

Analysis of Gradient Descent with Varying Step Sizes using Integral Quadratic Constraints

Ram Padmanabhan* and Peter Seiler†

Abstract

The framework of Integral Quadratic Constraints (IQCs) is used to perform an analysis of gradient descent with varying step sizes. Two performance metrics are considered: convergence rate and noise amplification. We assume that the step size is produced from a line search and varies in a known interval. Modeling the algorithm as a linear, parameter-varying (LPV) system, we construct a parameterized linear matrix inequality (LMI) condition that certifies algorithm performance, which is solved using a result for polytopic LPV systems. Our results provide convergence rate guarantees when the step size lies within a restricted interval. Moreover, we recover existing rate bounds when this interval reduces to a single point, i.e. a constant step size. Finally, we note that the convergence rate depends only on the condition number of the problem. In contrast, the noise amplification performance depends on the individual values of the strong convexity and smoothness parameters, and varies inversely with them for a fixed condition number.

1 Introduction

Convex optimization problems can be solved using numerous iterative algorithms, including gradient descent [1, 2], accelerated [3, 4] and proximal gradient methods [5, 6]. The use of control theoretic methods to analyze such algorithms has a rich history. Our focus is on using the notion of Integral Quadratic Constraints (IQCs) from robust control theory. IQCs were first introduced by Yakubovich [7] in the context of imposing quadratic constraints on an infinite-horizon control problem in a Lur’e system, and imposing multiple constraints via the S-procedure [8]. Megretski and Rantzer [9] unified the approach to robustness analysis of such systems using IQCs. Using the Kalman-Yakubovich-Popov (KYP) lemma [10], it was shown that verifying stability reduces to a linear matrix inequality (LMI) condition. Additional details can be found in [11, 12] with discussions connecting time- and frequency-domain IQCs [13] and discrete-time IQCs [14–16].

In [17], Drori and Teboulle developed a semidefinite program (SDP) approach to certify tight bounds for first-order optimization algorithms on strongly convex problems. Subsequent work by Lessard *et al.* [16] developed a unified framework using IQCs for the analysis of first-order algorithms on strongly convex functions. Importantly, while the approach in [17] scaled with the problem dimension, the use of IQCs circumvented this issue in [16]. A class of ρ -hard IQCs was introduced for characterizing convergence rates. Numerous IQCs for the gradient of strongly convex functions were presented, from simple sector bounds [4] to dynamic constraints using Zames-Falb IQCs [18]. Using these IQCs, certifying the convergence of first-order algorithms reduces to solving a small SDP independent of problem dimension. This framework has led to a significant body of subsequent work using IQCs to analyze and design optimization algorithms. This includes the analysis of the heavy-ball method [19], the biased stochastic gradient method [20], transient behavior of Nesterov’s accelerated method [21], analysis of non-strongly convex problems [22] and algorithm design [23–25]. A set of related case studies is available in [26].

Another important aspect of analyzing optimization algorithms is their robustness to noise, either in the iterate or in the gradient. Mohammadi *et al.* [27, 28] examine the variance in the iterate error when iterates are perturbed by additive white noise, for both gradient descent and Nesterov’s accelerated method. Based

*Department of Electrical and Computer Engineering and Coordinated Science Laboratory, University of Illinois Urbana-Champaign, Urbana, IL 61801, USA. Email: ramp3@illinois.edu

†Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI 48109, USA. Email: pseiler@umich.edu

on this, a set of tradeoffs between noise amplification and convergence are derived in [29] for strongly convex quadratic problems.

It must be noted that almost all prior work in this area assumes constant step sizes in the optimization algorithms. However, the use of constant step sizes usually requires additional knowledge of the problem, such as the Lipschitz constant L of the gradient. Varying step sizes, often based on a line search [30, 31] are commonly implemented in optimization algorithms. While these algorithms are computationally more intensive, they do not require any prior knowledge about function parameters and do not result in a reduction in performance. In this article, we focus on the analysis of gradient descent with a varying step size based on the framework developed in [16]. We assume a line search has been carried out, and produces a step size that lies in a restricted interval about $\alpha = \frac{1}{L}$, a popular choice for gradient descent on strongly convex functions. While constant step size algorithms can be represented as linear, time-invariant (LTI) systems, a varying step size indicates that the algorithm is now a linear, parameter-varying (LPV) system, where the step size is a parameter affecting the dynamics. IQCs have been used for the analysis of LPV systems, including obtaining their worst-case \mathcal{L}_2 gain [32, 33] and developing a robust synthesis algorithm [34].

In our work, we combine the framework developed in [16] with the LPV analysis approaches discussed above, for analyzing gradient descent with a varying step size. Using a result for polytopic LPV systems [35], we obtain convergence rates for the parameter-varying algorithm. In addition to the analysis of convergence rate, we use the noise amplification framework developed in [28] for constant step sizes, for our parameter-varying algorithm affected by gradient noise.

The remainder of this article is organized as follows. In Section 2, we present the two problems we examine, on convergence rate and noise amplification, as well as some background on IQCs. Our approach is described in Section 3 based on the characterization of the algorithm as an LPV system. This includes our main theorems, which extend prior constant step size results [16, 27, 28]. We also discuss how the implementation of the results can be simplified by tweaking the existing LMIs. Section 4 presents our main results. First, we show that our approach recovers known constant step size expressions when the step size interval reduces to a single point. Further, we derive analytical expressions that characterize convergence and noise amplification for varying step size gradient descent based on our formulation. Finally, numerical results that are obtained by solving the LMIs are discussed. Section 5 provides concluding remarks.

2 Preliminaries

2.1 Notation

Throughout this article, I_n denotes the $n \times n$ identity matrix and 0_n denotes the $n \times n$ matrix of zeros. \mathbb{Z}_+ denotes the set of all non-negative integers, and \mathbb{R}^n denotes the vector space of all n -tuples of real numbers. For a vector $x = [x_1, \dots, x_n]^T \in \mathbb{R}^n$, $\|x\|$ denotes the Euclidean norm. $A \otimes B$ denotes the Kronecker product of two matrices A and B . The Kronecker Delta function, denoted δ_j is defined as $\delta_j = 1$ if $j = 0$, and $\delta_j = 0$ otherwise.

2.2 Problem Formulation

Consider the following unconstrained optimization problem:

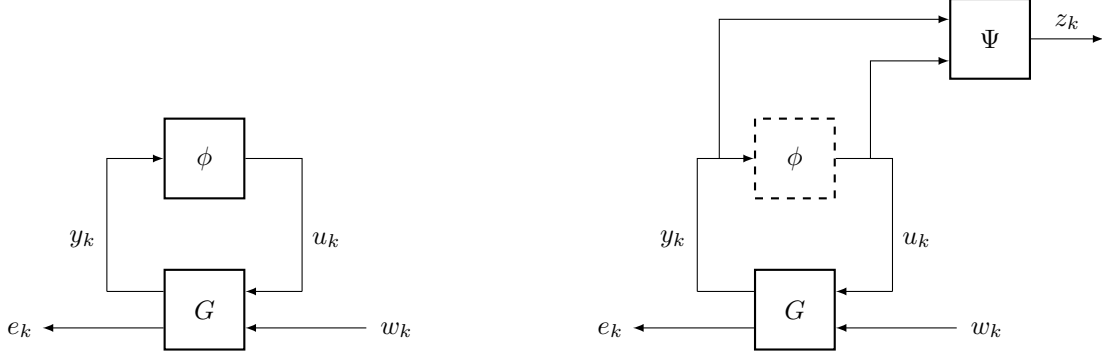
$$\min_{x \in \mathbb{R}^n} f(x), \quad (1)$$

where $f : \mathbb{R}^n \rightarrow \mathbb{R}$. Let $\nabla f(x)$ denote the gradient of f at x . We assume that f is m -strongly convex and ∇f is L -Lipschitz continuous, i.e., for all $x, y \in \mathbb{R}^n$,

$$f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{m}{2} \|y - x\|^2 \quad (2)$$

$$\text{and} \quad \|\nabla f(y) - \nabla f(x)\| \leq L \|y - x\|. \quad (3)$$

The condition number of the function f is $\kappa = \frac{L}{m}$. The class of functions f that satisfy (2) and (3) is denoted $\mathcal{S}(m, L)$. If $f \in \mathcal{S}(m, L)$ there exists a unique solution x^* to (1) [1, Chapter 4]. We assume, without loss of generality by a coordinate shift, that this optimal point occurs at the origin, i.e. $x^* = 0$.



(a) A linear system G in feedback with a nonlinear component ϕ , affected by white stochastic noise w_k . In gradient descent, $\phi = \nabla f$.

(b) Representation of the IQC framework. $z_k = \Psi(y_k, u_k)$, where Ψ is a stable LTI system in the most general case.

Figure 1: ϕ is the nonlinear component we wish to analyze, and is replaced by the constraints it imposes on the input-output pair (u, y) . These are written as constraints on z_k .

In this article, we consider the use of gradient descent with a time-varying step size α_k as an iterative approach to solve the problem (1). Initialized at some $x_0 \in \mathbb{R}^n$, gradient descent is characterized by the following update rule:

$$x_{k+1} = x_k - \alpha_k (\nabla f(x_k) + w_k). \quad (4)$$

Here, w_k is a perturbation used to model a noisy gradient. For the noise-free gradient method, $w_k = 0$ for all $k \in \mathbb{Z}_+$.

We examine two fundamental performance metrics associated with the use of gradient descent: (i) convergence rate and (ii) noise amplification. These metrics are defined as follows:

Definition 1 (Convergence Rate). The noise-free gradient method converges with a rate $\rho \in (0, 1)$ if there exists a constant $\beta > 0$ such that:

$$\|x_k\| \leq \beta \rho^k \|x_0\| \quad (5)$$

for all $x_0 \in \mathbb{R}^n$ and for all $k \in \mathbb{Z}_+$.

Definition 2 (Noise Amplification). Assume w_k is additive white stochastic noise with zero mean and identity covariance matrix, i.e. $\mathbb{E}[w_k] = 0$ and $\mathbb{E}[w_k w_l^T] = I \delta_{k-l}$. The noise amplification of the algorithm (4) is characterized by the metric γ , where:

$$\gamma \triangleq \left(\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=0}^N \mathbb{E} [\|x_k\|^2] \right)^{\frac{1}{2}} = \left(\lim_{N \rightarrow \infty} \mathbb{E} [\|x_N\|^2] \right)^{\frac{1}{2}}. \quad (6)$$

γ^2 is the steady-state variance of the iterate x_N , which is also the steady-state variance of the iterate error $x_N - x^*$ as the optimal solution is $x^* = 0$. The limit above exists when the update (4) is stable.

Our objective is to use the framework of Integral Quadratic Constraints (IQCs) to analyze these two performance metrics for gradient descent with a varying step size. We assume that α_k is a varying step size produced by some line search algorithm. Then, the two analysis problems we consider can be stated as follows:

Problem 1 (Convergence Rate). For the noise-free gradient method, find the smallest possible convergence rate ρ such that (5) is satisfied.

Problem 2 (Noise Amplification). For the gradient method affected by gradient noise (4), find the smallest possible bound on the steady-state iterate variance γ^2 , and subsequently the metric γ defined in (6).

2.3 Integral Quadratic Constraints

The use of Integral Quadratic Constraints (IQCs) is motivated by the fact that eq. (4) can be written as a linear update rule separated from the gradient, which is a nonlinear component:

$$\begin{aligned} x_{k+1} &= x_k - \alpha_k u_k + \alpha_k w_k, \\ y_k &= x_k, \\ u_k &= \nabla f(y_k). \end{aligned} \tag{7}$$

A general version of this configuration is shown in Fig. 1(a), where a linear system G is in feedback with a static, memoryless nonlinear function ϕ such that $u = \phi(y)$. The dynamics of the linear system can, in general, be represented as:

$$\begin{aligned} x_{k+1} &= A(\alpha_k)x_k + B_u(\alpha_k)u_k + B_w(\alpha_k)w_k, \\ \begin{bmatrix} y_k \\ e_k \end{bmatrix} &= \begin{bmatrix} C_y(\alpha_k) \\ C_e(\alpha_k) \end{bmatrix} x_k. \end{aligned} \tag{8}$$

In this configuration, e_k is a performance output. While G is linear, it may be time-varying or parameter-varying in which case the matrices A , B_u , B_w , C_y and C_e depend on the time index k or a parameter α_k .

Analyzing the general interconnection in Fig. 1(a) is not straightforward due to the nonlinearity ϕ . The IQC framework provides a convenient method to analyze this interconnection by replacing ϕ with the (usually quadratic) constraints it imposes on the input-output pair (u, y) . In Fig. 1(b), the general representation used for this framework is illustrated. Ψ is a stable, LTI system that generates an auxiliary sequence z_k from the sequences y_k and u_k , i.e. $z = \Psi(y, u)$. Throughout this article, we consider the simplest case where $\Psi = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \otimes I_n$, i.e. $z_k = \begin{bmatrix} y_k \\ u_k \end{bmatrix}$. Ψ is simply a static map in this case. We now define two classes of IQCs that are useful in analyzing optimization algorithms.

Definition 3. Consider the sequences y_k and u_k , and let $z_k = \begin{bmatrix} y_k \\ u_k \end{bmatrix}$. The nonlinear function $u_k = \phi(y_k)$ satisfies the pointwise IQC defined by M if for all $k \geq 0$,

$$z_k^T M z_k \geq 0. \tag{9}$$

Definition 4. Let $\rho > 0$ be given. Consider the sequences y_k and u_k , and let $z_k = \begin{bmatrix} y_k \\ u_k \end{bmatrix}$. The nonlinear function $u_k = \phi(y_k)$ satisfies the ρ -hard IQC defined by (M, ρ) if for all $T \geq 0$,

$$\sum_{k=0}^T \rho^{-2k} z_k^T M z_k \geq 0. \tag{10}$$

Note that if ϕ satisfies a pointwise IQC defined by M , it also satisfies the ρ -hard IQC defined by (M, ρ) for all $\rho < 1$. This fact is useful in deriving convergence rates. If $z \in \mathbb{R}^{n_z}$, then $M \in \mathbb{R}^{n_z \times n_z}$, and is symmetric. Typically, M is indefinite. Finally, more general classes of IQCs where Ψ is a dynamical system are defined in [16], which fall under the class of discrete-time Zames-Falb IQCs [18].

For the case of gradient descent, the nonlinearity $\phi = \nabla f$ satisfies the following IQC (adapted from [16]), and this is used throughout this article.

Lemma 1 (Sector IQC, [16]). Let $f \in \mathcal{S}(m, L)$ and $\phi = \nabla f$. Then, ∇f satisfies the pointwise IQC defined by:

$$M = \begin{bmatrix} -2mL & (L+m) \\ (L+m) & -2 \end{bmatrix} \otimes I_n, \tag{11a}$$

and the corresponding quadratic constraint is:

$$\begin{bmatrix} y_k \\ u_k \end{bmatrix}^T \begin{bmatrix} -2mL I_n & (L+m)I_n \\ (L+m)I_n & -2I_n \end{bmatrix} \begin{bmatrix} y_k \\ u_k \end{bmatrix} \geq 0. \tag{11b}$$

3 Analysis with Varying Step Sizes

In this section, we describe our approach to solve Problems 1 and 2 for gradient descent with a varying step size, using Integral Quadratic Constraints. This involves considering the algorithm (4) as a linear, parameter-varying (LPV) system. In Section 3.1, we present our main result on the convergence rate for the noise-free gradient method with a varying step size. In Section 3.2, we consider gradient descent affected by a noisy gradient and present our main result on noise amplification. Finally, Section 3.3 provides details on numerical implementations for our results.

A line search is one of the most common methods leading to an optimization algorithm with varying step sizes [2, 30, 31]. In our setup, we assume that a line search has been carried out, and produces a step size at each time step k , denoted α_k . Moreover, we assume that the step size satisfies the following condition:

$$\underline{\alpha} = \frac{1}{cL} \leq \alpha_k \leq \frac{c}{L} = \bar{\alpha}, \quad (12)$$

where $c \geq 1$ is some constant characterizing the interval $\mathcal{A} = [\underline{\alpha}, \bar{\alpha}]$. In other words, the step size α is not constant in time, but takes a value in this interval at each time step k . This interval is also geometrically centered about the popular step size choice $\alpha = 1/L$. While conditions on the step size after carrying out a Wolfe line search can be derived, the resulting bounds can be too conservative to be useful, particularly for high condition numbers.

Under the assumption (12), the gradient algorithm (4) can be written as an LPV system:

$$\begin{aligned} x_{k+1} &= Ax_k + B_u(\alpha_k)u_k + B_w(\alpha_k)w_k, \\ \begin{bmatrix} y_k \\ e_k \end{bmatrix} &= \begin{bmatrix} C_y \\ C_e \end{bmatrix} x_k, \end{aligned} \quad (13)$$

where

$$A = I_n; B_u(\alpha_k) = -\alpha_k I_n; B_w(\alpha_k) = \alpha_k I_n; C_y = C_e = I_n, \quad (14)$$

and the matrices B_u and B_w depend on the parameter $\alpha_k \in \mathcal{A}$.

3.1 Convergence Rate

We first consider the noise-free case, where $w_k = 0$ for all $k \in \mathbb{Z}_+$. The following theorem provides a method to characterize convergence rates for gradient descent with a varying step size satisfying (12), and is an extension of the main result in [16].

Theorem 1 (Problem 1, Convergence Rate). *Suppose the nonlinear function $u = \phi(y)$ satisfies a pointwise IQC defined by M as given in Definition 3. Suppose that there exists a positive definite matrix P , non-negative scalar λ , and scalar $\rho \in [0, 1]$ such that the following LMI is feasible for all $\alpha \in \mathcal{A}$:*

$$\begin{bmatrix} A^T P A - \rho^2 P & A^T P B_u(\alpha) \\ B_u^T(\alpha) P A & B_u^T(\alpha) P B_u(\alpha) \end{bmatrix} + \lambda \begin{bmatrix} C_y & 0_n \\ 0_n & I_n \end{bmatrix}^T M \begin{bmatrix} C_y & 0_n \\ 0_n & I_n \end{bmatrix} \preceq 0, \quad (15)$$

Then, for any x_0 we have:

$$\|x_k\| \leq \sqrt{\text{cond}(P)\rho^k} \|x_0\|. \quad (16)$$

Proof. The proof closely follows the arguments in [16, Theorem 4], and is omitted to conserve space. \blacksquare

3.2 Noise Amplification

Now consider the gradient algorithm affected by gradient noise, as written in (13) and (14). The following theorem characterizes the noise amplification metric γ for this algorithm, and extends the constant step size result in [27, 28].

Theorem 2 (Problem 2, Noise Amplification). *Suppose the nonlinear function $u = \phi(y)$ satisfies a pointwise IQC defined by M as given in Definition 3. Suppose that there exists a positive definite matrix P and a non-negative scalar λ such that the following LMI is feasible for all $\alpha \in \mathcal{A}$:*

$$\begin{bmatrix} A^T P A - P + C_e^T C_e & A^T P B_u(\alpha) \\ B_u^T(\alpha) P A & B_u^T(\alpha) P B_u(\alpha) \end{bmatrix} + \lambda \begin{bmatrix} C_y & 0_n \\ 0_n & I_n \end{bmatrix}^T M \begin{bmatrix} C_y & 0_n \\ 0_n & I_n \end{bmatrix} \preceq 0. \quad (17)$$

Then, the metric γ is bounded by:

$$\gamma \leq \sup_{\alpha \in \mathcal{A}} \left(\text{tr} \left(B_w^T(\alpha) P B_w(\alpha) \right) \right)^{\frac{1}{2}}. \quad (18)$$

Proof. Define a storage function $V(x_k) = x_k^T P x_k$. The LMI (17) must hold for $\alpha = \alpha_k$. Following similar steps to [28, Lemma 1], we can show that:

$$\frac{1}{N} \sum_{k=0}^N \mathbb{E}[\|e_k\|^2] \leq \frac{1}{N} \mathbb{E}[V(x_0) - V(x_{N+1})] + \text{tr} \left(B_w^T(\alpha_k) P B_w(\alpha_k) \right), \quad (19)$$

where $\mathbb{E}[\cdot]$ is over different realizations of w_k . Note that

$$\text{tr} \left(B_w^T(\alpha_k) P B_w(\alpha_k) \right) \leq \sup_{\alpha \in \mathcal{A}} \text{tr} \left(B_w^T(\alpha) P B_w(\alpha) \right).$$

Using this in (19), taking the limit as $N \rightarrow \infty$, using the definition of γ in (6) and noting that in our setup, the performance output e_k is the state x_k , the result (18) follows. \blacksquare

We note here that a common approach to reduce conservatism in parameter-varying LMI problems such as (15) and (17) is to use a parameter-dependent matrix $P(\alpha)$ instead of a constant matrix P . This is used to address bounded rates of variation of the parameter, as discussed in [35]. However, in our setup we do not assume any bounds on the rate of variation of α . Our only assumption is that α_k satisfies the constraint (12) at each time step k , but can vary at any rate between these quantities. Thus, we consider only a constant matrix P .

3.3 Numerical Implementation

We now discuss a few details on implementing the LMIs in (15) and (17). Each LMI is actually an infinite family of LMIs, as they must be satisfied for each $\alpha \in \mathcal{A}$. We can simplify this to a finite number of constraints using a result for parameterized LMIs [35]. Note how the LMIs (15) and (17) can be written as:

$$\begin{bmatrix} A^T \\ B_u^T(\alpha) \end{bmatrix} P \begin{bmatrix} A & B_u(\alpha) \end{bmatrix} + X(P, \lambda) \preceq 0, \quad (20)$$

for an appropriate matrix $X(P, \lambda)$ that is an affine function of the decision variables P and λ . By Schur complements, this is equivalent to:

$$\left[\begin{array}{c|c} X(P, \lambda) & \begin{matrix} A^T \\ B_u^T(\alpha) \end{matrix} \\ \hline \begin{matrix} A & B_u(\alpha) \end{matrix} & -P^{-1} \end{array} \right] \preceq 0. \quad (21)$$

Multiply the above equation on the left and right by the symmetric, block-diagonal matrix $\begin{bmatrix} I_{2n} & \mathbf{0} \\ \mathbf{0}^T & P \end{bmatrix}$ where $\mathbf{0}$ denotes the zero matrix of appropriate dimensions. Then, (21) is equivalent to:

$$\left[\begin{array}{c|c} X(P, \lambda) & \begin{matrix} A^T P \\ B_u^T(\alpha) P \end{matrix} \\ \hline \begin{matrix} P A & P B_u(\alpha) \end{matrix} & -P \end{array} \right] \preceq 0. \quad (22)$$

Note that (22) is affine in P and λ . Furthermore, in contrast to the form of the original LMIs (15), (17), the parameter α now enters affinely in the above expression since $B_u(\alpha) = -\alpha I_n$. Thus, it is sufficient to

check the constraints (15), (17) at the end points $\underline{\alpha}$ and $\bar{\alpha}$ [32,35]. In other words, there exist P and λ such that (22) holds for all $\alpha \in \mathcal{A}$ if and only if there exist P and λ such that (22) holds for $\alpha \in \{\underline{\alpha}, \bar{\alpha}\}$. This proves particularly useful in deriving analytical solutions of (15) and (17), and is a standard technique for polytopic LPV systems such as the ones we consider in (13) and (14).

Then, the convergence rate problem can be written as:

$$\begin{aligned} & \min_{\rho, P \succ 0, \lambda \geq 0} \rho^2 \\ \text{s.t.} \quad & \text{LMI (15) holds for } \underline{\alpha} \text{ and } \bar{\alpha} \end{aligned} \quad (23)$$

This is a bilinear problem due to the presence of the term $\rho^2 P$ in the LMI (15). However, it is quasiconvex and is known as a generalized eigenvalue problem [36]. While there exist special solvers for such problems, it can also be solved via bisection on ρ^2 , checking the feasibility of the constraints (23) at each step. The noise amplification problem can be written as:

$$\begin{aligned} & \min_{P \succ 0, \lambda \geq 0, \gamma > 0} \gamma^2 \\ \text{s.t.} \quad & \text{LMI (17) holds for } \underline{\alpha} \text{ and } \bar{\alpha}, \\ & \text{and } \gamma^2 \geq \text{tr}(B_w^T(\alpha) P B_w(\alpha)) \text{ at } \underline{\alpha} \text{ and } \bar{\alpha} \end{aligned} \quad (24)$$

This problem is an SDP and can be solved using freely available solvers.

4 Results

In this section, we present the results of our study on gradient descent with varying step sizes. In Section 4.1 we discuss the reduction to the constant step size case by setting the interval constant $c = 1$, and show that prior results on convergence rate and noise amplification are recovered. We also present some insights that are particular to the gradient noise setting, including a tradeoff between convergence rate and noise amplification for strongly convex functions. In Section 4.2 we first present analytical expressions for the convergence rate and noise amplification metric as functions of the condition number κ and the interval constant c , based on the parameter-varying LMIs (15) and (17). Finally, we present a set of numerical results based on Theorems 1 and 2.

4.1 Reduction to a Constant Step Size

When $c = 1$, the step size interval \mathcal{A} reduces to a single point, resulting in a constant step size α . The two LMIs (15) and (17) reduce to a single condition for $\alpha = \underline{\alpha} = \bar{\alpha}$.

For the convergence rate problem, substituting (14) and (11) in (15) when $c = 1$ produces the following:

$$\begin{bmatrix} (1 - \rho^2) & -\alpha \\ -\alpha & \alpha^2 \end{bmatrix} + \lambda \begin{bmatrix} -2mL & L + m \\ L + m & -2 \end{bmatrix} \preceq 0, \quad (25)$$

which is the same LMI obtained in [16]. This follows from a dimensionality reduction argument described therein. An analytical solution can be obtained using Schur complements:

$$\rho^* = \max\{|1 - \alpha m|, |1 - \alpha L|\} = \begin{cases} (1 - \alpha m), & \alpha \leq \frac{2}{L + m}, \\ (\alpha L - 1), & \alpha > \frac{2}{L + m} \end{cases}, \quad (26)$$

although constant step sizes larger than $\alpha = \frac{2}{L+m}$ are not common and do not provide an improvement in performance. For $\alpha = 1/L$, we note that the known convergence rate $\rho^* = 1 - \frac{1}{\kappa}$ is recovered using (26). Figure 2 illustrates the approximate number of iterations to convergence, given by $\frac{1}{(1-\rho^*)}$, as a function of the condition number κ for different values of the interval constant c . We discuss the $c = 1$ case here, and other values of c in Section 4.2. As we expect, the black curve for $c = 1$ coincides with the dashed red curve denoting the theoretical number of iterations $\frac{1}{(1-\rho^*)} = \kappa$ for the constant step size $\alpha = 1/L$. Setting $c = 1$

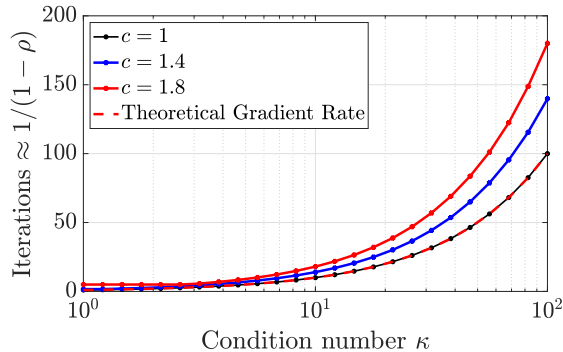


Figure 2: The approximate number of iterations to convergence as a function of the condition number for gradient descent with a varying step size characterized by c .

thus recovers the known constant step size result for the convergence of gradient descent, based on the family of LMIs in Theorem 1.

For the noise amplification problem, the constant step size case was discussed in [27, 28] using the LMI (17) for a single point α , under the iterate noise setting where $B_w = \sigma I_n$ and σ was the noise magnitude. Here, we consider a noisy gradient for which the dynamics are characterized by (14) and where $B_w(\alpha) = \alpha I_n$ when $c = 1$. In Theorem 2, from (18):

$$\gamma \leq (\text{tr}(B_w^T P B_w))^{\frac{1}{2}}, \quad (27)$$

where B_w is a constant for a given α . Following identical steps to [27], the metric γ satisfies the following bound:

$$\gamma \leq \gamma^* = \alpha \sqrt{\frac{n}{(1-\rho^2)}}, \quad (28)$$

where ρ is the convergence rate corresponding to the step size α , obtained from (26) and γ^* is the best upper bound on the noise amplification metric γ . For $\alpha = 1/L$, we have:

$$\gamma^* = \frac{1}{m} \sqrt{\frac{n}{2\kappa-1}} = \frac{1}{L} \sqrt{\frac{n}{2\kappa^{-1}-\kappa^{-2}}}, \quad (29)$$

and similarly if $\alpha = \frac{2}{L+m}$, we have:

$$\gamma^* = \frac{1}{m} \sqrt{\frac{n}{\kappa}} = \frac{1}{L} \sqrt{\frac{n}{\kappa^{-1}}}. \quad (30)$$

Notice how γ^* depends separately on m and L , and not just on the condition number κ . Further, note that this upper bound varies inversely with m and L for a fixed κ . Gradient noise amplification is thus an inherently different property from convergence rate, in that it depends individually on m and L and not just on κ . In particular, for a given condition number, γ^* can take different values based on the values of m and L .

Finally, these expressions are very closely related to the corresponding expressions in [27] for iterate noise amplification, and replacing α by $\sigma = 1$ in the above equations recovers the expressions obtained in [27]. Thus, setting $c = 1$ can recover known expressions for the noise amplification metric for constant step sizes. However, in the iterate noise setting, the best upper bound on γ no longer depends separately on m and L , and depends only on the condition number κ , unlike the gradient noise setting discussed above.

Figure 3 illustrates the variation of γ^* with condition number κ for different values of the interval constant c , and fixing one of the two parameters m and L . As before, we discuss the $c = 1$ case here, and other values of c in Section 4.2. When L is varied for particular values of m , the upper bound decreases as L and κ increase. This is shown in Figures 3(a) and (c). When m is varied for particular values of L , γ increases as m decreases and κ increases. This is shown in Figures 3(b) and (d). Furthermore, when m or L are made twice as large, the corresponding value of γ^* is half as large for all condition numbers, as seen by comparing Figure 3(a) with 3(b) and Figure 3(c) with 3(d). This illustrates that γ varies inversely with m and L for a given condition number.

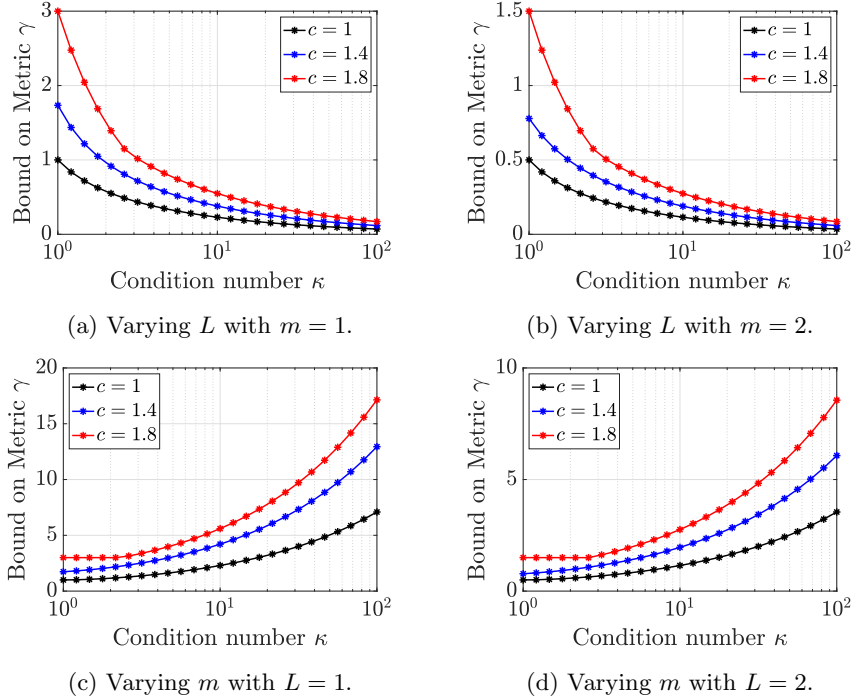


Figure 3: The upper bound on the noise amplification metric as a function of the condition number for gradient descent with a varying step size characterized by c .

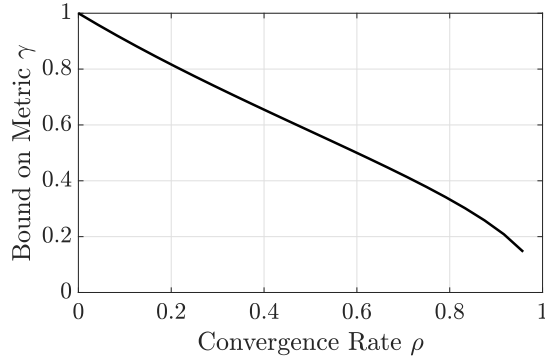


Figure 4: Tradeoff between noise amplification γ and convergence rate ρ , based on (31). A ‘faster’ algorithm has a larger value of metric γ , and is thus more sensitive to noise. In this figure, problem dimension $n = 1$ and strong convexity parameter $m = 1$.

We now discuss a tradeoff between convergence rate and noise amplification that arises from (26) and (28). A cursory examination of (28) seems to indicate that the noise amplification must increase as the convergence rate increases, i.e. worsens. However, note that γ^* also depends on the step size α , which is related to the convergence rate through (26). For $\alpha \leq \frac{2}{L+m}$, we have $\rho = (1 - \alpha m)$ or $\alpha = \frac{1-\rho}{m}$. Then,

$$\gamma^* = \frac{(1 - \rho)}{m} \sqrt{\frac{n}{(1 - \rho^2)}} = \frac{1}{m} \sqrt{\frac{n}{(1 + \rho)}}. \quad (31)$$

For $0 \leq \rho < 1$, γ^* decreases with increasing ρ as shown in Fig. 4. Thus, a choice of step size that results in faster convergence also is more sensitive to noise. This is closely related to a recent tradeoff result for accelerated momentum-based algorithms [29] on strongly convex quadratic problems. However, our result here applies to the more general class of strongly convex functions.

4.2 Results for a Varying Step Size

We now discuss our results for the convergence and noise amplification of gradient descent with a varying step size, based on Theorems 1 and 2. First, we present analytical expressions for the rate of convergence and noise amplification metric for varying step sizes, and then discuss a set of numerical results by solving the problems (23) and (24).

The following propositions provide analytical expressions for the convergence rate and noise amplification metric with varying step sizes in gradient descent. The proofs are provided in the appendix, and rely on the fact that it is sufficient to check the families of constraints (15) and (17) only at the end points $\underline{\alpha}$ and $\bar{\alpha}$ of the step size interval \mathcal{A} .

Proposition 1. The best upper bound on the convergence rate of varying step size gradient descent, based on the solution of the problem (23) for $1 \leq c < 2$ is given by:

$$\rho_{\text{varying}} = \begin{cases} c - 1, & \kappa \leq \frac{1}{c(2-c)} \\ 1 - \frac{1}{c\kappa}, & \kappa > \frac{1}{c(2-c)} \end{cases}, \quad (32)$$

depending on the interval constant c and the condition number κ as above.

Proposition 2. The noise amplification metric γ_{varying} for varying step size gradient descent, based on the solution of the SDP (24) for $1 \leq c < 2$ satisfies:

$$\gamma_{\text{varying}} \geq \gamma_{\text{varying}}^* = \begin{cases} \frac{c}{L} \sqrt{\frac{n}{2c-c^2}}, & \kappa \leq \frac{c}{2-c} \\ \frac{c}{m} \sqrt{\frac{n}{2c\kappa-c^2}}, & \kappa > \frac{c}{2-c} \end{cases}, \quad (33)$$

depending on the interval constant c and the condition number κ as above.

We now present a set of numerical results based on Theorems 1 and 2, which characterize how the convergence rate ρ and noise amplification metric γ vary with the condition number κ as well as the interval constant c . We choose three values of c to test: $c = 1$, $c = 1.4$ and $c = 1.8$. Note that $c = 1$ implies that the interval \mathcal{A} reduces to a point $\alpha_\kappa \equiv \alpha = 1/L$, a constant step size, which was discussed in Section 4.1.

In the convergence rate results presented here, we use $m = 1$ throughout. However, the results are unchanged for other values of m as the convergence rate depends only on the condition number κ as seen in (32). We now revisit Fig. 2, and discuss results for larger values of c . We first note that our results are consistent with the analytical expression derived in Proposition 1 for different values of c . As discussed earlier, setting $c = 1$ also recovers the theoretical number of iterations for convergence when $\alpha = 1/L$. For larger values of c , we observe that convergence is generally slower as the difference between $\underline{\alpha}$ and $\bar{\alpha}$ is larger, leading to a larger interval \mathcal{A} and introducing some conservatism in the rate bound. Further note that if $c \geq 2$, convergence would not be guaranteed based on our approach to solve the LMI in (15).

For the noise amplification results, we test different values of m and L , fixing one of the two parameters. These results are shown in Fig. 3. As with the case when $c = 1$, γ decreases with increasing L and m when the other parameter is fixed, for any value of c . Further, when κ is fixed and m or L are made twice as large, γ decreases by a factor of two, for any value of c . This can be seen by comparing Fig. 3(a) with Fig. 3(b), and comparing Fig. 3(c) with Fig. 3(d). Thus, γ varies inversely with m and L for a fixed κ , and continues to depend separately on these parameters. The value of γ increases as c increases, i.e. the algorithm is more sensitive to noise when the interval \mathcal{A} is larger. These results are also consistent with the bound derived in Proposition 2 for different values of c .

A few remarks are in order. In both Fig. 2 and Fig. 3, the guarantees for $c > 1$ are generally worse than the guarantees for $c = 1$, a constant step size. The results show that convergence is slower, and the algorithm is more sensitive to noise. This is primarily a consequence of the setup discussed in Section 3. The approach we discuss in Section 3.3 is inherently conservative, and the results demonstrate worse guarantees than the constant step size case, especially when the step size is allowed to vary over a larger set. However,

in practice, line search algorithms generally perform better (and not worse as predicted by our analysis) than their constant step size counterparts in terms of convergence or noise amplification. While the results do not accurately represent this, future work must focus on improving the methods discussed in Section 3. In particular, it is worth exploring how line search algorithms (such as a Wolfe line search) can be more accurately characterized so that the IQC approach may lead to improved guarantees. A difficulty with this, as mentioned in Section 3 is that the derivable range on the step size using a Wolfe line search can be much more conservative than the range used here, resulting in worse guarantees.

5 Concluding Remarks

In this article, we presented an analysis of gradient descent with varying step sizes, in terms of its convergence rate and noise amplification. Assuming a line search produces a step size in a given interval, the algorithm is modeled as an LPV system. Using a technique for polytopic LPV systems and building on prior work in the IQC framework, we construct SDPs to certify convergence rates and the steady-state variance in iterate error. Our condition provides a bound on the convergence rate when the step size is within a restricted set around $\alpha = 1/L$. Moreover, our condition recovers the corresponding gradient rate when the interval is a single point $1/L$. Further, the noise amplification metric depends on both parameters m and L individually, and varies inversely with them for a fixed condition number κ .

It is worth reiterating that the guarantees for a line search, i.e. for $c > 1$ are generally worse than those for a constant step size. This is not necessarily reflective of the practical performance of line search algorithms, and it is worth exploring how these algorithms can be more accurately characterized to obtain less conservative results. Another avenue for further work in this area is the use of dynamic IQCs or multiple IQCs that may reduce conservatism in the results. Numerous dynamic IQCs for convex functions (where Ψ in Section 2.3 is not simply a static map, but a dynamical system) are described in [16], and the use of multiple such IQCs may reduce conservatism in both the convergence rate bound and the noise amplification metric. Further avenues for future work include the analysis of time-varying step sizes in accelerated and stochastic gradient algorithms.

A Proof of Analytical Results

A.1 Proof of Proposition 1

First, note that (15) is a family of two LMIs at $\underline{\alpha}$ and $\bar{\alpha}$, based on the results in (20), (21) and (22). For a given α , we know from the constant step size case that the convergence rate satisfies (26). Let $\underline{\rho}$ and $\bar{\rho}$ be the convergence rate corresponding to the two step sizes $\underline{\alpha}$ and $\bar{\alpha}$. Then, the convergence rate from (15) can be written as:

$$\rho_{\text{varying}} = \max\{\underline{\rho}, \bar{\rho}\} \quad (34)$$

since this is the smallest step size for which (15) is feasible at both $\underline{\alpha}$ and $\bar{\alpha}$, with a common solution $P = 1$ which is discussed in [16]. Since $\underline{\alpha} = \frac{1}{cL} \leq \frac{2}{L+m}$, using (26):

$$\underline{\rho} = 1 - \frac{m}{cL} = 1 - \frac{1}{c\kappa}. \quad (35)$$

For $\bar{\rho}$, note that $\bar{\alpha} = \frac{c}{L} \geq \frac{2}{L+m}$ if $\kappa \leq \frac{c}{2-c}$ and $c < 2$. Then, using (26) with $c < 2$:

$$\bar{\rho} = \begin{cases} c-1, & \kappa \leq \frac{c}{2-c}, \\ 1 - \frac{c}{\kappa}, & \kappa > \frac{c}{2-c}. \end{cases} \quad (36)$$

Comparing (35) and (36), first note that $\frac{1}{c(2-c)} \leq \frac{c}{2-c}$ for $1 \leq c < 2$. When $\kappa \leq \frac{1}{c(2-c)} \leq \frac{c}{2-c}$, $\rho_{\text{varying}} = \max\{1 - \frac{1}{c\kappa}, c-1\} = c-1$ in this range of κ . Next, when $\frac{1}{c(2-c)} \leq \kappa \leq \frac{c}{2-c}$, $\rho_{\text{varying}} = \max\{1 - \frac{1}{c\kappa}, c-1\} = 1 - \frac{1}{c\kappa}$ in this range of κ . Finally, when $\kappa > \frac{c}{2-c}$, $\rho_{\text{varying}} = \max\{1 - \frac{1}{c\kappa}, 1 - \frac{c}{\kappa}\} = 1 - \frac{1}{c\kappa}$ for $c \geq 1$. Using

the above facts, we have:

$$\rho_{\text{varying}} = \begin{cases} c - 1, & \kappa \leq \frac{1}{c(2-c)} \\ 1 - \frac{1}{c\kappa}, & \kappa > \frac{1}{c(2-c)} \end{cases} \quad (37)$$

for $1 \leq c < 2$, which completes the proof of Proposition 1. \blacksquare

A.2 Proof of Proposition 2

From (24), we require $\gamma_{\text{varying}}^2 \geq \text{tr}(B_w^T(\alpha)PB_w(\alpha))$ at $\underline{\alpha}$ and $\bar{\alpha}$, where P is a common solution of the LMI (17) at these two points. Using the dimensionality reduction argument described in [16], we can restrict our search to those P of the form $P = P_0 \otimes I_n$, where P_0 is 1×1 . Substituting P and $B_w(\alpha)$, we require $\gamma_{\text{varying}}^2 \geq \alpha^2 n P_0$ at both $\underline{\alpha}$ and $\bar{\alpha}$, or

$$\gamma_{\text{varying}}^2 \geq \bar{\alpha}^2 n P_0 \quad (38)$$

since $\bar{\alpha} \geq \underline{\alpha}$. Let \underline{P}_0 and \bar{P}_0 be the two solutions at the points $\underline{\alpha}$ and $\bar{\alpha}$ respectively. Using (38), we then require:

$$\gamma_{\text{varying}}^2 \geq (\gamma_{\text{varying}}^*)^2 = \bar{\alpha}^2 n P_0^*, \quad P_0^* = \max\{\underline{P}_0, \bar{P}_0\}. \quad (39)$$

For a given α , the solution P_0 to (17) is given by:

$$P_0 = \begin{cases} \frac{1}{2\alpha m - \alpha^2 m^2}, & \alpha \leq \frac{2}{L+m}, \\ \frac{1}{2\alpha L - \alpha^2 L^2}, & \alpha \geq \frac{2}{L+m} \end{cases}, \quad (40)$$

which follows from similar arguments to [27]. Substituting $\underline{\alpha} = \frac{1}{cL}$ and $\bar{\alpha} = \frac{c}{L}$,

$$\underline{P}_0 = \frac{1}{\frac{2}{c\kappa} - \frac{1}{(c\kappa)^2}}, \quad (41)$$

$$\text{and } \bar{P}_0 = \begin{cases} \frac{1}{2c - c^2}, & \kappa \leq \frac{c}{2-c}, \\ \frac{1}{\frac{2c}{\kappa} - (\frac{c}{\kappa})^2}, & \kappa > \frac{c}{2-c} \end{cases}. \quad (42)$$

This follows from noticing $\bar{\alpha} = \frac{c}{L} \geq \frac{2}{L+m}$ if $\kappa \leq \frac{c}{2-c}$ and $c < 2$. Note that \underline{P}_0 and \bar{P}_0 are both positive since $1 \leq c < 2$ and $\kappa \geq 1$. Finally, it is easily shown that $\bar{P}_0 \geq \underline{P}_0$ for all κ and $c < 2$, and thus $P_0^* = \bar{P}_0$. Using (39) and simplifying,

$$\gamma_{\text{varying}} \geq \gamma_{\text{varying}}^* = \begin{cases} \frac{c}{L} \sqrt{\frac{n}{2c - c^2}}, & \kappa \leq \frac{c}{2-c} \\ \frac{c}{m} \sqrt{\frac{n}{2c\kappa - c^2}}, & \kappa > \frac{c}{2-c} \end{cases} \quad (43)$$

for $1 \leq c < 2$, which completes the proof of Proposition 2. \blacksquare

References

- [1] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, United Kingdom: Cambridge University Press, 2004.
- [2] J. Nocedal and S. J. Wright, *Numerical Optimization*, 2nd ed. New York, NY, USA: Springer Science + Business Media, 2006.
- [3] B. T. Polyak, "Some methods of speeding up the convergence of iteration methods," *USSR Computational Mathematics and Mathematical Physics*, vol. 4, no. 5, pp. 1–17, 1964.

- [4] Y. Nesterov, *Introductory Lectures on Convex Optimization*. New York, NY, USA: Springer Science + Business Media, 2004.
- [5] A. Beck, *First-Order Methods in Optimization*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2017.
- [6] N. Parikh and S. Boyd, “Proximal algorithms,” *Foundations and Trends in Optimization*, vol. 1, no. 3, pp. 123–231, 2014.
- [7] V. A. Yakubovich, “Nonconvex optimization problem: The infinite-horizon linear-quadratic control problem with quadratic constraints,” *Systems and Control Letters*, vol. 19, no. 1, pp. 13–22, Jul. 1992.
- [8] —, “S-procedure in nonlinear control theory,” *Vestnik Leningrad University*, vol. 4, pp. 73–93, 1971, (Russian).
- [9] A. Megretski and A. Rantzer, “System analysis via integral quadratic constraints,” *IEEE Transactions on Automatic Control*, vol. 42, no. 6, pp. 819–830, Jun. 1997.
- [10] A. Rantzer, “On the Kalman—Yakubovich—Popov lemma,” *Systems and Control Letters*, vol. 28, no. 1, pp. 7–10, Jun. 1996.
- [11] J. Veenman, C. W. Scherer, and H. Köroğlu, “Robust stability and performance analysis based on integral quadratic constraints,” *European Journal of Control*, vol. 31, pp. 1–32, Sep. 2016.
- [12] C. W. Scherer, “Dissipativity and integral quadratic constraints: Tailored computational robustness tests for complex interconnections,” *IEEE Control Systems Magazine*, vol. 42, no. 3, pp. 115–139, Jun. 2022.
- [13] P. Seiler, “Stability analysis with dissipation inequalities and integral quadratic constraints,” *IEEE Transactions on Automatic Control*, vol. 60, no. 6, pp. 1704–1709, Jun. 2015.
- [14] B. Hu, M. J. Lacerda, and P. Seiler, “Robustness analysis of uncertain discrete-time systems with dissipation inequalities and integral quadratic constraints,” *International Journal of Robust and Nonlinear Control*, vol. 27, no. 11, pp. 1940–1962, Jul. 2017.
- [15] R. Boczar, L. Lessard, and B. Recht, “Exponential convergence bounds using integral quadratic constraints,” in *2015 IEEE 54th Annual Conference on Decision and Control*, Osaka, Japan, Dec. 2015, pp. 7516–7521.
- [16] L. Lessard, B. Recht, and A. Packard, “Analysis and design of optimization algorithms via integral quadratic constraints,” *SIAM Journal on Optimization*, vol. 26, no. 1, pp. 57–95, 2016.
- [17] Y. Drori and M. Teboulle, “Performance of first-order methods for smooth convex minimization: A novel approach,” *Mathematical Programming*, vol. 145, pp. 451–482, Jun. 2013.
- [18] W. P. Heath and A. G. Willis, “Zames-Falb multipliers for quadratic programming,” in *Proceedings of the 44th IEEE Conference on Decision and Control*, Seville, Spain, Dec. 2005, pp. 963–968.
- [19] A. Badithela and P. Seiler, “Analysis of the heavy-ball algorithm using integral quadratic constraints,” in *2019 American Control Conference (ACC)*, Philadelphia, PA, USA, Jul. 2019, pp. 4081–4085.
- [20] B. Hu, P. Seiler, and L. Lessard, “Analysis of biased stochastic gradient descent using sequential semidefinite programs,” *Mathematical Programming*, vol. 187, pp. 383–408, May 2021.
- [21] H. Mohammadi, S. Samuelson, and M. R. Jovanović, “Transient growth of accelerated optimization algorithms,” *IEEE Transactions on Automatic Control*, vol. 68, no. 3, pp. 1823–1830, Mar. 2023.
- [22] M. Fazlyab, A. Ribeiro, M. Morari, and V. M. Preciado, “Analysis of optimization algorithms via integral quadratic constraints: Nonstrongly convex problems,” *SIAM Journal on Optimization*, vol. 28, no. 3, pp. 2654–2689, Sep. 2018.

- [23] B. Van Scoy, R. A. Freeman, and K. M. Lynch, “The fastest known globally convergent first-order method for minimizing strongly convex functions,” *IEEE Control Systems Letters*, vol. 2, no. 1, pp. 49–54, Jan. 2018.
- [24] S. Cyrus, B. Hu, B. Van Scoy, and L. Lessard, “A robust accelerated optimization algorithm for strongly convex functions,” in *2018 American Control Conference (ACC)*, Milwaukee, WI, USA, Jun. 2018, pp. 1376–1381.
- [25] L. Lessard and P. Seiler, “Direct synthesis of iterative algorithms with bounds on achievable worst-case convergence rate,” in *2020 American Control Conference (ACC)*, Denver, CO, USA, Jul. 2020, pp. 119–125.
- [26] L. Lessard, “The analysis of optimization algorithms: A dissipativity approach,” *IEEE Control Systems Magazine*, vol. 42, no. 3, pp. 58–72, Jun. 2022.
- [27] H. Mohammadi, M. Razaviyayn, and M. R. Jovanović, “Performance of noisy Nesterov’s accelerated method for strongly convex optimization problems,” in *2019 American Control Conference (ACC)*, Philadelphia, PA, USA, Jul. 2019, pp. 3426–3431.
- [28] H. Mohammadi, M. Razaviyayn, and M. R. Jovanović, “Robustness of accelerated first-order algorithms for strongly convex optimization problems,” *IEEE Transactions on Automatic Control*, vol. 66, no. 6, pp. 2480–2495, Jun. 2021.
- [29] —, “Tradeoffs between convergence rate and noise amplification for momentum-based accelerated optimization algorithms,” 2022. [Online]. Available: <https://arxiv.org/abs/2209.11920>
- [30] L. Armijo, “Minimization of functions having Lipschitz continuous first partial derivatives,” *Pacific Journal of Mathematics*, vol. 16, no. 1, pp. 1–3, Jan. 1966.
- [31] P. Wolfe, “Convergence conditions for ascent methods,” *SIAM Review*, vol. 11, no. 2, pp. 226–235, 1969.
- [32] H. Pfifer and P. Seiler, “Robustness analysis of linear parameter varying systems using integral quadratic constraints,” *International Journal of Robust and Nonlinear Control*, vol. 25, no. 15, pp. 2843–2864, Oct. 2015.
- [33] —, “Less conservative robustness analysis of linear parameter varying systems using integral quadratic constraints,” *International Journal of Robust and Nonlinear Control*, vol. 26, no. 16, pp. 3580–3594, Nov. 2016.
- [34] S. Wang, H. Pfifer, and P. Seiler, “Robust synthesis for linear parameter varying systems using integral quadratic constraints,” *Automatica*, vol. 68, pp. 111–118, Jun. 2016.
- [35] F. Wu, X. H. Yang, A. Packard, and G. Balas, “Induced \mathcal{L}_2 norm control for LPV systems with bounded parameter variation rates,” *International Journal of Robust and Nonlinear Control*, vol. 6, no. 9-10, pp. 983–998, Nov. 1996.
- [36] S. Boyd and L. El Ghaoui, “Method of centers for minimizing generalized eigenvalues,” *Linear Algebra and its Applications*, vol. 188, pp. 63–111, Jul. 1993.