

Massachusetts Institute of Technology

Department of Electrical Engineering and Computer Science

6.242, Fall 2004: MODEL REDUCTION \*

## Balanced Truncation<sup>1</sup>

This lecture introduces *balanced truncation* for LTI systems: an important projection model reduction method which delivers high quality reduced models by making an extra effort in choosing the projection subspaces.

### 5.1 The classical balanced truncation algorithm

This section describes the basic algorithm for balanced truncation of continuous time state space models.

#### 5.1.1 Motivation: removal of (almost) uncontrollable/unobservable modes

Removing an unobservable or an uncontrollable mode is an easy way of reducing the dimension of the state vector in a state space model. For example, system

$$\begin{aligned}\dot{x}_1 &= -x_1 + x_2 + f, \\ \dot{x}_2 &= -2x_2, \\ y &= x_1 + x_2,\end{aligned}$$

can be replaced by

$$\dot{x} = -x + f, \quad y = x$$

---

\*©A. Megretski, 2004

<sup>1</sup>Version of September 27, 2004

without changing its transfer function. In this state space model, with

$$A = \begin{bmatrix} -1 & 1 \\ 0 & -2 \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad C = [1 \quad 1], \quad D = 0,$$

the controllability matrix

$$M_c = [B \quad AB] = \begin{bmatrix} 1 & 0 \\ -1 & 0 \end{bmatrix}$$

satisfies  $pM_c = 0$  for  $p = [0 \quad 1]$ , and hence the variable  $px = x_2$  represents an *uncontrollable* mode. The removal of such mode can be viewed as a canonical projection model reduction

$$A \mapsto \hat{A} = UAV, \quad B \mapsto \hat{B} = UB, \quad C \mapsto \hat{C} = CV,$$

where the columns of  $V$  form a basis in the column range of  $M_c$  (which is the same as the null space of  $p$ ), and  $U$  can be selected quite arbitrarily subject to the usual constraint  $UV = I$ .

Strictly speaking, the example above cannot even be considered as “model reduction”, as the orders of the original and the projected systems are both equal to 1. A more interesting situation is represented by the perturbed model

$$\begin{aligned} \dot{x}_1 &= -x_1 + x_2 + f, \\ \dot{x}_2 &= -2x_2 + \epsilon f, \\ y &= x_1 + x_2, \end{aligned}$$

(same  $A, C, D$  but a modified  $B$ ), where  $\epsilon > 0$  is a parameter. Intuitively, one can expect that, when  $\epsilon > 0$  is small enough,

$$\dot{x} = -x + f, \quad y = x$$

is still a good reduced model. This expectation can be related to the fact that  $x_2$  is difficult to control by  $f$  when  $\epsilon > 0$  is small. One can say that the *mode*  $x_2 = px$ , which corresponds to the left (row) eigenvector of the  $A$  matrix ( $pA = -2p$  in this case), is almost uncontrollable, which can be seen directly from the transfer function

$$G(s) = \frac{s + 2 + \epsilon}{(s + 2)(s + 1)} = \frac{1 + \epsilon}{s + 1} - \frac{\epsilon}{s + 2},$$

which has a small coefficient at  $(s + 2)^{-1}$  in its partial fraction expansion.

One can attempt to introduce a measure of importance of an LTI system's mode (understood as a pole  $a$  of system's transfer function  $G(s)$ ) as something related to the absolute value of the coefficient with which  $(s - a)^{-1}$  enters the partial fraction expansion of  $G$ . However, it is rarely a good idea to base a model reduction algorithm solely on removal of "unimportant" system modes. For example, both modes  $a = 1 - \epsilon$  and  $a = 1 + \epsilon$  of the LTI system with transfer function

$$H(s) = \frac{1}{s + 1 - \epsilon} + \frac{1}{s + 1 + \epsilon},$$

where  $\epsilon > 0$  is small, are equally important, and none can be removed without introducing a huge model reduction error, despite the fact that a very good reduced model  $\hat{H}$  is given by

$$\hat{H}(s) = \frac{2}{s + 1}.$$

Balanced truncation is based on introducing a special joint measure of controllability and observability for *every* vector in the state space of an LTI system. Then, the reduced model is obtained by removing those components of the state vector which have the lowest importance factor in terms of this measure.

### 5.1.2 Observability measures

Consider a state space model

$$\dot{x}(t) = Ax(t) + Bf(t), \quad y(t) = Cx(t) + Df(t), \quad (5.1)$$

where  $A$  is an  $n$ -by- $n$  Hurwitz matrix (all eigenvalues have negative real part). When  $f(t) \equiv 0$  for  $t \geq 0$ , the value of the output  $y(t)$  at a given moment  $t$  is uniquely defined by  $x(0)$ , and converges to zero exponentially as  $t \rightarrow +\infty$ . Hence the integral

$$E_o = \int_0^{\infty} |y(t)|^2 dt,$$

measuring the "observable output energy" accumulated in the initial state, is a function of  $x(0)$ , i.e.  $E_o = E_o(x(0))$ . Moreover, since  $y(t)$  is a *linear* function of  $x(0)$ ,  $E_o$  will be a *quadratic form* with respect of  $x(0)$ , i.e.

$$E_o(x(0)) = x(0)'W_o x(0)$$

for some symmetric matrix  $W_o$ .

The quadratic form  $E_o = E_o(x(0))$  can be used as a *measure of observability* defined on the state space of system (5.1). In particular,  $E_o(x(0)) = 0$  whenever

$$M_o x(0) = [C; CA; \dots; CA^{n-1}]x(0) = 0.$$

The positive semidefinite matrix  $W_o$  of the quadratic form  $E_o$  is called the *observability Gramian*<sup>2</sup>. Since, by the definition,  $E_o(x(0)) \geq 0$  for all  $x(0)$ , the matrix  $W_o$  is always positive semidefinite. Indeed,  $W_o > 0$  whenever the pair  $(C, A)$  is observable.

It is important to notice that the word “system state” actually includes two different meanings. One meaning is that of a *primal* state, which, for a general model (5.1), is a column  $n$ -vector. For example, for a state space model

$$\dot{x}_1 = -x_1 + f, \quad \dot{x}_2 = -x_2 + f, \quad y = x_1 + x_2, \quad (5.2)$$

a primal state is a particular column vector value of  $x(t) \in \mathbf{R}^2$ , such as  $x(-1.2) = [1; -7]$ . Such column vector values are more precisely referred to as the *primal* states of (5.2).

On the other hand, one frequently refers to (5.2) as a “system with two states” (despite the fact that the set  $\mathbf{R}^2$  has an infinite number of elements), and, in this context,  $x_1 = x_1(t)$  and  $x_2 = x_2(t)$  can be referred to as the two states of the system. Let us call such states the *dual* states. For a general model (5.1), a dual state is a particular linear combination  $x_p = px(t)$  of the scalar components of  $x = x(t)$ , defined by a row  $n$ -vector  $p$ . For example, in (5.2), the dual states  $x_1 = x_1(t)$  and  $x_2 = x_2(t)$  are defined by row vectors  $p = [1 \ 0]$  and  $p = [0 \ 1]$  respectively.

Therefore, it is natural to ask for a definition of an observability measure of a given *dual* state of (5.1). It is defined as

$$E^o(p) = \inf_{x_0: px_0=1} E_o(x_0) \quad \text{for } p \neq 0,$$

i.e. as the minimal output energy which can be observed for  $t \geq 0$  when  $px(0) = 1$ . Note that infimum over an empty set equals plus infinity, hence  $E^o(0) = \infty$ . When the pair  $(C, A)$  is observable, and hence  $W_o > 0$  is invertible, the dual observability measure is given by

$$E^o(p) = \frac{1}{pW_o^{-1}p'} \quad \text{for } p \neq 0.$$

The following theorem is frequently utilized for computing  $W_o$  numerically.

**Theorem 5.1**  $W_o$  is the unique solution of the Lyapunov equation

$$W_o A + A' W_o + C' C = 0. \quad (5.3)$$

---

<sup>2</sup>Whether this should be spelled as “Grammian” or “Gramian”, is unclear

**Proof** Since system (5.1) is time invariant, the identity

$$x(0)'W_o x(0) = \int_0^\infty |Cx(\tau)|^2 d\tau$$

implies

$$x(t)'W_o x(t) = \int_t^\infty |Cx(\tau)|^2 d\tau.$$

Differentiating the second identity with respect to  $t$  at  $t = 0$  yields

$$2x(0)'W_o Ax(0) = -|Cx(0)|^2$$

for all  $x(0) \in \mathbf{R}^n$ . Comparing the coefficients on both sides of the quadratic identity yields (5.3). ■

Finding a numerical solution of (5.3) is not easy when  $n$  is about  $10^4$  and larger. In such situation, Theorem 5.1 can be used as a basis for finding an approximation of  $W_o$ .

It is important to understand that the observability measure alone should not be the only numerical test for choosing which states to eliminate in a model reduction procedure. Instead, a combination of observability *and* a controllability measures, to be introduced in the next subsection, should be used.

### 5.1.3 Controllability measures

Since  $A$  is a Hurwitz matrix, every input signal  $f = f(t)$  of finite energy, i.e. such that

$$\|f\|^2 = \int_{-\infty}^\infty |f(t)|^2 dt < \infty,$$

corresponds to a unique initial condition  $x(0)$  in (5.1) for which the corresponding solution  $x = x(t)$  satisfies  $x(t) \rightarrow 0$  as  $t \rightarrow -\infty$ . This solution is given by

$$x(t) = \int_0^\infty e^{A\tau} B f(t - \tau) d\tau,$$

where  $e^M$  denotes the *matrix exponent* of a square matrix  $M$ . One can say that input  $f = f(t)$  *drives* the system state from  $x(-\infty) = 0$  to  $x(0) = X(f(\cdot))$ .

Let  $p$  be a 1-by- $n$  row vector, so that the product  $px(t)$  is a dual state of (5.1) – a linear combination of components of the state space vector. The (dual) controllability measure  $E^c = E^c(p)$  is defined as the maximal value of  $|px(0)|^2$  which can be achieved by using an input  $f = f(t)$  of unit energy:

$$E^c(p) = \max\{|pX(f(\cdot))|^2 : \|f\| \leq 1\}.$$

Accordingly, the *primal* controllability measure  $E_c = E_c(x_0)$  is defined for a column vector  $x_0 \in \mathbf{R}^n$  as

$$E_c(x_0) = \inf_{p: px_0=1} E^c(p).$$

The following statement describes some basic properties of these controllability measures.

**Theorem 5.2** *Assuming that  $A$  is an  $n$ -by- $n$  Hurwitz matrix.*

(a)  $E^c(p) = pW_cp'$  is a quadratic form with the coefficient matrix

$$W_c = \int_0^\infty e^{At}BB'e^{A't}dt.$$

(b)  $W_c$  is the unique solution of the Lyapunov equation

$$AW_c + W_cA' = -BB'. \quad (5.4)$$

(c) A given state  $x_0 \in \mathbf{R}^n$  is reachable from zero if and only if  $E_c(x_0) > 0$  or, equivalently, the equation  $W_cp' = x_0$  has a solution  $p'$ . In this case  $E_c(x_0) = px_0$  is the minimum of  $\|f(\cdot)\|^2$  subject to  $X(f(\cdot)) = x_0$ .

**Proof** To prove (a), note that

$$\max_{\|f\| \leq 1} \int -\infty^\infty g(t)'f(t)dt = \|g\|,$$

hence

$$E^c(p) = \int_0^\infty |pe^{At}B|^2 dt = pW_cp'.$$

Statement (b) is actually a re-wording of Theorem 5.1, with  $C$  replaced by  $B'$ ,  $A$  replaced by  $A'$ ,  $W_o$  replaced by  $W_c$ , and  $x(0)$  replaced by  $p$ .

To prove (c), consider first the case when equation  $W_cp' = x_0$  has no solution  $p$ . Then there exists a row vector  $p_0$  such that  $p_0W_c = 0$  but  $p_0x_0 \neq 0$ . Here the equality means that  $|p_0X(f(\cdot))|^2$  equals zero for every finite energy signal  $f = f(t)$ . Since, from the inequality,  $|p_0x_0|^2 > 0$ , the state  $x_0$  is not reachable from zero.

Now assume that  $x_0 = W_cp'$  for some  $p$ . Then  $\|f\|^2 \geq pW_cp' = px_0$  whenever  $x_0 = X(f(\cdot))$ . On the other hand, for

$$f(t) = \begin{cases} B'e^{-A't}p', & t \leq 0, \\ 0, & t > 0, \end{cases}$$

we have  $\|f\|^2 = px_0$  and  $x_0 = X(f(\cdot))$ . ■

When the pair  $(A, B)$  is controllable, and, hence,  $W_c > 0$ , the primal controllability measure  $E_c = E_c(x_0)$  can be expressed as

$$E_c(x_0) = \frac{1}{x_0' W_c^{-1} x_0} \text{ for } x_0 \neq 0.$$

#### 5.1.4 Joint controllability and observability measures

The joint controllability and observability measures are defined as *products* of the corresponding controllability and observability measures:

$$E_{oc}(x_0) = E_o(x_0)E_c(x_0), \quad E^{oc}(p) = E^o(p)E^c(p).$$

For controllable and observable systems  $W_c$  and  $W_o$  are positive definite, and the formulae can be simplified to

$$E_{oc}(x_0) = \frac{x_0' W_o x_0}{x_0' W_c^{-1} x_0} (x_0 \neq 0), \quad E^{oc}(p) = \frac{p W_c p'}{p W_o^{-1} p'} (p \neq 0).$$

For model reduction purposes, we are interested in finding a subspace of primal state vectors for which the *minimum* of the joint controllability and observability measure over all non-zero elements is *maximal*. A basis in this subspace will yield columns of a projection matrix  $V$ . Similarly, we are interested in finding a subspace of dual state vectors for which the *minimum* of the joint controllability and observability measure over all non-zero elements is *maximal*. A basis in this subspace will yield rows of a projection matrix  $V$ .

The following theorem can be used in finding such  $V$  and  $U$ .

**Theorem 5.3** *Let  $W_c = L_c L_c'$  and  $W_o = L_o' L_o$  be the Choleski decompositions of the controllability and observability Gramians. Let*

$$\rho_1 \geq \cdots \geq \rho_r > \rho_{r+1} \geq \cdots \geq \rho_n \geq 0$$

*be the ordered eigenvalues of  $L_c' W_o L_c$  (possibly repeated). Let  $\psi_1, \dots, \psi_r$  be the corresponding first  $r$  normalized eigenvectors of  $L_c' W_o L_c$ , i.e.*

$$L_c' W_o L_c \psi_i = \rho_i \psi_i, \quad |\psi_i|^2 = 1, \quad \psi_i' \psi_k = 0 \text{ for } i \neq k.$$

*Let  $\sigma_i = \rho_i^{1/2}$ .*

(a)  $\rho_1 \geq \dots \geq \rho_r$  are also the eigenvalues of  $L_o W_c L'_o$ , and the corresponding normalized row eigenvectors can be defined by

$$\phi_i = \sigma_i^{-1} \psi'_i L'_c L'_o \quad (i = 1, \dots, r).$$

(b) The set of all linear combinations of vectors  $L_c \psi_i$  is the only  $r$ -dimensional linear subspace  $\mathcal{V}$  in  $\mathbf{R}^n$  such that  $E_{oc}(v) \geq \rho_r$  for every  $v \in \mathcal{V}$ .

(c) The set of all linear combinations of row vectors  $\phi_i L_o$  is the only  $r$ -dimensional linear subspace  $\mathcal{U}$  of row  $n$ -vectors such that  $E^{oc}(u) \geq \rho_k$  for every  $u \in \mathcal{U}$ .

(d)  $UV = I_r$ , where

$$V = L_c \begin{bmatrix} \psi_1 \sigma_1^{-1/2} & \dots & \psi_r \sigma_r^{-1/2} \end{bmatrix}, \quad U = \begin{bmatrix} \sigma_1^{-1/2} \phi_1 \\ \vdots \\ \sigma_r^{-1/2} \phi_r \end{bmatrix} L_o.$$

The proof of the theorem is obtained by inspection.

### 5.1.5 Balanced truncation

The projection model reduction algorithm which uses the projection matrices  $U, V$  described in the previous subsection is called *balanced truncation*. The “balancing” terminology comes from the following trivial observation.

**Theorem 5.4** *Assume that system (5.1) is both controllable and observable. Then  $W_c = L_c L'_c > 0$  and  $W_o = L'_o L_o > 0$  are positive definite. Moreover, if  $\Psi$  is the orthogonal matrix which columns are the ordered eigenvectors of  $L'_c W_o L_c$ , and  $\Sigma$  is the diagonal matrix with numbers  $\sigma_i$  on the diagonal, then the linear state transformation  $x = Sz$ , with  $S = L_c \Psi \Sigma^{-1/2}$ , yields an equivalent state space model with the coefficient matrices*

$$\bar{A} = S^{-1} A S, \quad \bar{B} = S^{-1} B, \quad \bar{C} = C S,$$

for which both the controllability and observability Gramians equal  $\Sigma$ .

State space models for which  $W_o = W_c$  are equal diagonal matrices with ordered positive diagonal entries  $\sigma_k$  are called *balanced*. The numbers  $\sigma_k$  are called the (Hankel) *singular numbers* of the corresponding system. The (canonical) method of balanced truncation is based on removing the states of a balanced realization which correspond to singular numbers below a certain threshold. Indeed, a practical implementation does not have to involve a calculation of a complete balanced model: only the projection matrices  $U, V$  are necessary.



### 5.1.6 Example

Let

$$G(s) = \frac{1}{s+1-\epsilon} + \frac{1}{s+1+\epsilon},$$

where  $\epsilon > 0$  is a small parameter. A state space model is

$$A = \begin{bmatrix} -1+\epsilon & 0 \\ 0 & -1-\epsilon \end{bmatrix}, \quad B = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad C = [1 \quad 1], \quad D = 0.$$

The controllability Gramians are given by

$$W_o = W_c = \frac{1}{2} \begin{bmatrix} \frac{1}{1-\epsilon} & 1 \\ 1 & \frac{1}{1+\epsilon} \end{bmatrix}.$$

The Hankel singular numbers are the eigenvalues of  $W_o = W_c$ , and equal

$$\sigma_{1,2} = \frac{1}{2} \frac{1 \pm \sqrt{1-\epsilon^2 + \epsilon^4}}{1-\epsilon^2}.$$

The corresponding eigenvectors are

$$v_{1,2} = \begin{bmatrix} 1 \\ 2\sigma_{1,2} - \frac{1}{1-\epsilon} \end{bmatrix}.$$

The dominant eigenvector defines

$$V \approx \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}, \quad U \approx \begin{bmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \end{bmatrix}.$$

The resulting reduced system is approximately given by

$$\hat{A} \approx -1, \quad \hat{B} \approx \sqrt{2}, \quad \hat{C} \approx \sqrt{2}.$$

## 5.2 Properties of balanced truncation

In this section, basic properties of reduced models obtained via balanced truncation (exact or approximate) are discussed.

### 5.2.1 Approximate Gramians

The two most expensive phases of numerical calculations associated with the canonical method of balanced truncation are finding the Gramians  $W_o, W_c$ , and calculating the dominant eigenvectors  $\psi_i$  of  $L_i'W_oL_c$ . for systems of large dimensions (more than  $10^4$  states) finding the Gramians exactly becomes difficult.

As a viable alternative, lower and upper bounds of the Gramians can be used to provide *provably* reliable results. Here by *lower bounds of Gramians*  $W_o, W_c$  we mean positive semidefinite matrices  $W_o^-, W_c^-$  for which the inequalities

$$W_o \geq W_o^-, \quad W_c \geq W_c^-$$

are guaranteed. The definition of upper bounds will be more strict: by *upper bounds of Gramians*  $W_o, W_c$  defined by the Lyapunov equalities

$$W_oA + A'W_o = -C'C, \quad AW_c + W_cA' = -BB',$$

where  $A$  is a Hurwitz matrix, we mean solutions  $W_o^+, W_c^+$  of the corresponding Lyapunov inequalities

$$W_o^+A + A'W_o^+ \leq -C'C, \quad AW_c^+ + W_c^+A' \leq -BB'.$$

These inequalities imply that  $W_o^+ \geq W_o$  and  $W_c^+ \geq W_c$ , but the inverse implication is not always true.

The following simple observation can be used to produce lower bounds of the Gramians.

**Theorem 5.5** *Let  $A$  be an  $n$ -by- $n$  Hurwitz matrix. Let  $F$  be an  $n$ -by- $m$  matrix. For  $s \in \mathbf{C}$  with  $\text{Re}(s) > 0$  define*

$$a = a(s) = (sI_n - A)^{-1}(\bar{s}I_n + A), \quad b = b(s) = \sqrt{2\text{Re}(s)}(sI_n - A)^{-1}B.$$

*Then*

- (a)  *$a$  is a Schur matrix (all eigenvalues strictly inside the unit disc);*
- (b) *an  $n$ -by- $n$  matrix  $P$  is a solution of the “continuous time” Lyapunov equation*

$$AP + PA' = -BB'$$

*if and only if it is a solution of the “discrete time” Lyapunov equation*

$$P = aPa' + bb';$$

(c) the matrix  $P$  from (b) is the limit

$$P = \lim_{k \rightarrow \infty} P_k,$$

where the symmetric matrices  $P_0 \leq P_1 \leq P_2 \leq \dots \leq P$  are defined by

$$P_0 = 0, \quad P_{k+1} = a(s_k)P_k a(s_k)' + b(s_k)b(s_k)',$$

and  $\{s_k\}$  is a sequence of complex numbers contained in a compact subset of the open right half plane.

The theorem reduces finding a lower bound of a solution of a Lyapunov equation to solving systems of linear equations (used to produce  $b(s_k)$  and the products  $a(s_k)P_k a(s_k)'$ ).

Calculation of an upper bound of a Gramian could be more tricky. One approach to finding such upper bounds relies on having a valid *energy function* for  $A$  available.

Indeed, assume that  $Q = Q'$  satisfies

$$AQ + QA' < 0.$$

If  $W_c^-$  is a good lower bound of the controllability Gramian  $W_c$ , defined by

$$AW_c + W_c A' = -BB',$$

then

$$AW_c^- + W_c^- A' \approx -BB',$$

and hence

$$A(W_c^- + \epsilon Q) + (W_c^- + \epsilon Q)A' \leq -BB'$$

for some  $\epsilon > 0$  (which will, hopefully, be small enough). Then  $W_c^-$  and  $W_c^- + \epsilon Q$  are a lower and an upper bound for the controllability Gramian  $W_c$ .

### 5.2.2 Lower bounds for model reduction error

A major benefit of doing balanced truncation is given by a *lower bound* of the error of *arbitrary* reduced models (not only those produced via balanced truncation).

**Theorem 5.6** *Let  $W_o^- = F_o' F_o$  and  $W_c^- = F_c' F_c$  be lower bounds of the observability and controllability Gramians  $W_o, W_c$  of a stable LTI model  $G = G(s)$ . Let*

$$\sigma_1^- \geq \sigma_2^- \geq \dots \geq 0$$

be the ordered singular numbers of  $F_o F_c$ . Then  $\sigma_k^-$  is a lower bound for the  $k$ -th Hankel singular number  $\sigma_k = \sigma_k(G)$  of the system, and

$$\|G - \hat{G}\|_\infty \geq \sigma_k^-$$

for every system  $\hat{G}$  of order less than  $k$ .

**Proof** Let  $\mathcal{Z}_k$  denote the subspace spanned by the  $k$  dominant eigenvectors of  $F_c' W_o^- F_c$ , i.e.

$$|F_o F_c z| \geq \sigma_k^- |z| \quad \forall z \in \mathcal{Z}_k.$$

Since  $W_c \geq F_c F_c'$ , every vector  $F_c z$  lies in the range of  $W_c$ , and  $q' F_c z \leq |z|^2$  whenever  $F_c z = W_c q$ . Hence every state  $x(0) = F_c z$  can be reached from  $x(-\infty) = 0$  using a minimal energy input  $f = f_z(t)$  (depending linearly on  $z$ ) of energy not exceeding  $|z|^2$ . On the other hand, every state  $x(0) = F_c z$  with  $z \in \mathcal{Z}_k$  will produce at least  $|F_o F_c z|^2 \geq (\sigma_k^-)^2 |z|^2$  of output energy. Since  $\hat{G}$  is a linear system of order less than  $k$ , there exists at least one non-zero  $z \in \mathcal{Z}_k$  for which the input  $f = f_z(t)$  produces a zero state  $\hat{x}(0) = 0$  at zero time. Then, assuming the input is zero for  $t > 0$ , the error output energy is at least  $(\sigma_k^-)^2 |z|^2$ . Since the testing input energy is not larger than  $|z|^2 > 0$ , this yields an energy gain of  $(\sigma_k^-)^2$ , which means that  $\|G - \hat{G}\|_\infty \geq \sigma_k^-$ . ■

### 5.2.3 Upper bounds for balanced truncation errors

The result from the previous subsection states that, for a stable LTI system  $G$ , no method can produce a reduced model  $\hat{G}$  of order less than  $k$  such that the H-Infinity error  $\|G - \hat{G}\|_\infty$  is less than the  $k$ -th singular number  $\sigma_k = \sigma_k(G)$ . The statement is easy to apply, since lower bounds  $\sigma_k^-$  of  $\sigma_k$  can be computed by using lower bounds of system Gramians.

The following theorem gives an *upper bound* of model reduction error for the *exact* implementation of the balanced truncation method.

**Theorem 5.7** *Let  $\sigma_1 > \sigma_2 > \dots > \sigma_h$  be the ordered set of different Hankel singular numbers of a stable LTI system  $G$ . Let  $\hat{G}$  be the reduced model obtained by removing the states corresponding to singular numbers not larger than  $\sigma_k$  from a balanced realization of  $G$ . Then  $\hat{G}$  is stable, and satisfies*

$$\|G - \hat{G}\|_\infty \leq 2 \sum_{i=k}^h \sigma_i.$$

The utility of Theorem 5.7 in practical calculations of H-Infinity norms of model reduction errors is questionable: an *exact* calculation of the H-Infinity norm is possible at about the same cost, and the upper bound itself can be quite conservative. Nevertheless, the theorem provides an important reassuring insight into the potential of balanced truncation: since the singular numbers of exponentially stable LTI systems decay exponentially, the upper bound of Theorem 5.7 is not expected to be much larger than the lower bound.

For example, for a system with singular numbers  $\sigma_i = 2^{-i}$ , a  $k$ th order reduced model cannot have quality better than  $2^{-k-1}$ , and exact balanced truncation is guaranteed to provide quality of at least  $2^{-k+1}$ .

The proof of Theorem 5.7 is based on estimating the quality of balanced truncation in the case when only the states of a balanced realization corresponding to the smallest Hankel singular number are removed, which is done in the following technical lemma.

**Lemma 5.1** *Let  $W = W' > 0$  be a positive definite symmetric  $n$ -by- $n$  matrix satisfying the Lyapunov equalities*

$$WA + A'W = -C'C, \quad AW + WA' = -BB'. \quad (5.5)$$

Assume that  $W, A, B, C$  can be partitioned as

$$W = \begin{bmatrix} P & 0 \\ 0 & \sigma I_r \end{bmatrix}, \quad A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad B = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}, \quad C' = \begin{bmatrix} C'_1 \\ C'_2 \end{bmatrix},$$

where  $A_{22}$  is an  $r$ -by- $r$  matrix, and matrices  $B_2, C'_2$  have  $r$  rows. Then

(a) *transfer matrices*

$$G(s) = C(sI_n - A)^{-1}B, \quad G_1(s) = C_1(sI_{n-r} - A_{11})^{-1}B_1$$

*are stable;*

(b) *the Lyapunov equalities*

$$PA_{11} + A'_{11}P = -C'_1C_1, \quad A_{11}P + PA'_{11} = -B_1B'_1$$

*are satisfied;*

(c)  $\|G - G_1\|_\infty \leq 2\sigma$ .

**Proof** It is sufficient to consider the case when the dimension  $m$  of  $f = f(t)$  equals the dimension  $k$  of  $y = y(t)$ . (If  $m < k$ , add zero columns to  $B$ , if  $m > k$ , add zero rows to  $C$ .) First note that re-writing (5.5) in terms of the blocks  $A_{ik}, B_i, C_i$  yields

$$PA_{11} + A'_{11}P = -C'_1C_1, \quad (5.6)$$

$$PA_{12} + \sigma A'_{21} = -C'_1C_2, \quad (5.7)$$

$$\sigma(A_{22} + A'_{22}) = -C'_2C_2, \quad (5.8)$$

$$A_{11}P + PA'_{11} = -B_1B'_1, \quad (5.9)$$

$$\sigma A_{12} + PA'_{21} = -B_1B'_2, \quad (5.10)$$

$$\sigma(A_{22} + A'_{22}) = -B_2B'_2. \quad (5.11)$$

Note (5.8) together with (5.11) implies that  $C'_2 = B_2\theta$  for some unitary matrix  $\theta$ . Also, (5.6) and (5.9) prove (b).

To prove (a), note that for every complex eigenvector  $v \neq 0$  of  $A$ ,  $Av = sv$  for some  $s \in \mathbf{C}$ , multiplication of the first equation in (5.5) by  $v'$  on the left and by  $v$  on the right yields

$$2\operatorname{Re}(s)v'Wv = -|Cv|^2.$$

Hence either  $\operatorname{Re}(s) < 0$  or  $\operatorname{Re}(s) = 0$  and  $Cv = 0$ . Hence all unstable modes of  $A$  are unobservable, and  $G = G(s)$  has no unstable poles. The same proof applies to  $G_1$ , since  $A_{11}$  satisfies similar Lyapunov equations.

To prove (c), consider the following state space model of the error system  $G - G_1$ :

$$\begin{aligned} \dot{x}_1 &= A_{11}x_1 + A_{12}x_2 + B_1f, \\ \dot{x}_2 &= A_{21}x_1 + A_{22}x_2 + B_2f, \\ \dot{x}_3 &= A_{11}x_3 + B_1f, \\ e &= C_1x_1 + c_2x_2 - C_1x_3. \end{aligned}$$

It would be sufficient to find a positive definite quadratic form  $V(x) = x'Hx$  such that

$$\psi(t) = 4\sigma^2|f(t)|^2 - |e(t)|^2 - \frac{dV(x(t))}{dt} \geq 0$$

for all solutions of system equations. Indeed, such Lyapunov function  $V$  can be readily presented, though there is no easy way to describe the intuitive meaning of its format:

$$V(x) = \sigma^2(x_1 + x_3)'P^{-1}(x_1 + x_3) + (x_1 - x_3)'P(x_1 - x_3) + 2\sigma|x_2|^2.$$

To streamline the derivation, introduce the shortcut notation

$$z = x_1 + x_3, \quad \Delta = x_1 - x_3, \quad \delta = C_1\Delta, \quad u = \sigma^{-1}B'_2x_2, \quad q = B'_1P^{-1}z.$$

The equations now take the form

$$\begin{aligned}\dot{\Delta} &= A_{11}\Delta + A_{12}x_2, \\ \dot{z} &= A_{11}z + A_{12}x_2 + 2B_1f, \\ \dot{x}_2 &= A_{22}x_2 + 0.5A_{21}(z + \Delta) + B_2f, \\ e &= C_1\Delta + \sigma\theta'u.\end{aligned}$$

We have

$$\begin{aligned}\psi &= 4\sigma^2|f|^2 - |C_1\Delta + \sigma\theta'u|^2 - 2\sigma^2z'P^{-1}(A_{11}z + A_{12}x_2 + 2B_1f) \\ &\quad - 2\Delta'P(A_{11}\Delta + A_{12}x_2) - 4\sigma x_2'[A_{22}x_2 + 0.5A_{21}(z + \Delta) + B_2f] \\ &= 4\sigma^2|f|^2 - |C_1\Delta + \sigma\theta'u|^2 + \sigma^2|q|^2 - 4\sigma^2q'f + \\ &\quad |\delta|^2 + 2\sigma^2|u|^2 - 4\sigma^2u'f - 2z'[\sigma^2P^{-1}A_{12} + \sigma A'_{21}]x_2 - 2\Delta'[PA_{12} + \sigma A'_{21}]x_2 \\ &= 4\sigma^2|f|^2 - |C_1\Delta + \sigma\theta'u|^2 + \sigma^2|q|^2 - 4\sigma^2q'f + |\delta|^2 \\ &\quad + 2\sigma^2|u|^2 - 4\sigma^2u'f + 2\sigma^2q'u + 2\sigma\delta'\theta'u \\ &\geq 4\sigma^2|f|^2 - |C_1\Delta + \sigma\theta'u|^2 + |\delta|^2 + 2\sigma^2|u|^2 + 2\sigma\delta'\theta'u - \sigma^2|u - 2f|^2 \\ &= 0\end{aligned}$$

(the first identity transformation utilized (5.6), (5.9), (5.8), the second identity used (5.7), (5.10), the third (inequality) applied minimization with respect to  $q$ , the last (identity) depended on  $\theta$  being a unitary matrix.  $\blacksquare$ )

It is important to realize that the lemma remains valid when the original Lyapunov equalities are replaced by the corresponding Lyapunov inequalities (of course, the equalities in (b) will get replaced by inequalities as well). Indeed, Lyapunov inequalities

$$WA + A'W \leq -C'C, \quad AW + WA' \leq -BB'$$

are equivalent to Lyapunov equalities

$$WA + A'W = -\tilde{C}'\tilde{C}, \quad AW + WA' = -\tilde{B}\tilde{B}'$$

where

$$\tilde{B} = \begin{bmatrix} B & B_\delta \end{bmatrix}, \quad \tilde{C} = \begin{bmatrix} C \\ C_\delta \end{bmatrix},$$

and  $B_\delta, C_\delta$  are chosen appropriately. Note that  $G(s)$  is a left upper corner block in the transfer matrix

$$\tilde{G} = \tilde{C}(sI_n - A)^{-1}\tilde{B}.$$

Since H-Infinity norm of a transfer matrix is not smaller than H-Infinity norm of its block, applying Lemma 5.1 to the system defined by  $A, \tilde{B}, \tilde{C}$  yields the stated generalization.