > 0$ , running Algorithm 6 after  $K(\mu_0, a, \delta, \varepsilon_{\text{dist}})$   
527 outer iteration

$$\|W_{\mu_k, \epsilon_k} - W_G\|_2 \leq \varepsilon_{\text{dist}} \quad (23)$$

528 where

$$K(\mu_0, a, \delta, \varepsilon_{\text{dist}}) \geq \frac{1}{\ln(1/(1-\delta))} \max \left\{ \ln \frac{\mu_0}{\beta^2 a^2}, \ln \frac{72\mu_0}{a^2(1-(1/2)^{1/4})}, \ln \left( \frac{3(4-\delta)\mu_0}{\varepsilon_{\text{dist}}^2} \right), \frac{1}{2} \ln \left( \frac{46656\mu_0^2}{a^2 \varepsilon_{\text{dist}}^2} \right), \frac{1}{3} \ln \left( \frac{46656\mu_0^3}{a^4 \varepsilon_{\text{dist}}^2} \right) \right\}$$

529 The total gradient descent steps are

$$\begin{aligned} & \sum_{k=0}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{2(\mu_k(a^2+1) + 3a^2)(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\epsilon_k^2} \\ & \leq 2(\mu_0(a^2+1) + 3a^2) \left( \frac{1}{\beta^6 a^6} + \left( \max \left\{ \frac{3(4-\delta)}{\varepsilon_{\text{dist}}^2}, \frac{216}{a \varepsilon_{\text{dist}}}, \left( \frac{216}{a \varepsilon_{\text{dist}}} \right)^{2/3}, \frac{1}{\beta^2 a^2}, \frac{72}{(1-(1/2)^{1/4})a^2} \right\} \right)^3 \right) g_{\mu_0}(W_{\mu_0}^{\epsilon_0}) \\ & \lesssim O(\mu_0 a^2 + a^2 + \mu_0) \left( \frac{1}{\beta^6 a^6} + \frac{1}{\varepsilon_{\text{dist}}^6} + \frac{1}{a^3 \varepsilon_{\text{dist}}^3} + \frac{1}{a^2 \varepsilon_{\text{dist}}^2} + \frac{1}{a^6} \right) \end{aligned}$$

530 *Proof.* Upon substituting gradient flow with gradient descent, it becomes possible to only identify an  
531  $\epsilon$ -stationary point for  $g_\mu(W)$ . This modification necessitates specifying the stepsize  $\eta$  for gradient  
532 descent, as well as an updating rule for  $\mu$ . The adjustment procedure used can substantially influence  
533 the result of Algorithm 6. In this proof, we will impose limitations on the update scheme  $\mu_k$ , the  
534 stepsize  $\eta_k$ , and the tolerance  $\epsilon_k$  to ensure their effective operation within Algorithm 6. The approach  
535 employed for this proof closely mirrors that of the proof for Theorem 1 albeit with more careful  
536 scrutiny. In this proof, we will work out all the requirements for  $\mu, \epsilon, \eta$ . Subsequently, we will verify  
537 that our selection in Theorem 4 conforms to these requirements.

538 In the proof, we occasionally use  $\mu, \epsilon$  or  $\mu_k, \epsilon_k$ . When we employ  $\mu, \epsilon$ , it signifies that the given  
539 inequality or equality holds for any  $\mu, \epsilon$ . Conversely, when we use  $\mu_k, \epsilon_k$ , it indicates we are  
540 examining how to set these parameters for distinct iterations.

541 **Establish the Bound**  $y_{\mu, \epsilon}^{**} \geq \sqrt{\mu}$  First, let us consider  $r_\epsilon(\sqrt{\mu}; \mu) \leq 0$ , i.e.

$$r_\epsilon(\sqrt{\mu}; \mu) = \frac{a + \epsilon/\mu}{\sqrt{\mu}} - (a^2 + 1) - \frac{\mu(a - \epsilon/\mu)^2}{4\mu^2} \leq 0$$

542 This is always true when  $\mu > 4/a^2$ , and we require

$$\epsilon \leq 2\mu^{3/2} + a\mu - 2\sqrt{2a\mu^{5/2} - \mu^3 a^2} \quad \text{when } \mu \leq \frac{4}{a^2}$$

543 Now we name it condition 1.

**Condition 1.**

$$\epsilon \leq 2\mu^{3/2} + a\mu - 2\sqrt{2a\mu^{5/2} - \mu^3 a^2} \quad \text{when } \mu \leq \frac{4}{a^2}$$

544 Under the assumption that Condition 1 is satisfied. Since  $r_\epsilon(y; \mu)$  is increasing function with  
545 interval  $y \in [y_{\text{lb}, \epsilon}, y_{\text{ub}, \epsilon}]$ , and we know  $y_{\text{lb}, \epsilon} \leq \sqrt{\mu} \leq y_{\text{ub}, \epsilon}$  and based on Theorem 7(ii), we have  
546  $y_{\text{lb}, \epsilon} \leq y_{\mu, \epsilon}^{**} \leq y_{\text{ub}, \epsilon}$ ,  $r_\epsilon(\sqrt{\mu}; \mu) \leq r_\epsilon(y_{\mu, \epsilon}^{**}; \mu) = 0$ . Therefore,  $y_{\mu, \epsilon}^{**} \geq \sqrt{\mu}$ .

547 **Ensuring the Correct Solution Path via Gradient Descent** Following the argument when we  
 548 prove Theorem 1, we strive to ensure that the gradient descent, when initiated at  $(x_{\mu_k, \epsilon_k}, y_{\mu_k, \epsilon_k})$ , will  
 549 converge within the "correct"  $\epsilon_{k+1}$ -stationary point region (namely,  $\|\nabla g_{\mu_{k+1}}(W)\|_2 < \epsilon_{k+1}$ ) which  
 550 includes  $(x_{\mu_{k+1}}^*, y_{\mu_{k+1}}^*)$ . For this to occur, we necessitate that:

$$y_{\mu_{k+1}, \epsilon_{k+1}} \stackrel{(1)}{>} y_{\mu_{k+1}, \epsilon_{k+1}}^{**} \stackrel{(2)}{>} \sqrt{\mu_{k+1}} \stackrel{(3)}{\geq} (2\mu_k^2)^{1/3} \frac{(a + \epsilon_k/\mu_k)^{1/3}}{(a - \epsilon_k/\mu_k)^{2/3}} \stackrel{(4)}{>} y_{\mu_k, \epsilon_k}^* \stackrel{(5)}{>} y_{\mu_k, \epsilon_k} \quad (24)$$

551 Here (1), (5) are due to Corollary 3; (2) comes from the boundary we established earlier; (3) is  
 552 based on the constraints we have placed on  $\mu_k$  and  $\mu_{k+1}$ , which we will present as Condition 2  
 553 subsequently; (4) is from the Theorem 7(ii) and relationship  $y_{\mu_k, \epsilon_k}^* < y_{\text{lb}, \mu_k, \epsilon_k}$ . Also, from the  
 554 Lemma 9,  $\max_{\mu \leq \tau} x_{\mu, \epsilon}^{**} \leq \min_{\mu > 0} x_{\mu, \epsilon}^*$ . Hence, by invoking Lemma 6, we can affirm that our  
 555 gradient descent consistently traces the correct stationary point. Now we state condition to make it  
 556 happen,

**Condition 2.**

$$(1 - \delta)\mu_k \geq \mu_{k+1} \geq (2\mu_k^2)^{2/3} \frac{(a + \epsilon_k/\mu_k)^{2/3}}{(a - \epsilon_k/\mu_k)^{4/3}}$$

557 In this context, our requirement extends beyond merely ensuring that  $\mu_k$  decreases. We further  
 558 stipulate that it should decrease by a factor of  $1 - \delta$ . Next, we impose another important constraint

**Condition 3.**

$$\epsilon_k \leq \mu_k^{3/2}$$

559 **Updating Rules** Now we are ready to check our updating rules satisfy the conditions above

$$\begin{aligned} \epsilon_k &= \min\{\beta a \mu_k, \mu_k^{3/2}\} \\ \mu_{k+1} &= (2\mu_k^2)^{2/3} \frac{(a + \epsilon_k/\mu_k)^{2/3}}{(a - \epsilon_k/\mu_k)^{4/3}} \end{aligned}$$

560 **Check for Conditions** First, we check the condition 2. condition 2 requires

$$(1 - \delta)\mu_k \geq (2\mu_k^2)^{2/3} \frac{(a + \epsilon_k/\mu_k)^{2/3}}{(a - \epsilon_k/\mu_k)^{4/3}} \Rightarrow \mu_k \frac{(a + \epsilon_k/\mu_k)^2}{(a - \epsilon_k/\mu_k)^4} \leq \frac{(1 - \delta)^3}{4}$$

561 Note that  $\epsilon_k \leq \beta a \mu_k < a \mu_k$

$$\mu_k \frac{(a + \epsilon_k/\mu_k)^2}{(a - \epsilon_k/\mu_k)^4} \leq \mu_k \frac{(1 + \beta)^2}{(1 - \beta)^4} \frac{1}{a^2}$$

562 Therefore, once the following inequality is true, Condition 2 is satisfied.

$$\mu_k \frac{(1 + \beta)^2}{(1 - \beta)^4} \frac{1}{a^2} \leq \frac{(1 - \delta)^3}{4} \Rightarrow \mu_k \leq \frac{a^2 (1 - \delta)^3 (1 - \beta)^4}{4 (1 + \beta)^2}$$

563 Because  $\mu_k \leq \mu_0 \leq \frac{a^2 (1 - \delta)^3 (1 - \beta)^4}{4 (1 + \beta)^2}$  from the condition we impose for  $\mu_0$ . Consequently, Condition  
 564 2 is satisfied under our choice of  $\epsilon_k$ .

565 Now we focus on the Condition 1. Because  $\epsilon_k \leq a\beta\mu_k$ , if we can ensure  $a\beta\mu_k \leq 2\mu_k^{3/2} + a\mu_k -$   
 566  $2\sqrt{2a\mu_k^{5/2} - \mu_k^3 a^2}$  holds, then we can show Condition 1 is always satisfied.

$$\begin{aligned} a\beta\mu_k &\leq 2\mu_k^{3/2} + a\mu_k - 2\sqrt{2a\mu_k^{5/2} - \mu_k^3 a^2} \\ 2\sqrt{2a\mu_k^{5/2} - \mu_k^3 a^2} &\leq 2\mu_k^{3/2} + (1 - \beta)a\mu_k \\ 4(2a\mu_k^{5/2} - \mu_k^3 a^2) &\leq 4\mu_k^3 + (1 - \beta)^2 a^2 \mu_k^2 + 4(1 - \beta)a\mu_k^{5/2} \\ 0 &\leq 4(a^2 + 1)\mu_k^3 + (1 - \beta)^2 a^2 \mu_k^2 - 4(1 + \beta)a\mu_k^{5/2} \\ 0 &\leq 4(a^2 + 1)\mu_k - 4(1 + \beta)a\mu_k^{1/2} + (1 - \beta)^2 a^2 \quad \text{when } 0 \leq \mu_k \leq 4/a^2 \\ 0 &\leq \mu_k - \frac{(1 + \beta)a}{(a^2 + 1)} \mu_k^{1/2} + \frac{(1 - \beta)^2 a^2}{4(a^2 + 1)} \end{aligned}$$

567 We also notice that

$$\frac{(1+\beta)^2 a^2}{(a^2+1)^2} - 4 \frac{(1-\beta)^2 a^2}{4(a^2+1)} \leq 0 \Leftrightarrow \left( \frac{1+\beta}{1-\beta} \right)^2 \leq a^2 + 1$$

568 Because  $\left( \frac{1+\beta}{1-\beta} \right)^2 \leq (1-\delta)(a^2+1)$ , the inequality above always holds and this inequality implies  
569 that for any  $\mu_k \geq 0$

$$0 \leq \mu_k - \frac{(1+\beta)a}{(a^2+1)} \mu_k^{1/2} + \frac{(1-\beta)^2 a^2}{4(a^2+1)}$$

570 Therefore, Condition 2 holds. Condition 3 also holds because of the choice of  $\epsilon_k$ .

571 **Bound the Distance** Let  $c = 72/a^2$ , and assume that  $\mu$  satisfies the following

$$\mu \leq \min \left\{ \frac{1}{c} \left( 1 - (1/2)^{1/4} \right), \beta^2 a^2 \right\} \quad (25)$$

Note that when  $\mu$  satisfies (25), then  $\mu^{3/2} \leq \beta a \mu$ , so  $\epsilon = \mu^{3/2}$ .

$$\mu \leq \frac{1}{c} \left( 1 - (1/2)^{1/4} \right) = \frac{a^2}{72} \left( 1 - (1/2)^{1/4} \right) \leq \frac{a^2}{4}$$

572

$$\epsilon/\mu = \sqrt{\mu} \leq \frac{a}{2} \quad (26)$$

573 Then

$$\begin{aligned} t_\epsilon((a - \epsilon/\mu)(1 - c\mu); \mu) &= \frac{1}{1 - c\mu} - 1 - \frac{\mu(a + \epsilon/\mu)^2}{(\mu(a^2 + 1) + (a - \epsilon/\mu)^2(1 - c\mu)^2)^2} \\ &= \frac{c\mu}{1 - c\mu} - \frac{\mu(a + \epsilon/\mu)^2}{(\mu(a^2 + 1) + (a - \epsilon/\mu)^2(1 - c\mu)^2)^2} \\ &\geq c\mu - \mu \frac{(a + \epsilon/\mu)^2}{(a - \epsilon/\mu)^4(1 - c\mu)^4} \\ &\geq c\mu - \mu \frac{(a + a/2)^2}{(a - a/2)^4(1 - c\mu)^4} \\ &= \mu \left( c - \frac{36}{a^2(1 - c\mu)^4} \right) \\ &= \mu \left( \frac{72}{a^2} - \frac{36}{a^2(1 - c\mu)^4} \right) > 0 \end{aligned}$$

574 Then we know  $(a - \epsilon/\mu)(1 - c\mu) < x_{\mu, \epsilon}^*$ . Now we can bound the distance  $\|W_{\mu_k, \epsilon_k} - W_G\|$ , it is  
575 important to note that

$$\begin{aligned} \|W_{\mu_k, \epsilon_k} - W_G\| &= \sqrt{(x_{\mu_k, \epsilon_k} - a)^2 + (y_{\mu_k, \epsilon_k})^2} \\ &\leq \max \left\{ \sqrt{(x_{\mu_k, \epsilon_k}^* - a)^2 + (y_{\mu_k, \epsilon_k}^*)^2}, \sqrt{(x_{\mu_k, \epsilon_{k-}}^* - a)^2 + (y_{\mu_k, \epsilon_k}^*)^2} \right\} \end{aligned}$$

576 We use the fact that  $x_{\mu_k, \epsilon_k}^* < x_{\mu_k, \epsilon_k} < a$ ,  $x_{\mu_k, \epsilon_k} < x_{\mu_k, \epsilon_{k-}}^*$  and  $y_{\mu_k, \epsilon_k} < y_{\mu_k, \epsilon_k}^*$ . Next, we can  
577 separately establish bounds for these two terms. Due to (24),  $y_{\mu_k, \epsilon_k}^* < (2\mu_k^2)^{1/3} \frac{(a + \epsilon_k/\mu_k)^{1/3}}{(a - \epsilon_k/\mu_k)^{2/3}} =$   
578  $\sqrt{\mu_{k+1}}$  and  $(a - \epsilon_k/\mu_k)(1 - c\mu_k) < x_{\mu_k, \epsilon_k}^*$

$$\sqrt{(x_{\mu_k, \epsilon_k}^* - a)^2 + (y_{\mu_k, \epsilon_k}^*)^2} \leq \sqrt{\mu_{k+1} + (a - (a - \epsilon_k/\mu_k)(1 - c\mu_k))^2}$$

579 Given that if  $x_{\mu_k, \epsilon_{k-}}^* \leq a$ , then  $\sqrt{(x_{\mu_k, \epsilon_k}^* - a)^2 + (y_{\mu_k, \epsilon_k}^*)^2} \geq \sqrt{(x_{\mu_k, \epsilon_{k-}}^* - a)^2 + (y_{\mu_k, \epsilon_k}^*)^2}$ . There-  
580 fore, if  $x_{\mu_k, \epsilon_{k-}}^* \geq a$ , we can use the fact that  $x_{\mu_k, \epsilon_{k-}}^* \leq a + \frac{\epsilon_k}{\mu_k}$ . In this case,

$$\sqrt{(x_{\mu_k, \epsilon_{k-}}^* - a)^2 + (y_{\mu_k, \epsilon_k}^*)^2} \leq \sqrt{\mu_{k+1} + (\epsilon_k/\mu_k)^2} = \sqrt{\mu_{k+1} + \mu_k} \leq \sqrt{(2 - \delta)\mu_k}$$

581 As a result, we have

$$\|W_{\mu_k, \epsilon_k} - W_G\| \leq \max\{\sqrt{\mu_{k+1} + (a - (a - \epsilon_k/\mu_k)(1 - c\mu_k))^2}, \sqrt{(2 - \delta)\mu_k}\}$$

582

$$\begin{aligned} \mu_{k+1} + (a - (a - \epsilon_k/\mu_k)(1 - c\mu_k))^2 &\leq (1 - \delta)\mu_k + (ac\mu_k + \sqrt{\mu_k} - c\mu_k^{3/2})^2 \\ &\leq (1 - \delta)\mu_k + 3(a^2c^2\mu_k^2 + \mu_k + c^2\mu_k^3) \\ &= (4 - \delta)\mu_k + 3a^2c^2\mu_k^2 + 3c^2\mu_k^3 \end{aligned}$$

583

$$\begin{aligned} \|W_{\mu_k, \epsilon_k} - W_G\| &\leq \max\{\sqrt{\mu_{k+1} + (a - (a - \epsilon_k/\mu_k)(1 - c\mu_k))^2}, \sqrt{(2 - \delta)\mu_k}\} \\ &\leq \max\{\sqrt{(4 - \delta)\mu_k + 3a^2c^2\mu_k^2 + 3c^2\mu_k^3}, \sqrt{(2 - \delta)\mu_k}\} \\ &= \sqrt{(4 - \delta)\mu_k + 3a^2c^2\mu_k^2 + 3c^2\mu_k^3} \end{aligned}$$

584 Just let

$$(4 - \delta)\mu_k \leq (4 - \delta)(1 - \delta)^k \mu_0 \leq \frac{\varepsilon_{\text{dist}}^2}{3} \Rightarrow k \geq \frac{\ln(3(4 - \delta)\mu_0/\varepsilon_{\text{dist}}^2)}{\ln(1/(1 - \delta))} \quad (27)$$

$$3a^2c^2\mu_k^2 \leq 3a^2c^2(1 - \delta)^{2k} \mu_0^2 \leq \frac{\varepsilon_{\text{dist}}^2}{3} \Rightarrow k \geq \frac{\ln(46656\mu_0^2/(a^2\varepsilon_{\text{dist}}^2))}{2\ln(1/(1 - \delta))} \quad (28)$$

$$3c^2\mu_k^3 \leq 3c^2(1 - \delta)^{3k} \mu_0^3 \leq \frac{\varepsilon_{\text{dist}}^2}{3} \Rightarrow k \geq \frac{\ln(46656\mu_0^3/(a^4\varepsilon_{\text{dist}}^2))}{3\ln(1/(1 - \delta))} \quad (29)$$

585 We use the fact that  $\mu_k \leq (1 - \delta)^k \mu_0$ . In order to satisfy (25).

$$\mu_k \leq \mu_0(1 - \delta)^k \leq \frac{a^2}{72}(1 - (1/2)^{1/4}) \Rightarrow k \geq \frac{\ln \frac{72\mu_0}{a^2(1 - (1/2)^{1/4})}}{\ln \frac{1}{1 - \delta}} \quad (30)$$

$$\mu_k \leq \mu_0(1 - \delta)^k \leq \beta^2 a^2 \Rightarrow k \geq \frac{\ln(\mu_0/(\beta^2 a^2))}{\ln \frac{1}{1 - \delta}} \quad (31)$$

586 Consequently, running Algorithm 6 after  $K(\mu_0, a, \delta, \varepsilon_{\text{dist}})$  outer iteration

$$\|W_{\mu_k, \epsilon_k} - W_G\|_2 \leq \varepsilon_{\text{dist}}$$

587 where

$$K(\mu_0, a, \delta, \varepsilon_{\text{dist}}) \geq \frac{1}{\ln(1/(1 - \delta))} \max\left\{\ln \frac{\mu_0}{\beta^2 a^2}, \ln \frac{72\mu_0}{a^2(1 - (1/2)^{1/4})}, \ln\left(\frac{3(4 - \delta)\mu_0}{\varepsilon^2}\right), \frac{1}{2} \ln\left(\frac{46656\mu_0^2}{a^2\varepsilon^2}\right), \frac{1}{3} \ln\left(\frac{46656\mu_0^3}{a^4\varepsilon^2}\right)\right\}$$

588 By Lemma 6,  $k$  iteration of Algorithm 6 need the following step of gradient descent

$$\frac{2(\mu_k(a^2 + 1) + 3a^2)(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\epsilon_k^2}$$

589 Let  $\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})$  satisfy  $\mu_{\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \leq \beta^2 a^2 < \mu_{\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})-1}$ . Hence, the total number  
 590 of gradient steps required by Algorithm 6 can be expressed as follows:

$$\begin{aligned}
& \sum_{k=0}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{2(\mu_k(a^2 + 1) + 3a^2)(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\epsilon_k^2} \\
& \leq 2(\mu_0(a^2 + 1) + 3a^2) \left( \sum_{k=0}^{\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})-1} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\epsilon_k^2} + \sum_{k=\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\epsilon_k^2} \right) \\
& = 2(\mu_0(a^2 + 1) + 3a^2) \left( \sum_{k=0}^{\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})-1} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\beta^2 a^2 \mu_k^2} + \sum_{k=\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\mu_k^3} \right) \\
& \leq 2(\mu_0(a^2 + 1) + 3a^2) \left( \sum_{k=0}^{\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})-1} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\beta^6 a^6} + \sum_{k=\widehat{K}(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\mu_k^3} \right) \\
& \leq 2(\mu_0(a^2 + 1) + 3a^2) \left( \sum_{k=0}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\beta^6 a^6} + \sum_{k=0}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^3} \right) \\
& = 2(\mu_0(a^2 + 1) + 3a^2) \left( \frac{1}{\beta^6 a^6} + \frac{1}{\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^3} \right) \left( \sum_{k=0}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} (g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}})) \right) \\
& \leq 2(\mu_0(a^2 + 1) + 3a^2) \left( \frac{1}{\beta^6 a^6} + \frac{1}{\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^3} \right) \left( \sum_{k=0}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} (g_{\mu_k}(W_{\mu_k}^{\epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}}^{\epsilon_{k+1}})) \right) \\
& = 2(\mu_0(a^2 + 1) + 3a^2) \left( \frac{1}{\beta^6 a^6} + \frac{1}{\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^3} \right) (g_{\mu_0}(W_{\mu_0, \epsilon_0}) - g_{\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})+1}}(W_{\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})+1}}^{\epsilon_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})+1}})) \\
& \leq 2(\mu_0(a^2 + 1) + 3a^2) \left( \frac{1}{\beta^6 a^6} + \frac{1}{\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})}^3} \right) g_{\mu_0}(W_{\mu_0, \epsilon_0})
\end{aligned}$$

591 Note from (27) and (30), the following should holds

$$\mu_{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} = \min\left\{\frac{\varepsilon_{\text{dist}}^2}{3(4-\delta)}, \frac{a\varepsilon_{\text{dist}}}{216}, \left(\frac{a\varepsilon_{\text{dist}}}{216}\right)^{2/3}, \beta^2 a^2, \frac{a^2}{72}(1 - (1/2)^{1/4})\right\}$$

592 Therefore,

$$\begin{aligned}
& \sum_{k=0}^{K(\mu_0, a, \delta, \varepsilon_{\text{dist}})} \frac{2(\mu_k(a^2 + 1) + 3a^2)(g_{\mu_{k+1}}(W_{\mu_k, \epsilon_k}) - g_{\mu_{k+1}}(W_{\mu_{k+1}, \epsilon_{k+1}}))}{\epsilon_k^2} \\
& \leq 2(\mu_0(a^2 + 1) + 3a^2) \left( \frac{1}{\beta^6 a^6} + \left( \max\left\{\frac{3(4-\delta)}{\varepsilon_{\text{dist}}^2}, \frac{216}{a\varepsilon_{\text{dist}}}, \left(\frac{216}{a\varepsilon_{\text{dist}}}\right)^{2/3}, \frac{1}{\beta^2 a^2}, \frac{72}{(1 - (1/2)^{1/4})a^2}\right\} \right)^3 \right) g_{\mu_0}(W_{\mu_0}^{\epsilon_0})
\end{aligned}$$

## 594 D Additional Theorems and Lemmas

595 **Theorem 5** (Detailed Property of  $r(y; \mu)$ ). *For  $r(y; \mu)$  in (8), then*

596 (i) *For  $\mu > 0$ ,  $\lim_{y \rightarrow 0^+} r(y; \mu) = \infty$ ,  $r(\frac{a}{a^2+1}, \mu) < 0$*

597 (ii) *For  $\mu > 0$ ,  $r(\sqrt{\mu}, \mu) < 0$ .*

598 (iii) *For  $\mu > \frac{a^2}{4}$*

$$\frac{dr(y; \mu)}{dy} < 0$$

*For  $0 < \mu \leq \frac{a^2}{4}$*

$$\begin{cases} \frac{dr(y; \mu)}{dy} > 0 & y_{\text{lb}} < y < y_{\text{ub}} \\ \frac{dr(y; \mu)}{dy} \leq 0 & \text{Otherwise} \end{cases} \quad (32a)$$

$$(32b)$$

599 *where*

$$y_{\text{lb}} = \frac{(4\mu)^{1/3}}{2} (a^{1/3} - \sqrt{a^{2/3} - (4\mu)^{1/3}}) \quad y_{\text{ub}} = \frac{(4\mu)^{1/3}}{2} (a^{1/3} + \sqrt{a^{2/3} - (4\mu)^{1/3}})$$

600 *Moreover,*

$$y_{\text{lb}} \leq \sqrt{\mu} \leq y_{\text{ub}}$$

601 (iv) *For  $0 < \mu < \frac{a^2}{4}$ , let  $p(\mu) = r(y_{\text{ub}}, \mu)$ , then  $p'(\mu) < 0$  and there exist a unique solution to*  
 602  *$p(\mu) = 0$ , denoted as  $\tau$ . Additionally,  $\tau < \frac{a^2}{4}$ .*

603 (v) *There exists a  $\tau > 0$  such that,  $\forall \mu > \tau$ , the equation  $r(y; \mu) = 0$  has only one solution. At*  
 604  *$\mu = \tau$ , the equation  $r(y; \mu) = 0$  has two solutions, and  $\forall \mu < \tau$ , the equation  $r(y; \mu) = 0$*   
 605 *has three solutions. Moreover,  $\mu < \frac{a^2}{4}$ .*

606 (vi)  *$\forall \mu < \tau$ , the equation  $r(y; \mu) = 0$  has three solution, i.e.  $y_{\mu}^* < y_{\mu}^{**} < y_{\mu}^{***}$ .*

$$\frac{dy_{\mu}^*}{d\mu} > 0 \quad \frac{dy_{\mu}^{**}}{d\mu} > 0 \quad \frac{dy_{\mu}^{***}}{d\mu} < 0 \text{ and } \lim_{\mu \rightarrow 0} y_{\mu}^* = 0, \lim_{\mu \rightarrow 0} y_{\mu}^{**} = 0, \lim_{\mu \rightarrow 0} y_{\mu}^{***} = \frac{a}{a^2 + 1}$$

607 *Moreover,*

$$y_{\mu}^* < y_{\text{lb}} < \sqrt{\mu} < y_{\mu}^{**} < y_{\text{ub}} < y_{\mu}^{***}$$

608 **Theorem 6** (Detailed Property of  $t(x; \mu)$ ). *For  $t(x; \mu)$  in (9), then*

609 (i) *For  $\mu > 0$ ,  $\lim_{x \rightarrow 0^+} t(x; \mu) = \infty$ ,  $t(a, \mu) < 0$*

610 (ii) *If  $\mu < \left( \frac{a(\sqrt{a^2+1}-a)}{2(a^2+1)} \right)^2$  or  $\mu > \left( \frac{a(\sqrt{a^2+1}+a)}{2(a^2+1)} \right)^2$ , then  $t(\sqrt{\mu(a^2+1)}, \mu) < 0$ .*

611 (iii) *For  $\mu > \frac{a^2}{4(a^2+1)^3}$*

$$\frac{dt(x; \mu)}{dx} < 0$$

*For  $0 < \mu \leq \frac{a^2}{4(a^2+1)^3}$*

$$\begin{cases} \frac{dt(x; \mu)}{dx} > 0 & x_{\text{lb}} < x < x_{\text{ub}} \\ \frac{dt(x; \mu)}{dx} \leq 0 & \text{Otherwise} \end{cases} \quad (33a)$$

$$(33b)$$

612

where

$$x_{\text{lb}} = \frac{(4\mu a)^{1/3}(1 - \sqrt{1 - \frac{(4\mu)^{1/3}(a^2+1)}{a^{2/3}}})}{2} \quad x_{\text{ub}} = \frac{(4\mu a)^{1/3}(1 + \sqrt{1 - \frac{(4\mu)^{1/3}(a^2+1)}{a^{2/3}}})}{2}$$

613

Moreover,

$$x_{\text{lb}} \leq \sqrt{\mu(a^2 + 1)} \leq x_{\text{ub}}$$

614

(iv) For  $0 < \mu < \frac{a^2}{4(a^2+1)^3}$  and let  $q(\mu) = t(x_{\text{lb}}, \mu)$ , then  $q'(\mu) > 0$  and there exist a unique solution to  $q(\mu) = 0$ , denoted as  $\tau$  and  $\tau < \frac{a^2}{4(a^2+1)^3} \leq \frac{1}{27}$ .

615

616

(v) There exists a  $\tau > 0$  such that,  $\forall \mu > \tau$ , the equation  $t(x; \mu) = 0$  has only one solution. At  $\mu = \tau$ , the equation  $t(x; \mu) = 0$  has two solutions, and  $\forall \mu < \tau$ , the equation  $t(x; \mu) = 0$  has three solutions. Moreover,  $\tau < \frac{a^2}{4(a^2+1)^3} \leq \frac{1}{27}$

617

618

(vi)  $\forall \mu < \tau$ ,  $t(x; \mu) = 0$  has three stationary points, i.e.  $x_{\mu}^{***} < x_{\mu}^{**} < x_{\mu}^*$ .

619

$$\frac{dx_{\mu}^*}{d\mu} < 0 \quad \frac{dx_{\mu}^{***}}{d\mu} > 0 \quad \text{and} \quad \lim_{\mu \rightarrow 0} x_{\mu}^* = a, \lim_{\mu \rightarrow 0} x_{\mu}^{**} = 0, \lim_{\mu \rightarrow 0} x_{\mu}^{***} = 0$$

620

Besides,

$$\max_{\mu \leq \tau} x_{\mu}^{**} \leq \frac{a(\sqrt{a^2+1}-a)}{2\sqrt{a^2+1}} \quad \text{and} \quad \frac{a(\sqrt{a^2+1}+a)}{2\sqrt{a^2+1}} \leq \min_{\mu > 0} x_{\mu}^*$$

621

It also implies that  $t(\frac{a(\sqrt{a^2+1}-a)}{2\sqrt{a^2+1}}; \mu) \geq 0$  and  $\max_{\mu \leq \mu_0} x_{\mu}^{**} < \min_{\mu > 0} x_{\mu}^*$

622

**Lemma 7.** Algorithm 1 with input  $f = g_{\mu}(x, y)$ ,  $\mathbf{z}_0 = (x(0), y(0))$  where  $(x(0), y(0)) \in C_{\mu 3}$  in (41), then  $\forall t \geq 0$ ,  $(x(t), y(t)) \in C_{\mu 3}$ . Moreover,  $\lim_{t \rightarrow \infty} (x(t), y(t)) = (x_{\mu}^*, y_{\mu}^*)$

623

624

**Lemma 8.** For any  $(x, y) \in C_{\mu 3}$  in (41), and  $(x, y) \neq (x_{\mu}^*, y_{\mu}^*)$

$$g_{\mu}(x, y) > g_{\mu}(x_{\mu}^*, y_{\mu}^*)$$

625

**Theorem 7** (Detailed Property of  $r_{\epsilon}(y; \mu)$ ). For  $r_{\epsilon}(y; \mu)$  in (15), then

626

(i) For  $\mu > 0, \epsilon > 0$ ,  $\lim_{y \rightarrow 0^+} r_{\epsilon}(y; \mu) = \infty$ ,  $y(\frac{a}{a^2+1}, \mu) < 0$

(ii) For  $\mu > \frac{(a-\epsilon/\mu)^4}{4(a+\epsilon/\mu)^2}$ , then  $\frac{dr_{\epsilon}(y; \mu)}{dy} < 0$ . For  $0 < \mu \leq \frac{(a-\epsilon/\mu)^4}{4(a+\epsilon/\mu)^2}$

$$\begin{cases} \frac{dr_{\epsilon}(y; \mu)}{dy} > 0 & y_{\text{lb}, \mu, \epsilon} < y < y_{\text{ub}, \mu, \epsilon} \\ \frac{dr_{\epsilon}(y; \mu)}{dy} \leq 0 & \text{Otherwise} \end{cases} \quad (34a)$$

(34b)

627

where

$$y_{\text{lb}, \mu, \epsilon} = \frac{(4\mu)^{1/3}}{2} \left( \left( \frac{(a-\epsilon/\mu)^2}{a-\epsilon/\mu} \right)^{1/3} - \sqrt{\left( \frac{(a-\epsilon/\mu)^2}{a-\epsilon/\mu} \right)^{2/3} - (4\mu)^{1/3}} \right)$$

$$y_{\text{ub}, \mu, \epsilon} = \frac{(4\mu)^{1/3}}{2} \left( \left( \frac{(a-\epsilon/\mu)^2}{a-\epsilon/\mu} \right)^{1/3} + \sqrt{\left( \frac{(a-\epsilon/\mu)^2}{a-\epsilon/\mu} \right)^{2/3} - (4\mu)^{1/3}} \right)$$

628

Also,

$$y_{\text{lb}, \mu, \epsilon} \leq (2\mu^2)^{1/3} \frac{(a+\epsilon/\mu)^{1/3}}{(a-\epsilon/\mu)^{2/3}}$$

629

$$y_{\text{lb}, \mu, \epsilon} \leq \sqrt{\mu} \leq y_{\text{ub}, \mu, \epsilon}$$

630

**Theorem 8** (Detailed Property of  $r_{\beta}(y; \mu)$ ). For  $r_{\beta}(y; \mu)$  in (19), then



631 (i) For  $\mu > 0, \epsilon > 0, \lim_{y \rightarrow 0^+} r_\beta(y; \mu) = \infty$

(ii) For  $\mu > \frac{a^2(1-\beta)^4}{4(1+\beta)^2}$ , then  $\frac{dr_\beta(y; \mu)}{dy} < 0$ . For  $0 < \mu \leq \frac{a^2(1-\beta)^4}{4(1+\beta)^2}$

$$\begin{cases} \frac{dr_\beta(y; \mu)}{dy} > 0 & y_{\text{lb}, \mu, \beta} < y < y_{\text{ub}, \mu, \beta} \\ \frac{dr_\beta(y; \mu)}{dy} \leq 0 & \text{Otherwise} \end{cases} \quad (35a)$$

$$\quad (35b)$$

632 where

$$y_{\text{lb}, \mu, \beta} = \frac{(4\mu)^{1/3}}{2} \left( \frac{a(1-\beta)^2}{1+\beta} \right)^{1/3} \left( 1 - \sqrt{1 - \frac{(4\mu)^{1/3}}{a^{2/3}} \left( \frac{1+\beta}{(1-\beta)^2} \right)^{2/3}} \right)$$

$$y_{\text{ub}, \mu, \beta} = \frac{(4\mu)^{1/3}}{2} \left( \frac{a(1-\beta)^2}{1+\beta} \right)^{1/3} \left( 1 + \sqrt{1 - \frac{(4\mu)^{1/3}}{a^{2/3}} \left( \frac{1+\beta}{(1-\beta)^2} \right)^{2/3}} \right)$$

633 Also,

$$y_{\text{lb}, \mu, \beta} \leq \frac{(4\mu)^{2/3} (1+\beta)^{1/3}}{2a^{1/3} (1-\beta)^{2/3}}$$

634

$$y_{\text{lb}, \mu, \beta} \leq \sqrt{\mu} \leq y_{\text{ub}, \mu, \beta}$$

635 **Theorem 9** (Detailed Property of  $t_\beta(x; \mu)$ ). For  $t_\beta(x; \mu)$  in (20), then

636 (i) For  $\mu > 0, \lim_{x \rightarrow 0^+} t_\beta(x; \mu) = \infty, t_\beta(a; \mu) < 0$

637 (ii) For  $\mu > \frac{a^2}{4(a^2+1)^3} \frac{(\beta+1)^4}{(\beta-1)^2}$

$$\frac{dt_\beta(x; \mu)}{dx} < 0$$

For  $0 < \mu \leq \frac{a^2}{4(a^2+1)^3} \frac{(\beta+1)^4}{(\beta-1)^2}$

$$\begin{cases} \frac{dt_\beta(x; \mu)}{dx} > 0 & x_{\text{lb}, \mu, \beta} < x < x_{\text{ub}, \mu, \beta} \\ \frac{dt_\beta(x; \mu)}{dx} \leq 0 & \text{Otherwise} \end{cases} \quad (36a)$$

$$\quad (36b)$$

638 where

$$x_{\text{lb}, \mu, \beta} = \frac{1}{2} \left( \frac{4a\mu(1+\beta)^2}{1-\beta} \right)^{1/3} \left( 1 - \sqrt{1 - \frac{(4\mu)^{1/3}(a^2+1)}{a^{2/3}} \left( \frac{1-\beta}{(1+\beta)^2} \right)^{2/3}} \right)$$

$$x_{\text{ub}, \mu, \beta} = \frac{1}{2} \left( \frac{4a\mu(1+\beta)^2}{1-\beta} \right)^{1/3} \left( 1 + \sqrt{1 - \frac{(4\mu)^{1/3}(a^2+1)}{a^{2/3}} \left( \frac{1-\beta}{(1+\beta)^2} \right)^{2/3}} \right)$$

639 (iii) If  $0 < \beta < \frac{\sqrt{(a^2+1)}-1}{\sqrt{(a^2+1)}+1}$ , then there exists a  $\tau_\beta > 0$  such that,  $\forall \mu > \tau_\beta$ , the equation

640  $r_\beta(x; \mu) = 0$  has only one solution. At  $\mu = \tau_\beta$ , the equation  $r_\beta(x; \mu) = 0$  has two

641 solutions, and  $\forall \mu < \tau_\beta$ , the equation  $r_\beta(x; \mu) = 0$  has three solutions. Moreover,  $\mu <$

642  $\frac{a^2}{4(a^2+1)^3} \frac{(\beta+1)^4}{(\beta-1)^2}$ .

643 (iv) If  $0 < \beta < \frac{\sqrt{(a^2+1)}-1}{\sqrt{(a^2+1)}+1}$ , then  $\forall \mu < \tau_\beta$ ,  $t_\beta(x; \mu) = 0$  has three stationary points, i.e.

644  $x_{\mu, \beta}^{***} < x_{\mu, \beta}^{**} < x_{\mu, \beta}^*$ . Besides,

$$\max_{\mu \leq \tau_\beta} x_{\mu, \beta}^{**} \leq \frac{a((1-\beta)\sqrt{a^2+1} - \sqrt{(1-\beta)^2(a^2+1) - (\beta+1)^2})}{2\sqrt{a^2+1}}$$

$$\frac{a((1-\beta)\sqrt{a^2+1} + \sqrt{(1-\beta)^2(a^2+1) - (\beta+1)^2})}{2\sqrt{a^2+1}} \leq \min_{\mu > 0} x_{\mu, \beta}^*$$

645

*It implies that*

$$\max_{\mu \leq \tau_\beta} x_{\mu,\beta}^{**} < \min_{\mu > 0} x_{\mu,\beta}^*$$

646 **Lemma 9.** *Under the same setting as Corollary 3,*

$$\max_{\mu \leq \tau} x_{\mu,\epsilon}^{**} < \min_{\mu > 0} x_{\mu,\epsilon}^*$$

647 **E Technical Proofs**648 **E.1 Proof of Theorem 3**649 *Proof.* For the sake of completeness, we have included the proof here. Please note that this proof can  
650 also be found in [33].651 *Proof.* We use the fact that  $f$  is  $L$ -smooth function if and only if for any  $W, Y \in \text{dom}(f)$ 

$$f(W) \leq f(Y) + \langle \nabla f(Y), Y - W \rangle + \frac{L}{2} \|Y - W\|_2^2$$

652 Let  $W = W^{t+1}$  and  $Y = W^t$ , then using the updating rule  $W^{t+1} = W^t - \frac{1}{L} \nabla f(W^t)$ 

$$\begin{aligned} f(W^{t+1}) &\leq f(W^t) + \langle \nabla f(W^t), W^{t+1} - W^t \rangle + \frac{L}{2} \|W^{t+1} - W^t\|_2^2 \\ &= f(W^t) - \frac{1}{L} \|\nabla f(W^t)\|_2^2 + \frac{1}{2L} \|\nabla f(W^t)\|_2^2 \\ &= f(W^t) - \frac{1}{2L} \|\nabla f(W^t)\|_2^2 \end{aligned}$$

653 Therefore,

$$\min_{0 \leq t \leq n-1} \|\nabla f(W^t)\|_2^2 \leq \frac{1}{n} \sum_{t=0}^{n-1} \|\nabla f(W^t)\|_2^2 \leq \frac{2L(f(W^0) - f(W^n))}{n} \leq \frac{2L(f(W^0) - f(W^*))}{n}$$

654

$$\min_{0 \leq t \leq n-1} \|\nabla f(W^t)\|_2^2 \leq \frac{2L(f(W^0) - f(W^*))}{n} \leq \epsilon^2 \Rightarrow n \geq \frac{2L(f(W^0) - f(W^*))}{\epsilon^2}$$

655

□

656

□

657 **E.2 Proof of Theorem 5**658 *Proof.* (i) For any  $\mu > 0$ ,

$$\begin{aligned} \lim_{y \rightarrow 0^+} r(y; \mu) &= \lim_{y \rightarrow 0^+} \frac{a}{y} - \frac{a^2}{\mu} - (a^2 + 1) = \infty \\ r\left(\frac{a}{a^2 + 1}\right) &= -\frac{\mu a^2}{\left(\frac{a}{a^2 + 1}\right)^2 + \mu} < 0. \end{aligned}$$

(ii)

$$\begin{aligned} r(\sqrt{\mu}, \mu) &= \frac{a}{\sqrt{\mu}} - \frac{a^2}{4\mu} - (a^2 + 1) \\ &= -\frac{a^2}{4} \left( \frac{1}{\sqrt{\mu}} - \frac{2}{a} \right)^2 - a^2 < 0 \end{aligned}$$

(iii)

$$\begin{aligned}\frac{dr(y; \mu)}{dy} &= -\frac{a}{y^2} + \frac{4a^2\mu y}{(y^2 + \mu)^3} \\ &= \frac{4a^2\mu y^3 - a(y^2 + \mu)^3}{y^2(y^2 + \mu)^3} \\ &= \frac{a((4a\mu)^{2/3}y^2 + (4a\mu)^{1/3}y(y^2 + \mu) + (y^2 + \mu)^2)((4a\mu)^{1/3}y - y^2 - \mu)}{y^2(y^2 + \mu)^3}\end{aligned}$$

For  $\mu \geq \frac{a^2}{4}$ ,  $((4a\mu)^{1/3}y - y^2 - \mu) < 0 \Leftrightarrow \frac{dr(y; \mu)}{dy} < 0$ .  
 For  $\mu < \frac{a^2}{4}$ ,  $y_{\text{lb}} < y < y_{\text{ub}}$ ,  $((4a\mu)^{1/3}y - y^2 - \mu) > 0 \Leftrightarrow \frac{dr(y; \mu)}{dy} > 0$ . For  $\mu < \frac{a^2}{4}$ ,  
 $y < y_{\text{lb}}$  or  $y_{\text{ub}} < y$ ,  $((4a\mu)^{1/3}y - y^2 - \mu) \leq 0 \Leftrightarrow \frac{dr(y; \mu)}{dy} \leq 0$ .

Note that

$$\frac{dr(y; \mu)}{d\mu} = 0 \Leftrightarrow ((4a\mu)^{1/3}y - y^2 - \mu) = 0 \Leftrightarrow (4a\mu)^{1/3} = y + \frac{\mu}{y}$$

The intersection between line  $(4a\mu)^{1/3}$  and function  $y + \frac{\mu}{y}$  are exactly  $y_{\text{lb}}$  and  $y_{\text{ub}}$ , and  
 $y_{\text{lb}} < \sqrt{\mu} < y_{\text{ub}}$ .

(iv) Note that for  $0 < \mu < \frac{a^2}{4}$ ,

$$\frac{\partial r}{\partial \mu} = -a^2 \frac{y^2 - \mu}{(\mu + y^2)^3} \quad \text{and} \quad y_{\text{lb}} < \sqrt{\mu} < y_{\text{ub}}$$

then  $\frac{\partial r}{\partial \mu} \Big|_{y=y_{\text{ub}}} < 0$ . Let  $p(\mu) = r(y_{\text{ub}}, \mu)$ , because  $\frac{\partial r}{\partial y} \Big|_{y=y_{\text{ub}}} = 0$ , then

$$\frac{dp(\mu)}{d\mu} = \frac{dr(y_{\text{ub}}, \mu)}{d\mu} = \frac{\partial r}{\partial y} \Big|_{y=y_{\text{ub}}} \frac{dy_{\text{ub}}}{d\mu} + \frac{\partial r}{\partial \mu} \Big|_{y=y_{\text{ub}}} = \frac{\partial r}{\partial \mu} \Big|_{y=y_{\text{ub}}} < 0$$

Also note that when  $\mu = \frac{a^2}{4}$ ,  $y_{\text{ub}} = \sqrt{\mu}$ ,  $p(\mu) = r(y_{\text{ub}}, \mu) = r(\sqrt{\mu}, \mu) < 0$ , and also if  
 $\mu < \frac{a^2}{4}$ , then

$$y_{\text{ub}} < \frac{(4\mu)^{1/3}}{2} 2a^{1/3} = (4\mu a)^{1/3}$$

Thus,

$$\begin{aligned}r((4\mu a)^{1/3}, \mu) &= \frac{a}{(4\mu a)^{1/3}} - \frac{\mu a^2}{((4\mu a)^{2/3} + \mu)^2} - (a^2 + 1) \\ &= \frac{a}{(4\mu a)^{1/3}} - \frac{a^2}{(\mu)^{1/3}((4a)^{2/3} + \mu^{1/3})^2} - (a^2 + 1) \\ &> \frac{1}{\mu^{1/3}} \left( \frac{a}{(4a)^{1/3}} - \frac{a^2}{(4a)^{4/3}} \right) - (a^2 + 1)\end{aligned}$$

Because  $\frac{a}{(4a)^{1/3}} > \frac{a^2}{(4a)^{4/3}}$ , it is easy to see when  $\mu \rightarrow 0$ ,  $r((4\mu a)^{1/3}, \mu) \rightarrow \infty$ . We know  
 $r(y_{\text{ub}}, \mu) > r((4\mu a)^{1/3}, \mu) \rightarrow \infty$  as  $\mu \rightarrow 0$  because of the monotonicity of  $r(y; \mu)$  in  
 Theorem 5(iii). Combining all of these, i.e.

$$\frac{dp(\mu)}{d\mu} < 0, \quad \lim_{\mu \rightarrow 0^+} p(\mu) = \infty, \quad p\left(\frac{a^2}{4}\right) < 0$$

There exists a  $\tau < \frac{a^2}{4}$  such that  $p(\tau) = 0$

(v) From Theorem 5(iv), for  $\mu > \tau$ , then  $p(\mu) = r(y_{\text{ub}}, \mu) > 0$ , and for  $\mu = \tau$ , then  
 $p(\mu) = r(y_{\text{ub}}, \mu) = 0$ . For  $\mu < \tau$ , then  $p(\mu) = r(y_{\text{ub}}, \mu) < 0$ , combining Theorem  
 5(i), 5(iii), we get the conclusions.

676 (vi) By Theorem 5(v),  $\forall \mu < \tau$ , there exists three stationary points such that  $0 < y_\mu^* < y_{\text{lb}} <$   
 677  $\sqrt{\mu} < y_\mu^{**} < y_{\text{ub}} < y_\mu^{***}$ . Because  $\left. \frac{dr(y; \mu)}{dy} \right|_{y=y_{\text{lb}}} = \left. \frac{dr(y; \mu)}{dy} \right|_{y=y_{\text{ub}}} = 0$ , then

$$\left. \frac{dr(y; \mu)}{dy} \right|_{y=y_\mu^*} \neq 0, \quad \left. \frac{dr(y; \mu)}{dy} \right|_{y=y_\mu^{**}} \neq 0, \quad \left. \frac{dr(y; \mu)}{dy} \right|_{y=y_\mu^{***}} \neq 0$$

678 By implicit function theorem [14], for solution to equation  $r(y; \mu) = 0$ , there exists a  
 679 unique continuously differentiable function such that  $y = y(\mu)$  and satisfies  $r(y(\mu), \mu) = 0$ .  
 680 Therefore,

$$\frac{\partial r}{\partial \mu} = -a^2 \frac{y^2 - \mu}{(\mu + y^2)^3}, \quad \frac{\partial r}{\partial y} = -\frac{a}{y^2} + \frac{4a^2 \mu y}{(y^2 + \mu)^3}, \quad \frac{dy(\mu)}{d\mu} = -\frac{\partial r / \partial \mu}{\partial r / \partial y}$$

681 Therefore by Theorem 5(iii),

$$\left. \frac{dy}{d\mu} \right|_{y=y_\mu^*} > 0, \quad \left. \frac{dy}{d\mu} \right|_{y=y_\mu^{**}} > 0, \quad \left. \frac{dy}{d\mu} \right|_{y=y_\mu^{***}} < 0$$

682 Because  $\lim_{\mu \rightarrow 0^+} y_{\text{lb}} = \lim_{\mu \rightarrow 0^+} y_{\text{ub}} = 0$ , then  $\lim_{\mu \rightarrow 0^+} y_\mu^* = \lim_{\mu \rightarrow 0^+} y_\mu^{**} = 0$ . Let us  
 683 consider  $r(\frac{a}{a^2+1}(1-c\mu), \mu)$  where  $c = 32 \frac{(a^2+1)^3}{a^2}$  and  $\mu < \frac{1}{2c}$

$$\begin{aligned} & r\left(\frac{a}{a^2+1}(1-c\mu), \mu\right) \\ &= \frac{a}{\frac{a}{a^2+1}(1-c\mu)} - \frac{\mu a^2}{\left(\frac{a^2}{(a^2+1)^2}(1-c\mu)^2 + \mu\right)^2} - (a^2+1) \\ &= (a^2+1)\left(\frac{c\mu}{1-c\mu}\right) - \frac{\mu a^2}{\left(\frac{a^2}{(a^2+1)^2}(1-c\mu)^2 + \mu\right)^2} \\ &\geq c(a^2+1)\mu - \frac{\mu a^2}{\left(\frac{a^2}{(a^2+1)^2}(1-c\mu)^2\right)^2} \\ &= c(a^2+1)\mu - \frac{16(a^2+1)^4}{a^2}\mu \\ &= \frac{16(a^2+1)^4}{a^2}\mu > 0 \end{aligned}$$

684 By Theorem 5(iii), then  $\frac{a}{a^2+1}(1-c\mu) < y_\mu^{***}$ , then

$$\frac{a}{a^2+1} = \lim_{\mu \rightarrow 0^+} \frac{a}{a^2+1}(1-c\mu), \mu \leq \lim_{\mu \rightarrow 0^+} y_\mu^{***} \leq \frac{a}{a^2+1}$$

685 Consequently,

$$\lim_{\mu \rightarrow 0^+} y_\mu^{***} = \frac{a}{a^2+1}$$

686

□

### 687 E.3 Proof of Theorem 6

688 *Proof.* (i) For  $\mu > 0$ ,

$$\begin{aligned} \lim_{x \rightarrow 0^+} t(x; \mu) &= \lim_{x \rightarrow 0^+} \frac{a}{x} - \frac{a^2}{\mu(a^2+1)^2} - 1 = \infty \\ t(a, \mu) &= -\frac{\mu a^2}{(\mu(a^2+1) + a^2)^2} < 0 \end{aligned}$$

(ii)

$$t(\sqrt{\mu(a^2+1)}, \mu) = \frac{a}{\sqrt{a^2+1}} \frac{1}{\sqrt{\mu}} - \frac{a^2}{4\mu(a^2+1)^2} - 1$$

689 If  $t(\sqrt{\mu(a^2+1)}, \mu) = 0$ , then

$$\frac{1}{\sqrt{\mu}} = 2 \frac{(a^2+1)^{3/2}}{a} \pm 2(a^2+1) \Rightarrow \mu = \left( \frac{a(\sqrt{a^2+1} \mp a)}{2(a^2+1)} \right)^2$$

690 so when  $\mu < \left( \frac{a(\sqrt{a^2+1}-a)}{2(a^2+1)} \right)^2$  or  $\mu > \left( \frac{a(\sqrt{a^2+1}+a)}{2(a^2+1)} \right)^2$ , then  $t(\sqrt{\mu(a^2+1)}, \mu) < 0$

(iii)

$$\begin{aligned} & \frac{dt(x, \mu)}{dx} \\ &= -\frac{a}{x^2} + \frac{4\mu a^2 x}{(\mu(a^2+1) + x^2)^3} \\ &= \frac{4\mu a^2 x^3 - a(\mu(a^2+1) + x^2)^3}{x^2(\mu(a^2+1) + x^2)^3} \\ &= \frac{a((\mu(a^2+1) + x^2)^2 + (\mu(a^2+1) + x^2)(4\mu a)^{1/3}x + (4\mu a)^{2/3}x^2)((4\mu a)^{1/3}x - \mu(a^2+1) - x^2)}{x^2(\mu(a^2+1) + x^2)^3} \end{aligned}$$

691 For  $\mu > \frac{a^2}{4(a^2+1)^3}$ , then  $(4\mu a)^{1/3}x - \mu(a^2+1) - x^2 < 0 \Leftrightarrow \frac{dt(x, \mu)}{dx} < 0$ . For  $\mu < \frac{a^2}{4(a^2+1)^3}$ ,  
 692 and  $x_{lb} < x < x_{ub}$ , then  $(4\mu a)^{1/3}x - \mu(a^2+1) - x^2 > 0 \Leftrightarrow \frac{dt(x, \mu)}{dx} > 0$ . For  $\mu < \frac{a^2}{4(a^2+1)^3}$ ,  
 693  $x < x_{lb}$  or  $x > x_{ub}$ ,  $(4\mu a)^{1/3}x - \mu(a^2+1) - x^2 < 0 \Leftrightarrow \frac{dt(x, \mu)}{dx} < 0$ .

694 We use the same argument as before to show that

$$x_{lb} < \sqrt{\mu(a^2+1)} < x_{ub}$$

695 (iv) Note that for  $0 < \mu < \frac{a^2}{4(a^2+1)^3}$

$$\frac{\partial t}{\partial \mu} = -a^2 \frac{x^2 - \mu(a^2+1)}{(\mu(a^2+1) + x^2)^3} \quad \text{and} \quad x_{lb} < \sqrt{\mu(a^2+1)} < x_{ub}$$

696 then  $\frac{\partial t}{\partial \mu} \Big|_{x=x_{lb}} > 0$ . Let  $q(\mu) = t(x_{lb}, \mu)$ , because  $\frac{\partial t}{\partial x} \Big|_{x=x_{lb}} = 0$ , then

$$\frac{dq(\mu)}{d\mu} = \frac{dt(x_{lb}, \mu)}{d\mu} = \frac{\partial t}{\partial x} \Big|_{x=x_{lb}} \frac{dx_{lb}}{d\mu} + \frac{\partial t}{\partial \mu} \Big|_{x=x_{lb}} = \frac{\partial t}{\partial \mu} \Big|_{x=x_{lb}} > 0$$

697 Note that  $\mu = \frac{a^2}{4(a^2+1)^3}$ ,  $x_{ub} = x_{lb} = \frac{(4\mu a)^{1/3}}{2}$ ,  $t(\frac{(4\mu a)^{1/3}}{2}, \frac{a^2}{4(a^2+1)^3}) = \frac{a}{(4\mu a)^{1/3}} - 1 > 0$ .

698 When  $\mu < \left( \frac{a(\sqrt{a^2+1}-a)}{2(a^2+1)} \right)^2$ , then  $t(\sqrt{\mu(a^2+1)}, \mu) < 0$  by Theorem 6(ii). It implies that  
 699  $q(\mu) < 0$  when  $\mu \rightarrow 0^+$ . By Theorem 6(iii),  $q(\mu) = t(x_{lb}, \mu) < t(\sqrt{\mu(a^2+1)}, \mu) < 0$ .  
 700 Combining all of the theses, i.e.

$$\frac{dq(\mu)}{d\mu} > 0, \quad \lim_{\mu \rightarrow 0^+} q(\mu) < 0, \quad q\left(\frac{a^2}{4(a^2+1)^3}\right) > 0$$

701 There exists a  $\tau < \frac{a^2}{4(a^2+1)^3}$ ,  $q(\tau) = 0$ . Such  $\tau$  is the same as in Theorem 5(iv).

702 (v) We follow the same proof from the proof of Theorem 5(v).

703 (vi) By Theorem 6(v),  $\forall \mu < \mu_0$ , there exists three stationary points such that  $0 < x_{\mu}^{***} < x_{lb} <$

704  $x_{\mu}^{**} < x_{ub} < x_{\mu}^* < a$ . Because  $\frac{dt(x; \mu)}{dx} \Big|_{x=x_{lb}} = \frac{dt(x; \mu)}{dx} \Big|_{x=x_{ub}} = 0$ , then

$$\frac{dt(x; \mu)}{dx} \Big|_{x=x_{\mu}^*} \neq 0, \quad \frac{dt(x; \mu)}{dx} \Big|_{x=x_{\mu}^{**}} \neq 0, \quad \frac{dt(x; \mu)}{dx} \Big|_{x=x_{\mu}^{***}} \neq 0$$

By implicit function theorem [14], for solutions to equation  $t(x; \mu) = 0$ , there exists a unique continuously differentiable function such that  $x = x(\mu)$  and satisfies  $t(x(\mu), \mu) = 0$ . Therefore,

$$\frac{dx}{d\mu} = -\frac{\partial t / \partial \mu}{\partial t / \partial x} = a^2 \frac{\frac{x^2 - \mu(a^2 + 1)}{(\mu(a^2 + 1) + x^2)^3}}{-\frac{a}{x^2} + \frac{4\mu a^2 x}{(\mu(a^2 + 1) + x^2)^3}}$$

Therefore, by Theorem 6(iii)

$$\left. \frac{dx}{d\mu} \right|_{x=x_\mu^*} < 0 \quad \left. \frac{dx}{d\mu} \right|_{x=x_\mu^{**}} > 0$$

Because  $0 < x_\mu^{***} < x_{lb} < x_\mu^{**} < x_{ub}$  and  $\lim_{\mu \rightarrow 0^+} x_{lb} = \lim_{\mu \rightarrow 0^+} x_{ub} = 0$ .

$$\lim_{\mu \rightarrow 0} x_\mu^{**} = \lim_{\mu \rightarrow 0} x_\mu^{***} = 0$$

Let us consider  $t(a(1 - c\mu), \mu)$  where  $c = \frac{32}{a^2}$  and  $\mu < \frac{1}{2c}$

$$\begin{aligned} & t(a(1 - c\mu); \mu) \\ &= \frac{a}{a(1 - c\mu)} - \frac{\mu a^2}{(\mu(a^2 + 1) + a^2(1 - c\mu)^2)^2} - 1 \\ &= \frac{c\mu}{1 - c\mu} - \frac{\mu a^2}{(\mu(a^2 + 1) + a^2(1 - c\mu)^2)^2} \\ &\geq c\mu - \frac{\mu a^2}{(a^2(1 - c\mu)^2)^2} \\ &\geq c\mu - \frac{16}{a^2}\mu > 0 \end{aligned}$$

By Theorem 6(iii). It implies

$$a(1 - c\mu) \leq x_\mu^*$$

taking  $\mu \rightarrow 0^+$  on both side,

$$a = \lim_{\mu \rightarrow 0^+} a(1 - c\mu) \leq \lim_{\mu \rightarrow 0^+} x_\mu^* \leq a$$

Hence,  $\lim_{\mu \rightarrow 0} x_\mu^* = a$ .

When  $\mu = \tau$ , because  $t(x_{lb}; \mu) = 0$  and  $x_{ub} > \sqrt{\mu(a^2 + 1)} > x_{lb}$ ,  $t(x; \mu)$  is increasing function between  $[x_{lb}, x_{ub}]$  then  $t(\sqrt{\mu(a^2 + 1)}; \mu) > t(x_{lb}; \mu) = 0$ . Moreover,  $t(\sqrt{\mu(a^2 + 1)}; \mu)$ ,  $x_{lb}$  and  $x_\mu^{**}$  are continuous function w.r.t  $\mu$ ,  $\exists \delta > 0$  which is really small, such that  $\mu = \tau - \delta$  and  $t(\sqrt{\mu(a^2 + 1)}; \mu) > 0$ ,  $t(x_{lb}, \mu) < 0$  (by Theorem 6(iv)) and  $x_\mu^{**} > x_{lb}$ , hence  $\left. \frac{dx}{d\mu} \right|_{x=x_\mu^{**}} < 0$ . It implies when  $\mu$  decreases, then  $x_\mu^{**}$  increases. This relation holds until  $x_\mu^{**} = \sqrt{\mu(a^2 + 1)}$

$$\begin{aligned} & t(x_\mu^{**}, \mu) = t(\sqrt{\mu(a^2 + 1)}, \mu) = 0 \\ & \Rightarrow \mu = \left( \frac{a(\sqrt{a^2 + 1} - a)}{2(a^2 + 1)} \right)^2 \end{aligned}$$

and  $\sqrt{\mu(a^2 + 1)} = \frac{a(\sqrt{a^2 + 1} - a)}{2\sqrt{a^2 + 1}}$ . Note that when  $\mu < \left( \frac{a(\sqrt{a^2 + 1} - a)}{2(a^2 + 1)} \right)^2$ ,  $t(\sqrt{\mu(a^2 + 1)}, \mu) < 0$ , it implies that  $x_\mu^{**} > \sqrt{\mu(a^2 + 1)}$  and  $\left. \frac{dx}{d\mu} \right|_{x=x_\mu^{**}} > 0$ , thus decreasing  $\mu$  leads to decreasing  $x_\mu^{**}$ . We can conclude

$$\max_{\mu \leq \tau} x_\mu^{**} \leq \frac{a(\sqrt{a^2 + 1} - a)}{2\sqrt{a^2 + 1}}$$

723 Note that  $\forall \mu$  s.t.  $\left(\frac{a(\sqrt{a^2+1}-a)}{2(a^2+1)}\right)^2 < \mu < \tau$ ,  $x_\mu^{**} < \left(\frac{a(\sqrt{a^2+1}-a)}{2(a^2+1)}\right)^2$ , so  
 724  $t\left(\left(\frac{a(\sqrt{a^2+1}-a)}{2(a^2+1)}\right)^2, \mu\right) \geq 0$ .

725 Note that when  $\mu > \frac{a^2}{a^2+1}$ , i.e.  $(x_\mu^*)^2 \geq \mu(a^2+1)$  then

$$\frac{dx}{d\mu}\bigg|_{x=x_\mu^*} > 0$$

726 It implies that when  $\mu$  decreases,  $x_\mu^*$  also decreases. It holds true until  $x_\mu^* = \sqrt{\mu(a^2+1)}$ .  
 727 The same analysis can be applied to  $x_\mu^*$  like above, we can conclude that

$$\min_{\tau} x_\mu^* = \frac{a(\sqrt{a^2+1}+a)}{2\sqrt{a^2+1}}$$

728 Hence

$$\max_{\mu \leq \tau} x_\mu^{**} \leq \frac{a(\sqrt{a^2+1}-a)}{2\sqrt{a^2+1}} < \frac{a(\sqrt{a^2+1}+a)}{2\sqrt{a^2+1}} \leq \min_{\mu > 0} x_\mu^*$$

729

□

#### 730 E.4 Proof of Theorem 7,8 and 9

731 *Proof.* The proof is similar to the proof of Theorem 5 and Theorem 6.

□

#### 732 E.5 Proof of Lemma 1

*Proof.*

$$\nabla^2 g_\mu(x, y) = \begin{pmatrix} \mu + y^2 & 2xy \\ 2xy & \mu(a^2+1) + x^2 \end{pmatrix}$$

733 Let  $\lambda_1(\nabla^2 g_\mu(x, y)), \lambda_2(\nabla^2 g_\mu(x, y))$  be the eigenvalue of matrix  $\nabla^2 g_\mu(x, y)$ , then

$$\begin{aligned} & \lambda_1(\nabla^2 g_\mu(x, y)) + \lambda_2(\nabla^2 g_\mu(x, y)) \\ &= \text{Tr}(\nabla^2 g_\mu(x, y)) = \mu + y^2 + \mu(a^2+1) + x^2 > 0 \end{aligned}$$

734 Now we calculate the product of eigenvalue

$$\begin{aligned} & \lambda_1(\nabla^2 g_\mu(x, y)) \cdot \lambda_2(\nabla^2 g_\mu(W)) \\ &= \det(\nabla^2 g_\mu(W)) \\ &= (\mu + y^2)(\mu(a^2+1) + x^2) - 4x^2y^2 \\ &= \frac{\mu a}{x} \frac{\mu a}{y} - 4x^2y^2 > 0 \\ &\Leftrightarrow \left(\frac{a\mu}{2}\right)^{2/3} > xy \\ &\Leftrightarrow \left(\frac{a\mu}{2}\right)^{2/3} > \frac{a\mu}{y^2 + \mu} y \\ &\Leftrightarrow y + \frac{\mu}{y} > (4a\mu)^{1/3} \end{aligned}$$

735 Note that for  $(x_\mu^*, y_\mu^*), (x_\mu^{***}, y_\mu^{***})$ , they satisfy (11a) and (11b), this fact is used in third equality and  
 736 second “ $\Leftrightarrow$ ”. By (32b), we know  $\lambda_1(\nabla^2 g_\mu(x, y)) \cdot \lambda_2(\nabla^2 g_\mu(x, y)) > 0$  for  $(x_\mu^*, y_\mu^*), (x_\mu^{***}, y_\mu^{***})$ ,  
 737 and  $\lambda_1(\nabla^2 g_\mu(x, y)) \cdot \lambda_2(\nabla^2 g_\mu(x, y)) < 0$  for  $(x_\mu^{**}, y_\mu^{**})$ , then

$$\begin{aligned} & \lambda_1(\nabla^2 g_\mu(x, y)) > 0, \lambda_2(\nabla^2 g_\mu(x, y)) > 0 \quad \text{for } (x_\mu^*, y_\mu^*), (x_\mu^{***}, y_\mu^{***}) \\ & \lambda_1(\nabla^2 g_\mu(x, y)) < 0 \text{ or } \lambda_2(\nabla^2 g_\mu(x, y)) < 0 \quad \text{for } (x_\mu^{**}, y_\mu^{**}) \end{aligned}$$

739 and

$$\nabla g_\mu(x, y) = 0$$

740 Then  $(x_\mu^*, y_\mu^*), (x_\mu^{***}, y_\mu^{***})$  are locally minima,  $(x_\mu^{**}, y_\mu^{**})$  is saddle point for  $g_\mu(W)$ .

□

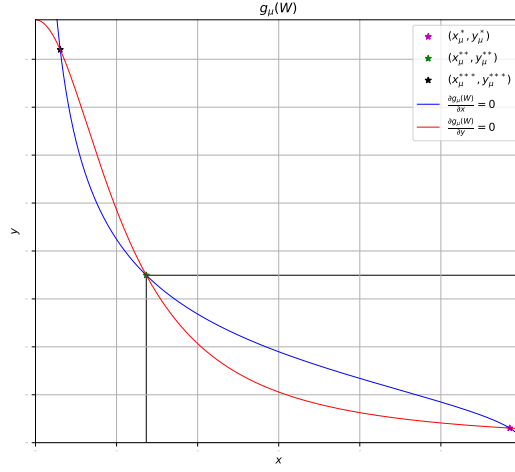


Figure 8: Stationary points when  $\mu < \tau$

*Proof.* Let us define the functions as below

$$\begin{cases} y_{\mu 1}(x) = \sqrt{\mu \left( \frac{a-x}{x} \right)} & 0 < x \leq a \\ y_{\mu 2}(x) = \frac{\mu a}{\mu(a^2 + 1) + x^2} & 0 < x \leq a \end{cases} \quad (37a)$$

$$\quad (37b)$$

$$\begin{cases} x_{\mu 1}(y) = \frac{\mu a}{y^2 + \mu} & 0 < y < \frac{a}{a^2 + 1} \\ x_{\mu 2}(y) = \sqrt{\mu \left( \frac{a}{y} - (a^2 + 1) \right)} & 0 < y < \frac{a}{a^2 + 1} \end{cases} \quad (38a)$$

$$\quad (38b)$$

742 with simple calculations,

$$y_{\mu 1} \geq y_{\mu 2} \Leftrightarrow t(x; \mu) \geq 0 \Leftrightarrow x \in (0, x_{\mu}^{***}] \cup [x_{\mu}^{**}, x_{\mu}^*]$$

743 and

$$x_{\mu 1} \geq x_{\mu 2} \Leftrightarrow r(y; \mu) \leq 0 \Leftrightarrow y \in [y_{\mu}^*, y_{\mu}^{**}] \cup [y_{\mu}^{***}, \frac{a}{a^2 + 1})$$

744 Here we divide  $B_{\mu}$  into three parts,  $C_{\mu 1}, C_{\mu 2}, C_{\mu 3}$

$$C_{\mu 1} = \{(x, y) | x_{\mu}^{**} < x \leq x_{\mu}^*, y_{\mu 1} < y < y_{\mu}^{**}\} \cup \{(x, y) | x_{\mu}^* < x \leq a, y_{\mu 2} < y < y_{\mu}^{**}\} \quad (39)$$

$$C_{\mu 2} = \{(x, y) | x_{\mu}^{**} < x \leq x_{\mu}^*, 0 \leq y < y_{\mu 2}\} \cup \{(x, y) | x_{\mu}^* < x \leq a, 0 \leq y < y_{\mu 1}\} \quad (40)$$

$$C_{\mu 3} = \{(x, y) | x_{\mu}^{**} < x \leq x_{\mu}^*, y_{\mu 2} \leq y \leq y_{\mu 1}\} \cup \{(x, y) | x_{\mu}^* < x \leq a, y_{\mu 1} \leq y \leq y_{\mu 2}\} \quad (41)$$

745 Also note that

$$\begin{aligned} \forall (x, y) \in C_{\mu 1} &\Rightarrow \frac{\partial g_{\mu}(x, y)}{\partial x} > 0, \frac{\partial g_{\mu}(x, y)}{\partial y} > 0 \\ \forall (x, y) \in C_{\mu 2} &\Rightarrow \frac{\partial g_{\mu}(x, y)}{\partial x} < 0, \frac{\partial g_{\mu}(x, y)}{\partial y} < 0 \end{aligned}$$

746 The gradient flow follows

$$\begin{pmatrix} x'(t) \\ y'(t) \end{pmatrix} = - \begin{pmatrix} \frac{\partial g_{\mu}(x(t), y(t))}{\partial x} \\ \frac{\partial g_{\mu}(x(t), y(t))}{\partial y} \end{pmatrix} = -\nabla g_{\mu}(x(t), y(t))$$



747 then

$$\forall (x, y) \in C_{\mu 1} \Rightarrow \begin{pmatrix} x'(t) \\ y'(t) \end{pmatrix} < 0, \quad \|\nabla g_{\mu}\| > 0 \quad (42)$$

$$\forall (x, y) \in C_{\mu 2} \Rightarrow \begin{pmatrix} x'(t) \\ y'(t) \end{pmatrix} > 0, \quad \|\nabla g_{\mu}\| > 0 \quad (43)$$

748 Note that  $\|\nabla g_{\mu}\|$  is not diminishing and bounded away from 0. Let us consider the  $(x(0), y(0)) \in$   
 749  $C_{\mu 1}$ , since  $\nabla g_{\mu}(x, y) \neq 0$ ,  $-\nabla g_{\mu}(x, y) < 0$  in (42) and boundness of  $C_{\mu 1}$ , it implies there exists a  
 750 finite  $t_0 > 0$  such that

$$(x(t_0), y(t_0)) \in \partial C_{\mu 1}, (x(t), y(t)) \in C_{\mu 1} \text{ for } 0 \leq t < t_0$$

751 where  $\partial C_{\mu 1}$  is defined as

$$\partial C_{\mu 1} = \{(x, y) | x_{\mu}^{**} < x \leq x_{\mu}^*, y = y_{\mu 1}\} \cup \{(x, y) | x_{\mu}^* < x \leq a, y = y_{\mu 2}\} \subseteq C_{\mu 3}$$

752 For the same reason, if  $(x(0), y(0)) \in C_{\mu 2}$ , there exists a finite time  $t_1 > 0$ ,

$$(x(t_0), y(t_0)) \in \partial C_{\mu 2}, (x(t), y(t)) \in C_{\mu 2} \text{ for } 0 \leq t < t_1$$

753 where  $\partial C_{\mu 2}$  is defined as

$$\partial C_{\mu 2} = \{(x, y) | x_{\mu}^{**} < x \leq x_{\mu}^*, y = y_{\mu 2}\} \cup \{(x, y) | x_{\mu}^* < x \leq a, y = y_{\mu 1}\} \subseteq C_{\mu 3}$$

754 then by lemma 7,  $\lim_{t \rightarrow \infty} (x(t), y(t)) = (x_{\mu}^*, y_{\mu}^*)$ . □

### 755 E.7 Proof of Lemma 3

756 *Proof.* This is just a result of the Theorem 5. □

### 757 E.8 Proof of Lemma 5

758 *Proof.* Note that

$$\nabla^2 g_{\mu}(W) = \begin{pmatrix} \mu + y^2 & 2xy \\ 2xy & \mu(a^2 + 1) + x^2 \end{pmatrix} = \begin{pmatrix} \mu & 0 \\ 0 & \mu(a^2 + 1) \end{pmatrix} + \begin{pmatrix} y^2 & 2xy \\ 2xy & x^2 \end{pmatrix}$$

759 Let  $\|\cdot\|_{\text{op}}$  is the spectral norm, and it satisfies triangle inequality

$$\begin{aligned} \|\nabla^2 g_{\mu}(W)\|_{\text{op}} &\leq \left\| \begin{pmatrix} \mu & 0 \\ 0 & \mu(a^2 + 1) \end{pmatrix} \right\|_{\text{op}} + \left\| \begin{pmatrix} y^2 & 2xy \\ 2xy & x^2 \end{pmatrix} \right\|_{\text{op}} \\ &= \mu(a^2 + 1) + \left\| \begin{pmatrix} y^2 & 2xy \\ 2xy & x^2 \end{pmatrix} \right\|_{\text{op}} \end{aligned}$$

760 The spectral norm of the second term in area A is bounded by

$$\max_{(x, y) \in A} \frac{(x^2 + y^2) + \sqrt{(x^2 + y^2)^2 + 12x^2y^2}}{2} \leq \frac{2a^2 + \sqrt{4a^4 + 12a^4}}{2} = 3a^2$$

761 We use  $x^2 \leq a^2, y^2 \leq a^2$  in the inequality. Therefore,

$$\|\nabla^2 g_{\mu}(W)\|_{\text{op}} \leq 3a^2 + \mu(a^2 + 1)$$

762 Also, according to [5, 33], for any  $f$ , if  $\nabla^2 f$  exists, then  $f$  is  $L$  smooth if and only if  $|\nabla^2 f|_{\text{op}} \leq L$ .

763 With this, we conclude the proof. □

### 764 E.9 Proof of Lemma 7

765 *Proof.* First we prove  $\forall t \geq 0, (x(t), y(t)) \in C_{\mu 3}$ , because if  $(x(t), y(t)) \notin C_{\mu 3}$ , then there exists a  
 766 finite  $t$  such that

$$(x(t), y(t)) \in \partial C_{\mu 3}$$

767 where  $\partial C_{\mu 3}$  is the boundary of  $C_{\mu 3}$ , defined as

$$\partial C_{\mu 3} = \{(x, y) | y = y_{\mu 1}(x) \text{ or } y = y_{\mu 2}(x), x_{\mu}^{**} < x \leq a\}$$

768 W.L.O.G, let us assume  $(x(0), y(0)) \in \partial C_{\mu 3}$  and  $(x(0), y(0)) \neq (x_\mu^*, y_\mu^*)$ . Here are four different  
769 cases,

$$\nabla g_\mu(x(t), y(t)) = \begin{cases} \begin{pmatrix} = 0 \\ > 0 \end{pmatrix} & \text{if } y(0) = y_{\mu 1}(x(0)), x_\mu^{**} < x(0) < x_\mu^* \\ \begin{pmatrix} = 0 \\ < 0 \end{pmatrix} & \text{if } y(0) = y_{\mu 1}(x(0)), x_\mu^* < x(0) \leq a \\ \begin{pmatrix} < 0 \\ = 0 \end{pmatrix} & \text{if } y(0) = y_{\mu 2}(x(0)), x_\mu^{**} < x(0) < x_\mu^* \\ \begin{pmatrix} > 0 \\ = 0 \end{pmatrix} & \text{if } y(0) = y_{\mu 2}(x(0)), x_\mu^* < x(0) \leq a \end{cases}$$

770 This indicates that  $-\nabla g_\mu(x(t), y(t))$  are pointing to the interior of  $C_{\mu 3}$ , then  $(x(t), y(t))$  can not  
771 escape  $C_{\mu 3}$ . Here we can focus our attention in  $C_{\mu 3}$ , because  $\forall t \geq 0, (x(t), y(t)) \in C_{\mu 3}$ . For  
772 Algorithm 1,

$$\frac{df(z_t)}{dt} = \nabla f(z_t) \dot{z}_t = -\|\nabla f(z_t)\|_2^2$$

773 In our setting,  $\forall (x, y) \in C_{\mu 3}$

$$\begin{cases} \nabla g_\mu(x, y) \neq 0 & (x, y) \neq (x_\mu^*, y_\mu^*) \\ \nabla g_\mu(x, y) = 0 & (x, y) = (x_\mu^*, y_\mu^*) \end{cases}$$

774 so

$$\frac{dg_\mu(x(t), y(t))}{dt} = \begin{cases} -\|\nabla g_\mu\|_2^2 < 0 & (x, y) \neq (x_\mu^*, y_\mu^*) \\ -\|\nabla g_\mu\|_2^2 = 0 & (x, y) = (x_\mu^*, y_\mu^*) \end{cases}$$

775 Plus,  $(x_\mu^*, y_\mu^*)$  is the unique stationary point of  $g_\mu(W)$  in  $C_{\mu 3}$ . By lemma 8

$$g_\mu(x, y) > g_\mu(x_\mu^*, y_\mu^*) \quad (x, y) \neq (x_\mu^*, y_\mu^*)$$

776 By Lyapunov asymptotic stability theorem [28], and applying it to gradient flow for  $g_\mu(x, y)$  in  $C_{\mu 3}$ ,  
777 we can conclude  $\lim_{t \rightarrow \infty} (x(t), y(t)) = (x_\mu^*, y_\mu^*)$ .  $\square$

## 778 E.10 Proof of Lemma 8

779 *Proof.* For any  $(x, y) \in C_{\mu 3}$  in 41, and  $(x, y) \neq (x_\mu^*, y_\mu^*)$ , in Algorithm 7. W.L.O.G, we can assume  
780  $x \in (x_\mu^{**}, x_\mu^*)$ , the analysis details can also be applied to  $x \in (x_\mu^*, a)$ . It is obvious that  $\tilde{x}_j < \tilde{x}_{j+1}$   
781 and  $\tilde{y}_{j+1} < \tilde{y}_j$ . Also,  $\lim_{j \rightarrow \infty} (\tilde{x}_j, \tilde{y}_j) = (x_\mu^*, y_\mu^*)$ . Otherwise either  $\tilde{x}_j \neq x_\mu^*$  or  $\tilde{y}_j \neq y_\mu^*$  hold,  
782 Algorithm 7 continues until  $\lim_{j \rightarrow \infty} (\tilde{x}_j, \tilde{y}_j) = \lim_{j \rightarrow \infty} (y_{\mu 2}(\tilde{y}_j), x_{\mu 1}(\tilde{x}_j))$ , i.e.  $(\tilde{x}_j, \tilde{y}_j)$  converges  
783 to  $(x_\mu^*, y_\mu^*)$ .

784 Moreover, note that for any  $j = 0, 1, \dots$

$$g_\mu(\tilde{x}_{j-1}, \tilde{y}_{j-1}) > g_\mu(\tilde{x}_{j-1}, \tilde{y}_j) > g_\mu(\tilde{x}_j, \tilde{y}_j)$$

785 Because

$$g_\mu(\tilde{x}_{j-1}, \tilde{y}_{j-1}) - g_\mu(\tilde{x}_{j-1}, \tilde{y}_j) = \frac{\partial g_\mu(\tilde{x}_{j-1}, \tilde{y})}{\partial y}(\tilde{y}_{j-1} - \tilde{y}_j) \quad \text{where } \tilde{y} \in (\tilde{y}_j, \tilde{y}_{j-1})$$

786 Note that

$$\frac{\partial g_\mu(\tilde{x}_{j-1}, \tilde{y})}{\partial y} > 0 \Rightarrow g_\mu(\tilde{x}_{j-1}, \tilde{y}_{j-1}) > g_\mu(\tilde{x}_{j-1}, \tilde{y}_j)$$

787 By the same reason,

$$g_\mu(\tilde{x}_{j-1}, \tilde{y}_j) > g_\mu(\tilde{x}_j, \tilde{y}_j)$$

788 By Lemma 1,  $(x_\mu^*, y_\mu^*)$  is local minima, and there exists a  $r_\mu > 0$  and any  $\{(x, y) \mid \|(x, y) -$   
789  $(x_\mu^*, y_\mu^*)\|_2 \leq r_\mu\}$ ,  $g_\mu(x, y) > g_\mu(x_\mu^*, y_\mu^*)$ . Since  $\lim_{j \rightarrow \infty} (\tilde{x}_j, \tilde{y}_j) = (x_\mu^*, y_\mu^*)$ , there exists a  $J > 0$   
790 such that  $\forall j > J$ ,  $\|(\tilde{x}_j, \tilde{y}_j) - (x_\mu^*, y_\mu^*)\|_2 \leq r_\mu$ , combining them all

$$g_\mu(x, y) > g_\mu(\tilde{x}_j, \tilde{y}_j) > g_\mu(x_\mu^*, y_\mu^*)$$

791

792  $\square$

---

**Algorithm 7:** Path goes to  $(x_\mu^*, y_\mu^*)$ 

---

**Input:**  $(x, y) \in C_{\mu 3}, x_{\mu 1}(y), y_{\mu 2}(x)$  as (38a),(37b)**Output:**  $\{(\tilde{x}_j, \tilde{y}_j)\}_{j=0}^\infty$ 

```
1  $(\tilde{x}_0, \tilde{y}_0) \leftarrow (x, y)$ 
2 for  $j = 1, 2, \dots$  do
3    $\tilde{y}_j \leftarrow y_{\mu 2}(\tilde{x}_{j-1})$ 
4    $\tilde{x}_j \leftarrow x_{\mu 1}(\tilde{y}_{j-1})$ 
5 end
```

---

**E.11 Proof of Lemma 4**

*Proof.* From the proof of Theorem 1, any any scheduling for  $\mu_k$  satisfies following will do the job

$$(2/a)^{2/3} \mu_{k-1}^{4/3} \leq \mu_k < \mu_{k-1}$$

Note that in Algorithm 4, we have  $\hat{a} = \sqrt{4(\mu_0 + \varepsilon)} < a$ , then it is obvious

$$(2/a)^{2/3} \mu_{k-1}^{4/3} < (2/\hat{a})^{2/3} \mu_{k-1}^{4/3}$$

The same analysis for Theorem 1 can be applied here.  $\square$

**E.12 Proof of Lemma 6**

*Proof.* By the Theorem 3 and Lemma 5 and the fact that  $A_{\mu, \epsilon}^1$  is  $\mu$ -stationary point region, we use the same argument as proof of Lemma 7 to demonstrate the gradient descent will never go to  $A_{\mu, \epsilon}^2$ .  $\square$

**E.13 Proof of Lemma 9**

*Proof.* By Theorem 9(iv)

$$\max_{\mu \leq \tau_\beta} x_{\mu, \beta}^{**} \leq \min_{\mu > 0} x_{\mu, \beta}^*$$

We also know from the proof of Corollary 3,  $x_{\mu, \epsilon}^{**} < x_{\mu, \beta}^{**}$  and  $x_{\mu, \beta}^* < x_{\mu, \epsilon}^*$ . Consequently,

$$\max_{\mu \leq \tau_\beta} x_{\mu, \epsilon}^{**} \leq \min_{\mu > 0} x_{\mu, \epsilon}^*$$

Because  $\tau_\beta > \tau$ , so

$$\max_{\mu \leq \tau} x_{\mu, \epsilon}^{**} \leq \max_{\mu \leq \tau_\beta} x_{\mu, \epsilon}^{**} \leq \min_{\mu > 0} x_{\mu, \epsilon}^*$$

$\square$

**E.14 Proof of Corollary 1**

*Proof.* Note that

$$\frac{a^2}{4(a^2 + 1)^3} \leq \frac{1}{27} \quad a > 0$$

when  $a > \sqrt{\frac{5}{27}}$ , then  $\frac{a^2}{4} > \mu_0 = \frac{1}{27} \geq \frac{a^2}{4(a^2 + 1)^3}$ , it satisfies condition in Lemma 4, we obtain the same result.  $\square$

**E.15 Proof of Corollary 2**

*Proof.* Use Theorem 5(vi) and Theorem 6(vi).  $\square$

## 811 E.16 Proof of Corollary 3

812 *Proof.* It is easy to know that

$$r_\beta(y; \mu) > r_\epsilon(y; \mu) > r(y; \mu)$$

813 and

$$t_\beta(x; \mu) < t_\epsilon(x; \mu) < t(x; \mu)$$

814 and when  $\mu < \tau$ , there are three solutions to  $r(y; \mu) = 0$  by Theorem 5. Also, we know from  
815 Theorem 7, 8

$$\lim_{y \rightarrow 0^+} r_\epsilon(y; \mu) = \infty \quad \lim_{y \rightarrow 0^+} r_\beta(y; \mu) = \infty$$

816 Note that when  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq a^2 + 1$

$$r_\beta(\sqrt{\mu}; \mu) = \frac{a(1+\beta)}{\sqrt{\mu}} - (a^2 + 1) - \frac{a^2(1-\beta)^2}{4\mu} \leq 0 \quad \forall \mu > 0$$

817 Therefore,

$$0 \geq r_\beta(\sqrt{\mu}; \mu) > r_\epsilon(\sqrt{\mu}; \mu) > r(\sqrt{\mu}; \mu)$$

818 Also, we know that for  $y_{\text{ub}}$  defined in Theorem 5(iii), we know  $r(y_{\text{ub}}; \mu) > 0$  from Theorem 5(iv).  
819 Therefore,

$$r_\beta(y_{\text{ub}}; \mu) > r_\epsilon(y_{\text{ub}}; \mu) > r(y_{\text{ub}}; \mu) > 0$$

820 Besides,  $\sqrt{\mu} < y_{\text{ub}}$ . By monotonicity of  $r_\beta(y; \mu)$  and  $r_\epsilon(y; \mu)$  from the Theorem 7(ii) and Theorem  
821 8(ii), it implies that there are at least two solutions to  $r_\beta(y; \mu)$  and  $r_\epsilon(y; \mu)$ . From the geometry  
822 of  $r_\beta(y; \mu), r_\epsilon(y; \mu), r(y; \mu)$  and  $t_\beta(x; \mu), t_\epsilon(x; \mu), t(x; \mu)$ , it is trivial to know that  $x_{\mu, \epsilon}^* \leq x_\mu^*$ ,  
823  $y_{\mu, \epsilon}^* \geq y_\mu^*, x_{\mu, \epsilon}^{**} \geq x_\mu^{**}, y_{\mu, \epsilon}^* \leq y_\mu^{**}$ .

824 Finally, for every point  $(x, y) \in A_{\mu, \epsilon}^1$ , there exists a pair  $\epsilon_1, \epsilon_2$ , each satisfying  $|\epsilon_1| \leq \epsilon$  and  $|\epsilon_2| \leq \epsilon$ ,  
825 such that  $(x, y)$  is the solution to

$$x = \frac{\mu a + \epsilon_1}{\mu + y^2} \quad y = \frac{\mu a + \epsilon_2}{x^2 + \mu(a^2 + 1)}$$

826 We can repeat the same analysis above to show that  $x_{\mu, \epsilon}^* \leq x, y_{\mu, \epsilon}^* \geq y$ . Applying the same logic  
827 to  $\forall (x, y) \in A_{\mu, \epsilon}^2$ , we find  $x_{\mu, \epsilon}^{**} \geq x, y_{\mu, \epsilon}^* \leq y$ . Thus,  $(x_\mu^*, y_\mu^*)$  is the extreme point of  $A_{\mu, \epsilon}^1$  and  
828  $(x_\mu^{**}, y_\mu^{**})$  is the extreme point of  $A_{\mu, \epsilon}^2$ , we get the results.  $\square$

## 829 F Experiments Details

830 In this section, we present experiments to validate the global convergence of Algorithm 6. Our  
831 goal is twofold: First, we aim to demonstrate that irrespective of the starting point, Algorithm 6  
832 using gradient descent consistently returns the global minimum. Second, we contrast our updating  
833 scheme for  $\mu_k, \epsilon_k$  as prescribed in Algorithm 6 with an arbitrary updating scheme for  $\mu_k, \epsilon_k$ . This  
834 comparison illustrates how inappropriate setting of parameters in gradient descent could lead to  
835 incorrect solutions.

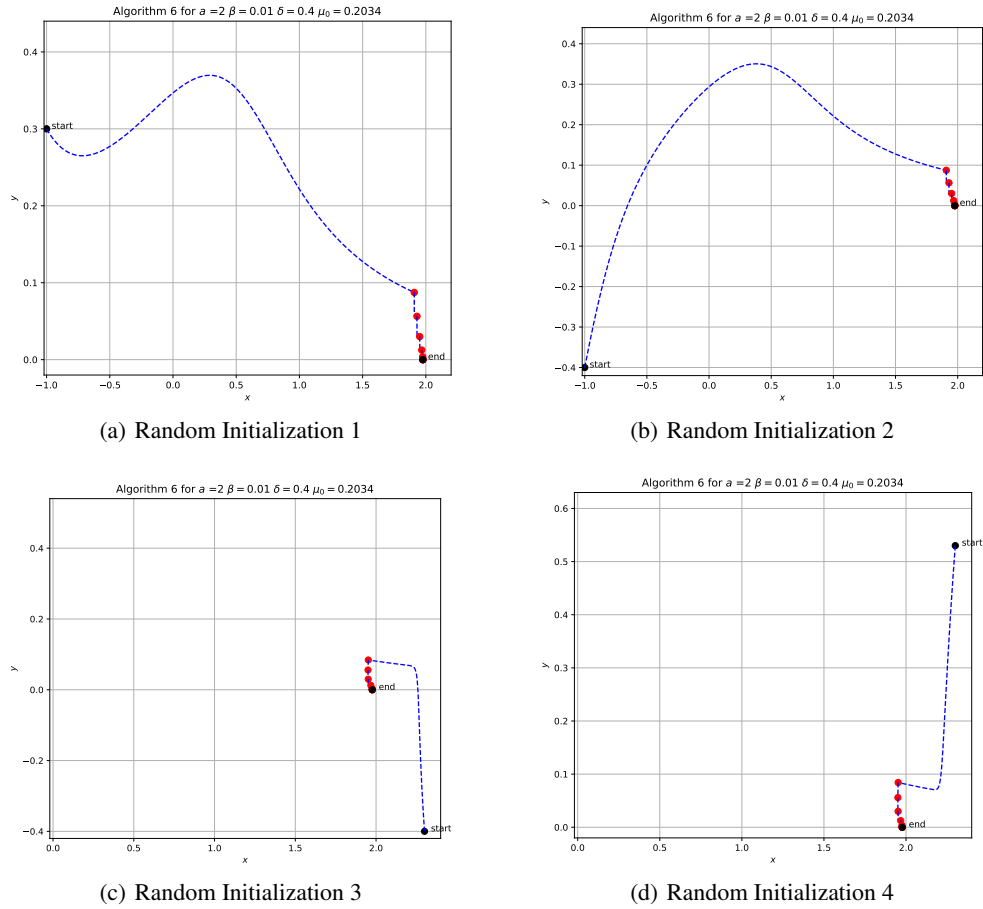
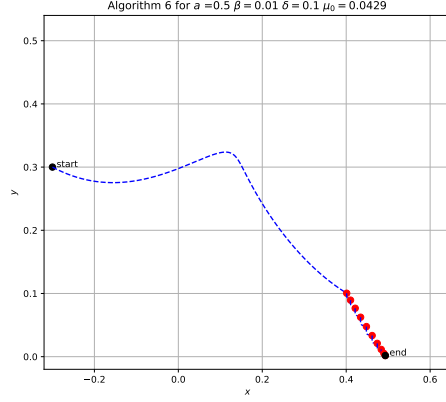
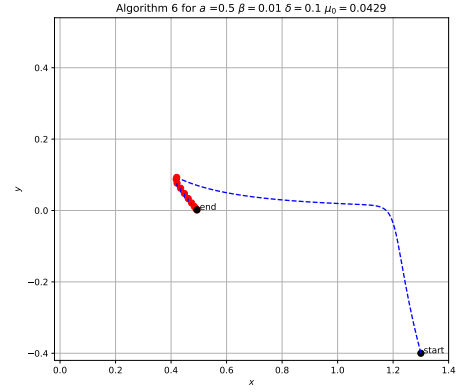


Figure 9: Trajectory of the gradient descent path with the different initializations for  $a = 2$ . We observe that regardless of the initialization, Algorithm 6 always converges to the global minimum.

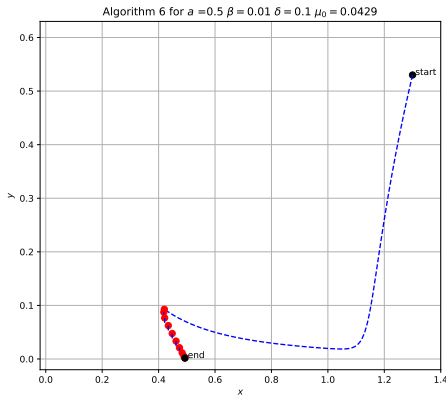
$$\text{Initial } \mu_0 = \frac{a^2 (1-\delta)^3 (1-\beta)^4}{4 (1+\beta)^2}$$



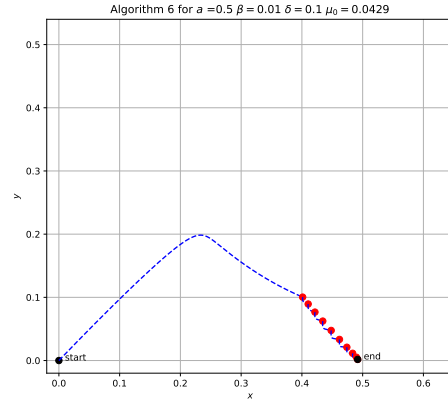
(a) Random Initialization 1



(b) Random Initialization 2



(c) Random Initialization 3



(d) Random Initialization 4

Figure 10: Trajectory of the gradient descent path with the different initializations for  $a = 0.5$ . We observe that regardless of the initialization, Algorithm 6 always converges to the global minimum.

$$\text{Initial } \mu_0 = \frac{a^2}{4} \frac{(1-\delta)^3(1-\beta)^4}{(1+\beta)^2}$$

837 **F.2 Wrong Specification of  $\delta$  Leads to Spurious Local Optimal**

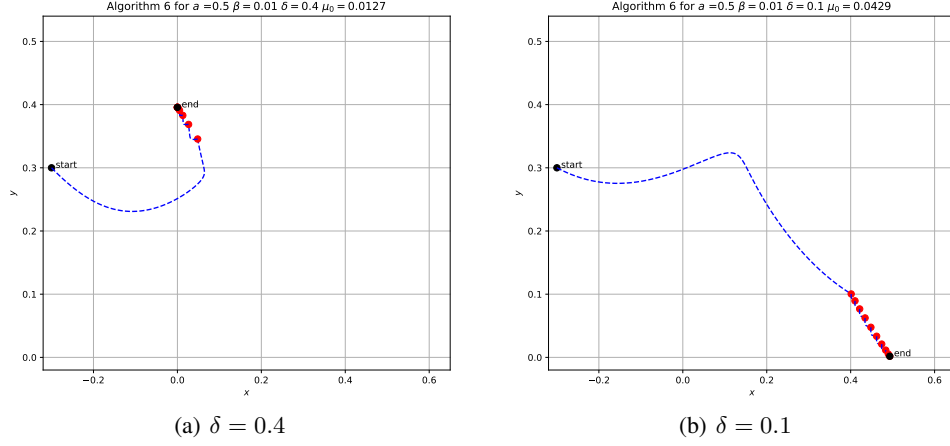


Figure 11: Trajectory of the gradient descent path for two difference  $\delta$ . Left:  $\beta$  violates requirement  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq (1-\delta)(a^2+1)$  in Theorem 4, leading to spurious local minimum. Right:  $\beta$  follows requirement  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq (1-\delta)(a^2+1)$  in Theorem 4, leading to global minimum. Initial  $\mu_0 = \frac{a^2}{4} \frac{(1-\delta)^3(1-\beta)^4}{(1+\beta)^2}$

838 **F.3 Wrong Specification of  $\beta$  Leads to Incorrect Solution**

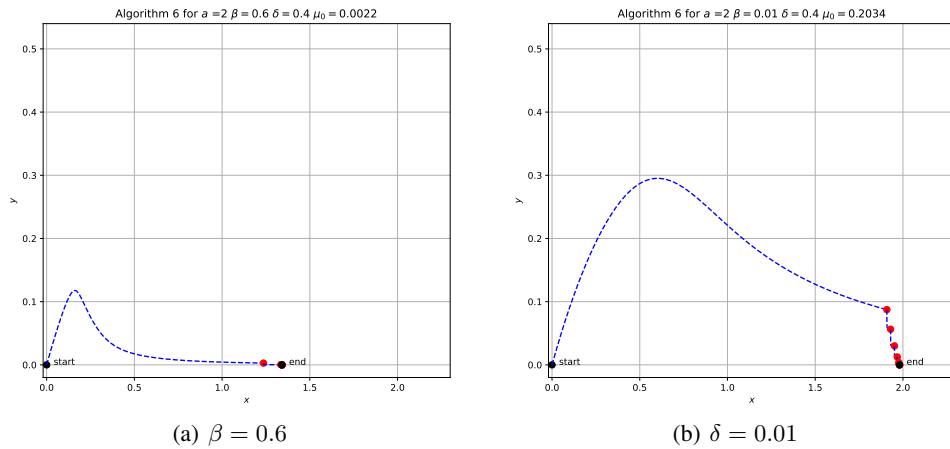


Figure 12: Trajectory of the gradient descent path for two difference  $\beta$ . Left:  $\beta$  violates requirement  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq (1-\delta)(a^2+1)$  in Theorem 4, leading to incorrect solution. Right:  $\beta$  follows requirement  $\left(\frac{1+\beta}{1-\beta}\right)^2 \leq (1-\delta)(a^2+1)$  in Theorem 4, leading to global minimum. Initial  $\mu_0 = \frac{a^2}{4} \frac{(1-\delta)^3(1-\beta)^4}{(1+\beta)^2}$

839 **F.4 Faster decrease of  $\mu_k$  Leads to Incorrect Solution**

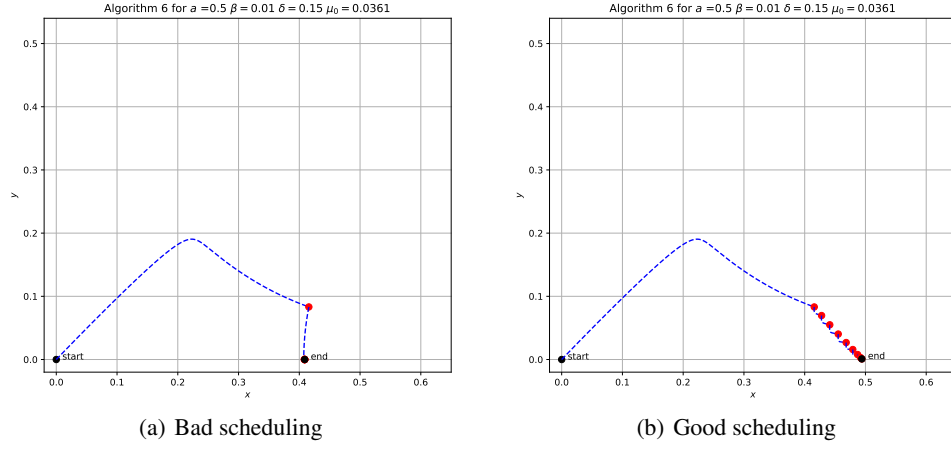


Figure 13: Trajectory of the gradient descent path for two difference update rules for  $\mu_k$  with the same initialization. Left: “Bad scheduling” uses a faster-decreasing scheme for  $\mu_k$ , leading to an incorrect solution, even a non-local optimal solution. Right: “Good scheduling” follows updating rule for  $\mu_k$  in Algorithm 6, leading to the global minimum. Initial  $\mu_0 = \frac{a^2}{4} \frac{(1-\delta)^3(1-\beta)^4}{(1+\beta)^2}$