

Lecture 2

Unifying the Analysis in Control and Optimization via Semidefinite Programs

Lecturer: Bin Hu, Date: 08/27/2020

In this lecture, we review some stability analysis tools in the controls literature, and then tailor them to analyze the convergence rates of some simple optimization methods. Hopefully you will be convinced that there are similarities between the analysis problems in control and optimization, and hence it is not too surprising that some analysis tools developed in the controls field can be applied to study large-scale optimization.

2.1 Stability Analysis in Control

2.1.1 Autonomous Systems and Internal Stability

Possibly the simplest system in the controls literature is the following so-called (autonomous) linear time-invariant system

$$\xi_{k+1} = A\xi_k \quad (2.1)$$

Here we consider discrete-time systems, and ξ_k is the state at time step k . Given the initial condition ξ_0 , then the sequence $\{\xi_k\}$ is completely determined by (2.1). One fundamental question control people usually ask is whether (2.1) is stable. The system (2.1) is internally stable if ξ_k converges to 0 given any arbitrary initial condition ξ_0 . Notice (2.1) just states that we have $\xi_k = A^k \xi_0$. There are four possible cases.

1. When A is Schur stable (or equivalently the spectral radius of A is smaller than 1), the term A^k converges to a zero matrix. For example, if $A = \begin{bmatrix} 0.9 & 0 \\ 0 & 0.9 \end{bmatrix}$, then $A^k = 0.9^k I$.
2. When the spectral radius of A is equal to 1 and all the Jordan blocks corresponding to eigenvalues with magnitude equal to 1 are 1×1 , A^k remains bounded for any k . This is the so-called marginal stability case where $A^k \xi_0$ remains bounded but may not converge to 0. For example, if $A = \begin{bmatrix} 1 & 0 \\ 0 & 0.9 \end{bmatrix}$, then $A^k = \begin{bmatrix} 1 & 0 \\ 0 & 0.9^k \end{bmatrix}$.
3. When the spectral radius of A is equal to 1 and all the Jordan blocks corresponding to eigenvalues with magnitude equal to 1 are not 1×1 , A^k becomes unbounded. For example, if $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$, then $A^k = \begin{bmatrix} 1 & k \\ 0 & 1 \end{bmatrix}$.

4. When the spectral radius of A is larger than 1, A^k also becomes unbounded. For example, if $A = \begin{bmatrix} 1.1 & 0 \\ 0 & 1.1 \end{bmatrix}$, then $A^k = 1.1^k I$.

Therefore, it is straightforward to verify that (2.1) is stable if and only if the spectral radius of A is strictly less than 1. The spectral radius condition only works for such linear time-invariant (LTI)¹ system. It is hard to extend such conditions for time-varying or nonlinear systems. Alternatively, one can also formulate necessary and sufficient stability conditions for (2.1) using semidefinite programs.

Theorem 2.1. *The system (2.1) is internally stable if and only if there exists a positive definite matrix P such that*

$$A^T P A - P < 0 \quad (2.2)$$

Here the inequality holds in the definite sense (so what we really mean here is that the matrix $(A^T P A - P)$ needs to be a negative definite matrix).

Proof: We will only show sufficiency since this direction can be generalized for time-varying or nonlinear systems. If (2.2) holds, then there exists a sufficiently small positive number $\varepsilon > 0$ such that $A^T P A - P \leq -\varepsilon P$ which can be rewritten as

$$A^T P A - (1 - \varepsilon)P \leq 0.$$

Therefore we can left and right multiply both sides of the above inequality with ξ_k^T and ξ_k and obtain

$$(A\xi_k)^T P (A\xi_k) - (1 - \varepsilon)\xi_k^T P \xi_k \leq 0$$

We have $\xi_{k+1} = A\xi_k$ and the above inequality is equivalent to $\xi_{k+1}^T P \xi_{k+1} \leq (1 - \varepsilon)\xi_k^T P \xi_k$. By induction, we have

$$\xi_k^T P \xi_k \leq (1 - \varepsilon)^k \xi_0^T P \xi_0$$

Since P is positive definite, we have $\xi_k^T P \xi_k \geq \lambda_{\min}(P)\|\xi_k\|^2$ where $\lambda_{\min}(P)$ is the smallest eigenvalue of P and is a positive number. Finally we have

$$\|\xi_k\|^2 \leq (1 - \varepsilon)^k c \quad (2.3)$$

where $c = \frac{\xi_0^T P \xi_0}{\lambda_{\min}(P)}$. We know $0 \leq 1 - \varepsilon < 1$ and hence $\|\xi\|$ converges to 0 as k goes to ∞ . We have established the internal stability of (2.1).

The proof for necessity relies on the LTI assumption and is omitted here. ■

¹This just means A is a constant matrix and does not change over time.

How to use the condition (2.2)? The testing condition (2.2) leads to a semidefinite program (or equivalently linear matrix inequality) problem. Given A , the left side of (2.2) is linear in P . One just needs to search such positive definite P satisfying the matrix inequality condition in (2.2). Numerically this can be done using semidefinite programming solvers. In the controls field, many analysis and design conditions are formulated as linear matrix inequality (LMI) conditions, and (2.2) is one of the simplest. We will see more such LMI conditions later.

Lyapunov functions. The proof of Theorem 2.1 relies on constructing the Lyapunov function $V(\xi) = \xi^T P \xi$. A physical interpretation for this function is that it measures how much energy is stored in the system. This function is nonnegative for all ξ and is zero at the $\xi = 0$ (which is the fixed point of (2.1)). In addition, it is radially unbounded. In the above proof, we have shown $V(\xi_{k+1}) \leq (1 - \varepsilon)V(\xi_k)$. So we have shown that the internal energy of the system is decreased at every step and eventually the minimum energy is attained at the fixed point. Lyapunov arguments can be applied in many cases and provide a powerful unified framework for stability analysis. We will learn more about this approach later.

Advantages of (2.2). It is emphasized that people do not really use (2.2) when testing the stability of (2.1). A more efficient approach is to look at the spectral radius of A directly. However, (2.2) can be extended to time-varying/nonlinear systems which one cannot apply the spectral radius arguments to analyze. For example, consider the so-called linear parameter-varying (LPV) system described by the following state space model:

$$\xi_{k+1} = A(\zeta_k)\xi_k \quad (2.4)$$

where the matrix A becomes a function of some scheduling parameter ζ_k . The parameter ζ can be measured at every step but we do not know it in advance. We do know how A depends on the value of ζ_k . Now we cannot come to a conclusion of the internal stability of (2.4) by just looking at the spectral radius of A for all ζ . However, the Lyapunov argument still works. We can show (2.4) is internally stable if there exists a positive definite matrix P such that

$$A(\zeta)^T P A(\zeta) - (1 - \varepsilon)P \leq 0, \quad \forall \zeta \quad (2.5)$$

The proof is almost identical to the proof of Theorem 2.1. We left and right multiply both sides of the above inequality with ξ_k^T and ξ_k and obtain $V(\xi_{k+1}) \leq (1 - \varepsilon)V(\xi_k)$ which immediately leads to the desired conclusion. Here we do not have necessity. If there is no solution for (2.5), it is still possible that (2.4) is internally stable. The numerical implementation of (2.5) is tricky since it has to be satisfied for all ζ . A heuristic is to grid ζ and then the infinite dimensional LMI condition (2.5) is approximated by a finite dimensional condition on the grid of ζ . This approach does introduce some numerical errors. It is also worth mentioning that sometimes we allow P to depend on the parameter ζ and this leads to the so-called parameter-dependent Lyapunov functions which can reduce the conservatism in the stability analysis. Now we give similar stability results for two more types of systems.

- LTV system: Now we consider an LTV system $\xi_{k+1} = A_k \xi_k$ where we do know how A explicitly depends on k . This system is internally stable if there exists a positive definite matrix P such that

$$A_k^\top P A_k - (1 - \varepsilon)P \leq 0, \forall k \quad (2.6)$$

Again, the proof is based on Lyapunov arguments. Define $V(\xi) = \xi^\top P \xi$. We left and right multiply both sides of the above inequality with ξ_k^\top and ξ_k and obtain $V(\xi_{k+1}) \leq (1 - \varepsilon)V(\xi_k)$ which immediately leads to the desired conclusion. Again the LMI condition here is infinite dimensional. One may need to solve this condition analytically. Similar to the LPV case, we can allow P to depend on k and formulate a less conservative LMI condition. If there exist a sequence of positive definite matrices P_k such that $P_k \geq cI, \forall k$ for some positive c and

$$A_k^\top P_{k+1} A_k - (1 - \varepsilon)P_k \leq 0, \forall k$$

then the LTV system is stable. The proof is based on defining a time-varying Lyapunov function as $V(\xi_k) = \xi_k^\top P_k \xi_k$. We can left and right multiply both sides of the above LMI with ξ_k^\top and ξ_k and obtain $V(\xi_{k+1}) = \xi_{k+1}^\top P_{k+1} \xi_{k+1} \leq (1 - \varepsilon)\xi_k^\top P_k \xi_k = (1 - \varepsilon)V(\xi_k)$ which immediately leads to the desired conclusion.

- Jump systems: Consider the system $\xi_{k+1} = A_{i_k} \xi_k$ where $\{i_k\}$ itself is a stochastic process. This system is mean square stable if $\mathbb{E}\|\xi_k\|^2$ converges to 0 given any initial conditions. Notice the state matrix depends on the jump parameter i_k . For simplicity, we assume $\{i_k\}$ is an I.I.D process sampled from a finite set $\{1, 2, \dots, n\}$. Suppose $\Pr(i_k = i) = p_i$. Again, we can use Lyapunov arguments to obtain stability conditions in the form of LMIs. This jump system is mean square stable if there exists a positive definite matrix P such that

$$\sum_{i=1}^n p_i A_i^\top P A_i - P < 0 \quad (2.7)$$

The proof is similar to the LTI system case but we need to use a little bit probability theory. The above LMI ensures $\sum_{i=1}^n p_i A_i^\top P A_i - (1 - \varepsilon)P \leq 0$ for some sufficiently small positive ε . Again we define $V(\xi_k) = \xi_k^\top P \xi_k$. A key relation is $\mathbb{E}[V(\xi_{k+1}) | \xi_k] = \sum_{i=1}^n p_i \xi_k^\top A_i^\top P A_i \xi_k$.² Therefore, we can left and right multiply both sides of the LMI with ξ_k^\top and ξ_k and obtain $\mathbb{E}[V(\xi_{k+1}) | \xi_k] = \sum_{i=1}^n p_i \xi_k^\top A_i^\top P A_i \xi_k \leq (1 - \varepsilon)\xi_k^\top P \xi_k = (1 - \varepsilon)V(\xi_k)$. Then we can take the full expectation and iterate the resultant inequality to establish the mean square stability. Notice when $n = 1$, the condition (2.7) just recovers the standard LMI condition for the LTI system. One can also allow i_k to be sampled from a Markov chain, and that leads to the so-called Markov jump linear system (MJLS). The stability conditions for MJLS are in the form of coupled LMIs, which are more complicated than (2.7). We skip the details here.

²To be more precise, the conditional expectation should be taken on \mathcal{F}_k which is the σ -algebra at k . We avoid such mathematical machinery here. Just think that if ξ_k is known, then the only source for randomness is i_k and we just average V_k based on the distribution of i_k .

There are many other types of linear dynamical systems including periodic systems that can be handled by similar Lyapunov arguments. We will not cover all of them. The key message is that time-invariance is not required by Lyapunov theory.

Convergence rate. Inequality (2.3) in the above proof actually gives an exponential convergence rate $\sqrt{1 - \varepsilon}$ which quantifies how fast $\|\xi_k\|$ approaches 0. An LTI system is either exponentially stable or not stable. Actually one can modify the LMI condition (2.2) to test whether (2.1) converges at a given testing rate or not. If there exists a positive definite matrix P such that

$$A^\top P A - \rho^2 P \leq 0 \quad (2.8)$$

then the system (2.1) converges at the exponential rate ρ , i.e. $\|\xi_k\| \leq c\rho^k$ where c is a constant. This can be proved using a similar Lyapunov argument.

Continuous-time results. There are similar results for continuous-time systems. For example, consider a continuous-time LTI system $\dot{\xi} = A\xi$. This system is internally stable if there exists a positive definite P such that $A^\top P + P A < 0$. The proof is similar. There exists a sufficiently small ε such that $A^\top P + P A \leq -2\varepsilon P$. We can define a quadratic Lyapunov function $V(\xi) = \xi^\top P \xi$ and then $\dot{V} = \dot{\xi}^\top P \xi + \xi^\top P \dot{\xi} = \xi^\top (A^\top P + P A) \xi \leq -2\varepsilon V$. Therefore, $V \leq C e^{-2\varepsilon t}$ and ξ decays at an exponential rate. Similar results can be derived for LPV or LTV systems.

2.1.2 Taking Inputs into Accounts: Input-Output Gain

In the controls field, we study how inputs can be used to change the behavior of the system. Built upon the autonomous system model (2.1), now we introduce more general state-space models for dynamical systems. Let a dynamic system G be governed by a linear state-space model, which is described by the following recursive iteration:

$$\begin{aligned} \xi_{k+1} &= A\xi_k + Bu_k \\ y_k &= C\xi_k + Du_k \end{aligned} \quad (2.9)$$

where $\xi_k \in \mathbb{R}^{n_\xi}$, $u_k \in \mathbb{R}^{n_u}$, $y_k \in \mathbb{R}^{n_y}$, $A \in \mathbb{R}^{n_\xi \times n_\xi}$, $B \in \mathbb{R}^{n_\xi \times n_u}$, $C \in \mathbb{R}^{n_y \times n_\xi}$, and $D \in \mathbb{R}^{n_y \times n_u}$. At each step k , the variables ξ_k , u_k , and y_k are referred to as the state, input, and output of the system G . When the initial condition ξ_0 is given, the state $\{\xi_k\}$ and the output $\{y_k\}$ will be completely determined by the input sequence $\{u_k\}$.

Block diagram. In the controls field, block diagrams are widely used. The input-output relationship of the dynamical system G can be described by the following block diagram.

The above block diagram just states (u, y) satisfies $y = G(u)$ when one views the dynamical system G (with some fixed initial condition) as an input-output map.

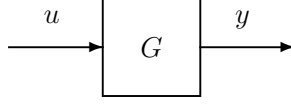


Figure 2.1. The Block-Diagram for a Dynamic System G

Clearly one can set $u = 0$ and study the internal stability of the resultant autonomous system. We have already talked about this type of analysis. Another important question is how the input u_k will affect the output y_k . A useful tool for answering such questions is the following LMI condition.

Theorem 2.2. *If there exists a positive semidefinite matrix P such that*

$$\begin{bmatrix} A^\top P A - P & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} + \begin{bmatrix} C & D \end{bmatrix}^\top \begin{bmatrix} I & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \begin{bmatrix} C & D \\ 0 & I \end{bmatrix} \leq 0 \quad (2.10)$$

then for any ξ_0 and arbitrary input sequence $\{u_k\}$, the system (2.9) satisfies the following bound with any N

$$\sum_{k=0}^N \|y_k\|^2 \leq \gamma^2 \sum_{k=0}^N \|u_k\|^2 + \xi_0^\top P \xi_0 \quad (2.11)$$

Proof: Based on the condition (2.10), we have

$$\begin{bmatrix} \xi_k \\ u_k \end{bmatrix}^\top \left(\begin{bmatrix} A^\top P A - P & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} + \begin{bmatrix} C & D \end{bmatrix}^\top \begin{bmatrix} I & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \begin{bmatrix} C & D \\ 0 & I \end{bmatrix} \right) \begin{bmatrix} \xi_k \\ u_k \end{bmatrix} \leq 0 \quad (2.12)$$

Notice we have $\xi_{k+1}^\top P \xi_{k+1} = (A\xi_k + Bu_k)^\top P (A\xi_k + Bu_k) = \begin{bmatrix} \xi_k \\ u_k \end{bmatrix}^\top \begin{bmatrix} A^\top P A & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} \begin{bmatrix} \xi_k \\ u_k \end{bmatrix}$.

Therefore, we have

$$\begin{bmatrix} \xi_k \\ u_k \end{bmatrix}^\top \begin{bmatrix} A^\top P A - P & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} \begin{bmatrix} \xi_k \\ u_k \end{bmatrix} = \xi_{k+1}^\top P \xi_{k+1} - \xi_k^\top P \xi_k$$

Similarly, we have

$$\|y_k\|^2 - \gamma^2 \|u_k\|^2 = (C\xi_k + Du_k)^\top (C\xi_k + Du_k) - \gamma^2 \|u_k\|^2 = \begin{bmatrix} \xi_k \\ u_k \end{bmatrix}^\top \begin{bmatrix} C & D \\ 0 & I \end{bmatrix}^\top \begin{bmatrix} I & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \begin{bmatrix} C & D \\ 0 & I \end{bmatrix} \begin{bmatrix} \xi_k \\ u_k \end{bmatrix}$$

Consequently, (2.12) just leads to

$$\xi_{k+1}^\top P \xi_{k+1} - \xi_k^\top P \xi_k + \|y_k\|^2 - \gamma^2 \|u_k\|^2 \leq 0 \quad (2.13)$$

Since P is positive semidefinite, we know $\xi_{N+1}^\top P \xi_{N+1} \geq 0$. We can directly sum the above inequality from $k = 0$ to N to finish the proof of Theorem 2.2. ■

Interpretations of γ . The smaller γ is, the more stable G is subject to the input u . Therefore, γ is a measure for input-output stability. Many control problems including tracking and disturbance rejection can be formulated as optimization problems whose objectives are minimizing such input-output gain γ . The famous \mathcal{H}_∞ control is based on this idea.

How to use the condition (2.10)? When (A, B, C, D) are given, the condition (2.10) is linear in P and γ^2 . Therefore, minimizing γ^2 subject to the constraints (2.10) and $P \geq 0$ can also be done via semidefinite programs.

Extensions. Similar analysis can be performed for time-varying/stochastic systems.

- LPV systems: Now we consider the LPV system

$$\begin{aligned}\xi_{k+1} &= A(\zeta_k)\xi_k + B(\zeta_k)u_k \\ y_k &= C(\zeta_k)\xi_k + D(\zeta_k)u_k\end{aligned}\tag{2.14}$$

where the matrices (A, B, C, D) depend on the scheduling parameter ζ_k . If there exists a positive semidefinite matrix P such that

$$\begin{bmatrix} A(\zeta)^\top P A(\zeta) - P & A(\zeta)^\top P B(\zeta) \\ B(\zeta)^\top P A(\zeta) & B(\zeta)^\top P B(\zeta) \end{bmatrix} + \begin{bmatrix} C(\zeta) & D(\zeta) \end{bmatrix}^\top \begin{bmatrix} I & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \begin{bmatrix} C(\zeta) & D(\zeta) \\ 0 & I \end{bmatrix} \leq 0, \forall \zeta$$

then for any ξ_0 and arbitrary input sequence $\{u_k\}$, the system (2.14) satisfies the input-output bound $\sum_{k=0}^N \|y_k\|^2 \leq \gamma^2 \sum_{k=0}^N \|u_k\|^2 + \xi_0^\top P \xi_0$ for any N . The proof is almost identical. Define $V(\xi) = \xi^\top P \xi$. We left and right multiply both sides of the LMI condition with $\begin{bmatrix} \xi_k^\top & u_k^\top \end{bmatrix}$ and $\begin{bmatrix} \xi_k \\ u_k \end{bmatrix}$ and obtain $V(\xi_{k+1}) - V(\xi_k) + \|y_k\|^2 - \gamma^2 \|u_k\|^2 \leq 0$ which immediately leads to the desired conclusion. Again, the numerical implementation of the LMI relies on griding heuristics. We may allow P to depend on the parameter ζ to reduce the conservatism in the analysis. However, the use of such parameter-dependent Lyapunov functions further increases the computational cost.

- LTV system: We can use similar Lyapunov arguments to obtain the following input-output analysis condition for LTV systems

$$\begin{bmatrix} A_k^\top P_{k+1} A_k - P_k & A_k^\top P_{k+1} B_k \\ B_k^\top P_{k+1} A_k & B_k^\top P_{k+1} B_k \end{bmatrix} + \begin{bmatrix} C_k & D_k \end{bmatrix}^\top \begin{bmatrix} I & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \begin{bmatrix} C_k & D_k \\ 0 & I \end{bmatrix} \leq 0, \forall k$$

Detailed derivations are omitted.

- Jump systems: Consider the following jump system

$$\begin{aligned}\xi_{k+1} &= A_{i_k} \xi_k + B_{i_k} u_k \\ y_k &= C_{i_k} \xi_k + D_{i_k} u_k\end{aligned}\tag{2.15}$$

where $\{i_k\}$ is the jump parameter sampled from a finite set $\{1, 2, \dots, n\}$ in an I.I.D. manner. Suppose $\Pr(i_k = i) = p_i$. If there exists a positive semidefinite matrix P such that

$$\sum_{i=1}^n \left(p_i \begin{bmatrix} A_i^\top P A_i - P & A_i^\top P B_i \\ B_i^\top P A_i & B_i^\top P B_i \end{bmatrix} + p_i \begin{bmatrix} C_i & D_i \\ 0 & I \end{bmatrix}^\top \begin{bmatrix} I & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \begin{bmatrix} C_i & D_i \\ 0 & I \end{bmatrix} \right) \leq 0$$

then the system (2.15) satisfies $\sum_{k=0}^N \mathbb{E} \|y_k\|^2 \leq \gamma^2 \sum_{k=0}^N \mathbb{E} \|u_k\|^2 + \mathbb{E} \xi_0^\top P \xi_0$ for any N . The proof is again based on standard Lyapunov arguments. (Verify this yourself!) We can see the same trick has been applied again and again to obtain all these different results.

An interesting variant of Theorem 2.2. Suppose the input u_k satisfies a bound $\|u_k\| \leq U$ where U is a constant. Then we can modify Theorem 2.3 to show ξ_k converges to a ball around 0. The result is formally stated as follows.

Theorem 2.3. *If there exists a positive semidefinite matrix P such that*

$$\begin{bmatrix} A^\top P A - \rho^2 P & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \leq 0 \quad (2.16)$$

then for any ξ_0 and arbitrary input sequence $\{u_k\}$ satisfying $\|u_k\| \leq U$, the following inequality holds

$$\xi_k^\top P \xi_k \leq \rho^{2k} \xi_0^\top P \xi_0 + \frac{\gamma^2 U^2}{1 - \rho^2} \quad (2.17)$$

Proof: Based on the LMI condition, we have

$$\begin{bmatrix} \xi_k \\ u_k \end{bmatrix}^\top \left(\begin{bmatrix} A^\top P A - \rho^2 P & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \right) \begin{bmatrix} \xi_k \\ u_k \end{bmatrix} \leq 0 \quad (2.18)$$

Notice we have $\xi_{k+1}^\top P \xi_{k+1} = (A\xi_k + Bu_k)^\top P (A\xi_k + Bu_k) = \begin{bmatrix} \xi_k \\ u_k \end{bmatrix}^\top \begin{bmatrix} A^\top P A & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} \begin{bmatrix} \xi_k \\ u_k \end{bmatrix}$.

Therefore, we have

$$\begin{bmatrix} \xi_k \\ u_k \end{bmatrix}^\top \begin{bmatrix} A^\top P A - \rho^2 P & A^\top P B \\ B^\top P A & B^\top P B \end{bmatrix} \begin{bmatrix} \xi_k \\ u_k \end{bmatrix} = \xi_{k+1}^\top P \xi_{k+1} - \rho^2 \xi_k^\top P \xi_k$$

Eventually we have

$$\xi_{k+1}^\top P \xi_{k+1} \leq \rho^2 \xi_k^\top P \xi_k + \gamma^2 \|u_k\|^2 \leq \rho^2 \xi_k^\top P \xi_k + \gamma^2 U^2 \quad (2.19)$$

We can iterate the above inequality to obtain the desired conclusion. ■

As a matter of fact, if $A^\top P A - \rho^2 P < 0$, then there exists a sufficiently large γ that (2.16) holds. This can be proved using Schur complement lemma.

Continuous-time results. There are similar results for continuous-time systems. For example, consider a continuous-time LTI system

$$\begin{aligned}\dot{\xi} &= A\xi + Bu \\ y &= C\xi + Du\end{aligned}$$

If there exists a positive semidefinite matrix P such that

$$\begin{bmatrix} A^\top P + PA & PB \\ B^\top P & 0 \end{bmatrix} + \begin{bmatrix} C & D \\ 0 & I \end{bmatrix}^\top \begin{bmatrix} I & 0 \\ 0 & -\gamma^2 I \end{bmatrix} \begin{bmatrix} C & D \\ 0 & I \end{bmatrix} \leq 0$$

then we have

$$\int_{t=0}^T \|y(t)\|^2 dt \leq \gamma^2 \int_{t=0}^T \|u(t)\|^2 dt + \xi_0^\top P \xi_0$$

A key fact used in the proof is $\dot{V} = \dot{\xi}^\top P \xi + \xi^\top P \dot{\xi} = \begin{bmatrix} \xi \\ u \end{bmatrix}^\top \begin{bmatrix} A^\top P + PA & PB \\ B^\top P & 0 \end{bmatrix} \begin{bmatrix} \xi \\ u \end{bmatrix}$.

2.2 Convergence Analysis in Optimization

In this section, we will apply the linear system tools reviewed in the last section to analyze some simple iterative algorithms in optimization. Hopefully this convinces you that there is some similarity between analyzing a control system and analyzing an optimization algorithm.

- **Example 1: Optimization of a Quadratic Function.** Suppose we want to minimize a quadratic function:

$$\min_{x \in \mathbb{R}^p} g(x) = \frac{1}{2} x^\top Q x + q^\top x + r \quad (2.20)$$

where $Q > 0$. Clearly we have $\nabla f(x) = Qx + q$. Suppose x^* is the optimal point, i.e. $Qx^* + q = 0$. Since the function is strongly convex ($Q > 0$), there exists a unique x^* for this problem. To find such an optimal point, we can start from an arbitrary initial point x_0 and then run an iterative gradient-based algorithm. The simplest algorithm we can run is the gradient descent method $x_{k+1} = x_k - \alpha \nabla f(x_k)$. Since $\nabla f(x) = Qx + q$, the gradient descent method is equivalent to the following linear system:

$$x_{k+1} = x_k - \alpha(Qx_k + q) = (I - \alpha Q)x_k - \alpha q = (I - \alpha Q)x_k + \alpha Qx^*$$

which can be rewritten as $x_{k+1} - x^* = (I - \alpha Q)(x_k - x^*)$. Hence we can denote $\xi_k = x_k - x^*$. Analyzing how fast x_k converges to x^* is equivalent to analyzing how fast ξ_k converges to 0. Hence if there exists a positive definite P such that

$$(I - \alpha Q)^\top P (I - \alpha Q) - \rho^2 P \leq 0,$$

then the gradient descent method is guaranteed to converge at the rate ρ for the above quadratic optimization problem. Now we assume f is L -smooth and m -strongly convex, i.e. $mI \leq Q \leq LI$. Then we can just use $P = I$ to recover the standard convergence rate of the gradient descent method. We can also consider more complicated algorithms such as the Heavy-ball method and the Nesterov's accelerated method. The heavy-ball method applies the following update:

$$x_{k+1} = x_k - \alpha \nabla f(x_k) + \beta(x_k - x_{k-1}) = x_k - \alpha(Qx_k + q) + \beta(x_k - x_{k-1})$$

which can be rewritten as

$$\begin{bmatrix} x_{k+1} - x^* \\ x_k - x^* \end{bmatrix} = \begin{bmatrix} (1 + \beta)I - \alpha Q & -\beta I \\ I & 0 \end{bmatrix} \begin{bmatrix} x_k - x^* \\ x_{k-1} - x^* \end{bmatrix}$$

Hence the Heavy-ball method becomes an LTI system with

$$\xi_k = \begin{bmatrix} x_k - x^* \\ x_{k-1} - x^* \end{bmatrix}, \quad A = \begin{bmatrix} (1 + \beta)I - \alpha Q & -\beta I \\ I & 0 \end{bmatrix}$$

The Nesterov's method applies the following update:

$$\begin{aligned} x_{k+1} &= x_k - \alpha \nabla f((1 + \beta)x_k - \beta x_{k-1}) + \beta(x_k - x_{k-1}) \\ &= x_k - \alpha(Q((1 + \beta)x_k - \beta x_{k-1}) + q) + \beta(x_k - x_{k-1}) \end{aligned}$$

which can be rewritten as

$$\begin{bmatrix} x_{k+1} - x^* \\ x_k - x^* \end{bmatrix} = \begin{bmatrix} (1 + \beta)(I - \alpha Q) & -\beta(I - \alpha Q) \\ I & 0 \end{bmatrix} \begin{bmatrix} x_k - x^* \\ x_{k-1} - x^* \end{bmatrix}$$

Hence the Heavy-ball method becomes an LTI system with

$$\xi_k = \begin{bmatrix} x_k - x^* \\ x_{k-1} - x^* \end{bmatrix}, \quad A = \begin{bmatrix} (1 + \beta)(I - \alpha Q) & -\beta(I - \alpha Q) \\ I & 0 \end{bmatrix}$$

You will be asked to discuss how to apply the LTI system theory to recover standard convergence rates of the above three methods at the beginning of Lecture 3. Be prepared!

- **Example 2: Inaccurate gradients.** Suppose the gradient information is corrupted as $\nabla f(x_k) + u_k$ where we only know some upper bounds on the norm of u_k . Can we apply Theorem 2.3 to show that the gradient descent method still converges to a ball around x^* ? The answer is yes! The gradient descent method becomes

$$x_{k+1} - x^* = (I - \alpha Q)(x_k - x^*) - \alpha u_k.$$

In this case, x_k does not converge to x^* exactly, but it is going to stay within a ball around x^* . Theorem 2.3 can be directly applied!

There are many other examples. For example, SGD on ridge regression becomes an MJLS. The policy evaluation in reinforcement learning can be viewed as an LTI system. A question is how to extend the above analysis to general settings involving more general cost functions and stochastic algorithms. We will talk about this in the next two lectures.