

**A Thick-Restarted Krylov Subspace Projection Method for Model Order Reduction**

By

EFREM B. RENSI

B.S. (San Jose State University, San Jose) 2006

THESIS

Submitted in partial satisfaction of the requirements for the degree of

MASTER OF SCIENCE

in

Mathematics

in the

OFFICE OF GRADUATE STUDIES

of the

UNIVERSITY OF CALIFORNIA

DAVIS

Approved:

---

Roland Freund

---

Robert Guy

---

Zhaojun Bai

Committee in Charge

2013

# Contents

<b>Contents</b>	<b>2</b>
<b>1 Introduction</b>	<b>6</b>
What is model order reduction?	6
1.1 Content overview	6
Chapter 2	7
Chapter 3	7
Chapter 4	8
Chapter 5	8
Additional contributions	9
<b>2 The model and its transfer-function</b>	<b>10</b>
2.1 The LTI-system model	10
2.2 System transfer-function	11
2.2.1 Moments	13
Moment-matching	13
2.2.2 Shift-invert representation of the transfer-function	14
Moment representation	15
Region of convergence for moment-matching	15
2.2.3 Invariant-subspaces	16
Generalized-eigenvalues and invariance	16
2.2.4 Pole-Residue representation	18
Poles and residues of the standard transfer-function formulation	18
Poles and residues from shifted transfer-function formulation	20
2.2.5 Pole-weight	21
Total system-mass	22
2.3 Reduced-order transfer-function via projection	24
<b>3 Order-reduction via Krylov-subspace projection</b>	<b>26</b>
3.1 Krylov-subspace projection methods	26
3.1.1 The Krylov subspace	27
3.2 The Arnoldi process	28
3.2.1 Complexity of Arnoldi (with MGS orthogonalization)	29
3.2.2 ROM size vs. construction cost	30
3.2.3 The Arnoldi relation	31

The remaining candidate-vector . . . . .	31
3.2.4 Approximate eigenvalues from Arnoldi relation . . . . .	32
3.3 Implicit vs. explicit Ritz-values and vectors . . . . .	33
Implicit Ritz-values . . . . .	33
Explicit Ritz-values . . . . .	34
3.4 Moment-matching property of Krylov-subspace projected ROM . . . . .	35
3.4.1 Moment matching of the implicitly-projected ROM . . . . .	36
3.4.2 Moment matching of the explicitly-projected ROM . . . . .	37
<b>4 Interpolation-point selection and resulting projection-bases . . . . .</b>	<b>41</b>
4.1 Multiple point moment-matching . . . . .	41
4.1.1 Merging bases . . . . .	42
4.1.2 Interpolation-point translation . . . . .	43
4.2 Complex expansion-points . . . . .	46
4.2.1 Producing a real basis for a complex Krylov-subspace . . . . .	47
Ruhe's method . . . . .	48
4.2.2 Using a real inner-product . . . . .	48
4.3 Equivalent-real formulations . . . . .	50
4.3.1 Equivalence of split complex and equivalent-real subspaces . . . . .	51
4.3.2 Reduced-order models via equivalent-real formulations . . . . .	54
via explicit projection . . . . .	54
via implicit projection . . . . .	55
<b>5 The band-Arnoldi process and a proposed thick-restarted variant . . . . .</b>	<b>56</b>
5.1 Band-Arnoldi algorithm . . . . .	56
5.1.1 Candidates/residual term . . . . .	58
5.1.2 Deflation term . . . . .	58
5.1.3 Residual norms . . . . .	60
5.2 Thick-restarting the Band-Arnoldi process . . . . .	62
5.2.1 Implicit restart . . . . .	63
Disclaimer . . . . .	63
Implicit restart . . . . .	64
With a general nearly-invariant subspace . . . . .	64
5.2.2 Intermediate ROM from a thick-restarted band-Arnoldi process . . . . .	68
5.3 Proposing a new model-reduction method . . . . .	70
5.4 Results . . . . .	73
5.4.1 ex308 . . . . .	74
ex308 Benchmarks . . . . .	75
ex308 thick-restart example 1 . . . . .	81
ex308 thick-restart example 2 . . . . .	88
5.4.2 ex1841 . . . . .	90
ex1841 test1 . . . . .	92
<b>References . . . . .</b>	<b>96</b>

## **Abstract**

The major contribution of this dissertation is a new thick-restarted Krylov method that allows for re-starting the process at different interpolation points while preserving moment matching and a controllable degree of numerical linear dependence of the resulting basis. It reduces computational cost by recycling a few vectors it from cycle to restarted cycle (as opposed to all of them), for each different interpolation-point. The algorithms comprising the method are outlined and their efficiency is demonstrated by applying it to selected test data sets.

## Acknowledgments

Thanks to my supportive parents Cornelia & Giuseppe Rensi, and the rest of my family!

Maria “Bem” Cayco!!!! who has been my friend, confidant, and mentor for years, helped me make it through my undergraduate years at San Jose State, helped me not shoot myself in the foot with graduate school applications.

My office-mate and colleague Yuji Nakatsukasa for his intricate knowlegde of numerical linear algebra and snacks from Kim’s market, Boba tea, and great company while sharing an office for 5 years!

Thanks to the Math departement staff, especially Sylvia Davis and Tina Deneena.

The Davis Bike Collective!

# Chapter 1

## Introduction

### What is model order reduction?

Model order-reduction (MOR), also known as model-reduction, is the process by which a mathematical model characterized by some number of parameters  $N$ , referred to here as the unreduced model (URM), is approximated by a so-called reduced-order model (ROM) of size  $n$ .

### 1.1 Content overview

We address the question of whether augmenting a Band-Krylov process with Ritz-vectors could be usefully implemented as a tool in the development of more efficient Krylov-subspace projection methods for model order-reduction. In the first three chapters we establish how that makes sense from a theoretical standpoint. Chapters 2 and 3 introduce the basic concepts required as building blocks for a theory of model order-reduction via Krylov-subspace projection, and are probably optional reading for those familiar with standard Krylov-subspace methods for model-reduction, except that we introduce the notion of *pole-weight*, which differs slightly from the notion of dominant poles that appears in current model-reduction literature.

Chapters 4 and 5 use the analytical framework outlined in prior chapters to argue the potential usefulness of thick-restarting the band-Krylov algorithm for more efficient model-reduction methods and contain the bulk of the novel research presented here.

## Chapter 2

Chapter 2 introduces the matrix-valued system transfer-function

$$\mathcal{H} : \mathbb{C} \rightarrow \mathbb{C}^{m \times p} \quad (1.1)$$

of an  $m \times p$  LTI descriptor-system, which describes a model with  $m$  inputs and  $p$  outputs. An example of such a model is a signal processor, which takes an input signal and outputs a modified response signal. For a given complex-frequency  $s$ , the  $i, j$ -th component of  $\mathcal{H}(s)$  gives the system's response from output-terminal  $j$  when  $s$  is input into input-terminal  $i$ . We would like to create a ROM such that its transfer-function  $\tilde{\mathcal{H}}$  approximates the URM transfer-function  $\mathcal{H}$ , so in some sense the problem of model reduction is that of approximating the function (1.1) over a given subset of  $\mathbb{C}$ .

The transfer function can be represented (approximated) as a quotient of two polynomials, or as a Taylor-series, and either representation is taken with respect to an expansion-point, also referred to as an interpolation-point. Approximating (1.1) with some number of Taylor series terms is known as moment-matching and we talk about that in §2.2.1.

The transfer function can also be represented as a rational function, which for  $m = p = 1$  is the quotient

$$\mathcal{H}(s) = \frac{P(s)}{Q(s)}$$

of two polynomials, where zeros of  $Q(s)$  and  $P(s)$  are known as poles and zeros of  $\mathcal{H}(s)$ . One way to approximate (1.1) is via its poles and zeros; its poles are of particular interest in model reduction. Some poles are more important than others, and we introduce a measure of this importance called pole-weight.

## Chapter 3

In chapter 3 we begin a discussion of Krylov subspace based model-reduction methods in earnest. We introduce the Krylov subspace, and the Arnoldi algorithm, which is the most basic Krylov method used for subspace-projection based MOR. Before its application to MOR, the Arnoldi

algorithm was primarily used to find eigenvalues/vectors of large matrices. We address how this is related to model-reduction. Location and convergence of approximate eigenvalues found by an iterative Krylov algorithm like Arnoldi are an essential part of how we determine poles of the system transfer-function.

We define two different kinds of ROM transfer function that can be constructed with a given Krylov-subspace basis: the ROM via *implicit projection* (or implicitly-projected ROM) and that via *explicit projection*. The former is cheaper to construct but is not desirable for use in applications because it can have “bad” poles. It is useful for analysis of the emerging approximate model while we are constructing it. The explicit projection is the ROM generally used for applications but generally is not constructed until the process is done.

## Chapter 4

Chapter 4 addresses some of the difficulties that arise with use of multiple complex interpolation-points, as opposed to one real point. The reason for using a single real expansion-point is to avoid the computational cost of complex arithmetic. The definition of the transformation (4.5) that provides a translation of residual space from one interpolation-point to another in §4.1.2 may be new. It turns out that we did not use it for our work but it may be useful to others.

In §4.2.2 we introduce a possible improvement that reduces the cost of orthogonalizing complex Krylov-subspace vectors by half, and is fairly simple to implement. As far as we know this is a novel contribution to the field.

## Chapter 5

Chapter 5 introduces the Band-Arnoldi algorithm, which is a generalization of the Arnoldi algorithm that allows for iterating simultaneously with several candidate vectors while remaining essentially a single-vector iterative process. The Band-Arnoldi algorithm seems not to be as popular as its *block* variant that iterates an entire block of vectors on each step, but it allows for a simpler and more intuitive restart-scheme. Restarting a Krylov process is done generally to reduce computational costs of orthogonalizing each iterate vector against every previous one, and we can take advantage



of this to restart at a different interpolation-point. The “thick” modifier indicates that on each cycle we re-use some of the information from the previous cycle, but not all of it as with full-orthogonalization.

Exploration of thick-restarting the Band-Arnoldi process was the original purpose of our research that led to this document, and is the thesis of our work. We show how thick-restarting the Band-Arnoldi algorithm can be incorporated into a model-reduction method. It is debatable whether the method presents an improvement over existing methods in general, but we do show via a few examples that smaller models can be achieved than using a single interpolation point, and that thick-restarting presents an improvement in efficiency over a multiple-point method using full-orthogonalization with minimal drawbacks.

### **Additional contributions**

An innovation that may not be obvious from this document is that we developed a library of matlab routines that implement the restarted band-Arnoldi method and will be freely available for future researchers to build on. We also developed a set of visualization tools that are nice to look at, and in our opinion, make model reduction more appealing as a field of study to visually-oriented students. In particular, the surface plot of a transfer-function. Plotting the transfer function as a 3-dimensional surface does not seem like a particularly new idea but there is surprisingly little on-line evidence that anyone has bothered to do it for models with more than one or two poles. Some of the transfer function plots are quite nice landscapes that are also functional in that they convey pole information better than the 2D frequency-response gain plots typically presented in current model-reduction publications. An example of this is Figure 5.8.

## Chapter 2

# The model and its transfer-function

### 2.1 The LTI-system model

Our basic problem is to approximate the ( $N$ -dimensional) linear, time-invariant (LTI) descriptor-system

$$\begin{aligned} \mathbf{E} \frac{dx}{dt} &= \mathbf{A}x + \mathbf{B}u \\ y &= \mathbf{C}^T x \end{aligned} \tag{2.1}$$

with a system that is realized by smaller matrices  $\mathbf{A}_n$ ,  $\mathbf{E}_n$ ,  $\mathbf{B}_n$ , and  $\mathbf{C}_n$ . The collection of constant matrices  $(\mathbf{A}, \mathbf{E}, \mathbf{B}, \mathbf{C})$  is called a realization of the model, which we sometimes call the unreduced model (URM).

Matrices  $\mathbf{E}$  and  $\mathbf{A}$  are singular in general. We only assume that  $\mathbf{A} - s\mathbf{E}$  is invertible for all  $s \in \mathbb{C}$ , except for a finite set of so-called eigenvalues.

$\mathbf{B} \in \mathbb{R}^{N \times m}$  and  $\mathbf{C} \in \mathbb{R}^{N \times p}$  are matrices that condition the input signal  $u(t) \in \mathbb{R}^m$  and output signal  $y(t) \in \mathbb{R}^p$ . If  $p = m = 1$  then (2.1) is a single-input, single-output (SISO) system; its response (output) is a scalar-valued function. If the dimension of  $\mathbf{B}$  and  $\mathbf{C}$  are both greater than one then (2.1) is a multi-input, multi-output (MIMO) system.

We assume that the order- $N$  system (2.1) is too large to work with and we want a model that behaves like (2.1), but with significantly reduced state-space dimension  $n$ . For example, it may not be necessary for our approximate model to have unobservable states, or uncontrollable

ones.

Suppose that, for some reason, we believe restricting the state space of the model (2.1) to an  $n$ -dimensional subspace  $\mathcal{K}$ , will yield a good reduced-order model. If  $V$  is an orthogonal basis of  $\mathcal{K}$  then the reduced-order model (ROM) obtained via orthogonal projection onto  $\mathcal{K}$  is a new descriptor system

$$\begin{aligned} \mathbf{E}_n \frac{d\tilde{x}}{dt} &= \mathbf{A}_n \tilde{x} + \mathbf{B}_n u \\ \tilde{y} &= \mathbf{C}_n^T \tilde{x}, \end{aligned} \tag{2.2}$$

with orthogonal projections

$$\mathbf{A}_n = V^T \mathbf{A} V, \quad \mathbf{E}_n = V^T \mathbf{E} V, \quad \mathbf{C}_n = V^T \mathbf{C}, \quad \mathbf{B}_n = V^T \mathbf{B}, \tag{2.3}$$

where  $\tilde{x}(t) \in \mathbb{C}^n$  is the state of the reduced-order system such that  $V\tilde{x}(t)$  approximates the state of the unreduced model. The  $p$  output(s)  $\tilde{y}(t) \in \mathbb{R}^p$  approximate  $y(t) \in \mathbb{R}^p$  from (2.1), given the same  $m$  input(s)  $u(t) \in \mathbb{R}^m$ , and ideally  $\|y - \tilde{y}\|$  is small.

The reduced-order model (2.2), (2.3) is called the *explicitly-projected* ROM, because in the computational setting we must actually compute (2.3). Ultimately the desired ROM is in this form.

## 2.2 System transfer-function

The transfer-function is a direct relationship between input and output of the model in the frequency domain. If we temporarily ignore the state of the model and view it simply as a mapping of an input signal  $u$ , to an output signal  $y$ , the system (2.1) acts as a system-function  $y = h(u)$ . The transfer-function is obtained by applying the Laplace transform (eg.  $X(s) = \mathcal{L}\{x(t)\}$ ) to (2.1) and assuming a zero initial condition  $X(0) = 0$ , which yields the algebraic equations

$$\begin{aligned} s\mathbf{E}X &= \mathbf{A}X + \mathbf{B}U, \\ Y &= \mathbf{C}^T X. \end{aligned}$$

Then  $Y(s) = \mathcal{H}(s)U(s)$ , where

$$\mathcal{H}(s) = \mathbf{C}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \quad (2.1)$$

is the transfer-function over  $\mathbb{C}$ . Note that  $\mathcal{H}(s)$  is defined only if the matrix pencil  $(\mathbf{A}, \mathbf{E})$  is regular, meaning that the matrix  $\mathbf{A} - s\mathbf{E}$  is invertible for all but a finite set of eigenvalues.

For a general MIMO transfer-function (2.1) where  $\mathbf{C} = \begin{bmatrix} \mathbf{c}_1 & \mathbf{c}_2 & \dots & \mathbf{c}_p \end{bmatrix}$  and  $\mathbf{B} = \begin{bmatrix} \mathbf{b}_1 & \mathbf{b}_2 & \dots & \mathbf{b}_m \end{bmatrix}$ , we can consider (2.1) to be  $mp$  scalar-valued SISO (single input single output) transfer-functions

$$\mathcal{H}_{ij}(s) = \mathbf{c}_i^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{b}_j \in \mathbb{C},$$

and for example in the  $2 \times 2$  case we have

$$\mathcal{H}(s) = \begin{bmatrix} \mathcal{H}_{11}(s) & \mathcal{H}_{12}(s) \\ \mathcal{H}_{21}(s) & \mathcal{H}_{22}(s) \end{bmatrix}.$$

For a general MIMO transfer-function  $\mathcal{H}(s) = \mathbf{C}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{B}$ ,

$$[\mathcal{H}(s)]^H = \mathbf{B}^T (s\mathbf{E}^T - \mathbf{A}^T)^{-1} \mathbf{C},$$

which implies that  $\mathcal{H}_{ij}(s) = \overline{\mathcal{H}_{ji}(s)}$ , or more importantly,

$$|\mathcal{H}_{ij}(s)| = |\mathcal{H}_{ji}(s)|,$$

so it suffices to consider  $\mathcal{H}_{ij}(s)$  only for  $i < j$ .

### 2.2.1 Moments

The transfer-function is a rational function, and thus can be represented by a Taylor series about an expansion-point  $\sigma \in \mathbb{C}$ , having the general form

$$\mathcal{H}(s) = \sum_{j=0}^{\infty} (s - \sigma)^j \mathcal{H}^{(j)}, \quad \text{or equivalently,} \quad \mathcal{H}(s + \sigma) = \sum_{j=0}^{\infty} s^j \mathcal{H}^{(j)} \quad (2.2)$$

where the Taylor coefficient

$$\mathcal{H}^{(j)} = \frac{1}{j!} \left. \frac{d^j \mathcal{H}}{ds^j} \right|_{s=\sigma} \quad (2.3)$$

is called the  $j$ -th *moment* of the transfer-function about  $\sigma$ .

### Moment-matching

Krylov-subspace projection methods boast moment-matching properties. The reduced-order model transfer-function implied by a Krylov-subspace method is guaranteed to share a number of terms of the Taylor series about one or several points, with that of the full unreduced model.

Suppose the URM (unreduced model) transfer-function expressed as a Taylor series about  $\sigma$  is

$$\mathcal{H}(s) = \mathcal{H}^{(0)} + (s - \sigma)\mathcal{H}^{(1)} + (s - \sigma)^2\mathcal{H}^{(2)} + \dots + (s - \sigma)^{n-1}\mathcal{H}^{(n-1)} + \dots$$

A reduced-order model (ROM) whose transfer-function can be written as

$$\hat{\mathcal{H}}(s) = \hat{\mathcal{H}}^{(0)} + (s - \sigma)\hat{\mathcal{H}}^{(1)} + (s - \sigma)^2\hat{\mathcal{H}}^{(2)} + \dots + (s - \sigma)^{n-1}\hat{\mathcal{H}}^{(n-1)} + \dots$$

where

$$\hat{\mathcal{H}}^{(j)} = \mathcal{H}^{(j)} \quad \text{for } j = 0, 1, 2, \dots, n-1$$

is said to match  $n$ -moments about  $\sigma$ .

Moments can be matched about any number of expansion-points; also called interpolation-points.

### 2.2.2 Shift-invert representation of the transfer-function

Moment matching properties of Krylov-subspace methods are accomplished via the the following reformulation of the transfer-function (2.1).

Let  $\sigma \in \mathbb{C}$  be a point for which  $\sigma \mathbf{E} - \mathbf{A}$  is invertible. Then

$$\begin{aligned}\mathcal{H}(s) &= \mathbf{C}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \\ &= \mathbf{C}^T (\sigma\mathbf{E} - \mathbf{A} + (s - \sigma)\mathbf{E})^{-1} \mathbf{B} \\ &= \mathbf{C}^T (I - (s - \sigma)\mathbf{H})^{-1} \mathbf{R}\end{aligned}\tag{2.4}$$

where<sup>1</sup>

$$\mathbf{H} := (\mathbf{A} - \sigma\mathbf{E})^{-1}\mathbf{E} \quad \text{and} \quad \mathbf{R} := (\sigma\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}.\tag{2.5}$$

(2.4) is sometimes called the shifted transfer-function formulation, with shift  $\sigma$ , although it does not depend on  $\sigma$ . The shift matters when we consider the ROM transfer-function that approximates (2.4).

The generally non-sparse  $\mathbf{H} = \mathbf{H}(\sigma) \in \mathbb{C}^{N \times N}$  is a sort of operator or multiplier that acts on  $\mathbf{R} = \mathbf{R}(\sigma) \in \mathbb{C}^{N \times p}$ .  $\mathbf{H}$  is dense in general and is rarely if ever explicitly formed. We only need a way to obtain matrix-vector products  $\mathbf{H}v$  for vectors  $v \in \mathbb{C}^N$ .  $\mathbf{H}$  and  $\mathbf{R}$  are the building blocks for the moments of the transfer-function about  $\sigma$ .

The shifted transfer-function representation (2.4) can alternatively be considered the transfer-function for the shifted descriptor system

$$\begin{aligned}\mathbf{H} \frac{dx}{dt} &= (I - \sigma\mathbf{H})x + \mathbf{R}u \\ y &= \mathbf{C}^T x,\end{aligned}\tag{2.6}$$

which is equivalent to (2.1) for any  $\sigma$  such that  $\mathbf{A} - \sigma\mathbf{E}$  is invertible. This is notable because some order-reduction schemes work by replacing  $\mathbf{H}$ ,  $\mathbf{R}$ , and  $\mathbf{C}$  with reduced-order approximations  $\tilde{\mathbf{H}} = V^T \mathbf{H} V$ ,  $\tilde{\mathbf{R}}_n = V^T \mathbf{R}$ , and  $\mathbf{C}_n = V^T \mathbf{C}$ . We will call such a ROM *implicitly* projected on to

---

<sup>1</sup>To verify (2.4), note that  $I = (\sigma\mathbf{E} - \mathbf{A})^{-1}(\sigma\mathbf{E} - \mathbf{A})$ .

span  $V$ , as opposed to the *explicitly* projected model (2.2). A model obtained via implicit-projection is not equivalent to (2.2) in general and is undesirable for some applications, but is much cheaper to construct.

### Moment representation

We now express transfer-function moments about  $\sigma$  in terms of  $\mathbf{H}$  and  $\mathbf{R}$ . Via Neumann series expansion (power series for matrices) re-write (2.4) as

$$\begin{aligned}\mathcal{H}(s) &= \mathbf{C}^T \left( \sum_{j=0}^{\infty} (s - \sigma)^j \mathbf{H}^j \right) \mathbf{R} \\ &= \sum_{j=0}^{\infty} (s - \sigma)^j \mathbf{C}^T \mathbf{H}^j \mathbf{R}.\end{aligned}\tag{2.7}$$

The moments  $\mathcal{H}^{(j)}$  from (2.2) are specified exactly in (2.7):

$$\mathcal{H}^{(j)} = \mathbf{C}^T \mathbf{H}^j \mathbf{R}.\tag{2.8}$$

### Region of convergence for moment-matching

The power (Taylor) series representation seems to imply that (2.7) is only valid for  $s$  in a disc of radius  $1/\|\mathbf{H}\|_{op}$  around  $\sigma$ , where the operator norm

$$\|\mathbf{H}\|_{op} = \sup_{v \neq 0} \left\{ \frac{\|\mathbf{H}v\|}{\|v\|} \right\}.$$

Then certainly  $\|\mathbf{H}\|_{op} \geq |\lambda_1|$ , where  $\lambda_1$  is the largest eigenvalue of  $\mathbf{H}$ . Equivalently,  $\|\mathbf{H}\|_{op} \geq 1/|(\mu_1 - \sigma)|$  where  $\mu_1$  is the closest pole to  $\sigma$ . Thus, the region of convergence for (2.7) is the largest disc centered at  $\sigma$  that does not contain a pole (see §2.2.4). The closer  $\sigma$  is to a pole of the transfer-function, the smaller region of convergence we theoretically have for moment-matching about  $\sigma$ . In practice Krylov-subspace methods are observed to converge well outside of the theoretical region of convergence.

### 2.2.3 Invariant-subspaces

Given a transformation  $\mathbf{H} : \mathbb{C}^N \rightarrow \mathbb{C}^N$ , a subspace  $\mathcal{Q}$  of  $\mathbb{C}^N$  is called **H**-invariant if

$$\mathbf{H}\mathcal{Q} \subseteq \mathcal{Q}. \quad (2.9)$$

The span of a set of eigenvectors of **H** is an invariant subspace under **H**, and an invariant subspace (with respect to **H**) always has a basis consisting of eigenvectors of **H**.

If  $Q$  is a basis for  $\mathcal{Q}$  then

$$\mathbf{H}Q = QT \quad (2.10)$$

for some matrix  $T \in \mathbb{C}^{\ell \times \ell}$ . If the basis vectors are eigenvectors  $Z$  then (2.10) becomes

$$\mathbf{H}Z = \Lambda Z,$$

where  $\Lambda = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_\ell\}$  is a diagonal matrix of eigenvalues associated with the vectors

$$Z = \begin{bmatrix} z_1 & z_2 & \cdots & z_\ell \end{bmatrix}.$$

If the basis  $Q = \begin{bmatrix} u_1 & u_2 & \cdots & u_\ell \end{bmatrix}$  for the **H**-invariant-subspace

$$\mathcal{Q} = \text{span} \begin{bmatrix} u_1 & u_2 & \cdots & u_\ell \end{bmatrix} = \text{span} \begin{bmatrix} z_1 & z_2 & \cdots & z_\ell \end{bmatrix}$$

is orthonormal then we call vectors  $u_j$  Schur-vectors and sometimes call  $\mathcal{Q}$  a Schur space. Also, (2.10) is called a Schur-decomposition, and  $T$  is upper triangular with eigenvalues  $\lambda_j$  associated with  $z_j$  along its diagonal. A Schur decomposition is often preferred over an eigen-decomposition because it is easier to compute and Schur vectors are more numerically stable.

### Generalized-eigenvalues and invariance

An eigenvalue of matrix pencil  $(\mathbf{A}, \mathbf{E})$ , called a generalized eigenvalue, is a  $\mu \in \mathbb{C}$  such that  $(\mathbf{A} - \mu\mathbf{E})z = 0$  has nonzero solutions  $z \neq 0 \in \mathbb{C}^N$ , which are called right-eigenvectors since they multiply from the right.



The notion of invariance under a general operator  $\mathbf{H}$  extends to that of a matrix pencil. The subspace  $\mathcal{Q} = \text{span } Z$  is called invariant, or deflating, [32] with respect to  $(\mathbf{A}, \mathbf{E})$  if

$$\dim(\mathbf{A}\mathcal{Q} + \mathbf{E}\mathcal{Q}) \leq \dim \mathcal{Q}. \quad (2.11)$$

For a regular matrix pencil  $(\mathbf{A}, \mathbf{E})$  and any  $\sigma \in \mathbb{C}$  that is not an eigenvalue of  $(\mathbf{A}, \mathbf{E})$  it is shown in [12] that  $(\mathbf{A}, \mathbf{E})$ -invariance is equivalent to  $\mathbf{H}$ -invariance for

$$\mathbf{H} = (\mathbf{A} - \sigma \mathbf{E})^{-1} \mathbf{E}, \quad (2.12)$$

which happens to be the shift-invert operator defined by (2.5).

Then  $\mathbf{H}(\sigma) = (\mathbf{A} - \sigma \mathbf{E})^{-1} \mathbf{E}$  has the same eigenvectors (regardless of  $\sigma$ ) as the pencil  $(\mathbf{A}, \mathbf{E})$ . An eigenvalue  $\lambda$  of  $\mathbf{H}$  and its corresponding eigenvalue  $\mu$  of  $(\mathbf{A}, \mathbf{E})$  are related by

$$\lambda = \frac{1}{\mu - \sigma}, \quad \mu = \sigma + \frac{1}{\lambda} \quad (2.13)$$

and they share the same eigenvector. That is,

$$\begin{aligned} \mathbf{H}z = \lambda z & \quad \Longleftrightarrow \quad \mathbf{A}z = \mu \mathbf{E}z \\ & = \left( \frac{1}{\mu - \sigma} \right) z \quad \Longleftrightarrow \quad = \left( \sigma + \frac{1}{\lambda} \right) \mathbf{E}z \end{aligned}$$

To see this, observe that

$$\begin{aligned} (\mu \mathbf{E} - \mathbf{A}) &= [(\mu - \sigma) \mathbf{E} + (\sigma \mathbf{E} - \mathbf{A})] \\ &= [(\mu - \sigma) \underbrace{(\sigma \mathbf{E} - \mathbf{A})^{-1} \mathbf{E}}_{-\mathbf{H}} + I] \\ &= [(\mu - \sigma) \mathbf{H} - I] \\ &= \left[ \mathbf{H} - \left( \frac{1}{\mu - \sigma} \right) I \right] \\ &= (\mathbf{H} - \lambda I). \end{aligned} \quad (2.14)$$

In other words, invariant-subspaces under the shifted operator  $\mathbf{H}(\sigma)$  are shift-invariant. This is useful because if we decide to change the shift  $\sigma$ , any previously discovered invariant-subspace will be still be invariant under the new operator.

#### 2.2.4 Pole-Residue representation

Recall the transfer-function

$$\mathcal{H}(s) = \mathbf{C}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{B}, \quad (2.1)$$

which includes the linear matrix pencil  $(\mathbf{A}, \mathbf{E})$ . Poles of the transfer-function (2.1) are values  $\mu \in \mathbb{C} \cup \infty$  such that  $\|\mathcal{H}(\mu)\| = \infty$ . Poles of  $\mathcal{H}(s)$  are eigenvalues of the matrix pencil  $(\mathbf{A}, \mathbf{E})$ , but their significance is determined by  $\mathbf{B}$  and  $\mathbf{C}$ .

The  $\mathbf{A}$  and  $\mathbf{E}$  that arise from Modified Nodal Analysis [16] circuit representations are singular in general. In particular,  $\mathbf{E}$  is singular, which means  $(\mathbf{A}, \mathbf{E})$  has eigenvalues at  $\infty$  at least some of which are associated with the null space of  $\mathbf{E}$ .

#### Poles and residues of the standard transfer-function formulation

Let us assume that  $(\mathbf{A}, \mathbf{E})$  has a full eigen-decomposition

$$\mathbf{A}\mathbf{Z} = \mathbf{E}\mathbf{Z}\mathcal{M} \quad \text{and} \quad \mathbf{A}^T\mathbf{W} = \mathbf{E}^T\mathbf{W}\mathcal{M}$$

where  $\mathbf{Z}, \mathbf{W} \in \mathbb{C}^{N \times N}$  represent the right and left eigenspaces of  $(\mathbf{A}, \mathbf{E})$ , respectively, and  $\mathcal{M}$  is the diagonal matrix of eigenvalues  $\mu_j \in \mathbb{C} \cup \{\infty\}$ . Note that since  $\mathbf{A}$  and  $\mathbf{E}$  are real, every eigenvalue is real, infinite, or is one of a complex conjugate pair.

The left and right eigenvectors are orthogonal, and can be scaled<sup>2</sup> so that  $\mathbf{W}^H \mathbf{E} \mathbf{Z} = \mathbf{I}$ . Then

---

<sup>2</sup>This scaling is not generally the default for eigensolver algorithms.

$$\begin{aligned}
\mathcal{H}(s) &= \mathbf{C}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \\
&= \mathbf{C}^T (W^{-H} (sI - \mathcal{M}) Z^{-1})^{-1} \mathbf{B} \\
&= \mathbf{C}^T Z (sI - \mathcal{M})^{-1} W^H \mathbf{B} \\
&= \mathbf{C}^T Z \left( \sum_{j=1}^q \frac{1}{s - \mu_j} \right) W^H \mathbf{B}.
\end{aligned} \tag{2.15}$$

The pole (eigen) decomposition (2.15) of the transfer-function suggests that it can be approximated using eigen-pairs of  $(\mathbf{A}, \mathbf{E})$ .

For the sake of expressing (2.15) without clutter we assumed all eigenvalues are simple and finite. Actually, multiple eigenvalues are possible and eigenvalues at infinity are inevitable because  $\mathbf{E}$  is singular. The eigenspace associated with infinite eigenvalues is the nullspace of  $\mathbf{E}$ .

Recall the left and right hand sides  $\mathbf{C}^T Z \in \mathbb{C}^{p \times N}$  and  $W^H \mathbf{B} \in \mathbb{C}^{N \times m}$  from (2.15) and consider the partitions

$$\mathbf{C}^T Z = \begin{bmatrix} \hat{c}_1 & \hat{c}_2 & \cdots & \hat{c}_N \end{bmatrix}$$

and

$$\mathbf{B}^T W = \begin{bmatrix} \hat{b}_1 & \hat{b}_2 & \cdots & \hat{b}_N \end{bmatrix}$$

into  $N$  columns.

Then we can express the the pole-residue form of transfer-function (2.1) as

$$\mathcal{H}(s) = \sum_{\mu_j = \infty} \hat{c}_j \hat{b}_j^T + \sum_{\mu_j \neq \infty} \frac{\hat{c}_j \hat{b}_j^T}{s - \mu_j}. \tag{2.16}$$

Note that  $\hat{c}_j$  and  $\hat{b}_j$  are scalars if this is a SISO model. Generally,  $\hat{c}_j \hat{b}_j^T \in \mathbb{C}^{p \times m}$ . The transpose  $\hat{b}_j^T$  is in fact a transpose and not a conjugate-transpose, even if  $\hat{b}_j$  is complex-valued.

Our necessary assumption  $W^H \mathbf{E} Z = I$  is not the default scaling for eigensolvers in practice. In that case, we must consider the scaling factor

$$\xi_j = 1/w_j^H \mathbf{E} z_j$$

for  $j = 1, 2, \dots, q$  and (2.16) generalizes to

$$\mathcal{H}(s) = \sum_{\mu_j = \infty} \xi_j \hat{c}_j \hat{b}_j^T + \sum_{\mu_j \neq \infty} \xi_j \frac{\hat{c}_j \hat{b}_j^T}{s - \mu_j}. \quad (2.17)$$

### Poles and residues from shifted transfer-function formulation

For model reduction it is often favorable to work with the so-called  $\sigma$ -shifted transfer-function

$$\mathcal{H}(s) = \mathbf{C}^T (I - (s - \sigma)\mathbf{H})^{-1} \mathbf{R}, \quad (2.4)$$

rather than the standard formulation (2.1).

Assume that  $\mathbf{H}$  is diagonalizable with right-eigenbasis  $Z$ , so that

$$\mathbf{H}Z = Z\Lambda$$

for a  $N \times N$  diagonal matrix  $\Lambda = \begin{bmatrix} \lambda_1 & \lambda_2 & \dots & \lambda_N \end{bmatrix}$  of eigenvalues.

Then

$$\begin{aligned} \mathcal{H}(s) &= \mathbf{C}^T (I - (s - \sigma)\mathbf{H})^{-1} \mathbf{R} \\ &= \mathbf{C}^T (Z(I - (s - \sigma)\Lambda)Z^{-1})^{-1} \mathbf{R} \\ &= (\mathbf{C}^T Z) \Delta(s) (Z^{-1} \mathbf{R}) \end{aligned} \quad (2.18)$$

where  $\Delta(s) = (I - (s - \sigma)\Lambda)^{-1}$  is a diagonal matrix with diagonal entries  $\delta_j(s) = 1 - (s - \sigma)\lambda_j$ , or equivalently

$$\delta_j(s) = \begin{cases} \frac{\sigma - \mu_j}{s - \mu_j}, & \mu_j \neq \infty \\ 1, & \mu_j = \infty. \end{cases} \quad (2.19)$$

Then for

$$\mathbf{C}^T Z = \begin{bmatrix} \hat{f}_1 & \hat{f}_2 & \dots & \hat{f}_N \end{bmatrix} \quad \text{and} \quad (Z^{-1} \mathbf{R})^T = \begin{bmatrix} \hat{g}_1 & \hat{g}_2 & \dots & \hat{g}_N \end{bmatrix},$$

we have

$$\mathcal{H}(s) = \sum_j \frac{\hat{f}_j \hat{g}_j^T}{1 - (s - \sigma)\lambda_j} \quad (2.20)$$

$$\begin{aligned} &= \sum_{\lambda_j=0} \hat{f}_j \hat{g}_j^T + \sum_{\lambda_j \neq 0} \frac{\sigma - \mu_j}{s - \mu_j} \hat{f}_j \hat{g}_j^T \\ &= \sum_j \delta_j(s) \hat{f}_j \hat{g}_j^T, \end{aligned} \quad (2.21)$$

where  $\delta_j(s)$  is from (2.19). The transpose  $\hat{g}_j^T$  is in fact a standard (not-conjugated) transpose, even if  $\hat{g}_j$  is complex-valued. Both  $\hat{f}_j$  and  $\hat{g}_j^T$  are scalars in the SISO case. Note that a zero eigenvalue  $\lambda_j = 0$  of  $\mathbf{H}$  corresponds to an infinite  $\mu_j = \infty$  pole (eigenvalue of  $(\mathbf{A}, \mathbf{E})$ ).

### 2.2.5 Pole-weight

Poles of the transfer-function  $\mathcal{H}(s) = \mathbf{C}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{B}$  are values  $\mu \in \mathbb{C} \cup \infty$  such that  $\|\mathcal{H}(\mu)\| = \infty$ . Poles of  $\mathcal{H}(s)$  are eigenvalues of the matrix pencil  $(\mathbf{A}, \mathbf{E})$ , but their significance is determined by  $\mathbf{B}$  and  $\mathbf{C}$ . Pole dominance is a notion of a pole's influence on the transfer-function frequency response  $\mathcal{H}(i\omega)$  on an interval of the  $\Im$ -axis (or all of it). Pole-residue formulations (2.17) and (2.21) suggest a hierarchy of poles' importance for approximation.

We will define a measure of pole-dominance, which we will call its *mass* or *weight* with respect to the frequency response domain  $i[\omega_1, \omega_2] \subset \mathbb{C}$ . It is similar to the modal dominance index (MDI) of [1], but it considers a pole's influence over the frequency response domain rather than all of the positive  $\Im$ -axis, and it does not blow-up for poles on or near the  $\Im$ -axis.

If we take the norm of (2.21) over the interval  $i[\omega_1, \omega_2]$  of interest on the  $\Im$ -axis, we have

$$\|\mathcal{H}(i\omega)\|_\infty \leq \sum_j \|\delta_j(i\omega)\|_\infty \|\hat{f}_j\|_1 \|\hat{g}_j\|_1, \quad (2.22)$$

which is a sum of positive numbers, each one associated with a pole  $\mu_j$ . We call this positive number the weight of the pole. A relatively large pole-weight

$$\gamma_j = \|\delta_j(i\omega)\|_\infty \|\hat{f}_j\|_1 \|\hat{g}_j\|_1 \quad (2.23)$$

indicates that  $\mu_j$  is a so-called dominant pole.

For a SISO model,  $\hat{f}_j$  is a scalar so  $\|\hat{f}_j\|_1 = |\hat{f}_j|$  and it represents the weighting of the pole  $\mu_j$  by the left-hand multiplier  $\mathbf{C} = \mathbf{c}$  of the transfer-function (2.1).

A MIMO system has  $p$  such left-multipliers in the form of  $\mathbf{C} = \begin{bmatrix} \mathbf{c}_1 & \mathbf{c}_2 & \cdots & \mathbf{c}_p \end{bmatrix}$  and each one has an associated element in the column vector  $\hat{f}_j = (\hat{f}_{1j}, \hat{f}_{2j}, \dots, \hat{f}_{pj}) \in \mathbb{C}^p$ . By summing them we get an overall sense of how much  $\mu_j$  is favored by  $\mathbf{C}$ . Thus we use the 1-norm (column-sum)

$$\|\hat{f}_j\|_1 = \sum_i |\hat{f}_{ij}|.$$

The reasoning for using  $\|\hat{g}_j\|_1$  in (2.22) is similar.

The scalar-valued function  $\delta_j(i\omega)$  represents the influence of pole  $\mu_j$  on the system frequency response via its proximity to the segment  $i[\omega_1, \omega_2]$  of interest, and we take its maximum value

$$\begin{aligned} \|\delta_j(i\omega)\|_\infty &= \max_{\omega \in [\omega_1, \omega_2]} |\delta_j(i\omega)| \\ &= \begin{cases} \frac{|\sigma - \mu_j|}{\min\{|\mu_j - \omega_1|, |\mu_j - \omega_2|, |\Re(\mu_j)|\}}, & \mu_j \neq \infty \\ 1, & \mu_j = \infty \end{cases} \end{aligned} \quad (2.24)$$

over that interval as a sense of its over-all influence on the transfer-function in that region. The value  $\min\{|\mu_j - \omega_1|, |\mu_j - \omega_2|, |\Re(\mu_j)|\}$  is merely the distance of  $\mu_j$  to the segment  $i[\omega_1, \omega_2]$ , as illustrated in figure 2.1. Conjugate pairs must be considered together, so when determining weight we actually consider  $\Re(\mu_j) + i|\Im(\mu_j)|$ , rather than  $\mu_j$ . That way, each member of the pair gets assigned the same weight.

### Total system-mass

The right-hand-side of (2.22) (i.e.  $\sum_j \gamma_j$ ), is the total mass of the system (2.1) with respect to  $i[\omega_1, \omega_2]$ , and we attempted to use it as a measure of ROM convergence. The combined weight of a few dominant poles often comprises most of a system's total mass.

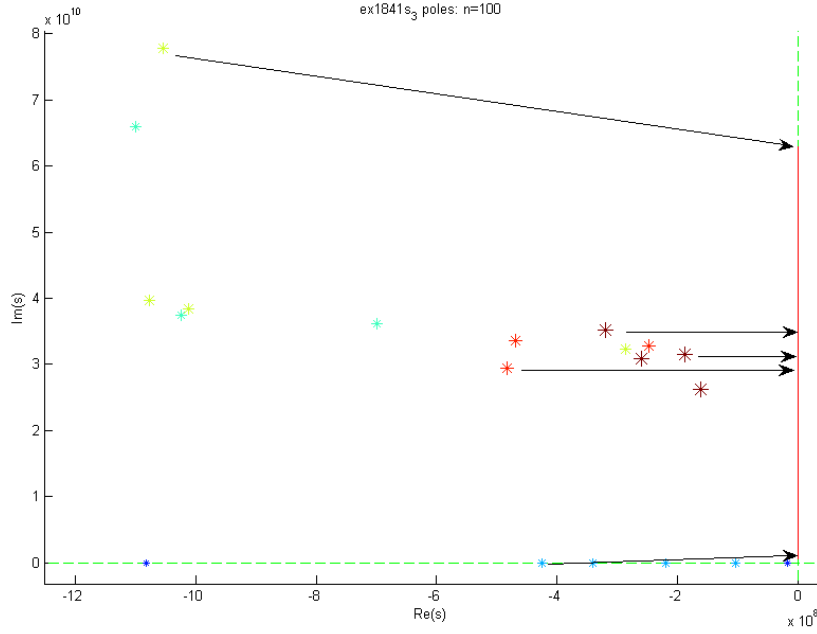


Figure 2.1: The red (solid) segment on the  $\Im$ -axis is the segment  $i(\omega_0, \omega_1)$  of interest. In this case it looks like it extends to the origin but it does not. The arrows indicate how we define a pole's distance to the segment, which we use to determine the pole's weight. Conjugate pairs must be considered together, so when determining weight we always consider  $\Re(\mu) + i|\Im(\mu)|$ , rather than  $\mu$ . That way, each member of the pair gets assigned the same weight.

## 2.3 Reduced-order transfer-function via projection

The URM (unreduced model) transfer-function formulations via explicit (2.1) and implicit (2.4) projection exist only in theory for large applications, which can be on the order of  $10^9$  at the time of this writing. The two transfer-functions (2.1) and (2.4) are mathematically equivalent. Subspace projected ROM transfer-functions can be obtained via *explicit* projection, or *implicit* projection, yielding two forms similar to (2.1) and (2.4) that are *not* mathematically equivalent but converge to the URM (unreduced) transfer-function (2.1), (2.4) as the projection subspace approaches  $\mathbb{C}^N$ . We will refer to the two formulations as the implicit, and explicit transfer-functions associated with a given subspace for projection.

Let  $V \in \mathbb{C}^{N \times n}$  be a matrix with orthogonal columns that form a basis of our projection subspace  $\mathcal{K}$ . If we make the orthogonal projections

$$\mathbf{A}_n := V^T \mathbf{A} V, \quad \mathbf{E}_n := V^T \mathbf{E} V, \quad \mathbf{C}_n := V^T \mathbf{C}, \quad \mathbf{B}_n := V^T \mathbf{B}, \quad (2.3)$$

of realization  $(\mathbf{A}, \mathbf{E}, \mathbf{B}, \mathbf{C})$  on to  $\mathcal{K}$ ,<sup>3</sup> the **explicitly projected** model  $(\mathbf{A}_n, \mathbf{E}_n, \mathbf{B}_n, \mathbf{C}_n)$  has transfer-function

$$\hat{\mathcal{H}}_n(s) = \mathbf{C}_n^T (s \mathbf{E}_n - \mathbf{A}_n)^{-1} \mathbf{B}_n. \quad (2.1)$$

The explicitly projected model has the property that the projected ROM (2.1) of a stable system is stable if the projection basis  $V$  is real-valued, which is not true in general of the implicit projected transfer-function, described next.

Some iterative subspace methods, notably Krylov-subspace methods, use the operator  $\mathbf{H}$  and operand  $\mathbf{R}$  to construct a basis  $V$  and a projected operator matrix  $\tilde{\mathbf{H}} \in \mathbb{C}^{n \times n}$  such that

$$\tilde{\mathbf{H}} = V^T \mathbf{H} V.$$

This permits what we will call the ROM transfer-function via implicit-projection, or implicit

---

<sup>3</sup>assuming the matrix pencil  $(\mathbf{A}_n, \mathbf{E}_n)$  is regular



transfer-function.

$$\tilde{\mathcal{H}}_n(s) = \mathbf{C}_n^H \left( I - (s - \sigma)\tilde{\mathbf{H}} \right)^{-1} \tilde{\boldsymbol{\rho}}_n, \quad (2.2)$$

where

$$\mathbf{C}_n := \mathbf{W}^H \mathbf{C} \quad \text{and} \quad \tilde{\boldsymbol{\rho}}_n := \mathbf{V}^T \mathbf{R}.$$

(2.2) is the projected ROM analog to the “shifted” transfer-function (2.4).

We say (2.2) is a transfer-function via implicit-projection because it exists without a projected realization  $(\mathbf{A}_n, \mathbf{E}_n, \mathbf{B}_n, \mathbf{C}_n)$  from (2.3). The reduced model implied by (2.2) is actually

$$\begin{aligned} \tilde{\mathbf{H}} \frac{d\tilde{x}}{dt} &= (I + \sigma\tilde{\mathbf{H}})\tilde{x} + \tilde{\boldsymbol{\rho}}_n u \\ \hat{y} &= \mathbf{C}_n^T \tilde{x}, \end{aligned} \quad (2.3)$$

which is not equivalent to the explicitly projected system

$$\begin{aligned} \mathbf{E}_n \frac{dz}{dt} &= \mathbf{A}_n z + \mathbf{B}_n u \\ \hat{y} &= \mathbf{C}_n^T z, \end{aligned} \quad (2.2)$$

given the same basis  $V$ , unless  $V$  spans  $\mathbb{C}^N$  (i.e.  $n = N$ ) in which case they are both equivalent to the original system (2.1).

Even if the original system is passive and/or stable, the implicitly projected ROM is not necessarily passive or stable which makes it less suitable as a ROM. Some efforts [30, 15, 17] have been made to remedy this situation, but these methods tend to sacrifice moment-matching properties in the process. The reduced models produced are still pretty good, but do not have proven error bounds. As a result the explicitly projected model is generally preferred in practice, although it requires quite a bit more computation to carry out the projections (2.3).

## Chapter 3

# Order-reduction via Krylov-subspace projection

### 3.1 Krylov-subspace projection methods

Recall two ways to express the system transfer-function  $\mathcal{H} : \mathbb{C} \rightarrow \mathbb{C}^{m \times p}$

$$\mathcal{H}(s) = \mathbf{C}^T (s\mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \quad (2.1)$$

$$= \mathbf{C}^T (I - (s - \sigma)\mathbf{H})^{-1} \mathbf{R} \quad (2.4)$$

with

$$\mathbf{H} = (\mathbf{A} - \sigma\mathbf{E})^{-1}\mathbf{E} \quad \text{and} \quad \mathbf{R} = (\sigma\mathbf{E} - \mathbf{A})^{-1}\mathbf{B}, \quad (2.5)$$

where (2.1) is the standard formulation and (2.4) is the so-called  $\sigma$ -shifted formulation, and  $\mathbf{H}$  sometimes called a shift-inverse operator. We approximate  $\mathcal{H}(s)$  by approximating  $\mathbf{H}$  and  $\mathbf{R}$ , since it is known that successive applications of  $\mathbf{H}$  to  $\mathbf{R}$ , as in  $\mathbf{H}\mathbf{R}, \mathbf{H}^2\mathbf{R}, \dots$ , creates progressively better approximations to the spectrum of  $\mathbf{H}$ .

A Krylov-subspace method iteratively constructs a basis but we often think of the progression in terms of a sequence of converging eigenvalues of  $\mathbf{H}$ , starting with the largest. In our case  $\mathbf{H}$  is a

shift-and-invert operator, so large eigenvalues  $\lambda$  of  $\mathbf{H}$  are eigenvalues

$$\mu = \sigma + 1/\lambda$$

of  $(\mathbf{A}, \mathbf{E})$  which are closest to  $\sigma$ . Since eigenvalues  $\mu$  of  $(\mathbf{A}, \mathbf{E})$  are poles of the transfer-function  $\mathcal{H}(s)$ , we speak of “poles converging” as progress towards an accurate ROM. This is consistent with the notion of a Taylor series giving a better approximation near  $\sigma$  with each additional moment.

A fundamental feature (and drawback) of Krylov-subspace methods for model order-reduction is that for a given shift  $\sigma$  there is only one way for the method to progress: from poles  $\mu$  closest to  $\sigma$  (i.e. smallest  $|\mu - \sigma|$ ), to those farthest away. This presents a problem because often only a few dominant poles (§2.2.5) influence the transfer-function over the segment of interest  $i[\omega_0, \omega_1]$  on the  $\Im$ -axis. For example, suppose poles  $\mu_1$  and  $\mu_2$  are dominant poles but are separated (in distance from  $\sigma$ ) by several insignificant poles, as in

$$|\mu_1 - \sigma| > |\mu_3 - \sigma| > |\mu_4 - \sigma| > \cdots > |\mu_\ell - \sigma| > |\mu_2 - \sigma|.$$

Then, using a straightforward Krylov process, after  $\mu_1$  converges,  $\mu_3, \dots, \mu_\ell$  must all converge before  $\mu_2$  does, and all of the associated vector information is added to the projection basis  $V$ , creating a larger than necessary model. One way around this is to change the shift  $\sigma$  after some number of iterations. It should also be noted that information about significant transfer-function *zeros* may be included with that of insignificant poles, so convergence of dominant poles is a somewhat dubious indicator of approximate model convergence. Ideally, both dominant pole and zero information should be considered and at the time of this writing there are no Krylov methods that do this.

### 3.1.1 The Krylov subspace

Krylov-subspace projection methods are developed out of the power iteration with  $\mathbf{H}$  acting on  $\mathbf{R} = \begin{bmatrix} \mathbf{r}_1 & \mathbf{r}_2 & \cdots & \mathbf{r}_m \end{bmatrix}$ . The goal of subspace projection-based model order-reduction is to obtain a basis  $V$  of a subspace  $\mathcal{K}$  of  $\mathbb{C}^N$  on which to project our system realization in order to obtain a reduced model. The Krylov-subspace is the ideal subspace to use for projection because reduced-

order models obtained via Krylov-subspaces have moment-matching properties (§2.2.1). Reduced-order models obtained by projection onto non-Krylov-subspaces do not have moment-matching properties in general. The  $n$ -th *Krylov-subspace* induced by  $\mathbf{H}$  and a vector  $\mathbf{r}$  is

$$\mathcal{K}_n(\mathbf{H}, \mathbf{r}) = \text{span} \{ \mathbf{r}, \mathbf{H}\mathbf{r}, \mathbf{H}^2\mathbf{r}, \dots, \mathbf{H}^{n-1}\mathbf{r} \}. \quad (3.1)$$

For the  $n$ -th *block* Krylov-subspace

$$\mathcal{K}_n(\mathbf{H}, \mathbf{R}) = \text{span} \{ \mathbf{R}, \mathbf{H}\mathbf{R}, \mathbf{H}^2\mathbf{R}, \dots, \mathbf{H}^{n-1}\mathbf{R} \} \quad (3.2)$$

$n$  is the dimension of the subspace, not the number  $\eta$  of powers of  $\mathbf{H}$  that are involved.

Krylov-subspace projection methods for MOR come out of methods for finding eigenvalues, so we will first address general projection methods and Krylov projection methods for eigensolving, then delve into MOR.

## 3.2 The Arnoldi process

An  $n$ -iteration cycle of Arnoldi's algorithm constructs a basis  $V$  for the Krylov subspace  $\mathcal{K}_n(\mathbf{H}, \mathbf{r})$ , and the projected operator

$$\tilde{\mathbf{H}} = V^T \mathbf{H} V,$$

which is an upper-Hessenberg matrix. An upper-Hessenberg matrix such as

$$\mathbf{H}_4 = \begin{bmatrix} h_{11} & h_{12} & h_{13} & h_{14} \\ h_{21} & h_{22} & h_{23} & h_{24} \\ & h_{32} & h_{33} & h_{34} \\ & & h_{43} & h_{44} \end{bmatrix}$$

is like an upper-triangular matrix, but with nonzero entries on the 1st subdiagonal.

The Arnoldi process (Algorithm 1) basically performs a power iteration and orthogonalizes each iterate against previous ones, thus producing an orthonormal basis for (3.1). The most costly part

---

**Algorithm 1: ARNOLDI**

---

**Input:**  $\mathbf{r} \in \mathbb{C}^N$ ,  $\mathbf{H} \in \mathbb{C}^{N \times N}$  (or some way to compute  $\mathbf{H}v$  for  $v \in \mathbb{C}^N$ )  
**Output:** orthonormal  $V \in \mathbb{C}^{N \times n}$ , Upper Hessenberg  $\tilde{\mathbf{H}} \in \mathbb{C}^{n \times n}$  where  $\text{span } V = \mathcal{K}_n(\mathbf{H}, \mathbf{r})$ ,  
and  $\tilde{\mathbf{H}} = V^T \mathbf{H} V$

```
1  $r_0 := \mathbf{r}$ 
2  $v_1 := r_0 / \|r_0\|_2$ 
3 for  $k = 1$  to  $n$  do
4    $r_k := \mathbf{H}v_k$ 
5   for  $j = 1$  to  $k$  do      % Make  $r_k$  orthogonal to previous  $\{v_1, v_2, \dots, v_k\}$ 
6      $h_{jk} := v_k^H r_k$ 
7      $r_k := r_k - h_{jk}v_j$ 
8   if  $\|r_k\|_2 \neq 0$  then
9      $h_{j+1,k} := \|r_k\|_2$ 
10     $v_{k+1} := r_k / \|r_k\|_2$ 
11  else exit  $k$ -loop
12 return  $V = [v_1 \ v_2 \ \dots \ v_n]$ ,  $\hat{v}_n = v_{n+1}h_{n+1,n}$ ,  $\tilde{\mathbf{H}} = [h_{ij}]$ 
```

---

of the algorithm is the matrix-vector product (line 4), followed by the orthogonalization part (lines 5-7). Algorithm 1 uses Modified Gram-Schmidt for orthogonalization but there are variants of the Arnoldi process that use other orthogonalizing methods. A notable alternative uses Householder reflectors making for a more stable and more costly method.

### 3.2.1 Complexity of Arnoldi (with MGS orthogonalization)

Take an  $n$ -iteration Arnoldi cycle with a general matrix  $\mathbf{H}$ : there are  $n$  matrix-vector products  $\mathbf{H}v_k$  (line 4), each requiring  $N^2$  scalar multiplications (flops). With Modified Gram-Schmidt (MGS) as the orthogonalization process, we have  $1 + 2 + \dots + n = n(n+1)/2$  inner-products (line 6) and an equal number of AXPYs<sup>1</sup> (line 7), each requiring  $N$  flops. Note that the  $k$ -th step of Arnoldi requires  $kN$  flops for orthogonalization. The process takes longer for each iterate, eventually grinding to a crawl if  $N$  is large. The total cost of an  $n$ -iteration cycle of Arnoldi method is roughly

$$nN^2 + n^2N$$

---

<sup>1</sup> $\alpha X + Y$  operations where  $\alpha$  is a scalar and  $X$  and  $Y$  are vectors.

flops:  $nN^2$  flops for matrix-vector products (sometimes called matvecs), and  $n^2N$  flops for orthogonalization. Clearly the computational cost of Arnoldi is dominated by matvecs. It should be noted that the  $\sigma$ -shifted inverse operator  $\mathbf{H} = (\mathbf{A} - \sigma\mathbf{E})^{-1}\mathbf{E}$  used for model reduction is not a general, dense matrix. The “matrix-vector product”

$$\mathbf{H}v = Q [U^{-1}L^{-1}\mathbf{B}(Pv)]$$

is actually implemented as a pair of sparse triangular solves requiring at most

$$2 \text{nnz}(U) + 2 \text{nnz}(L) \leq 2N(N+1)$$

flops, where  $\text{nnz}(T) \leq N(N+1)/2$  is the number of nonzero entries of an  $N \times N$  triangular matrix  $T$ . The computation of these “matvecs” is still  $\mathcal{O}(N^2)$  so we don’t commit a major crime by viewing the operation as a matrix-vector product, as long as we consider the one-time cost of sparse  $LU = P\mathbf{H}Q$  factorization.

The Arnoldi algorithm requires  $(n+1)N$  units of storage for the basis vectors  $v_j \in V$ , which is also an issue in large applications.

We consider the ROM size  $n$  to be negligible in comparison to the order  $N$  of the full model. Computation and storage cost are major issues when  $N$  is large. For a model of size  $n$  we have no choice but to compute  $n$  applications of  $\mathbf{H}$ , each  $\mathcal{O}(N^2)$ . Restarted Krylov methods attempt to make the process more computationally manageable by reducing the amount of orthogonalization. Since latter iterations require the most computation, the idea is to start over at a certain point.

### 3.2.2 ROM size vs. construction cost

Ultimately the goal of model order-reduction is to produce a small, accurate model (we would like to minimize  $n$  and model approximation error<sup>2</sup>), but the time taken to construct the model needs to be considered as well. In some applications, once a ROM is constructed it gets used repeatedly for several computations.

---

<sup>2</sup>one measure of this is  $\|\mathcal{H} - \mathcal{H}_n\|$  in some norm.

Next, consider the case where producing the model (or several models) is itself the major expense. We may, for example, only need to solve the system once and the original system of order  $N$  is just too large to solve. Maybe a new ROM needs to be generated at every step in some sequence. ROM construction efficiency is where Krylov methods excel. This distinction is important because it sets the context when discussing the best model reduction method for a particular application.

### 3.2.3 The Arnoldi relation

An  $n$ -step cycle of the Arnoldi process with  $\mathbf{H}$  and  $\mathbf{r}$  yields the so-called Arnoldi relation

$$\begin{aligned}\mathbf{H}V &= \begin{bmatrix} V & v_{n+1} \end{bmatrix} \begin{bmatrix} \leftarrow & \tilde{\mathbf{H}} & \rightarrow \\ 0 & \cdots & h_{n+1,n} \end{bmatrix} \\ &= V\tilde{\mathbf{H}} + h_{n+1,n}v_{n+1}e_n^T \\ &= V\tilde{\mathbf{H}} + \hat{v}_ne_n^T\end{aligned}\tag{3.1}$$

where  $V \in \mathbb{C}^{N \times n}$  is the orthogonal basis matrix for  $\mathcal{K}_n(\mathbf{H}, \mathbf{r})$ , starting with  $v_1 = \mathbf{r}/\|\mathbf{r}\|_2$ , and the upper Hessenberg matrix  $\tilde{\mathbf{H}} \in \mathbb{C}^{n \times n}$  is the Petrov-Galerkin projection of  $\mathbf{H}$  on to that space (also known as the Arnoldi matrix), and can be considered a reduced-order spectral approximation to  $\mathbf{H}$ , because eigenvalues of  $\tilde{\mathbf{H}}$  approximate those of  $\mathbf{H}$ . The largest eigenvalues of  $\tilde{\mathbf{H}}$  are the most accurate, a property inherited from the power iteration.

#### The remaining candidate-vector

The last ( $n$ -th) candidate-vector  $\hat{v}_n = h_{n+1,n}v_{n+1}$  of algorithm 1 is a notable quantity because it represents the error of the approximation  $\mathbf{H}V \approx V\tilde{\mathbf{H}}$  (of  $V$  to an invariant-subspace). One expects the sequence to decrease in general, since  $\|\hat{v}_n\|$  is zero for  $n \geq d(\mathbf{H}, \mathbf{r}) \leq N$ , but it is not monotonically decreasing. From a model reduction standpoint, a satisfactory model can be obtained without having  $n$  large enough to make  $\|\hat{v}_n\|$  small. This is because  $\|\hat{v}_n\|$  represents the amount of new spectral information of  $\mathbf{H}$  discovered on the  $n$ -th step, after previously discovered directions  $v_j, j = 1, 2, \dots, n$  have been subtracted off. A rapidly decreasing  $\|\hat{v}_n\|$  indicates that

further iterations with  $\mathbf{H}$  are not producing much new spectral information. Recall that eigenvalues of  $\mathbf{H}$  correspond to poles of the transfer-function. As long as new poles are being discovered (starting from  $\sigma$  and moving outward),  $\hat{v}_n$  is rich with information. The poles being discovered may or may not be significant in the sense of frequency response approximation, however.

### 3.2.4 Approximate eigenvalues from Arnoldi relation

The Arnoldi relation (3.1) (or Krylov relation in general) is what makes Krylov-subspaces ideal for subspace projection eigenvalue methods. Compared to  $\mathbf{H}$ , the matrix  $\tilde{\mathbf{H}}$  is small enough for its eigenvalue decomposition

$$\tilde{\mathbf{H}}W = W\Lambda.$$

Ritz-values (eigenvalues  $\lambda$  of  $\tilde{\mathbf{H}}$ ) are approximate eigenvalues of the large operator  $\mathbf{H}$ , and long<sup>3</sup> Ritz-vectors  $Vw \in \mathcal{K}_n(\mathbf{H}, \mathbf{r})$  are the associated approximate eigenvectors of  $\mathbf{H}$ .

Left multiplying (3.1) with  $W$  yields

$$\begin{aligned} \mathbf{H}VW &= V\tilde{\mathbf{H}}W + \hat{v}_n e_n^T W \\ &= VW\Lambda + \hat{v}_n e_n^T W \end{aligned} \tag{3.2}$$

$$\mathbf{H}Z = Z\Lambda + \hat{v}_n e_n^T W \tag{3.3}$$

so for  $z_j = Vw_j$ ,

$$\begin{aligned} \mathbf{H}z_j &= \lambda z_j + \xi_j \hat{v}_n \\ &= \lambda z_j + \xi_j h_{n+1,n} v_{n+1} \end{aligned}$$

where  $\xi_j = (e_n^T W)_j = W_{nj} \in \mathbb{C}$  is the  $j$ -th entry of the bottom ( $n$ -th) row of  $W$ . Here we see that every Ritz residual-vector

$$\mathbf{H}z_j - \lambda_j z_j = \xi_j \hat{v}_n$$

given by (3.3) is a scalar multiple of the residual-vector  $\hat{v}_n$ . Assuming  $\|z_j\|_2 = 1$ , the Arnoldi-

---

<sup>3</sup>vectors  $w \in W$  are sometimes called short Ritz-vectors.



relation (3.1) thus implies a simple formulation

$$\text{rr}_j = \frac{\|\mathbf{H}z_j - \lambda_j z_j\|_2}{\|\lambda_j z_j\|_2} = \frac{|\xi_j|}{|\lambda_j|} \|\hat{v}_n\|_2 = \frac{|\xi_j|}{|\lambda_j|} |h_{n+1,n}| \quad (3.4)$$

for the relative residual-errors of the Ritz-values/vectors. Ritz-pairs with low associated relative residual norms (3.4) are good approximations to eigenvalues/vectors of  $\mathbf{H}$  if they are well-conditioned. This is because (3.4) actually indicates that  $(\lambda_j, z_j)$  is an exact eigen-pair of the perturbed matrix  $\mathbf{H} + \mathcal{E}$  where the norm of the perturbation  $\|\mathcal{E}\| = \|\hat{v}_n\|$ . Rearrangement of the Arnoldi relation (3.1) reveals

$$\mathbf{H}V - \hat{v}_n e_n^T = (\mathbf{H} - \hat{v}_n v_n^T)V = V\tilde{\mathbf{H}}.$$

If an eigenvalue of  $\mathbf{H}$  is highly sensitive to perturbation (is badly conditioned) then a low or zero residual-norm (3.4) could be misleading. We can avoid this situation by noting that the largest eigenvalues of  $\mathbf{H}$  converge first. The relative residual errors are likely to be accurate for Ritz values of largest magnitude, which we expect to converge first. A very small eigenvalue of  $\mathbf{H}$  with small relative-residual error may be suspect.

### 3.3 Implicit vs. explicit Ritz-values and vectors

It should be noted that although eigenvalues  $\lambda$  of  $\mathbf{H} = (\mathbf{A} - \sigma \mathbf{E})^{-1} \mathbf{E}$  and  $\mu$  of  $(\mathbf{A}, \mathbf{E})$  are related by  $\lambda = 1/(\mu - \sigma)$  and share common eigenvectors, the same cannot be said for approximate eigenvalues  $\hat{\lambda}$  of  $\tilde{\mathbf{H}} = V^T \mathbf{H} V$  and  $\hat{\mu}$  of  $(\mathbf{A}_n, \mathbf{E}_n) = (V^T \mathbf{A} V, V^T \mathbf{E} V)$ .

The reader should note that there are two different sets of approximations to the spectrum of  $(\mathbf{A}, \mathbf{E})$ , both implied by projection on to  $\mathcal{K}_n(\mathbf{A}, \mathbf{R})$  via basis  $V$ :

**Implicit Ritz-values** The set of implicit Ritz-values

$$\left\{ 1/\hat{\lambda} + \sigma \mid \hat{\lambda} \in \sigma(\tilde{\mathbf{H}}) \right\} \quad (3.1)$$

of  $(A, E)$  with respect to  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ , where

$$\tilde{\mathbf{H}} = V^T \mathbf{H} V = V^T (\mathbf{A} - \sigma \mathbf{E})^{-1} \mathbf{E} V$$

is a byproduct of constructing  $V$  by  $n$  steps of the Arnoldi algorithm.

**Explicit Ritz-values** The set of of explicit Ritz-values

$$\{ \hat{\mu} \in \sigma(\mathbf{A}_n, \mathbf{E}_n) \} \quad (3.2)$$

of  $(A, E)$  with respect to  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$  where

$$(\mathbf{A}_n, \mathbf{E}_n) = (V^T \mathbf{A} V, V^T \mathbf{E} V),$$

is not implied by the Arnoldi process and must be computed. (3.1) and (3.2) are not equal in general but are related in that they both converge to the same spectrum  $\sigma(\mathbf{A}, \mathbf{E})$ . Note that both sets of approximate eigenvalues are dependent on  $\sigma$ ; we expect eigenvalue approximations closer to  $\sigma$  to be more accurate for both (3.1) and (3.2), because they both result from projection on to the Krylov-subspace  $\mathcal{K}_n(\mathbf{A}, \mathbf{R})$ , where  $\mathbf{H} = \mathbf{H}(\sigma)$  and  $\mathbf{R} = \mathbf{R}(\sigma)$ .

The values of the associated approximate eigenvectors are not dependent on  $\sigma$ . Only the order in which they converge depends on  $\sigma$ . Eigenvectors associated with (3.1) and with (3.2) are not equal in general, but sufficiently converged vectors are nearly equal.

When converged, implicit and explicit eigen-pairs are nearly identical. We consider approximate eigenvalues/vectors coming from explicit and implicit computation to be interchangeable if they are near  $\sigma$  and have low relative residual-error norm (3.4). Thus, if an eigen-pair  $(\hat{\lambda}_j, w_j)$  of  $\tilde{\mathbf{H}}$  is converged, then we can expect that  $(\sigma + 1/\hat{\lambda}_j, V w_j)$  is a converged eigen-pair of  $(\mathbf{A}_n, \mathbf{E}_n)$  with about the same order of error of approximation to an eigen-pair of  $(\mathbf{A}, \mathbf{E})$ .

The reason we care about both sets of approximate eigenvalues/vectors is that implicit (3.1) Ritz-values/vectors are far cheaper to compute than the explicit variety (3.2), but the explicit formulation (2.2) is the end goal of explicit-projection-based MOR. Un-converged poles of the implicitly

projected model transfer-function can and often do have positive real-part, which is unfavorable for ROM applications. These eigenvalue approximations all move to the left half of the complex plane as they converge to their final resting values, but as long as there are any implicit Ritz-values  $1/\hat{\lambda} + \sigma$  with positive real part, the implicitly projected model (2.2) is possibly unstable and not attractive for model order-reduction.<sup>4</sup> Implicitly obtained eigen-information is useful feedback to gauge and possibly direct progress of an adaptive method. Some MOR methods, typically called *restarted* methods including [15, 24, 17, 2], have been developed which attempt to purge subspace components associated with “bad” (destabilizing) or otherwise unwanted eigenvalues from the constructed basis  $V$ , but they destroy moment-matching properties and introduce other problems. Explicitly projected eigenvalues  $\hat{\mu}$  of  $(\mathbf{A}_n, \mathbf{E}_n)$  are always (for any  $n$ ) well-behaved as long as the projection basis  $V$  is real, and as long as  $\mathcal{K}_n(\mathbf{A}(\sigma), \mathbf{R}(\sigma)) \subseteq \text{span } V$ , the explicitly projected ROM on to  $V$  is guaranteed to be of matrix-Padé-type (match moments) with respect to  $\sigma$ .

### 3.4 Moment-matching property of Krylov-subspace projected ROM

We are going to prove that a reduced-order model implied by orthogonal projection (via one orthonormal basis  $V = \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix}$ ) on to a Krylov-subspace matches  $l$  moments about  $\sigma$ , where  $l$  is the block-degree of the Krylov-subspace

$$\text{span} \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} = \mathcal{K}_l(\mathbf{H}, \mathbf{R}) = \text{span} \begin{bmatrix} \mathbf{R} & \mathbf{H}\mathbf{R} & \mathbf{H}^2\mathbf{R} & \cdots & \mathbf{H}^{l-1}\mathbf{R} \end{bmatrix}.$$

We will show this for implicitly projected ROMs (2.2) in Theorem 1, and explicitly projected ROMs (2.1) in Theorem 2. It is interesting to note that both implicitly projected and explicitly projected ROMs match the same moments about an expansion-point. Significant differences in the two ROM approximations are present away from expansion-point(s)  $\sigma$ , but near  $\sigma$  they are approximations of the same order.

Recall the URM (unreduced model) transfer-function  $\mathcal{H}(s) = \mathbf{C}^T(\mathbf{A} - s\mathbf{E})^{-1}\mathbf{B}$  of LTI descriptor system (2.1), and its equivalent shift-invert formulation  $\mathcal{H}(s) = \mathbf{C}^T(I - (s - \sigma)\mathbf{H})^{-1}\mathbf{R}$  with shift

---

<sup>4</sup>Implicitly projected ROMs, such as those produced by PVL [7] often work fine in many practical applications despite being unstable, but they are currently unpopular.

$\sigma$ , or

$$\mathcal{H}(s + \sigma) = \mathbf{C}^T (I - s\mathbf{H})^{-1} \mathbf{R} \quad (3.1)$$

$$= \sum_{j=0}^{\infty} s^j \mathcal{H}^{(j)} \quad (2.2)$$

where (2.2) is the Taylor series expansion of  $\mathcal{H}(s)$  about  $\sigma \in \mathbb{C}$ . The  $j$ -th moment  $\mathcal{H}^{(j)}$  was shown in §2.2.1 to be

$$\mathcal{H}^{(j)} = \mathbf{C}^T \mathbf{H}^j \mathbf{R}. \quad (3.2)$$

### 3.4.1 Moment matching of the implicitly-projected ROM

The implicitly projected ROM transfer-function

$$\tilde{\mathcal{H}}(s + \sigma) = \mathbf{C}_n^T (I - s\tilde{\mathbf{H}})^{-1} \tilde{\boldsymbol{\rho}}_n \quad (3.3)$$

is defined via projection of (3.1), as

$$\tilde{\mathbf{H}} = V^T \mathbf{H} V, \quad \mathbf{C}_n = V^T \mathbf{C}, \quad \tilde{\boldsymbol{\rho}}_n = V^T \mathbf{R} \quad (3.4)$$

rather than by projecting the system realization  $(\mathbf{A}, \mathbf{E}, \mathbf{B}, \mathbf{C})$ , hence its specification as the transfer-function for an *implicitly* projected model. Moments of (3.3) about  $\sigma$  are given as  $\tilde{\mathcal{H}}^{(j)} = \mathbf{C}_n^T \tilde{\mathbf{H}}_n^j \tilde{\boldsymbol{\rho}}_n$ .

**Theorem 1.** *Suppose the span of an orthonormal basis  $V \in \mathbb{R}^{N \times n}$  contains the Krylov-subspace  $\mathcal{K}_l(\mathbf{H}, \mathbf{R})$  of block-degree  $l$  for some  $l \leq n$ . Then moments of  $\tilde{\mathcal{H}}^{(j)}(\sigma)$  of the implicitly projected ROM transfer-function (2.2), (3.3) and moments  $\mathcal{H}^{(j)}(\sigma)$  of the URM transfer-function (2.1) about  $\sigma$  are related by*

$$\tilde{\mathcal{H}}^{(j)} = \mathbf{C}_n^T \tilde{\mathbf{H}}_n^j \tilde{\boldsymbol{\rho}}_n = \mathbf{C}^T \mathbf{H}^j \mathbf{R} = \mathcal{H}^{(j)} \quad (3.5)$$

for  $j = 0, 1, \dots, l - 1$ .

*Proof.* The theorem follows from left-applying  $\mathbf{C}^T$  to

$$\mathbf{H}^j \mathbf{R} = V \tilde{\mathbf{H}}^j \tilde{\boldsymbol{\rho}}_n \quad (3.6)$$

for  $j = 0, 1, \dots, l-1$ , which we will show by induction.

For  $j = 0$ , (3.6) follows from (3.4). Now assume (3.6) holds for some  $j \in \{0, 1, \dots, l-2\}$ .

Applying  $\mathbf{H}$  to (3.6) yields

$$\begin{aligned} \mathbf{H}(\mathbf{H}^j \mathbf{R}) &= \mathbf{H}^{j+1} \mathbf{R} = \mathbf{H}(V \tilde{\mathbf{H}}^j \tilde{\boldsymbol{\rho}}_n) \\ &= V \tilde{\mathbf{H}} \tilde{\mathbf{H}}^j \tilde{\boldsymbol{\rho}}_n, \quad \text{since } \mathbf{H}V = V \tilde{\mathbf{H}} \\ &= V \tilde{\mathbf{H}}^{j+1} \tilde{\boldsymbol{\rho}}_n. \end{aligned}$$

□

### 3.4.2 Moment matching of the explicitly-projected ROM

Proof of moment-matching for the explicitly projected ROM (2.2) transfer-function is a little more involved than Theorem 1 for the implicitly projected model. It is included as Theorem 2. The proof is adapted from [9, proposition 6 and theorem 7].

Recall the explicitly-projected ROM (2.2) with transfer-function

$$\hat{\mathcal{H}}_n(s) = \mathbf{C}_n^T (s \mathbf{E}_n - \mathbf{A}_n)^{-1} \mathbf{B}_n, \quad (2.1)$$

where the system realization  $(\mathbf{A}, \mathbf{E}, \mathbf{B}, \mathbf{C})$  is said to be *explicitly* projected as

$$\mathbf{A}_n := V^T \mathbf{A} V, \quad \mathbf{E}_n := V^T \mathbf{E} V, \quad \mathbf{C}_n := V^T \mathbf{C}, \quad \mathbf{B}_n := V^T \mathbf{B}.$$

Moments of the ROM transfer-function (2.1) are

$$\hat{\mathcal{H}}^{(j)} = \mathbf{C}_n^T \hat{\mathbf{H}}^j \hat{\boldsymbol{\rho}}_n,$$

where the structures

$$\widehat{\mathbf{H}} := (\mathbf{A}_n - \sigma \mathbf{E}_n)^{-1} \mathbf{E}_n \quad \text{and} \quad \widehat{\boldsymbol{\rho}}_n := (\sigma \mathbf{E}_n - \mathbf{A}_n)^{-1} \mathbf{B}_n. \quad (3.7)$$

are analogous to the shift-invert operator and start-block

$$\mathbf{H} := (\mathbf{A} - \sigma \mathbf{E})^{-1} \mathbf{E}, \quad \mathbf{R} := (\sigma \mathbf{E} - \mathbf{A})^{-1} \mathbf{B} \quad (2.5)$$

of the unreduced model (2.1). The proof of Theorem 1 depended on  $\widetilde{\mathbf{H}} = V^T \mathbf{H} V$ , which we do not have for the explicitly projected ROM. In general,  $\widehat{\mathbf{H}} \neq V^T \mathbf{H} V$ . However, for an appropriate choice of  $F_n$ ,

$$\widehat{\mathbf{H}} = V^T F_n V$$

implies that (2.1) matches  $l$  moments.

**Theorem 2.** *Suppose the span of an orthonormal basis  $V \in \mathbb{R}^{N \times n}$  contains the Krylov-subspace  $\mathcal{K}_l(\mathbf{H}, \mathbf{R})$  of block-degree  $l$  for some  $l \leq n$ , and let*

$$F_n := V(\mathbf{A}_n - \sigma \mathbf{E}_n)^{-1} V^T \mathbf{E}. \quad (3.8)$$

*Then for  $j \leq l \leq n$ , the  $j$ -th moment  $\widehat{\mathcal{H}}^{(j)}$  of the explicitly projected ROM transfer-function (2.1) and moment  $\mathcal{H}^{(j)}$  of the unreduced model (2.1) are related by*

$$\widehat{\mathcal{H}}^{(j)} = \mathbf{C}^T F_n^i \mathbf{R} \quad \text{for } i = 0, 1, \dots \quad (3.9)$$

$$= \mathcal{H}^{(i)} \quad \text{for } i = 0, 1, \dots, j-1. \quad (3.10)$$

*Proof.* First we show (3.9). Since  $\text{span } \mathbf{H}^i \mathbf{R} \subseteq \mathcal{K}_l(\mathbf{H}, \mathbf{R})$  for  $i = 0, 1, \dots, j$  and  $\mathcal{K}_l(\mathbf{H}, \mathbf{R}) \subseteq \text{span } V$ , for each  $i = 1, 2, \dots, j$  there is a matrix  $X_i$  such that

$$\mathbf{H}^{i-1} \mathbf{R} = V X_i. \quad (3.11)$$

Recall that  $\mathbf{R} = (\sigma \mathbf{E} - \mathbf{A})^{-1} \mathbf{B}$ . Then for  $i = 1$ ,

$$\mathbf{B} = (\sigma \mathbf{E} - \mathbf{A}) \mathbf{R} = (\sigma \mathbf{E} \mathbf{V} - \mathbf{A} \mathbf{V}) X_1,$$

which when left-multiplied by  $V^T$  results in

$$V^T \mathbf{B} = V^T (\sigma \mathbf{E} \mathbf{V} - \mathbf{A} \mathbf{V}) X_1$$

$$\mathbf{B}_n = (\sigma \mathbf{E}_n - \mathbf{A}_n) X_1.$$

Then

$$X_1 = (\sigma \mathbf{E}_n - \mathbf{A}_n)^{-1} \mathbf{B}_n = \hat{\boldsymbol{\rho}}_n. \quad (3.12)$$

Right-multiplying (3.8) with  $V$  gives  $F_n V = V(\mathbf{A}_n - \sigma \mathbf{E}_n)^{-1} \mathbf{E}_n = V \hat{\mathbf{H}}$ , and by induction on  $i$ ,

$$F_n^i V = V \hat{\mathbf{H}}^i \quad \text{for } i = 0, 1, \dots \quad (3.13)$$

Then moments of the ROM transfer-function

$$\begin{aligned} \hat{\mathcal{H}}^{(i)} &= \mathbf{C}_n^T \hat{\mathbf{H}}^i \hat{\boldsymbol{\rho}}_n = \mathbf{C}^T V \hat{\mathbf{H}}^i \hat{\boldsymbol{\rho}}_n \\ &= \mathbf{C}^T (F_n^i V) X_1 \quad \text{by (3.12) and (3.13)} \\ &= \mathbf{C}^T F_n^i \mathbf{R} \quad \text{by (3.11) with } i = 1, \end{aligned}$$

which is (3.9).

Proof of (3.10) is implied by

$$\mathbf{H}^j \mathbf{R} = F_n^j \mathbf{R} \quad \text{for } i = 0, 1, \dots, j-1 \quad (3.14)$$

which we show by induction on  $i$ . (3.14) is trivial for  $i = 0$ . Now assume (3.14) is satisfied for some  $i \in \{0, 1, j-2\}$ . We will show that

$$F_n^{i+1} \mathbf{R} = \mathbf{H}^{i+1} \mathbf{R}$$

as follows:

$$((\mathbf{A} - \sigma \mathbf{E})^{-1} \mathbf{E})(F_n^i \mathbf{R}) = \mathbf{H}(\mathbf{H}^i \mathbf{R}) = \mathbf{H}^{i+1} \mathbf{R} = V X_{i+2} \quad (3.15)$$

where the rightmost expression follows from (3.11). Left-multiplying (3.15) with  $V^T(\mathbf{A} - \sigma \mathbf{E})$  yields

$$(V^T \mathbf{E})(F_n^i \mathbf{R}) = (V^T(\mathbf{A} - \sigma \mathbf{E})V)X_{i+2} = (\mathbf{A}_n - \sigma \mathbf{E}_n)X_{i+2}. \quad (3.16)$$

Then

$$\begin{aligned} F_n^{i+1} \mathbf{R} &= F_n(F_n^i \mathbf{R}) \\ &= \left( V(\mathbf{A}_n - \sigma \mathbf{E}_n)^{-1} V^T \mathbf{E} \right) (F_n^i \mathbf{R}) \\ &= V(\mathbf{A}_n - \sigma \mathbf{E}_n)^{-1} (V^T \mathbf{E})(F_n^i \mathbf{R}) \\ &= V X_{i+2}, \quad \text{by (3.16)} \\ &= \mathbf{H}^{i+1} \mathbf{R}, \quad \text{by (3.11)} \end{aligned}$$

which proves (3.14). Applying  $\mathbf{C}^T$  yields (3.10), i.e.

$$\widehat{\mathcal{H}}^{(j)} = \mathbf{C}^T F_n^j \mathbf{R} = \mathbf{C}^T \mathbf{H}^j \mathbf{R} = \mathcal{H}^{(j)}$$

□



## Chapter 4

# Interpolation-point selection and resulting projection-bases

### 4.1 Multiple point moment-matching

Theorems 1 and 2 can be extended to imply moment-matching about any number of expansion-points if the projection subspace contains the appropriate Krylov-subspaces. Much of the pioneering rational interpolation research, notably the rational-Lanczos method [11] (and [13]) for model order-reduction was done by Grimme in the mid and late 1990s. It is somewhat based on Ruhe’s Rational-Krylov [29, 27] eigenvalue method and formalization, and possibly Olsson’s [22]. Of particular interest are [14] and [6], both of which discuss interpolation-point selection. We refer the reader to those sources for the details of point selection.

[19, 20, 18, 8]) are more recent multi-point rational-interpolation methods. Also a Jacobi-Davidson MOR method [3]. Lee, Chu, and Feng’s RAMAO/AORA method (Rational Arnoldi Method with Adaptive Order selection/ Adaptive-Order Rational-Arnoldi) [19, 20] breaths new life into an adaptive point-selection method introduced by [11], based on the sequence of ROM *moment-errors* implied by the sequence of candidate-vectors  $\hat{v}_k$  of the Arnoldi process §3.2.

In rational-Krylov method literature, the shifts/interpolation-points are usually denoted by  $\sigma$  and we will adopt that convention for this chapter. Suppose that for each  $j = 1, 2, \dots, \tau$  the

Krylov-subspace

$$\mathcal{K}_j = \mathcal{K}_{n_j}(\mathbf{H}_j, \mathbf{R}_j) = \text{span} \begin{bmatrix} \mathbf{R}_j & \mathbf{H}_j \mathbf{R}_j & \mathbf{H}_j^2 \mathbf{R}_j & \cdots & \mathbf{H}_j^{l_j-1} \mathbf{R}_j \end{bmatrix}$$

of dimension  $n_j$ , induced by

$$\mathbf{H}_j := \mathbf{H}(\sigma_j) = (\mathbf{A} - \sigma_j \mathbf{E})^{-1} \mathbf{E} \quad \text{and} \quad \mathbf{R}_j := \mathbf{R}(\sigma_j) = (\sigma_j \mathbf{E} - \mathbf{A})^{-1} \mathbf{B}$$

is contained in the span of  $V$ , so that

$$\mathcal{K}_1 \cup \mathcal{K}_2 \cup \cdots \cup \mathcal{K}_\tau \subseteq \text{span } V \tag{4.1}$$

Then the ROM implied by orthogonal projection on to  $\text{span } V$  matches  $l_j$  moments about interpolation-point  $\sigma_j$  for each  $j = 1, 2, \dots, \tau$ .

#### 4.1.1 Merging bases

There are several ways to produce a basis for the composite space (4.1). The naive method suggested by our previous discussion of single-point Krylov methods is to use  $\tau$  consecutive runs of a basic Krylov method like Algorithm 1 (Arnoldi), each producing an orthonormal basis  $V_j$  for which

$$\text{span } V_j = \mathcal{K}_j,$$

and then somehow putting the bases together into  $V = \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix}$ , where

$$\text{span} \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix} = \text{span } V_1 \cup \text{span } V_2 \cup \cdots \cup \text{span } V_\tau. \tag{4.2}$$

For the general application of rational-interpolation we assume that complex interpolation-points are used. As will be discussed in §4.2 we are typically required to split the basis vectors into  $\Re$  and  $\Im$  parts. For this reason it is not necessary for the individual bases  $V_j$  to be orthogonal to one-another, or even linearly-independent. However, an  $\mathbf{H}_j$ -invariant-subspace  $\mathcal{Y}$  contained in

$\mathcal{K}_{l_j}(\mathbf{H}_j, \mathbf{R}_j)$  is also  $(\mathbf{A}, \mathbf{E})$ -invariant and thus has global significance.

The naive approach of producing bases for  $\mathcal{K}_{n_j}(\mathbf{H}(\sigma_j), \mathbf{R}(\sigma_j))$  separately and combining them in a post-processing step is inefficient because there is a significant degree of overlap between spaces. Recall that an invariant-subspace under  $\mathbf{H}(\sigma)$  is independent of the expansion-point (shift)  $\sigma$ . Suppose  $\mathbf{H}_1 = \mathbf{H}(\sigma_1)$  and  $\mathbf{H}_2 = \mathbf{H}(\sigma_2)$ . Then an invariant-subspace  $\mathcal{Y} \subset \text{span } V_1$  under  $\mathbf{H}_1$  is also invariant under  $\mathbf{H}_2$ .

It would be wasteful to spend computational effort re-discovering  $(\mathbf{A}, \mathbf{E})$ -invariant-subspace while computing a basis for  $\mathcal{K}_{l_j}(\mathbf{H}_j, \mathbf{R}_j)$ , if we already discovered it while constructing the basis  $V_{j-1}$  for  $\mathcal{K}_{l_{j-1}}(\mathbf{H}_{j-1}, \mathbf{R}_{j-1})$ . Traditional rational-interpolation methods for MOR such as rational-Lanczos [11] avoid this issue by doing full, Arnoldi-style orthogonalization of an iterate against every previous vector, generating one orthogonal basis  $V$  for (4.2) directly. Considering that for complex-valued interpolation-points we will have to split the basis (4.2) anyway, thus losing any orthogonality, doing full orthogonalization of every complex basis-vector is overkill.

#### 4.1.2 Interpolation-point translation

Suppose we have an approximate eigen-pair  $(\lambda, z)$  of  $\mathbf{H}_1 = \mathbf{H}(\sigma_1)$  so that

$$\mathbf{H}_1 z = \lambda z + r$$

and we would like to know its residual under another member  $\mathbf{H}_2 = \mathbf{H}(\sigma_2)$  of the shift-invert operator family

$$\left\{ \mathbf{H}(\sigma) = (\mathbf{A} - \sigma \mathbf{E})^{-1} \mathbf{E} \right\}. \quad (4.3)$$

Then

$$\mathbf{H}_2 z = T \mathbf{H}_1 z = \lambda T z + T r, \quad (4.4)$$

where  $T$  is the transformation

$$\begin{aligned} T(\sigma_1, \sigma_2) &:= (\mathbf{A} - \sigma_2 \mathbf{E})^{-1}(\mathbf{A} - \sigma_1 \mathbf{E}) \\ &= (\sigma_2 - \sigma_1) \mathbf{H}_2 + I. \end{aligned} \tag{4.5}$$

Note that  $T\mathbf{H}_1 = \mathbf{H}_2$  and  $T\mathbf{R}_1 = \mathbf{R}_2$ .

*Proof of (4.5).* The expression (4.5) comes from observing that

$$\begin{aligned} \tilde{v} &= Tv = (\mathbf{A} - \sigma_2 \mathbf{E})^{-1}(\mathbf{A} - \sigma_1 \mathbf{E})v, \\ (\mathbf{A} - \sigma_2 \mathbf{E})\tilde{v} &= (\mathbf{A} - \sigma_1 \mathbf{E})v \\ \mathbf{A}\tilde{v} - \sigma_2 \mathbf{E}\tilde{v} &= \mathbf{A}v - \sigma_1 \mathbf{E}v \\ \mathbf{A}(\tilde{v} - v) - \sigma_2 \mathbf{E}(\tilde{v} - v) &= (\sigma_2 - \sigma_1) \mathbf{E}v \\ \tilde{v} - v &= (\sigma_2 - \sigma_1)(\mathbf{A} - \sigma_2 \mathbf{E})^{-1} \mathbf{E}v \\ &= (\sigma_2 - \sigma_1) \mathbf{H}_2 v. \end{aligned}$$

□

The translation transformation (4.5) must be invertible, with

$$(T(\sigma_1, \sigma_2))^{-1} = T(\sigma_2, \sigma_1) = (\sigma_1 - \sigma_2) \mathbf{H}_1 + I. \tag{4.6}$$

For easier notation let  $\Delta = \sigma_2 - \sigma_1$ , so that  $T = T(\sigma_1, \sigma_2) = \Delta \mathbf{H}_2 + I$  and  $T^{-1} = -\Delta \mathbf{H}_1 + I$ .

Then (4.5) and (4.6) imply that

$$(\Delta \mathbf{H}_2 + I)(-\Delta \mathbf{H}_1 + I) = (-\Delta \mathbf{H}_1 + I)(\Delta \mathbf{H}_2 + I) = I,$$

from which it follows that

$$\mathbf{H}_2 \mathbf{H}_1 = \frac{\mathbf{H}_2 - \mathbf{H}_1}{\sigma_2 - \sigma_1} \tag{4.7}$$

and for  $\sigma_2 \neq \sigma_1$ ,

$$\mathbf{H}_1 \mathbf{H}_2 = \mathbf{H}_2 \mathbf{H}_1. \quad (4.8)$$

(4.8) shows that the set (4.3) of operators, commutes. (4.7) implies that

$$\frac{d}{d\sigma} \mathbf{H}(\sigma) = (\mathbf{H}(\sigma))^2.$$

It might interest the reader to observe that the shift-invert transfer-function representation

$$\mathcal{H}(s) = \mathbf{C}^T (I - (s - \sigma) \mathbf{H}(\sigma))^{-1} \mathbf{R}(\sigma), \quad (4.9)$$

defined about the interpolation-point  $\sigma$  but independent of  $\sigma$ , involves a transformation of the form (4.5), since

$$I - (s - \sigma) \mathbf{H} = (\sigma - s) \mathbf{H} + I = T(s, \sigma).$$

Then we may re-write (4.9) as

$$\begin{aligned} \mathcal{H}(s) &= \mathbf{C}^T T(\sigma, s) \mathbf{R}(\sigma) \\ &= \mathbf{C}^T \mathbf{R}(s) \\ &= \mathbf{C}^T (\mathbf{A} - s \mathbf{E})^{-1} \mathbf{B}. \end{aligned}$$

Now back to Ritz-residual translation. Recall from (4.4) that for eigen-pair  $(\lambda, z)$  of  $\mathbf{H}_1$  we have

$$\mathbf{H}_2 z = T \mathbf{H}_1 z = \lambda T z + T r$$

with the translation  $T = T(\sigma_1, \sigma_2) = \Delta \mathbf{H}_2 + I$  from (4.5) and  $\Delta = \sigma_2 - \sigma_1$ .

Then

$$\begin{aligned} \mathbf{H}_2 z &= \lambda (\Delta \mathbf{H}_2 + I) z + T r \\ &= \lambda z + \lambda \Delta \mathbf{H}_2 z + T r, \end{aligned}$$

so

$$(1 - \lambda\Delta)\mathbf{H}_2 z = \lambda z + Tr. \quad (4.10)$$

Define the scaling factor

$$\zeta := \frac{1}{1 - \lambda\Delta} = \frac{\mu - \sigma_1}{\mu - \sigma_2}$$

where  $\mu = 1/\lambda + \sigma_1$  is the approximate eigenvalue of  $(\mathbf{A}, \mathbf{E})$  associated with  $\lambda$ . Then

$$\begin{aligned} \mathbf{H}_2 z - (\zeta\lambda)z &= \zeta Tr \\ &= \zeta(\Delta\mathbf{H}_2 + I)r \end{aligned} \quad (4.11)$$

## 4.2 Complex expansion-points

If a shift  $\sigma_j \in \mathbb{C}$  is not strictly-real then the basis  $V_j \in \mathbb{C}^{N \times n_j}$  of its associated Krylov-subspace is also complex. Although Grimme discusses interpolation point selection in some depth in [14], properties of a ROM obtained via projection with a complex basis have not been fully explored. They are generally avoided in part due to the extra computation and storage required for complex arithmetic.

It should be noted however that using  $\sigma \in \mathbb{R}$  is only half as efficient as it appears to be and  $\sigma \in \mathbb{C}$  with  $\Re(\sigma) \neq 0$  is potentially twice as efficient as it appears. That is because the system pencil  $(\mathbf{A}, \mathbf{E})$  is real and its complex eigenvalues are conjugate pairs. If  $\sigma \in \mathbb{R}$ , eigenvalues of  $\mathbf{H} = (\mathbf{A} - \sigma\mathbf{E})^{-1}\mathbf{E}$  must converge pairwise, so each conjugate-pair of converged vectors provide only one piece of complex spectral information. Eigenvalues of  $\mathbf{H}$  for complex  $\sigma$  do not converge in pairs, but each converged eigenvalue  $\lambda$  implies that the pole  $\mu = \sigma + 1/\lambda$  and its conjugate  $\bar{\mu}$  is converged. For reasons discussed next we generally split a complex basis into  $\Re$  and  $\Im$  parts, so there is not as much difference between real and complex interpolation-points  $\sigma$  as there seems.

Real bases are preferred because the system (2.1) realization  $(\mathbf{A}, \mathbf{E}, \mathbf{B}, \mathbf{C})$  consists of real matrices; explicit-projection with a real basis yields a ROM characterized by  $(\mathbf{A}_n, \mathbf{E}_n, \mathbf{B}_n, \mathbf{C}_n)$  which is then also real and thus retains properties of the original model. One such property is symmetry of the transfer-function about the  $\Re$ -axis.

The typical procedure to obtain a real basis for a complex Krylov-subspace is to split the  $n$  vector basis  $V$  into  $2n$  real vectors  $v_j^{\mathbf{r}} = \Re(v_j)$  and  $v_j^{\mathbf{i}} = \Im(v_j)$ , forming the so-called split-basis  $V^* \in \mathbb{R}^{N \times 2n}$ , which spans the so-called *split-subspace*<sup>1</sup>

$$\begin{aligned}
& \text{span} \begin{bmatrix} v_1^{\mathbf{r}} & v_1^{\mathbf{i}} & v_2^{\mathbf{r}} & v_2^{\mathbf{i}} & \cdots & v_n^{\mathbf{r}} & v_n^{\mathbf{i}} \end{bmatrix} \\
&= \text{span } \mathcal{K}_\eta(\mathbf{H}, \mathbf{R}) \cup \mathcal{K}_\eta(\overline{\mathbf{H}}, \overline{\mathbf{R}}) \\
&= \text{span } \mathcal{K}_\eta(\mathbf{H}(\sigma), \mathbf{R}(\sigma)) \cup \mathcal{K}_\eta(\mathbf{H}(\bar{\sigma}), \mathbf{R}(\bar{\sigma})) \\
&= \mathcal{K}_\eta(\mathbf{H}, \mathbf{R})^*
\end{aligned} \tag{4.1}$$

of dimension  $\eta \leq 2n$ . The basis for a standard Krylov-subspace may have complex vectors but its span is generally considered over  $\mathbb{R}$ . The split complex Krylov-subspace admits a real basis but its span is over  $\mathbb{C}$ , so it should still be considered a complex space that contains  $\mathcal{K}_\eta(\mathbf{H}, \mathbf{R})$ .

Notice that the split Krylov-subspace (4.1) is the union of two Krylov-subspaces with complex conjugate shifts  $\sigma$  and  $\bar{\sigma}$  so we may consider a complex shift to be two shifts. Saad calls this idea “double-shifting” in [25], where it was first given significant analysis. Matching moments about a conjugate pair of points  $\sigma$  and  $\bar{\sigma}$  is not as advantageous so much as an unavoidable effect of requiring a real basis. Indeed, convergence of an eigenvalue  $\mu \in \mathbb{C}$  of  $(\mathbf{A}, \mathbf{E})$  with associated vector  $z$  is equivalent to convergence of the eigen-pair  $(\bar{\mu}, \bar{z})$  of  $(\mathbf{A}, \mathbf{E})$  as well. It would seem that the basis and the resulting ROM are potentially twice as large. This is true in theory, but in practice a complex quantity  $z = \alpha + i\beta \in \mathbb{C}$  is represented by two real quantities  $\alpha, \beta \in \mathbb{R}$  anyway. A complex ROM realization of order  $n$  is deceptively small because of the complex quantities involved. We avoid this ambiguity by always referring to the model size  $n$  as the number of vectors of the real projection basis  $V$ .

#### 4.2.1 Producing a real basis for a complex Krylov-subspace

If we use a shift  $\sigma$  with nonzero  $\Im$  part but use real basis for projection, we have no choice but to project on to a split-Krylov space of the form (4.1). One way to do that is to perform a run of the Arnoldi process (algorithm 1) in complex arithmetic with matrices  $\mathbf{H}(\sigma)$  and  $\mathbf{R}(\sigma)$ , yielding the

---

<sup>1</sup>this procedure is not novel, although only in this text do we call the resulting space a “split” subspace.

complex orthogonal basis  $V = \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix}$  and then split  $V$  into  $\Re$  and  $\Im$ , such as

$$\begin{bmatrix} v_1^{\Re} & v_1^{\Im} & v_2^{\Re} & v_2^{\Im} & \cdots & v_n^{\Re} & v_n^{\Im} \end{bmatrix}. \quad (4.2)$$

But the set (4.2) is no longer orthogonal, and possibly linearly dependent. Orthogonalization of (4.2) requires up to  $(2n)^2 N$  flops.

### Ruhe's method

It seems that it would be ideal to create an orthogonal, real basis for (4.1) directly during the iterative process. Ruhe addresses this in [27], in the context of a single-vector general rational-Krylov method for eigenvalue finding. His method involves considering  $\Re$  and  $\Im$  parts of the power iterate separately, so that each iteration yields two new real vectors. Implemented as a modification of the Arnoldi algorithm, line 4 of Algorithm 1 becomes

$$a_k + ib_k = \mathbf{H}v_k,$$

and we orthogonalize vectors  $a_k$  and  $b_k$  separately. The problem with this method just described is that it is not clear what vector we should iterate with next in order to build a basis for (4.1). It is not clear whether the real basis produced by Ruhe's method spans a Krylov-subspace, let alone a basis for the split-Krylov-subspace (4.1), nor whether projection with this basis matches moments. Surprisingly there is neither much literature nor results on this topic with regards to model order-reduction, but further work in this area could be promising, as the typical split and re-orthogonalize procedure seems somewhat redundant.

#### 4.2.2 Using a real inner-product

Ruhe's method for generating a real basis yields unpredictable results, compared with working strictly in complex arithmetic followed by splitting the complex basis into  $\Re$  and  $\Im$  parts and reorthogonalizing as a post processing step. One way to cut complex-vector orthogonalization costs in half for a Gram-Schmidt based process during the the Krylov process is to use the



alternate real-valued inner-product that we denote by  $\langle \cdot, \cdot \rangle_{\mathbb{R}}$ .

For complex vectors  $v = a + ib$  and  $w = x + iy$ , define

$$\langle v, w \rangle_{\mathbb{R}} = x^T a + y^T b \in \mathbb{R}. \quad (4.3)$$

Consider “orthogonalization” using (4.3) instead of the standard complex Euclidean inner-product

$$v^H w = x^T a + y^T b + i(x^T b - y^T a). \quad (4.4)$$

The real inner-product (4.3) is cheaper to compute than (4.4). The basis  $V$  produced by a cycle of the Arnoldi process (algorithm 1) using this inner-product for orthogonalization is not real-valued, nor is it orthogonal in the Euclidean sense. We cannot use it to make orthogonal projections as in (2.3), and an  $n$ -dimensional basis  $V$  orthogonal with respect to (4.3) does not span  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ , in general. However, it satisfies (4.1); that is,

$$\text{span } V^* = \mathcal{K}_n(\mathbf{H}, \mathbf{R})^* = \text{span } \mathcal{K}_n(\mathbf{H}(\sigma), \mathbf{R}(\sigma)) \cup \mathcal{K}_n(\mathbf{H}(\bar{\sigma}), \mathbf{R}(\bar{\sigma})),$$

which works out because we must split and re-orthogonalize anyway. We call a set *equivalent-real* orthogonal if it is orthogonal with respect to (4.3), where the notion of “equivalent-realness” is explained next.

Constructing such an equivalent-real orthogonal basis for  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$  can be accomplished by replacing line 6

$$h_{jk} = v_k^H r_k$$

in algorithm 1 with

$$g_{jk} = \langle v_k, r_k \rangle_{\mathbb{R}}$$

where  $\langle \cdot, \cdot \rangle$  is defined by (4.3).

However, the matrix

$$\mathbf{G} = [g_{jk}] \neq V^H \mathbf{H} V \quad (4.5)$$

of orthogonalization coefficients is no longer an orthogonal-projection in the Euclidean sense which limits this idea's utility for model order-reduction.

### 4.3 Equivalent-real formulations

A Krylov process that orthogonalizes iterates with respect to (4.3) is effective for constructing a split-worthy basis because it is an Euclidean inner-product with respect to the “equivalent-real” Krylov-subspace  $\mathcal{K}_n(\ddot{\mathbf{H}}, \ddot{\mathbf{R}}) \subset \mathbb{R}^{2N}$  induced by the *equivalent-real*, or *realified* formulations

$$\ddot{\mathbf{H}} = \begin{bmatrix} \mathbf{H}^r & -\mathbf{H}^i \\ \mathbf{H}^i & \mathbf{H}^r \end{bmatrix} \quad \text{and} \quad \ddot{\mathbf{R}} = \begin{bmatrix} \mathbf{R}^r \\ \mathbf{R}^i \end{bmatrix}. \quad (4.1)$$

of  $\mathbf{H} = \mathbf{H}^r + i\mathbf{H}^i$  and  $\mathbf{R} = \mathbf{R}^r + i\mathbf{R}^i$  from (2.5). This idea is from [25, Sec. 5] and [5, ‘K1-formulation’]. There is a general definition and discussion of realified spaces as a pure-mathematics topic in [23].

A basis

$$\ddot{\mathbf{V}} = \begin{bmatrix} \ddot{v}_1 & \ddot{v}_2 & \cdots & \ddot{v}_n \end{bmatrix} \quad (4.2)$$

for the Krylov-subspace  $\mathcal{K}_n(\ddot{\mathbf{H}}, \ddot{\mathbf{R}})$  induced by the realified matrices (4.1) consists of vectors

$$\ddot{v} = \begin{bmatrix} \ddot{v}^t \\ \ddot{v}^b \end{bmatrix}$$

where we call  $\ddot{v}^t$  and  $\ddot{v}^b$  in  $\mathbb{R}^N$  the *top* and *bottom* parts of  $\ddot{v}$ . We define a split of the equivalent-real basis (4.2) as

$$\ddot{\mathbf{V}}_n^* := \begin{bmatrix} \ddot{v}_1^t & \ddot{v}_1^b & \ddot{v}_2^t & \ddot{v}_2^b & \cdots & \ddot{v}_n^t & \ddot{v}_n^b \end{bmatrix}, \quad (4.3)$$

which is analogous to the split (4.1) of set of complex vectors into  $\Re$  and  $\Im$ -parts.

The next result establishes that

$$\text{span } \ddot{\mathbf{V}}_n^* = \mathcal{K}_n(\mathbf{H}, \mathbf{R})^*.$$

That is, whether we construct a basis for a complex Krylov-subspace  $\mathcal{K}_\eta(\mathbf{H}, \mathbf{R})$  using complex arithmetic, or using real arithmetic with equivalent real forms  $\ddot{\mathbf{H}}$  and  $\ddot{\mathbf{R}}$ , splitting the basis yields the same spanning set.

We remind the reader that equivalent-real forms never need to be explicitly formed. They are only implied by the use of the inner-product (4.3).

#### 4.3.1 Equivalence of split complex and equivalent-real subspaces

**Lemma 1.** *Consider the equivalent-real formulations  $\ddot{\mathbf{H}}$  and  $\ddot{\mathbf{R}}$  of  $\mathbf{H}$  and  $\mathbf{R}$  as defined by (4.1). Then equivalent real formulation of  $\mathbf{H}^j \mathbf{R}$  is  $\ddot{\mathbf{H}}^j \ddot{\mathbf{R}}$  for any integer  $j = 0, 1, 2, \dots$ , i.e.*

$$(\mathbf{H}^j \mathbf{R})^* = \ddot{\mathbf{H}}^j \ddot{\mathbf{R}}. \quad (4.4)$$

Equivalently,

$$\ddot{\mathbf{H}}^j \ddot{\mathbf{R}} = \begin{bmatrix} \Re(\mathbf{H}^j \mathbf{R}) \\ \Im(\mathbf{H}^j \mathbf{R}) \end{bmatrix}. \quad (4.4^*)$$

*Proof.* Trivially for  $j = 0$  we have  $\ddot{\mathbf{R}} := \begin{bmatrix} \mathbf{R}^r & \mathbf{R}^i \end{bmatrix}^T$ . For  $j \geq 1$ , let  $K = \mathbf{H}^{j-1} \mathbf{R}$ . Then  $\widehat{K} = \begin{bmatrix} K^r & K^i \end{bmatrix}^T$  is the equivalent-real form of  $K = K^r + iK^i$ , so

$$\ddot{\mathbf{H}}^j \ddot{\mathbf{R}} = \ddot{\mathbf{H}} \widehat{K} = \begin{bmatrix} \mathbf{H}^r & -\mathbf{H}^i \\ \mathbf{H}^i & \mathbf{H}^r \end{bmatrix} \begin{bmatrix} K^r \\ K^i \end{bmatrix} = \begin{bmatrix} \mathbf{H}^r K^r - \mathbf{H}^i K^i \\ \mathbf{H}^r K^i + \mathbf{H}^i K^r \end{bmatrix}$$

is the equivalent real formulation of

$$\begin{aligned} \mathbf{H}^j \mathbf{R} &= \mathbf{H} K = (\mathbf{H}^r + i\mathbf{H}^i)(K^r + iK^i) \\ &= (\mathbf{H}^r K^r - \mathbf{H}^i K^i) + i(\mathbf{H}^r K^i + \mathbf{H}^i K^r) \end{aligned}$$

□

It follows as a corollary that the split-Krylov-subspaces (4.1) and (4.3) induced by each pair, are equal.

$$\mathcal{K}_n(\ddot{\mathbf{H}}, \ddot{\mathbf{R}})^* = \mathcal{K}_n(\mathbf{H}, \mathbf{R})^* \quad (4.5)$$

The inner-products (4.4) and (4.3) yield different notions of orthogonality of a complex basis and its equivalent-real counterpart, and ultimately incompatible spaces. The inner-product (4.3) implies a weaker orthogonality than (4.4): if two complex vectors  $v, w \in \mathbb{C}^N$  are orthogonal then it follows that their equivalent real forms  $\hat{v}, \hat{w} \in \mathbb{R}^{2N}$  are also orthogonal, but the converse is not true in general. A basis  $\ddot{V}$  of the block-Krylov-subspace  $\mathcal{K}_n(\ddot{\mathbf{H}}, \ddot{\mathbf{R}})$  cannot be identified with a basis of  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ : if we express each basis vector  $\ddot{v}_j$  as a complex vector

$$v_j = \ddot{v}_j^{\mathbf{t}} + i\ddot{v}_j^{\mathbf{b}},$$

the resulting set of complex vectors  $\{v_j\}$  will generally neither be orthogonal, nor will it span  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ . However, we are not interested in  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ , but rather its split variation  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})^*$ , which is why the result (4.5) of Lemma 4.4 is promising.

The norms implied by (4.4) and (4.3) for a complex vector  $v \in \mathbb{C}^N$  and its equivalent real form  $\ddot{v} \in \mathbb{R}^{2N}$  are equal:

$$\|\ddot{v}\|_2^2 = \ddot{v}^T \ddot{v} = v^H v = \langle v, v \rangle = \|v\|_2^2. \quad (4.6)$$

Lemma 4.4 establishes that complex and realified forms of  $\mathbf{H}$  and  $\mathbf{R}$  run for equal numbers of iterations induce the same split Krylov-subspace  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})^*$ . The next result establishes that the basis vectors produced by an iteration of the Arnoldi process advance the split-Krylov-subspace (4.1) in the same order, regardless of whether we use complex or equivalent-real formulation (i.e. complex formulation with inner-product (4.3)).

We show this for the simpler case that  $\mathbf{R} = \mathbf{r} \in \mathbb{C}^N$  is a single vector. The result of Theorem 3 can be extended to the more general *Band-Krylov* process which is applicable to MIMO model reduction.

**Theorem 3.** *Consider  $\mathbf{H}, \mathbf{r}$  from (2.5) and their realified formulations  $\ddot{\mathbf{H}}$  and  $\ddot{\mathbf{r}}$  defined by (4.1).*

Let  $V = \begin{bmatrix} v_1 & v_2 & \cdots & v_n \end{bmatrix}$  be the orthonormal basis implied by  $n$  Arnoldi iterations of  $\mathbf{H}$  with start vector  $\mathbf{r}$ , and let  $\ddot{V} = \begin{bmatrix} \ddot{v}_1 & \ddot{v}_2 & \cdots & \ddot{v}_n \end{bmatrix}$ , with  $\ddot{v}_j = \begin{bmatrix} \ddot{v}_j^{\mathbf{t}} & \ddot{v}_j^{\mathbf{b}} \end{bmatrix}^T$ , be the analogous vectors produced by Arnoldi iterations using  $\ddot{\mathbf{H}}$  and  $\ddot{\mathbf{r}}$ . Then there exist scalars  $\alpha, \beta \in \mathbb{R}$  and real vectors  $w, z \in \mathcal{K}_n(\mathbf{H}, \mathbf{r})^*$  such that

$$\Re(v_n) = \alpha \ddot{v}_n^{\mathbf{t}} + w \quad \text{and} \quad \Im(v_n) = \beta \ddot{v}_n^{\mathbf{b}} + z. \quad (4.7)$$

*Proof.* We will prove (4.7) by induction. For  $n = 1$  we have  $v_1 = \mathbf{r} / \|\mathbf{r}\|_2$ ,  $\ddot{v}_1 = \ddot{\mathbf{r}} / \|\ddot{\mathbf{r}}\|_2$ , where  $\|\ddot{\mathbf{r}}\|_2 = \|\mathbf{r}\|_2$  by (4.6), so  $\Re(v_1) = \ddot{v}_1^{\mathbf{t}}$  and  $\Im(v_1) = \ddot{v}_1^{\mathbf{b}}$ , trivially satisfying (4.7).

Now assume we have performed  $n \geq 1$  steps of the standard Arnoldi process to obtain an orthonormal basis  $V$  for  $\mathcal{K}_n(\mathbf{H}, \mathbf{r})$ , and assume a complex span for its split space  $\mathcal{K}_n(\mathbf{H}, \mathbf{r})^*$  (of dimension  $n$ ), so that

$$\mathcal{K}_n(\mathbf{H}, \mathbf{r}) \subseteq \mathcal{K}_n(\mathbf{H}, \mathbf{r})^* = \text{span} \begin{bmatrix} \tilde{v}_1 & \tilde{v}_2 & \cdots & \tilde{v}_\eta \end{bmatrix} \quad (4.8)$$

for real basis vectors  $\tilde{v}_j \in \mathbb{R}^N$ . On the  $n \geq 1$ -th step, the Arnoldi process with  $\mathbf{H}$  and  $\mathbf{r}$  computes scalar orthogonalization coefficients  $\{h_{jn}\}_{j=1}^n \subset \mathbb{C}$  and  $h_{n+1,n} \in \mathbb{R}$  such that

$$\begin{aligned} h_{n+1,n} v_{n+1} &= \mathbf{H} v_n - \sum_{j=1}^n h_{jn} v_j \\ &= \mathbf{H}^n \mathbf{r} + \sum_{j=1}^n c_j v_j, \quad c_j \in \mathbb{R} \\ &= \mathbf{H}^n \mathbf{r} + \sum_{j=1}^{\eta} d_j \tilde{v}_j, \quad d_j \in \mathbb{C}, \quad \text{by (4.8)} \\ &= \mathbf{H}^n \mathbf{r} + w_1 + i z_1, \quad w_1, z_1 \in \mathcal{K}_n(\mathbf{H}, \mathbf{r})^* \cap \mathbb{R}^N. \end{aligned} \quad (4.9)$$

Lemma 1 implies that we can re-write (4.9) in realified form as

$$h_{n+1,n} \begin{bmatrix} \Re(v_{n+1}) \\ \Im(v_{n+1}) \end{bmatrix} = \ddot{\mathbf{H}}^n \ddot{\mathbf{r}} + \begin{bmatrix} w_1 \\ z_1 \end{bmatrix}. \quad (4.10)$$

On the other hand, after  $n$  iterations Arnoldi iterations with  $\ddot{\mathbf{H}}$  and  $\ddot{\mathbf{r}}$  we have

$$\hat{h}_{n+1,n} \ddot{v}_{n+1} = \ddot{\mathbf{H}}^n \ddot{\mathbf{r}} - \sum_{j=1}^n \hat{h}_{jn} \ddot{v}_j,$$

so

$$\begin{aligned} \hat{h}_{n+1,n} \begin{bmatrix} \ddot{v}_{n+1}^{\mathbf{t}} \\ \ddot{v}_{n+1}^{\mathbf{b}} \end{bmatrix} &= \ddot{\mathbf{H}}^n \ddot{\mathbf{r}} + \sum_{j=1}^n \hat{c}_j \begin{bmatrix} \ddot{v}_j^{\mathbf{t}} \\ \ddot{v}_j^{\mathbf{b}} \end{bmatrix}, \quad \hat{c}_j \in \mathbb{R} \\ &= \ddot{\mathbf{H}}^n \ddot{\mathbf{r}} + \sum_{j=1}^{\eta} \begin{bmatrix} \hat{a}_j \tilde{v}_j \\ \hat{b}_j \tilde{v}_j \end{bmatrix}, \quad \hat{a}_j, \hat{b}_j \in \mathbb{R} \\ &= \ddot{\mathbf{H}}^n \ddot{\mathbf{r}} + \begin{bmatrix} w_2 \\ z_2 \end{bmatrix}, \end{aligned}$$

where  $w_2, z_2 \in \mathcal{K}_n(\mathbf{H}, \mathbf{r})^* \cap \mathbb{R}^N$ . □

Theorem 3 establishes that Arnoldi vectors generated using  $\ddot{\mathbf{H}}$  and  $\ddot{\mathbf{r}}$  yield basis vectors for  $\mathcal{K}_n(\mathbf{H}, \mathbf{r})^*$  in the same order as those obtained from  $\mathbf{H}$  and  $\mathbf{r}$ ; in fact, up to finite precision error they yield exactly the same basis.

### 4.3.2 Reduced-order models via equivalent-real formulations

#### via explicit projection

The explicitly projected ROM (2.2) using a basis (4.1) for the split-Krylov-subspace  $\mathcal{K}_\eta(\mathbf{H}, \mathbf{R})^*$  is equivalent regardless of the orthogonalization method used to construct it.

### via implicit projection

The matrix  $\mathbf{G}$  of orthogonalization coefficients from the equivalent-real Arnoldi process, (4.5), is a Rayleigh-quotient approximant to the equivalent-real operator

$$\ddot{\mathbf{H}} = \begin{bmatrix} \mathbf{H}^r & -\mathbf{H}^i \\ \mathbf{H}^i & \mathbf{H}^r \end{bmatrix}$$

of (4.1), and not the original complex-shifted operator  $\mathbf{H}$ . It is the projection

$$\mathbf{G} = \ddot{V}^T \ddot{\mathbf{H}} \ddot{V}$$

that would be formed by orthogonal projection using  $\ddot{V} \in \mathbb{R}^{2N \times \eta}$ , where  $\text{span } \ddot{V} = \mathcal{K}_n(\ddot{\mathbf{H}}, \ddot{\mathbf{R}})$  for the implied Krylov-subspace induced by equivalent-real forms (4.1). Thus, an implicitly projected ROM (2.2) is not so simple to characterize. For example, it is known that the spectrum

$$\Lambda(\ddot{\mathbf{H}}) = \Lambda(\mathbf{H}) \cup \Lambda(\overline{\mathbf{H}})$$

of  $\ddot{\mathbf{H}}$  contains spectral information for the complex-conjugate  $\overline{\mathbf{H}}$  which complicates matters because we are interested only in the spectrum of  $\mathbf{H}$ . This makes analyzing a ROM transfer-function via implicit-projection onto a realified Krylov-subspace non-trivial, but it might be a promising improvement if developed further.

## Chapter 5

# The band-Arnoldi process and a proposed thick-restarted variant

“Thick”-starting a Krylov-process that iterates with  $\mathbf{H}$  starting on  $\mathbf{R}$  means that we instead start on  $\begin{bmatrix} Z & \mathbf{R} \end{bmatrix}$  where  $Z$  is a basis of known Ritz-vectors of  $\mathbf{H}$ . The notion of a re-start for multi-point MOR is that once we have iterated enough with  $\mathbf{H}_1 = \mathbf{H}(\sigma_1)$  on  $\mathbf{R}_1 = \mathbf{R}(\sigma_1)$ , creating an approximation about  $\sigma_1$ , we can start over, iterating about  $\sigma_2$  with  $\mathbf{H}_2$  and  $\mathbf{R}_2$ , avoiding wasting computation on rediscovering invariant-subspace. Recall that all  $\mathbf{H}(\sigma)$  (over  $\sigma$  at which  $\mathbf{H}(\sigma)$  is defined) share invariant-subspace. If we determine a convergent-enough invariant-subspace  $Y_1$  of  $\mathbf{H}_1$ , then we can thick-restart the process with  $\mathbf{H}_2$  acting on  $\begin{bmatrix} Y_1 & \mathbf{R}_2 \end{bmatrix}$ .

The thesis of this document and primary contribution<sup>1</sup> to the field of model order-reduction is the proposal and testing of a thick-restarted Arnoldi-type algorithm for multi-point MOR, where the model to be reduced is assumed to be MIMO. We show that we can obtain smaller models of the same or better accuracy, at less computational cost, using this process.

### 5.1 Band-Arnoldi algorithm

First we address the band-Arnoldi algorithm and the band-Arnoldi relation, which are the basic computational and analytical engine for the MOR method. The *band*-Arnoldi process of [10] (2003)

---

<sup>1</sup>as far as we know



is included here as algorithm 2. An older (possibly the original) version of a band-Krylov process was given by [28].

Band-Arnoldi differs from a block-Krylov process (like block-Arnoldi [21]) in that it cycles through a “band” of candidate-vectors  $\begin{bmatrix} \hat{v}_n & \hat{v}_{n+1} & \dots & \hat{v}_{n+m_c} \end{bmatrix}$ , where  $m_c(n)$  is the current band size on the  $n$ -th iteration of the main loop. The initial band is the start block

$$\begin{bmatrix} \hat{v}_1 & \hat{v}_2 & \dots & \hat{v}_m \end{bmatrix} = \begin{bmatrix} \mathbf{r}_1 & \mathbf{r}_2 & \dots & \mathbf{r}_m \end{bmatrix} = \mathbf{R}.$$

On the  $j$ -th iteration of the algorithm, the candidate-vector  $\hat{v}_j$  either gets deflated or becomes Krylov basis-vector  $v_j$ , which is then advanced via Arnoldi iteration to be candidate-vector  $\hat{v}_{j+m_c}$ . If we deflate  $\hat{v}_j$  then the band size  $m_c$  is decremented. Since the algorithm proceeds as a continuous cycle rather than a block iteration, at any step  $n$  it is simpler to refer to the computed basis  $V \in \mathbb{C}^{N \times n}$  for  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ , where  $n$  is the dimension of the basis, rather than a block-degree.

The band-Arnoldi algorithm run for  $n$ -iterations with operator  $\mathbf{H}$  and start-block  $\mathbf{R}$  returns a basis  $V \in \mathbb{C}^{N \times n}$  for  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ , deflated vectors  $\dot{V} = \begin{bmatrix} \dot{v}_1 & \dot{v}_2 & \dots & \dot{v}_d \end{bmatrix}$ , remaining candidate-vectors  $\hat{V} = \begin{bmatrix} \hat{v}_{n+1} & \hat{v}_{n+2} & \dots & \hat{v}_{n+m_c} \end{bmatrix}$ , and Rayleigh-quotient  $\tilde{\mathbf{H}} = V^H \mathbf{H} V$  that satisfy the band-Arnoldi relation

$$\mathbf{H}V = V\tilde{\mathbf{H}} + \begin{bmatrix} (I - VV^H)\dot{V} & \hat{V} \end{bmatrix} F \quad (5.1)$$

where  $F = \begin{bmatrix} F_1 \\ F_2 \end{bmatrix}$ .  $\tilde{\mathbf{H}}$  is block-upper-Hessenberg (strictly upper-Hessenberg in the single vector setting  $m = 1$ ) with possibly non-zero columns in the typically zero region<sup>2</sup> corresponding to the deflated vectors  $\dot{V}$ .  $F_1$  and  $F_2$  are indexing matrices that position vectors  $\dot{v}_j$  and  $\hat{v}_j$  respectively into the  $n$  available positions of the  $N \times n$  block (5.1).

In addition, algorithm 2 returns  $\tilde{\rho}$ , and  $\tilde{\rho}^{\text{defl}}$  where

$$\mathbf{R} = V\tilde{\rho} + \tilde{\rho}^{\text{defl}}, \quad (5.2)$$

---

<sup>2</sup>lower co-Hessenberg?

and  $V^H \mathbf{R} = \tilde{\boldsymbol{\rho}}$ .

### 5.1.1 Candidates/residual term

$$\begin{aligned} \hat{V}F_2 &= \begin{bmatrix} 0 & 0 & \dots & 0 & \hat{V} \end{bmatrix} \in \mathbb{C}^{N \times n} \\ &= \begin{bmatrix} \hat{v}_{n+1} & \hat{v}_{n+2} & \dots & \hat{v}_{n+m_c} \end{bmatrix} \begin{bmatrix} \dots & 1 & & \\ \dots & & 1 & \\ \dots & & & \ddots \\ \dots & & & & 1 \end{bmatrix} \end{aligned}$$

is the residual term involving the band  $\hat{V}$  of candidate-vectors after the  $n$ -th iteration of the main loop. The matrix  $F_2 \in \{0, 1\}^{m_c \times n} = \begin{bmatrix} 0 & 0 & \dots & 0 & I \end{bmatrix}$  simplifies to  $e_n^T$  for the single vector iteration.  $F_2$  places  $\hat{V}$  in the last  $m_c$  of  $n$  positions. Note that  $V^H \hat{V} = 0$ .

### 5.1.2 Deflation term

$\dot{V}F_1 \in \mathbb{C}^{N \times n}$  is the zero or mostly-zero matrix  $\hat{V}^{\text{defl}}$  implied by algorithm 2 (band-Arnoldi). If no deflation or only exact deflation occurred then  $\dot{V} = 0$  and  $\dot{V}F_1$  is an  $N \times n$  matrix of zeros. If inexact deflation was performed on the  $j$ -th candidate-vector then  $j \in \mathcal{I}$  and  $\hat{v}_j^{\text{defl}} \neq 0$ . Negative or zero indices in  $\mathcal{I}$  correspond to deflations that happened within the start block  $\mathbf{R}$ . For example, if  $j - m \leq 0$  then  $\tilde{\boldsymbol{\rho}}_j^{\text{defl}} = \hat{v}_{j-m}^{\text{defl}}$ . We may similarly define  $F_0$  so that  $\tilde{\boldsymbol{\rho}}^{\text{defl}} = \dot{V}F_0$ .

As an example of a deflation matrix  $\hat{V}^{\text{defl}}$ , suppose  $d = 2$  vectors  $\dot{v}_1 = \hat{v}_2^{\text{defl}}$  and  $\dot{v}_2 = \hat{v}_5^{\text{defl}}$  were deflated at iterations  $m_c + 2$  and  $m_c + 5$  of a band-Arnoldi process of  $n = m_c + 10$  iterations, with band-size  $m_c$ . Then for standard basis vectors  $e_2, e_5 \in \mathbb{R}^{10}$ ,

$$\begin{aligned} \dot{V}F_1 = \hat{V}^{\text{defl}} &= \begin{bmatrix} 0 & \hat{v}_2^{\text{defl}} & 0 & 0 & \hat{v}_5^{\text{defl}} & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} 0 & \dot{v}_1 & 0 & 0 & \dot{v}_2 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} \dot{v}_1 & \dot{v}_2 \end{bmatrix} \begin{bmatrix} e_2^T \\ e_5^T \end{bmatrix}. \end{aligned} \tag{5.3}$$

The band-Arnoldi algorithm deflates a candidate-vector  $\hat{v}_j$  (i.e. removes it from further iterations) if  $\|\hat{v}_j\| \leq \text{dtol}$ <sup>3</sup> after orthogonalizing  $\hat{v}_j$  against  $V = \{v_1, v_2, \dots, v_j\}$ , which means it is almost linearly dependent with previous basis vectors. Algorithm 2 then sets  $\hat{v}_j^{\text{defl}} := \hat{v}_j$  and removes it as a candidate, and the current band size  $m_c$  is decreased by one.  $\hat{v}_j^{\text{defl}}$  is no longer used for iterations and basis vectors  $v_{j+1}, v_{j+2}, \dots, v_{n+m_c}$  are not made orthogonal to  $\hat{v}_j^{\text{defl}}$ . Then

$$V^H \hat{v}_j^{\text{defl}} = \begin{bmatrix} 0 & 0 & \dots & 0 & v_{j+1}^H \hat{v}_j^{\text{defl}} & v_{j+2}^H \hat{v}_j^{\text{defl}} & \dots & \hat{v}_j^H \hat{v}_j^{\text{defl}} \end{bmatrix}^T. \quad (5.4)$$

(5.4) implies that  $\|V^H \dot{v}\| \leq \|\dot{v}\| \leq \text{dtol}$ .

If no/exact deflation was performed,  $\tilde{\mathbf{H}}$  is strictly block-upper-Hessenberg, otherwise  $\tilde{\mathbf{H}}$  may have non-zero entries in the triangular region,  $\tilde{\mathbf{H}}_{\mathcal{E}} = V^H \dot{V}$ , below the 1st subdiagonal of  $\tilde{\mathbf{H}}$ . If an inexact deflation occurred on the  $j$ -th iteration, (5.4) is included in the Rayleigh-quotient  $\tilde{\mathbf{H}}$  as the  $j$ -th column of  $\tilde{\mathbf{H}}_{\mathcal{E}}$ . Then

$$\|\tilde{\mathbf{H}}_{\mathcal{E}}\| = \|V^H \dot{V}\|_F \leq \|\dot{V}\|_F \leq \text{dtol} \sqrt{d}, \quad (5.5)$$

and

$$\|(I - VV^H)\dot{V}\|_F \leq \|\dot{V}\|_F. \quad (5.6)$$

$\tilde{\boldsymbol{\rho}}^{\text{defl}}$  in (5.2) is also an all or mostly-zero matrix of very small norm, representing deflations that occurred during the first  $m$ -iterations, i.e. linear dependence within the start block  $\mathbf{R}$ .

---

<sup>3</sup>[26] suggests  $\text{dtol} = \sqrt{\epsilon}$ , where  $\text{machine-}\epsilon = 2^{-52} \approx 2.22\text{e-}16$  in double-precision (64-bit) floating point.

### 5.1.3 Residual norms

A similarity decomposition  $\tilde{\mathbf{H}}S = SU$  such as a Schur or eigenvalue-decomposition of the Rayleigh-quotient together with (5.1) and setting  $Y = VS$  gives the block residual-norm bound

$$\begin{aligned}\|\mathbf{H}Y - YU\|_F^2 &= \left\| \left( (I - VV^H)\dot{V}F_1 + \hat{V}F_2 \right) S \right\|_F^2 \\ &\leq \|\dot{V}\|_F^2 + \|\hat{V}\|_F^2 \\ &\leq (\text{dtol})^2 d + \|\hat{V}\|_F^2 \\ &= d\varepsilon_M + \|\hat{V}\|_F^2, \quad \text{for } \text{dtol} = \sqrt{\varepsilon_M}\end{aligned}$$

in the Frobenius-norm (entry-wise 2-norm).

Given  $W = \begin{bmatrix} w_1 & w_2 & \cdots & w_n \end{bmatrix}$  for the eigenvalue decomposition  $\tilde{\mathbf{H}}W = W\Lambda$  and Ritz-basis  $Z = VW$ , the residual for a Ritz-pair  $(\lambda_j, z_j)$ , is

$$\begin{aligned}\mathbf{H}z_j - z_j\lambda_j &= \begin{bmatrix} (I - VV^H)\dot{V} & \hat{V} \end{bmatrix} Fw_j \\ &= (I - VV^H)\hat{V}^{\text{defl}}w_j + \hat{V}F_2w_j\end{aligned}$$

For determining convergence of Ritz-pairs,

$$\begin{aligned}\|\mathbf{H}z_j - z_j\lambda_j\|_2^2 &\leq d\varepsilon_M + \|\hat{V}F_2w_j\|_2^2 \\ &= d\varepsilon_M + \|\hat{V}\tilde{w}_j\|_2^2.\end{aligned}\tag{5.7}$$

where  $\tilde{w}_j \in \mathbb{C}^{m_c}$  is the last  $m_c$  elements (rows) of  $w_j$ .

(5.7) suggests a few different ways to cheaply estimate the relative residual norm

$$\frac{\|\mathbf{H}z_j - z_j\lambda_j\|}{\|\lambda_j z_j\|}\tag{5.8}$$

for a Ritz-pair  $(\lambda_j, z_j)$ . We assume  $\|z_j\| = 1$ , so that  $\|\lambda_j z_j\| = |\lambda_j|$ .

Some methods estimate the relative residual-norm as  $\|\mathbf{H}z_j - z_j\lambda_j\|/\|\mathbf{H}z_j\|$  with an estimate of  $\|\mathbf{H}\|$  or with  $\|\tilde{\mathbf{H}}\|$ . We use  $|\lambda_j|$  because it is uncertain whether  $\|\mathbf{H}\|$  or  $\|\tilde{\mathbf{H}}\|$  are good estimates of

$\|\mathbf{H}\|$ .

Assuming  $d\varepsilon_M$  is negligible,

$$\left\| \hat{V} \right\| \frac{\|\tilde{w}_j\|}{|\lambda_j|} \quad (5.9)$$

is an estimate of relative residual norm, as is

$$\frac{1}{|\lambda|} \begin{bmatrix} \|\hat{v}_1\| & \|\hat{v}_2\| & \cdots & \|\hat{v}_{m_c}\| \end{bmatrix} \begin{bmatrix} |\tilde{w}_j^{(1)}| \\ |\tilde{w}_j^{(2)}| \\ \vdots \\ |\tilde{w}_j^{(m_c)}| \end{bmatrix}. \quad (5.10)$$

Both (5.9) and (5.10) are cheaper to compute than norms  $\|\hat{V}\tilde{w}_j\|$  of potentially large matrix-vector products, but (5.10) might be better if  $\hat{V}F_2$  has rank greater than one. Estimates (5.9) and (5.10) are equal for rank-1 residuals.

---

**Algorithm 2:** BAND-ARNOLDI

---

**Input:**  $\mathbf{H}$  and start-block  $\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \cdots \ \mathbf{r}_m]$ ,  
**Output:** basis  $V$  for  $\mathcal{K}_n(\mathbf{H}, \mathbf{R})$ , deflated vectors  $\hat{V}$ , candidate-vectors  $\hat{V}$ ,  $\tilde{\mathbf{H}}$ , and  $\tilde{\boldsymbol{\rho}}$  that satisfy (5.1)

```
1  $\hat{v}_i := \mathbf{r}_i$  for  $i = 1, 2, \dots, m$ 
2  $m_c := m$ 
3  $\mathcal{I} := \emptyset$ 
4 for  $n = 1$  to  $n_{max}$  do
5   while  $\|\hat{v}_n\|_2 < dtol \cdot \|\mathbf{H}\|_{est}$  do      % remove  $\hat{v}_n$  if necessary (deflation)
6      $\hat{v}_{n-m_c}^{\text{defl}} := \hat{v}_n$ 
7      $\mathcal{I} = \mathcal{I} \cup \{n - m_c\}$  % locations in  $\hat{V}^{\text{defl}}$  (or  $\tilde{\boldsymbol{\rho}}^{\text{defl}}$ ) that contain deflated
      vectors
8      $m_c := m_c - 1$  % we assume no early termination
9      $\hat{v}_j := \hat{v}_{j+1}$  for  $j = n, n+1, \dots, n+m_c-1$ 
10     $h_{n,n-m_c} := \|\hat{v}_n\|_2$ 
11     $v_n := \hat{v}_n / \|\hat{v}_n\|_2$ 
12    for  $j = n+1$  to  $n+m_c-1$  do      % Make candidates  $\{\hat{v}_1, \hat{v}_2, \dots, \hat{v}_n\}$  orthogonal
      to  $v_n$ 
13       $h_{n,j-m_c} := v_n^H \hat{v}_j$ 
14       $\hat{v}_j := \hat{v}_j - h_{n,j-m_c} v_n$ 
15     $\hat{v}_{m_c+n} := \mathbf{H} v_n$ 
16    for  $j = 1$  to  $n$  do      % Make  $\hat{v}_{m_c+n}$  orthogonal to previous  $\{v_1, v_2, \dots, v_n\}$ 
17       $h_{jn} := v_j^H \hat{v}_{m_c+n}$ 
18       $\hat{v}_{m_c+n} := \hat{v}_{m_c+n} - h_{jn} v_j$ 
19    for  $j \in \mathcal{I}$  do
20       $h_{nj} := v_n^H \hat{v}_j^{\text{defl}}$ 
21 return  $V, \tilde{\mathbf{H}}, \tilde{\boldsymbol{\rho}} = [h_{ij}]_{i=1,2,\dots,m}^{j=1-m,\dots,1,0}$ ,  $\hat{V}, \hat{V}^{\text{defl}}, \tilde{\boldsymbol{\rho}}^{\text{defl}} = [\hat{v}_j^{\text{defl}}], j = 1-m, \dots, 1, 0$ ,
```

---

## 5.2 Thick-restarting the Band-Arnoldi process

Suppose that via (5.9) or (5.10) we have identified  $\ell$  Ritz-pairs  $(\lambda_j, z_j)$  with residuals  $\gamma_j$  so that

$$\mathbf{H}z_j - \lambda z_j = \gamma_j$$

for  $j = 1, 2, \dots, \ell$ . If we run the band-Arnoldi algorithm with  $\mathbf{H}$  and start block  $\begin{bmatrix} Z & \mathbf{R} \end{bmatrix}$  where  $Z = \begin{bmatrix} z_1 & z_2 & \cdots & z_\ell \end{bmatrix}$  and  $\mathbf{R} = \begin{bmatrix} \mathbf{r}_1 & \mathbf{r}_2 & \cdots & \mathbf{r}_m \end{bmatrix}$ , then after  $\ell$  iterations, the constructed orthonormal

basis is  $V_\ell = \begin{bmatrix} v_1 & v_2 & \dots & v_\ell \end{bmatrix}$  for  $\text{span } Z$  where

$$u_{jj}v_j = (I - V_{j-1}V_{j-1}^H)z_j$$

for  $j = 1, 2, \dots, \ell$ .  $u_{jj} = \|z_j\|_2 \in \mathbb{R}$  is the diagonal element of the upper-triangular Rayleigh-quotient  $U$ , where  $Z = VU$  can be regarded as a QR factorization. Assuming unit Ritz-vectors, the next  $\ell$  candidate-vectors are

$$\begin{aligned} \hat{v}_{\ell+m+j} &= (I - V_jV_j^H)\mathbf{H}v_j \\ &= (I - V_jV_j^H)\mathbf{H}(I - V_{j-1}V_{j-1}^H)z_j \\ &= (I - V_jV_j^H)\mathbf{H}z_j, \quad \text{because } \mathbf{H}V_{j-1} \subset \text{span } V_j \\ &= (I - V_jV_j^H)(\lambda_j z_j + \gamma_j) \\ &= (I - V_jV_j^H)\gamma_j \quad \text{because } z_j \in \text{span } V_j. \end{aligned} \tag{5.1}$$

for  $j = 1, 2, \dots, \ell$ . Then the candidate-vectors at that point are

$$\hat{V}_\ell = \begin{bmatrix} (I - V_\ell V_\ell^H)\mathbf{R} & \hat{v}_{\ell+m+1} & \hat{v}_{\ell+m+2} & \dots & \hat{v}_{2\ell+m} \end{bmatrix}. \tag{5.2}$$

Note that  $(I - V_jV_j^H)\gamma_j = \gamma_j$  if the residual  $\gamma_j$  is orthogonal to  $\{z_1, z_2, \dots, z_j\}$ , which is the case when the basis  $Z$  of Ritz-vectors came from one cycle of a Krylov process such as algorithm 2. If we consider restarts however, the Ritz-vectors from each restart will have a different orthogonal residual, so the set of Ritz-vectors will not be orthogonal to the set of residuals in general.

### 5.2.1 Implicit restart

**Disclaimer** In this subsection we address the theory of an implicit restart for the band-Arnoldi process, but the reader should note that the method that we developed and tested for this document (in §5.3) does not use implicit restarts. For that method we explicitly augmented the start block with a growing orthonormal basis of kept Ritz-vectors and endured the cost of the extra  $\ell_j$  iterations on each cycle. This is because interpolation-point translation (see §4.1.2) complicates matters a bit

and implicit-restarting was not necessary for the relatively small test-models that we used. The computational cost saved is almost negligible unless the system is very large or we are restarting with a lot of vectors. Thus this section serves mostly as a reference for those who wish to develop such an implicit restart scheme themselves, and can be considered optional reading otherwise.

**Implicit restart** It would be wasteful to actually compute  $HV_\ell = V_\ell U + \hat{V}$  (re-multiplying  $\mathbf{H}$  by Ritz-vectors or known invariant-subspace) via algorithm 2 since we already know (5.2). The first  $\ell$  iterations of the restarted process can be carried out implicitly as with the implicit restarts used for single-vector methods. In that case we let  $Y := V_\ell$  be the orthonormal basis for the collection of Ritz-vectors  $Z$ , and let  $\hat{Y} := \hat{V}$  be the set of candidate-vectors defined by (5.1) and (5.2), thus pre-loading the basis and avoiding  $\ell$  redundant matrix-vector products.

#### With a general nearly-invariant subspace

To be more general let us assume we have a general band-Krylov-Schur relation (not orthogonal in general), with basis  $Y = \begin{bmatrix} y_1 & y_2 & \cdots & y_\ell \end{bmatrix} \in \mathbb{C}^{N \times \ell}$  for subspace  $\mathcal{Y}$ , and with basis  $\hat{Y} = \begin{bmatrix} \hat{y}_1 & \hat{y}_2 & \cdots & \hat{y}_\nu \end{bmatrix}$  (with  $\nu \leq \ell$ ) for the residual subspace so that

$$\begin{aligned} \mathbf{H}Y &= \begin{bmatrix} Y & \hat{Y} \end{bmatrix} \begin{bmatrix} U \\ B \end{bmatrix} \\ &= YU + \hat{Y}B. \end{aligned} \tag{5.3}$$

Then

$$\mathbf{H}Y = YU' + (I - YY^H)\hat{Y}B \tag{5.4}$$

with  $U' = U + Y^H\hat{Y}$  is an orthogonal relation. If  $U$  is upper triangular then  $U'$  is no longer upper triangular, but  $Y$  is usually chosen so that the norm  $\|\hat{Y}\|$  is small.

If we run band-Arnoldi with  $\mathbf{H}$  on start block  $\begin{bmatrix} Y & \mathbf{R} \end{bmatrix}$  instead of just  $\mathbf{R}$ , then after the first  $\ell$  iterations we have Arnoldi basis  $V = Y$  and candidate-vectors  $\hat{V}_\ell$  (from (5.2)) such that the span



of the first  $\ell$  of them,

$$\text{span} \left\{ \hat{v}_{\ell+m+1}, \hat{v}_{\ell+m+2}, \dots, \hat{v}_{2\ell+m} \right\}, \quad (5.5)$$

form a basis for  $\hat{Y}$ . If  $\|\hat{Y}\| \leq \text{dtol}\sqrt{\ell}$  then candidates resulting from processing the first  $\ell$  start vectors are negligibly small and will get deflated by the standard band-Arnoldi process in the next iteration(s).

If  $Y$  is not deflatable ( $\|\hat{Y}\| > \text{dtol}\sqrt{\ell}$ ) we can manually set all vectors of (5.5) to zero and continue band-Arnoldi iterations from step  $\ell + 1$ .

Assuming we let the process proceed without intervention, the resulting relation after  $n$  iterations for a total of  $n + \ell$ , is

$$\mathbf{H} \begin{bmatrix} Y & V \end{bmatrix} = \begin{bmatrix} Y & V \end{bmatrix} \begin{bmatrix} U' & G \\ \mathcal{E} & \hat{H} \end{bmatrix} + \begin{bmatrix} (I - YY^H)\hat{Y} \\ (I - VV^H)\hat{Y} & (I - VV^H)\dot{V} & \hat{V} \end{bmatrix} \begin{bmatrix} B \\ F_1 \\ F_2 \end{bmatrix} \quad (5.6)$$

where  $\dot{V}$ ,  $\hat{V}$ , and  $F$  are defined as in (5.1). The  $(\ell + n) \times (\ell + n)$  Rayleigh-quotient

$$\tilde{\mathbf{H}} := \begin{bmatrix} U' & G \\ \mathcal{E} & \hat{H} \end{bmatrix} \quad (5.7)$$

includes the block  $\mathcal{E} = V^H \hat{Y} B$ .

If we had manually set the first  $\ell$  candidates to zero then (5.6) becomes the much simpler

$$\mathbf{H} \begin{bmatrix} Y & V \end{bmatrix} \approx \begin{bmatrix} Y & V \end{bmatrix} \begin{bmatrix} U' & G \\ 0 & \hat{H} \end{bmatrix} + \begin{bmatrix} 0 & (I - VV^H)\dot{V}F_1 & \hat{V}F_2 \end{bmatrix} \quad (5.6^*)$$

In the thick-restart [31] and Krylov-Schur [33] schemes for single-vector iterations, they assume  $\|\mathcal{E}\| \leq \|\hat{Y}B\|$  is small enough to be negligible, so that  $U' \approx U$  is upper-triangular. Also, they restart with the residual from a previous cycle, which in our case would be  $\hat{Y}$ . In that case  $\hat{Y}$  would be included in  $V$  and  $\mathcal{E}$  would be zero. Those methods then take advantage of the upper-triangular structure of the leading  $\ell \times \ell$  principle submatrix  $U$  of (5.7). For this analysis we assume  $\hat{Y}$  gets

deflated so that  $V$  and  $\hat{Y}$  are not orthogonal.

Taking a full or partial similarity decomposition  $\tilde{\mathbf{H}}W = W\Lambda$ , the associated block  $Z = \begin{bmatrix} Y & V \end{bmatrix} W$  has residual

$$\mathbf{H}Z - Z\Lambda = \begin{bmatrix} (I - YY^H)\hat{Y} \\ (I - VV^H)\hat{Y} & (I - VV^H)\dot{V} & \hat{V} \end{bmatrix} \begin{bmatrix} B \\ F_1 \\ F_2 \end{bmatrix} W.$$

Let  $\hat{y}^T$ ,  $\dot{v}^T$ , and  $\hat{v}^T$  be the row vectors of norms of the columns of  $\hat{Y}$ ,  $\dot{V}$ , and  $\hat{V}$ , respectively. For example,

$$\hat{y}^T = \begin{bmatrix} \|\hat{y}_1\| & \|\hat{y}_2\| & \cdots & \|\hat{y}_\nu\| \end{bmatrix}.$$

Let

$$f^T := \begin{bmatrix} \hat{y}^T & \dot{v}^T & \hat{v}^T \end{bmatrix} \begin{bmatrix} B \\ F_1 \\ F_2 \end{bmatrix}.$$

For many applications  $\|\dot{V}\|_F \leq \text{dtol}\sqrt{d}$  is negligible, in which case  $f^T \approx \begin{bmatrix} \hat{y}^T B & \hat{v}^T F_2 \end{bmatrix}$ . The residual-norm  $\|\hat{Y}\|$  of  $Y$  is small in the sense that it represents a nearly-invariant-subspace to some degree, but it is not negligible in general.

For an individual Ritz-pair  $z_j = \begin{bmatrix} Y & V \end{bmatrix} w_j \in Z_2$  and  $\lambda_j \in \Lambda$ , a bound for residual-norm is

$$\|\mathbf{H}z_j - \lambda_j z_j\| \leq |f^T w_j|. \quad (5.8)$$

Note that

$$|f^T w_j| \geq \|\hat{Y}\|,$$

so our residual-norm estimate can never be better than  $\|\hat{Y}\|$ . We consider Ritz-vector  $z_j$  to be converged if the relative residual-norm

$$\text{rr}_j = \frac{|f^T w_j|}{\|\mathbf{H}\|_{\text{est}}} \leq \text{ctol}. \quad (5.9)$$

$\|\mathbf{H}\|_{\text{est}} = \max_v \|\mathbf{H}v\|$  is an estimated operator-norm of  $\mathbf{H}$  obtained during iterations. A value suggested by [26] is  $\text{ctol} = \sqrt{\epsilon}$ , which is the same value used for  $\text{dtol}$  in [10].

The bound used by ARPACK suggests that  $(\lambda_j, z_j)$  is converged if

$$|f^T w_j| \leq \max\{\epsilon_M \|\tilde{\mathbf{H}}\|, \text{ctol} \cdot |\lambda|\}.$$

### 5.2.2 Intermediate ROM from a thick-restarted band-Arnoldi process

The following describes how we can form an intermediate ROM via implicit-projection cheaply in order to observe convergence of the model.

Recall that the band-Arnoldi algorithm

$$\begin{bmatrix} V & \tilde{\mathbf{H}} & \tilde{\boldsymbol{\rho}} \end{bmatrix} = \mathbf{bArnoldi}(\mathbf{H}, \mathbf{R}) \quad (5.10)$$

where  $\tilde{\mathbf{H}} = V^H \mathbf{H} V$  and  $\tilde{\boldsymbol{\rho}} = V^H \mathbf{R}$ , is the ROM approximation via implicit-projection

$$\tilde{\mathcal{H}}(s) = (\mathbf{C}^T V) (I - (s - \sigma) \tilde{\mathbf{H}})^{-1} \underbrace{(V^H \mathbf{R})}_{\tilde{\boldsymbol{\rho}}} \quad (5.11)$$

to the URM transfer-function

$$\mathcal{H}(s) = \mathbf{C}^T (I - (s - \sigma) \mathbf{H})^{-1} \mathbf{R}.$$

For simplicity let us assume we will augment, or “thicken” the start block of the band-Arnoldi process with a basis  $Y$  for an *exactly*  $\mathbf{H}$ -invariant-subspace, so that  $\mathbf{H}Y = YU$ .

Then the thick-started process

$$\begin{bmatrix} V & \tilde{\mathbf{H}} & \tilde{\boldsymbol{\rho}} \end{bmatrix} = \mathbf{bArnoldi}(\mathbf{H}, \begin{bmatrix} Y & \mathbf{R} \end{bmatrix}) \quad (5.12)$$

yields

$$V = \begin{bmatrix} Y & V' \end{bmatrix}, \quad \tilde{\mathbf{H}} = \begin{bmatrix} U & G \\ & \tilde{\mathbf{H}}' \end{bmatrix}, \quad \text{and} \quad \tilde{\boldsymbol{\rho}} = \begin{bmatrix} Y & V' \end{bmatrix}^H \begin{bmatrix} Y & \mathbf{R} \end{bmatrix} = \begin{bmatrix} \tilde{\boldsymbol{\rho}}_1 & \tilde{\boldsymbol{\rho}}_2 \end{bmatrix}$$

Note that  $\tilde{\boldsymbol{\rho}}_2 = \begin{bmatrix} Y & V' \end{bmatrix}^H \mathbf{R}$ , so the implied ROM transfer-function

$$\tilde{\mathcal{H}}(s) = \left( \mathbf{C}^T \begin{bmatrix} Y & V' \end{bmatrix} \right) \left( I - (s - \sigma) \tilde{\mathbf{H}} \right)^{-1} \underbrace{\left( \begin{bmatrix} Y & V' \end{bmatrix}^H \mathbf{R} \right)}_{\tilde{\boldsymbol{\rho}}_2} \quad (5.13)$$

makes use of only  $\tilde{\boldsymbol{\rho}}_2$  rather than all of  $\tilde{\boldsymbol{\rho}}$ , and  $\tilde{\boldsymbol{\rho}}_1 = \begin{bmatrix} I \\ 0 \end{bmatrix}$  is left out.

### 5.3 Proposing a new model-reduction method

Our restarted ROM method can be described as a multiple-shift method with invariant-subspace recycling. We construct a basis for a Krylov-subspaces at a (potentially) different shift on every cycle. The thick-restart mechanism allows us to hold on to a growing basis  $Y$  of known converged  $(\mathbf{A}, \mathbf{E})$ -invariant-subspace, and continue orthogonalizing new basis vectors against  $Y$ . Assuming  $Y_0 = \{\}$  and  $\ell_j = \dim Y_{j-1}$ , the general process is

1. For  $j = 1, 2, \dots, \tau$ , use Krylov-operator  $\mathbf{H}_j = (\mathbf{A} - \sigma_j \mathbf{E})^{-1} \mathbf{E}$  and operand  $\mathbf{R}_j = (\sigma_j \mathbf{E} - \mathbf{A})^{-1} \mathbf{B}$ .
  - (a) **Expand**  $Y_{j-1}$  into basis  $\begin{bmatrix} Y_{j-1} & V_j \end{bmatrix}$  for  $\mathcal{Y}_{j-1} \cup \mathcal{K}_{n_j}(\mathbf{H}_j, \mathbf{R}_j)$  via  $\ell_j + m$  iterations of thick-start Arnoldi process.
  - (b) The computed quantities  $\tilde{\rho}_2 = \tilde{V}_j^H \mathbf{R}_j$  and  $\tilde{\mathbf{H}} = \tilde{V}_j^H \mathbf{H}_j \tilde{V}_j$  can be used to **observe** the ROM transfer-function

$$\tilde{\mathcal{H}}_j(s) = \mathbf{C}^T \tilde{V}_j (I - (s - \sigma_j) \tilde{\mathbf{H}})^{-1} \tilde{\rho}_2$$

of (5.13) via implicit-projection on to the basis  $\tilde{V}_j = \begin{bmatrix} Y_{j-1} & V_j \end{bmatrix}$ .

- (c) **Deflate**  $\begin{bmatrix} Y_{j-1} & V_j \end{bmatrix}$  into  $\begin{bmatrix} Y_{j-1} & Y'_j \end{bmatrix}$ , where  $Y'_j \subset \text{span } V_j$  is newly converged  $(\mathbf{A}, \mathbf{E})$ -invariant-subspace, and set  $Y_j = \begin{bmatrix} Y_{j-1} & Y'_j \end{bmatrix}$ .

2. Now we have a set of orthogonal blocks  $\{V_1, V_2, \dots, V_\tau\}$  (but not orthogonal to each other) such that

$$\text{span} \begin{bmatrix} V_1 & V_2 & \dots & V_\tau \end{bmatrix} = \bigcup_{j=1}^{\tau} \mathcal{K}_{n_j}(\mathbf{H}_j, \mathbf{R}_j).$$

Let  $\hat{V} = \text{span} \begin{bmatrix} V_1 & V_2 & \dots & V_\tau \end{bmatrix} \in \mathbb{R}^{N \times n}$  (most-likely split into  $\Re$  and  $\Im$  parts and re-orthogonalized). The explicitly-projected ROM realization is then

$$\mathbf{A}_n = \hat{V}^T \mathbf{A} \hat{V}, \quad \mathbf{E}_n = \hat{V}^T \mathbf{E} \hat{V}, \quad \mathbf{B}_n = \hat{V}^T \mathbf{B}, \quad \mathbf{C}_n = \hat{V}^T \mathbf{C},$$

and it has transfer-function

$$\hat{\mathcal{H}}(s) = \mathbf{C}_n^T (\mathbf{A}_n - s \mathbf{E}_n)^{-1} \mathbf{B}_n.$$

Step 1 of the above outline is given as algorithm 3.

---

**Algorithm 3:** EXPLICIT THICK-RESTARTED BAND-ARNOLDI CYCLE

---

**Input:** System realization  $(\mathbf{A}, \mathbf{E}, \mathbf{B}, \mathbf{C})$ , initial interpolation-point  $\sigma_1 \in \mathbb{C}$ .

- 1 Set  $Y_0 = \{\}, V_{\text{ROM}} = \{\}, m := \dim \mathbf{B}$
- 2 **for**  $j = 1, 2, \dots$  **do**
- 3     Set  $\ell_j := \dim Y_{j-1}$
- 4     Let  $\mathbf{H}_j := (\mathbf{A} - \sigma_j \mathbf{E})^{-1} \mathbf{E}$  and  $\mathbf{R}_j := (\sigma_j \mathbf{E} - \mathbf{A})^{-1} \mathbf{B}$
- 5     Compute  $(V, \hat{V}, \dot{V}, \tilde{\mathbf{H}}, \tilde{\rho}) := \mathbf{bArnoldi}(1 : \ell_j, \mathbf{H}_j, [Y_{j-1} \ \mathbf{R}_j])$ .  
       % manually set candidates resulting from processing  $Y_{j-1}$ , to zero.
- 6     Set:  $\hat{v}_i := 0$ , for  $i = \ell_j + m + (1, 2, \dots, \ell_j)$
- 7     Continue  $(V, \hat{V}, \dot{V}, \tilde{\mathbf{H}}, \tilde{\rho}) := \mathbf{bArnoldi}(\ell_j + 1 : n_j, \mathbf{H}_j, [Y_{j-1} \ \mathbf{R}_j])$ .
- 8     Set  $V_{\text{ROM}} := [V_{\text{ROM}} \ V_{\mathbf{R}_j}]$ , where  $V = [v_1 \ v_2 \ \dots \ v_\ell \ V_{\mathbf{R}_j}]$ .  
       % The  $j$ -th implicitly projected ROM transfer-function is given by (5.13).
- 9     Take eigen-decomposition  $\tilde{\mathbf{H}}W = W\Lambda$ . The corresponding poles are  $\mu_i = \sigma_j + 1/\lambda_i$ .  
       Convergence of a Ritz-pair  $(\lambda_i, z_i)$  where  $z_i = [Y \ V] w_i$  is given by (5.10).
- 10    Compute pole-weights  $\gamma_1, \gamma_2, \dots, \gamma_{n_j}$  as (2.23) and (2.24).
- 11    Let  $Z_j$  consist of converged Ritz-vectors and those with large relative-weight  $|\gamma_i|/\Sigma|\gamma_i|$ .
- 12    Let  $(Y_j, T_j) := QR([Y_{j-1} \ Z_j])$
- 13    Select new interpolation-point  $\sigma_{j+1}$ .

**Output:** Basis  $V_{\text{ROM}}$  for  $\bigcup_j \mathcal{K}_{n_j}(\mathbf{H}_j, \mathbf{R}_j)$ .

---

The explicit thick-restarted Band-Arnoldi algorithm is given as algorithm 3. It consists of restarting the band-Arnoldi algorithm (algorithm 2) with a basis of Ritz-vectors and setting to zero the candidate vectors resulting from processing those Ritz-vectors. We experimented with allowing the Ritz-vectors to be processed normally, but it requires more computation and generally resulted in a less accurate ROM for a given size. In practice (for large  $N$ ), we would not process  $Y_{j-1}$  explicitly perform (steps 5 and 6 of algorithm 2). We would perform an implicit-restart method like Krylov-Schur[33], by pre-loading  $V$  with the already orthogonal-basis  $Y_{j-1}$  and  $\tilde{\mathbf{H}}$  with

$$U_{j-1} = T_{j-1} \Lambda_j T_{j-1}^{-1}.$$

Selection of a new interpolation-point (line 13) is left up to whatever method the user chooses; given that we have fairly cheap access to pole distribution data for the implicit ROM transfer-function at any iteration, we assume a point-selection method will take advantage of that. An example of a simple adaptive method is to choose  $\sigma_{j+1}$  to be close to the location of the un-converged pole with largest weight. That would be something like

$$\sigma_{j+1} = \Im(\mu_\tau)$$

where  $\tau = \operatorname{argmax}_i |\gamma_i|$  the un-converged pole with largest weight.



## 5.4 Results

First we consider our two example models, `ex308` and `ex1841` approximated at a single interpolation-point. We recorded the number of iterations of bArnoldi required to reach a relative transfer-function error

$$\frac{\|\mathcal{H} - \tilde{\mathcal{H}}\|}{\|\mathcal{H}\|} \leq \text{tf\_tol} = 0.01, \quad (5.1)$$

at each of three points that are canonical in some sense. Those are a real point  $\pi 10^{10}$ , the  $\Im$ -point  $\pi i 10^{10}$  located roughly at the midpoint of the segment of interest, and the complex point  $(1 + i)\pi 10^{10}$ , shown in figure 5.1. The resulting ROM size in each case depends on the dimension of the split-Krylov-subspace

$$\mathcal{K}_{n'}(\mathbf{H}, \mathbf{R})^* = \mathcal{K}_n(\mathbf{H}, \mathbf{R}) \cup \mathcal{K}_n(\overline{\mathbf{H}}, \overline{\mathbf{R}}),$$

so the dimension  $n'$  of the ROM explicitly projected onto a real basis is no larger than  $n$ . If no deflation occurred during re-orthogonalization of conjugate parts of the projection basis, then  $n' = 2n$  and that was typical for our experiments.

We will give a count of floating-point operations (flops) for producing ROMS. We consider flop-counts to be scalar products in  $\mathbb{R}^n$ , so when complex arithmetic is being used ( $\sigma \notin \mathbb{R}$ ) we must multiply the count by 4. Band-Arnoldi run for  $n$ -iterations with a constant band-size of  $m_c = m$  requires approximately

$$\text{bA\_count}(n) = nN^2 + N(n)(n-1)/2 + Nmn$$

flops. That is  $nN^2$  flops for  $n$ -matrix-vector products,  $N(n)(n-1)/2$  flops for orthogonalization of  $1 + 2 + \dots + n$  basis vectors, and  $Nmn$  flops for orthogonalization of  $m$  candidate-vectors at each iteration.

We include the flop count for tests because we wish to reduce this number using restarts, even if the ROM dimension itself is not appreciably smaller. We would like that for  $l$  cycles of

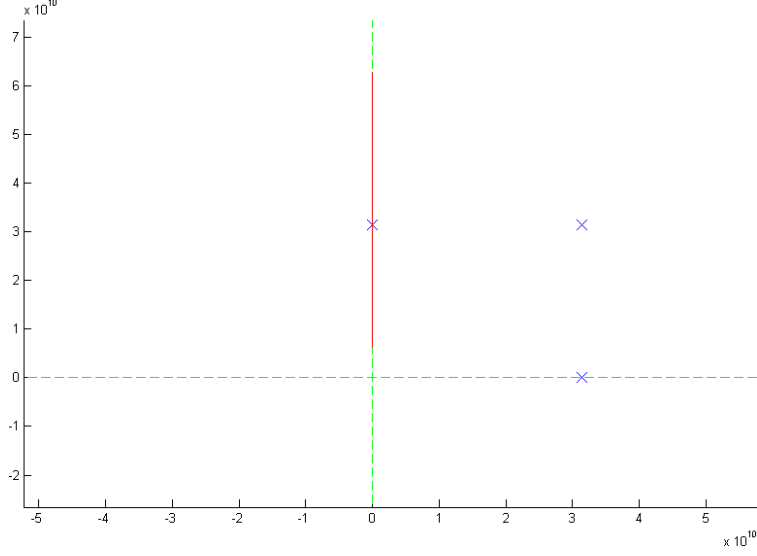


Figure 5.1: The three interpolation-points used for single point benchmarks. The segment of interest  $[10^9, 10^{10}]i$  on the  $\Im$ -axis, is highlighted. Note how small  $[0, 10^9]$  is, in comparison.

band-Arnoldi, each run for  $n_j$  iterations,

$$lM + \sum_j \text{bA\_count}(n_j) < M + \text{bA\_count} \left( \sum n_j \right)$$

$M$  represents the cost of factoring or (re)forming  $\mathbf{H}_j$  and  $\mathbf{R}_j$  which, for a restarted method, must be done  $l$  times (for each  $j = 1, 2, \dots, l$ ). It only needs to be done once for a single-point method. We do not have a value for  $M$  because it varies with the application. It may be negligibly small or prohibitively large, and depends on the size and sparsity of the model realization  $(\mathbf{A}, \mathbf{E}, \mathbf{B}, \mathbf{C})$ .

#### 5.4.1 ex308

ex308 is a  $2 \times 2$  MIMO model of a RCL circuit with 2 input and 2 output terminals, that comes from from a test problem for PEEC modelling of interconnect from IBM or Carnegie Mellon University. ex308 is characterized by many poles very near the  $\Im$ -axis, giving its transfer-function gain a spikey appearance.

$\sigma$	iterations ( $n$ )	ROM size ( $n'$ )	LI	rel-err	flops	figure
$\pi 10^{10}$	144	144	1	7.368e-05	16,920,288 + M	<b>5.4</b>
$i\pi 10^{10}$	71	142	0.992958	5.913e-4	30,177,840 + M	<b>5.5</b>
$(1+i)\pi 10^{10}$	97	194	0.793814	5.1713e-3	42,782,432 + M	<b>5.6</b>

Table 5.1: Benchmark data for **ex308**. flops is a count of real (in  $\mathbb{R}$ ), non-zero scalar products required for matrix-vector multiplication and inner-products.

### **ex308 Benchmarks**

Benchmark data for **ex308** is given in table **5.1**.

ROM size (projection basis dimension) is given as the dimension  $n'$  of the real basis  $V_{\text{ROM}}$  obtained by splitting and re-orthogonalizing  $\hat{V}$ . “LI” is a linear-independence measure defined as

$$\text{LI}(V_{\text{ROM}}) = \frac{\text{rank}_{\text{eff}}(V_{\text{ROM}})}{n}$$

where  $\text{rank}_{\text{eff}}$  is the “effective-rank” of  $V_{\text{ROM}}$  as determined by matlab. We expect the restarted method to produce a less “effectively” linearly-independent basis.

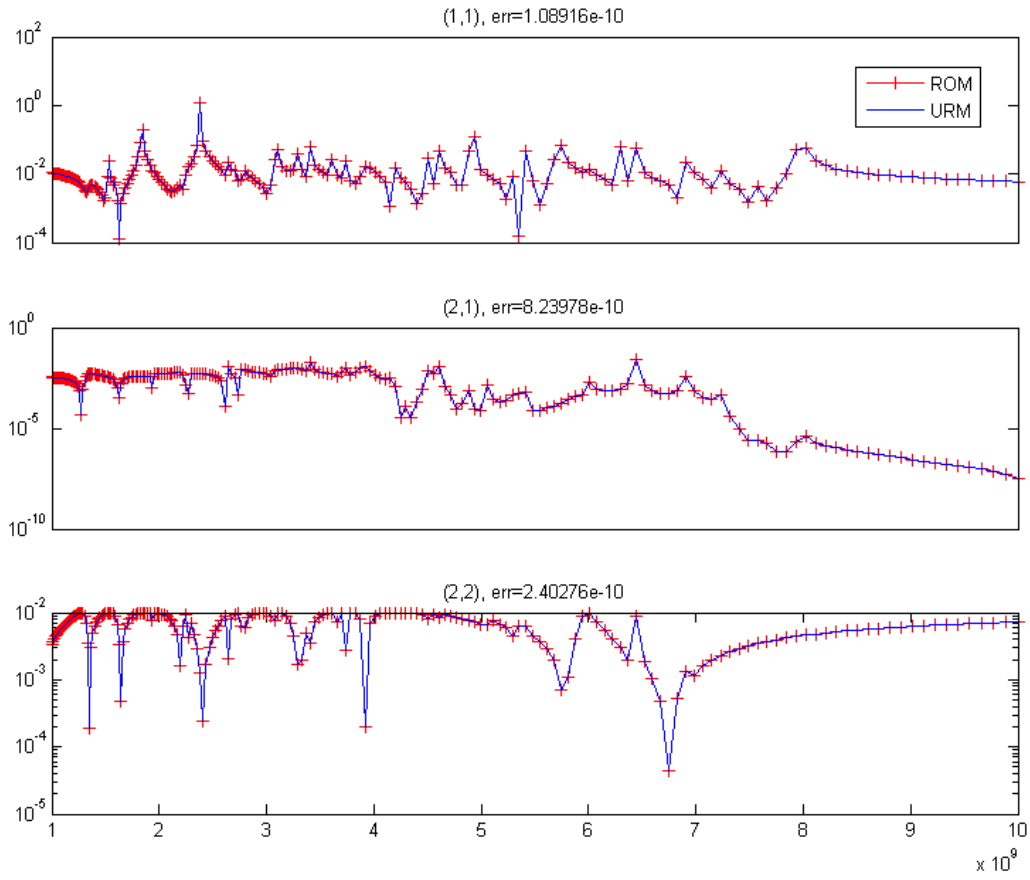


Figure 5.2: These are the three unique gain plots for **ex308**.

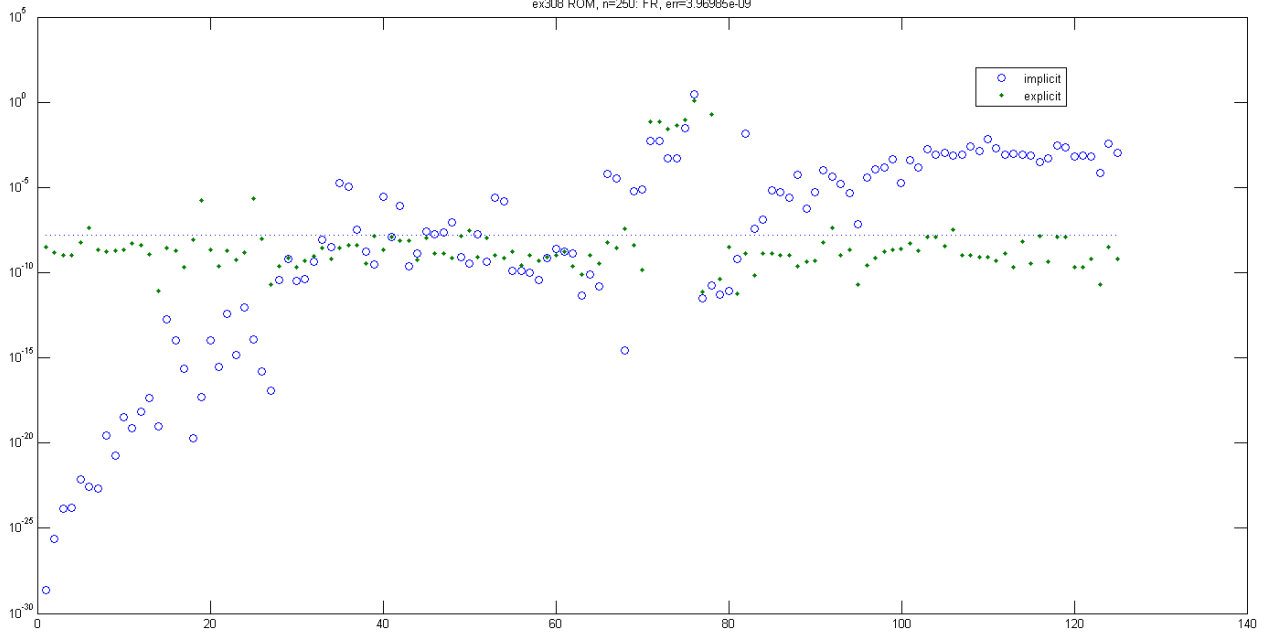


Figure 5.3: This is a plot of relative-residual errors for the 250 poles of an  $n = 250$  ROM (about  $\sigma = (1 + i)\pi 10^{10}$ ) of **ex308**. The circles are the poles derived from eigenvalues of  $\tilde{\mathbf{H}} = V^H \mathbf{H} V$  (the implicit ROM), and the dots are the eigenvalues of the explicitly projected matrix pencil  $(\mathbf{A}_n, \mathbf{E}_n) = (V^H \mathbf{A} V, V^H \mathbf{E} V)$  (the explicit ROM). These are different sets of poles for the most part, except that they converge to the same set of eigenvalues of  $(\mathbf{A}, \mathbf{E})$  as  $n$  increases. We can expect the two ROMs to share *converged* poles. In practice, only the implicit ROM poles (the circles) will be available because relative residual norms are cheap to compute for Ritz-values from  $\tilde{\mathbf{H}}$ . Computing eigen-pairs  $(\mu, z)$  of  $(\mathbf{A}_n, \mathbf{E}_n)$  would require an expensive explicit-projection and there is no cheap formula like (5.8) for the residual norm  $\|\mathbf{A}z - \mu \mathbf{E}z\|$ .

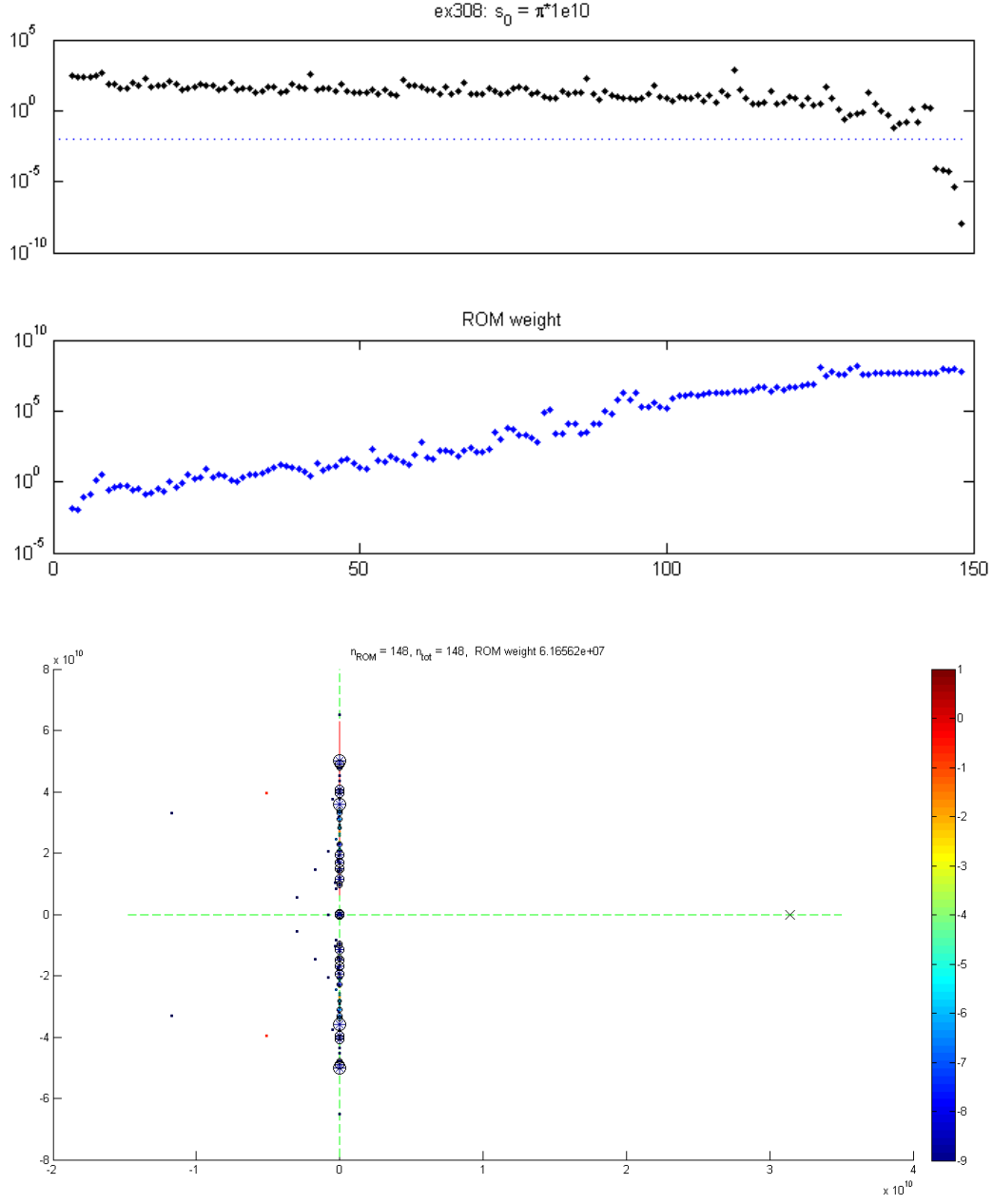


Figure 5.4: Transfer-function relative-error (5.1) (of explicitly projected ROM) and ROM weight (of the implicitly projected ROM) vs.  $n$  for **ex308**, at real interpolation-point  $s_0 = \pi 10^{10} \in \mathbb{R}$ , indicated by ‘x’. The dotted line in the first plot represents a relative-error of 0.01. **ex308** is characterized by a dense distribution of poles on or very near the  $\Im$ -axis, which is evident from the pole distribution of the implicitly defined ROM transfer-function at 148 iterations. That particular ( $n = 148$ ) ROM implies a ROM via *explicit* projection with relative transfer-function error  $\approx 10^{-8}$ .

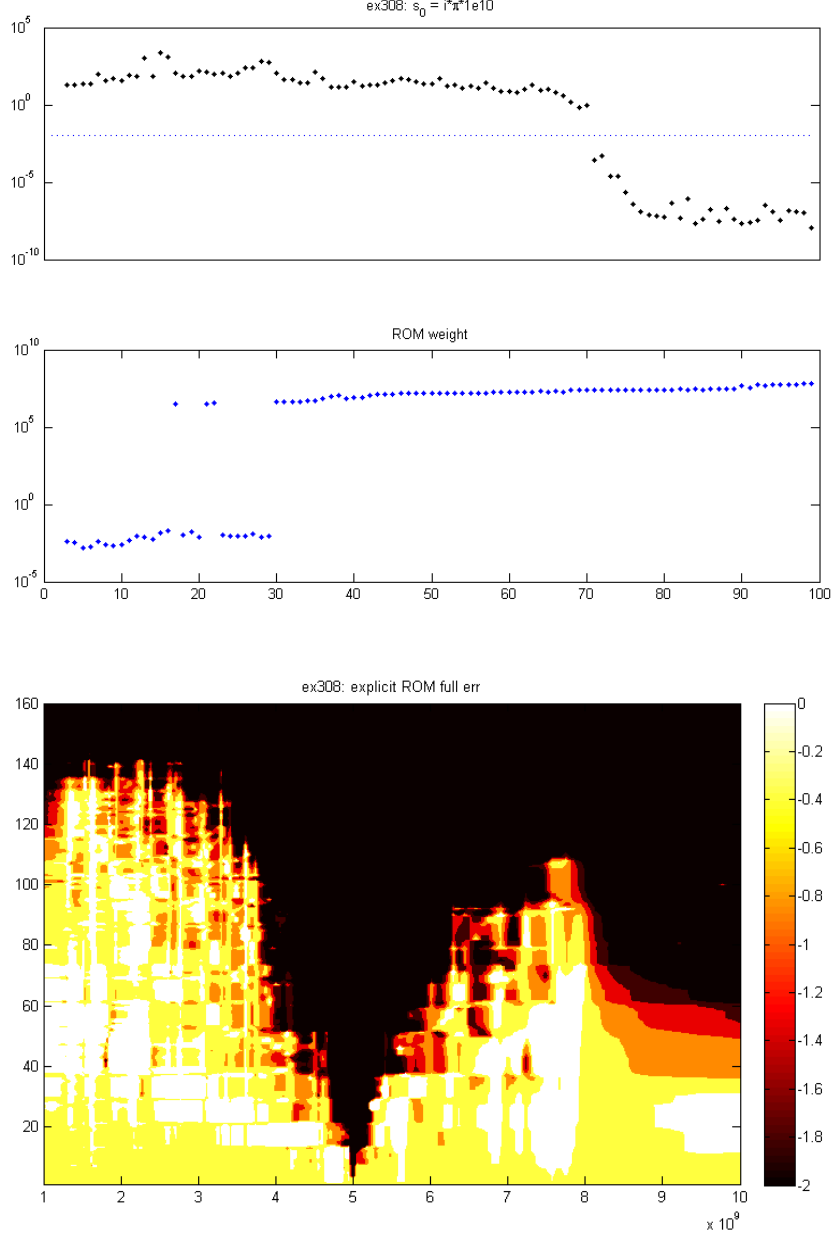


Figure 5.5: Transfer-function relative-error and ROM weight vs.  $n$  for `ex308`, at  $\Im$  interpolation-point  $s_0 = i\pi \cdot 10^{10} \in i\mathbb{R}$ . Unfortunately it appears that ROM weight is not a consistently reliable indicator of transfer-function convergence when using a single interpolation-point. In this example it looks like ROM weight converges after about 30 iterations and after that, only its distribution changes. It is also possible that there is one very dominant pole that appears at  $n = 30$  and it remains one pole as it converges to its resting position. The second plot is relative (explicit) ROM error over iterations  $1, 2, \dots, 160$ . This plot gives a sense of localized convergence of the transfer-function. Since the single interpolation-point is placed near the center of the segment of interest, we see that the transfer-function approximation is most accurate (dark region indicates rel-error is less than 0.01) near the center and convergence works outward from there.

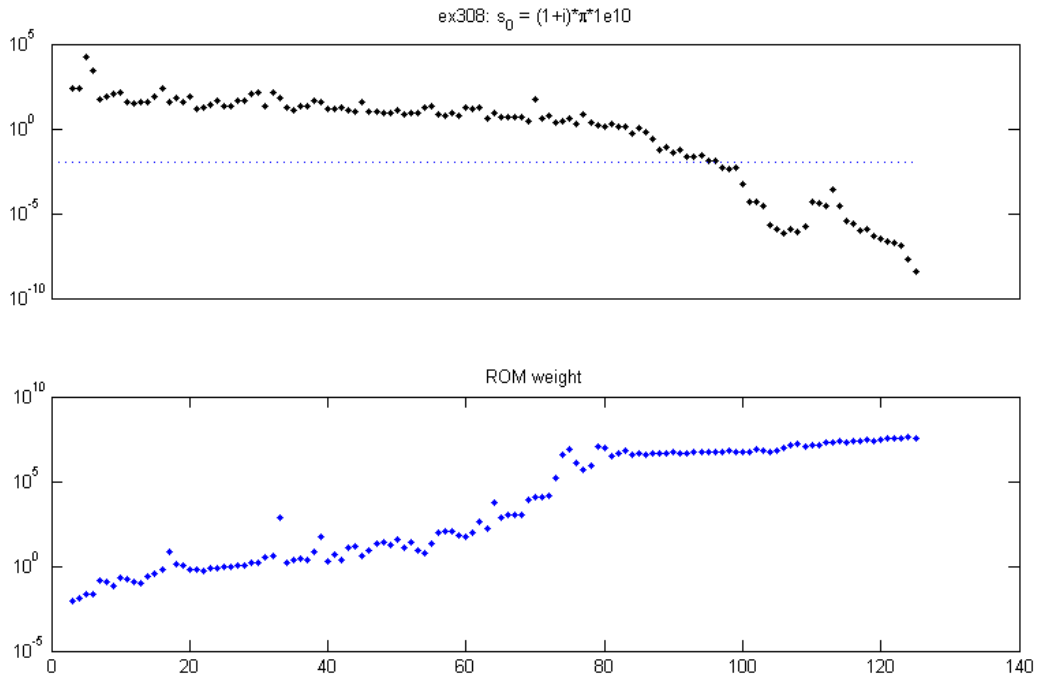


Figure 5.6: Transfer-function relative-error and ROM-weight vs.  $n$  for **ex308**, at  $\sigma = (1+i)\pi \cdot 10^{10}$ . This is much the way we would like the relationship between ROM-weight and transfer-function error to look. ROM-weight leveling-off would indicate transfer-function error convergence.



### ex308 thick-restart example 1

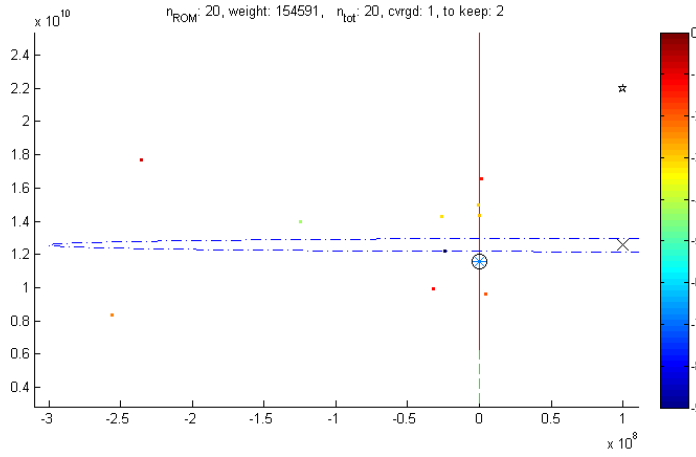
Here we show an example run of the thick-restarted band-Arnoldi process using three pre-set interpolation-points  $\sigma_j = 10^8 + 2\pi i \cdot 10^{10} p_j$  for  $p_j = 2, 3.5, 6$ . When scaled this way, the frequency range of interest  $p \in (1, 10)$  corresponds to  $s \in i(10^9, 10^{10})$  so our choices of  $p$  suggest convergence of the frequency-response at those localities first, and outward from there. We ran the algorithm for for  $n_j = 20, 25, 25$  iterations.

Converged Ritz-vectors (those with relative residual less than  $\text{ctol} = \sqrt{\epsilon} \approx 1.49\text{e-}8$ ) and those associated with dominant poles ( $\text{wt}_i / \sum \text{wt}_i \leq 0.05$ ) were recycled.

The resulting ROM required a total of 70 iterations (not including re-processing thick-restarting Ritz-vectors), was of size  $n' = 140$  and had a relative-error of  $7.14155\text{e-}05$ , making it compare favorably with the benchmark examples in table 5.1. It required  $30,217,264 + 3M$  flops, where  $M$  is the cost of creating  $\mathbf{H}_j$  and  $\mathbf{R}_j$ .

Execution of `test4('ex308', 1e8+2i*pi*1e9*[2 3.5 6], [20 25 25])` yields

```
cycle 1 expanding at s_0 = 2\pi i 10^9(0.0159155 + 2i), band_size = 2 + 0
... ROM: 20, n_tot: 20, converged: 1, keep: 2 weight: 154591
...updating thick-restart basis...dim Y = 2
```

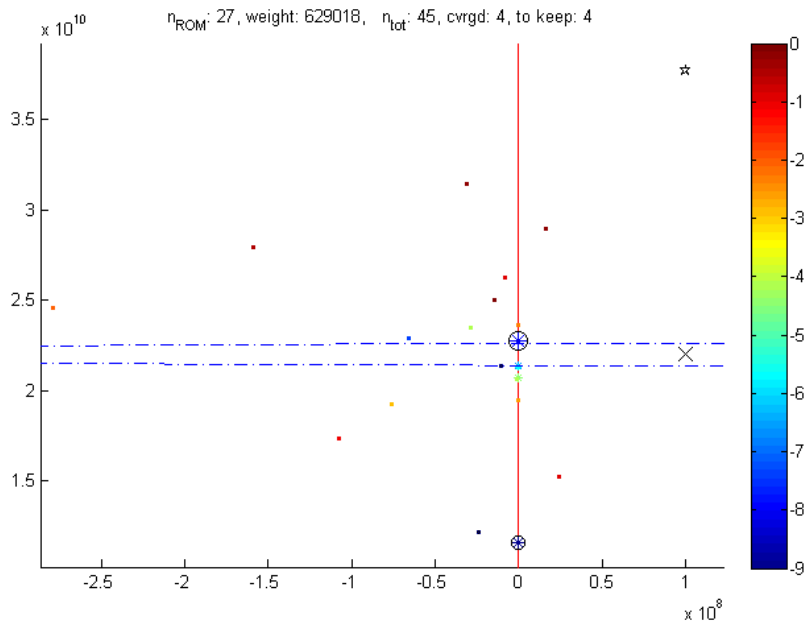


In the above plot, poles of the implicitly projected ROM of the first cycle are indicated by ‘\*’ symbols. The color of a pole indicates its degree of convergence. The interpolation-point is indicated by ‘x’, and a dashed-circle<sup>4</sup> around the interpolation-point indicates distance to the first

<sup>4</sup>elongated in the plot due to greatly asymmetric scaling. This is one reason why interpolation-points on or near the  $\Im$ -axis can result in the convergence of un-necessarily large ROMs.

converged pole. The ‘★’ symbol indicates placement of the next interpolation-point. Pole-size in the plot corresponds to weight. Some poles have circles around them, indicating that those will be kept for thick-restarting the next cycle. In this example we are lucky to have had a very dominant pole among the two that converged on the first cycle.

```
cycle 2 expanding at s_0 = 2\pi 10^9(0.0159155 + 3.5i), band_size = 2 + 2
v_defl(5)/H_est = 0 < 1.49012e-08, mc = 3
v_defl(5)/H_est = 0 < 1.49012e-08, mc = 2
... ROM: 27, n_tot: 45, converged: 4, keep: 4 weight: 629018
...updating thick-restart basis...dim Y = 4
```



```
cycle 3 expanding at s_0 = 2\pi 10^9(0.0159155 + 6i), band_size = 2 + 4
v_defl(7)/H_est = 0 < 1.49012e-08, mc = 5
v_defl(7)/H_est = 0 < 1.49012e-08, mc = 4
v_defl(7)/H_est = 0 < 1.49012e-08, mc = 3
v_defl(7)/H_est = 0 < 1.49012e-08, mc = 2
... ROM: 29, n_tot: 70, converged: 6, keep: 8 weight: 0.00604195
...updating thick-restart basis...dim Y = 8
```



	$\Re(\mu)$	$\Im(\mu)$	rr	wt	keep
1	-2.3746e+07	1.2186e+10i	3.6646e-11	0.0138714	1
2	4.8802e+01	1.1562e+10i	1.75228e-07	154560	1
3	-1.2457e+08	1.3997e+10i	5.15493e-05	0.0166039	0
4	-3.2363e+08	1.0214e+10i	0.000147586	0.0136769	0
5	-1.7281e+09	1.4688e+10i	0.00067386	0.0223279	0
6	-2.6226e+07	1.4295e+10i	0.000702401	0.0122022	0
7	-9.2396e+05	1.4978e+10i	0.00092006	14.3325	0
8	1.2915e+05	1.4334e+10i	0.00126677	1.344	0
9	-2.5567e+08	8.3294e+09i	0.00515879	0.0234752	0
10	4.1068e+06	9.6085e+09i	0.0136536	0.503765	0

(a) Cycle 1

	$\Re(\mu)$	$\Im(\mu)$	rr	wt	keep
1	-2.3746e+07	1.2186e+10i	0	0.0110811	1
2	1.9498e+02	1.1562e+10i	0	44474.2	1
3	-1.0463e+07	2.1391e+10i	3.12586e-09	0.151564	1
4	1.1927e-01	2.2714e+10i	7.41398e-09	567340	1
5	-6.6011e+07	2.2864e+10i	4.3831e-08	0.013322	0
6	-4.9619e+00	2.1332e+10i	2.98909e-07	11524.3	0
7	5.1626e+00	2.1293e+10i	2.11094e-06	834.36	0
8	-8.2289e+08	2.0635e+10i	1.2973e-05	0.0239924	0
9	-4.4901e+01	2.0682e+10i	3.32799e-05	4827.55	0
10	-2.8524e+07	2.3487e+10i	6.76036e-05	0.0287073	0

(b) Cycle 2

	$\Re(\mu)$	$\Im(\mu)$	rr	wt	keep
1	-2.3746e+07	1.2186e+10i	0	3.06212e-08	1
2	4.9552e+00	2.2714e+10i	0	6.69624e-09	1
3	-1.0463e+07	2.1391e+10i	0	5.41799e-08	1
4	-2.4655e+02	1.1562e+10i	0	1.49686e-07	1
5	-4.3476e+06	3.7641e+10i	5.19823e-23	1.98325e-08	1
6	-5.1591e+08	3.7457e+10i	2.15413e-14	2.57324e-07	1
7	-1.2248e+00	3.5908e+10i	3.04755e-08	4.08767e-07	0
8	-3.5711e+00	3.9570e+10i	4.79573e-07	5.11497e-08	0
9	-2.8146e+07	4.0486e+10i	3.09744e-06	8.83886e-08	0
10	-3.6394e+04	3.9670e+10i	3.63502e-06	3.98871e-09	0

(c) Cycle 3

Table 5.2: The 10 Ritz-poles of lowest relative-residual for each cycle. One thing to note is that pole-weight does not stay consistent from cycle to cycle for this test. Two converged poles (indicated by shaded rows in the tables) appear to change dominance quite drastically. This could be due to the way we define pole-weight.

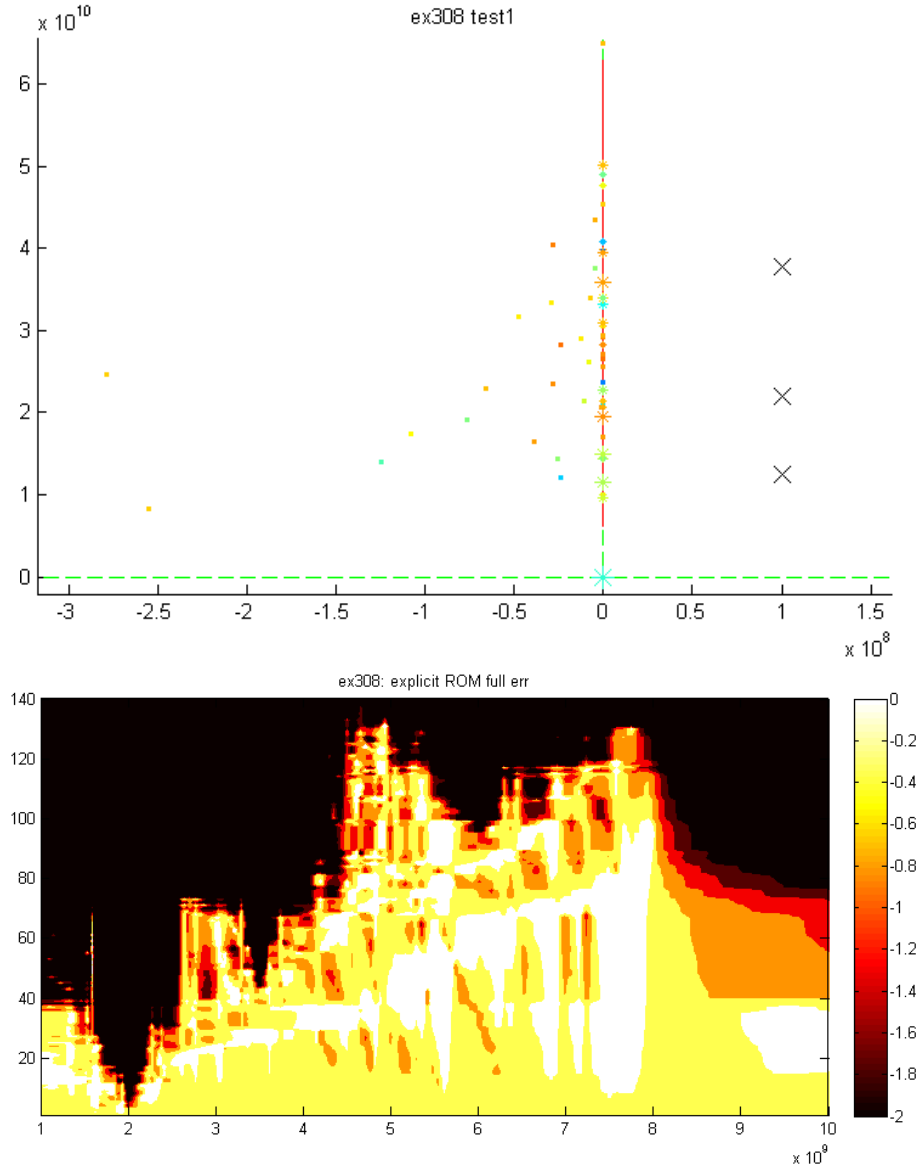


Figure 5.7: The pole distribution for the explicitly-projected transfer-function and interpolation-points are in the first plot. The second plot, of local transfer-function-error over the frequency range of interest evolving with inclusion of basis vectors in  $V_{\text{ROM}}$  reflects expansion about the points  $p = 2, 3.5$ , and  $6$ . It appears that a smaller and more accurate ROM could have been constructed with fewer iterations at  $p = 2$  and more iterations at  $p = 6$ , or possibly two more interpolation-points at  $p = 5$  and  $8$ . We found the  $1\text{e}8$  offset from the  $\Im$ -axis to yield good results in numerous test-runs of the process for this example.

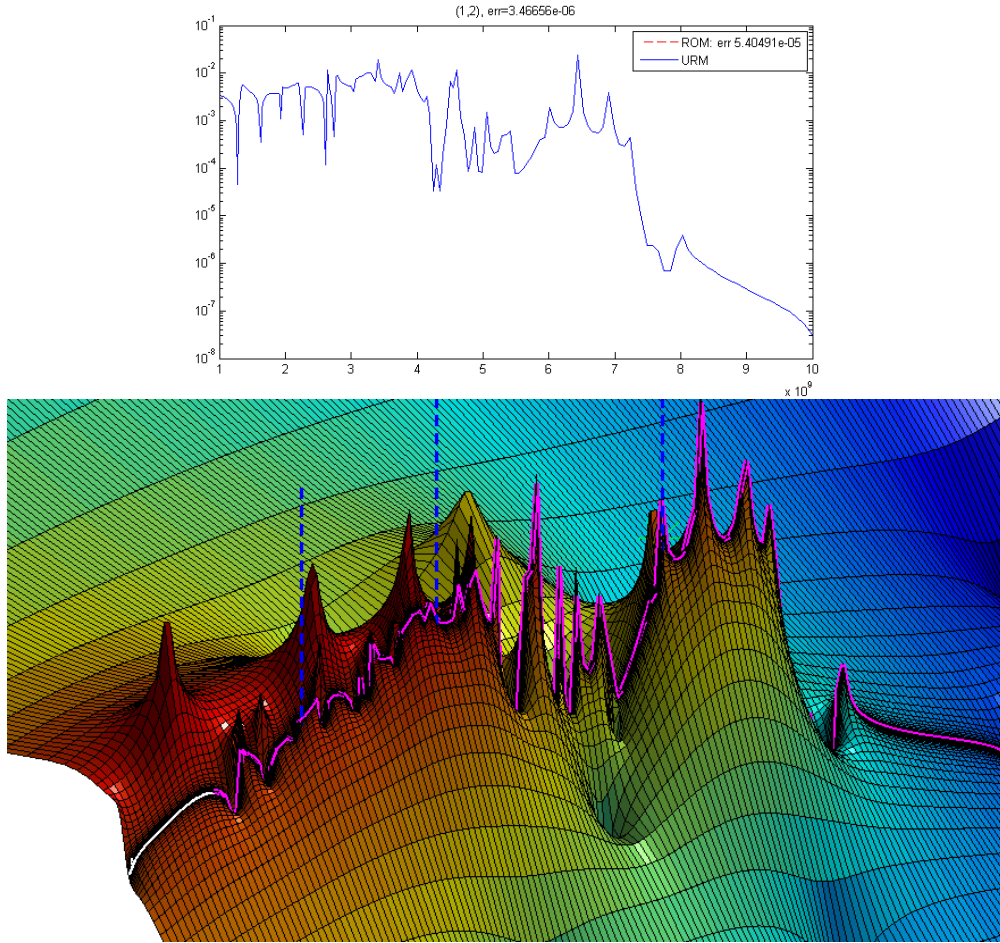


Figure 5.8: Plots of the (1,2) component of the frequency-response and transfer-function surface for example 1. It looks like the interpolation-point placement was almost ideal for `ex308`, in the sense that the points are near centers of pole-mass. Note that the interpolation-points are actually offset `1e8` into the  $\Re$  half-plane, but the scale of the surface plot is such that they appear to be on the segment of interest.

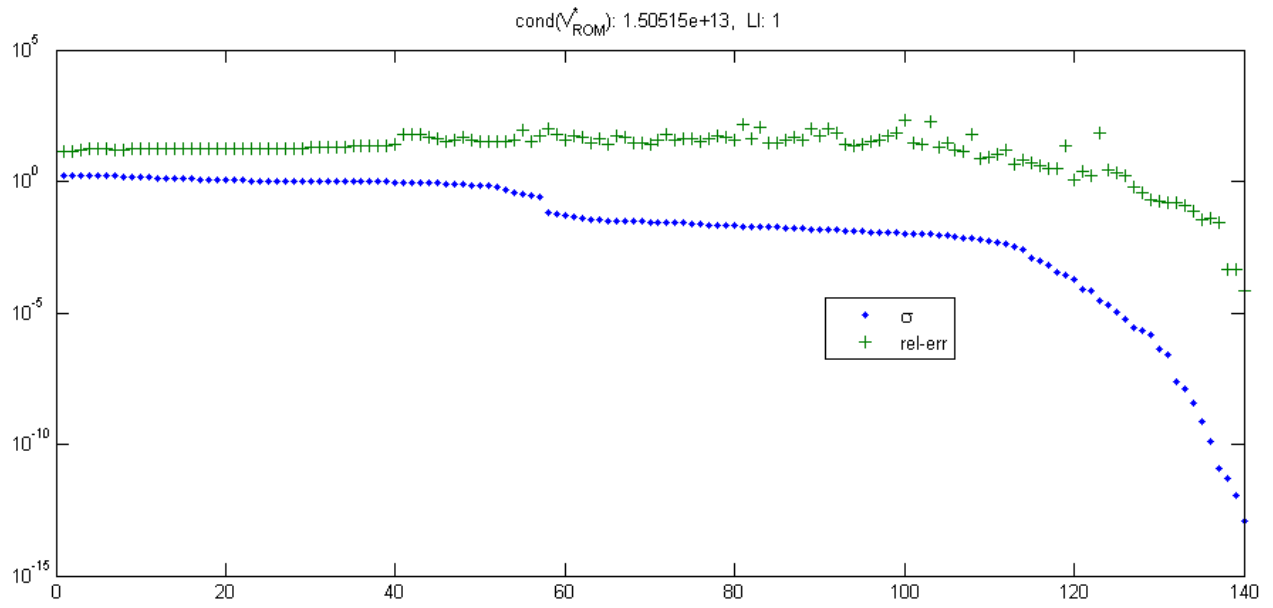


Figure 5.9: Here we looked at relative-error for explicitly-projected ROM transfer-functions for successively larger bases  $U_n$  for  $n = 1, 2, \dots, n' = 140$ , where  $U$  is the basis of left singular-vectors of  $\hat{V}$ , and  $\sigma$  is the corresponding singular value.

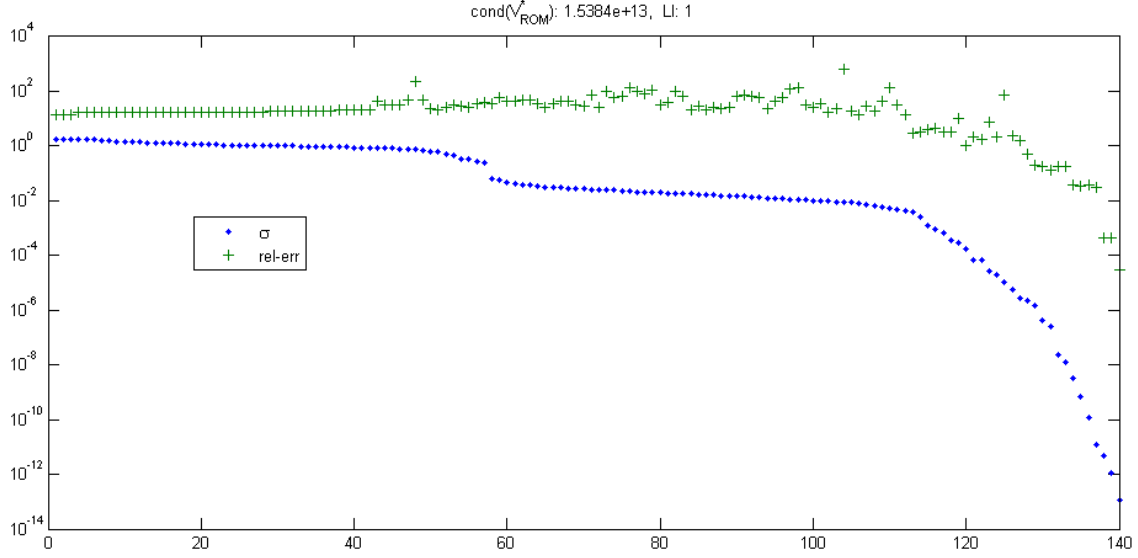


Figure 5.10: Transfer-function relative error for **ex308** test 2, plotted in order of decreasing singular values. The condition number of the split basis in this case is  $1.5384\text{e}13$ , which is not much different from that from test 1, indicating that keeping Ritz-vectors did improve linear independence of the projection basis, but not by much.

### **ex308 thick-restart example 2**

The next test is similar to the previous one except that we did not keep any Ritz-vectors from cycle to cycle. In this case the model of size  $n' = 140$  had about the same accuracy at with a relative transfer-function error of  $2.84772\text{e}-05$ . It was negligibly cheaper to construct at  $27,707,680 + 3M$  flops. There was very little overlap between Krylov-subspaces in this case, possibly because of the low number of iterations we performed.



	$\Re(\mu)$	$\Im(\mu)$	<b>rr</b>	<b>wt</b>	<b>keep</b>
1	-2.3746e+07	1.2186e+10i	3.6646e-11	0.0138714	0
2	4.8802e+01	1.1562e+10i	1.75228e-07	154560	0
3	-1.2457e+08	1.3997e+10i	5.15493e-05	0.0166039	0
4	-3.2363e+08	1.0214e+10i	0.000147586	0.0136769	0
5	-1.7281e+09	1.4688e+10i	0.00067386	0.0223279	0
6	-2.6226e+07	1.4295e+10i	0.000702401	0.0122022	0
7	-9.2396e+05	1.4978e+10i	0.00092006	14.3325	0
8	1.2915e+05	1.4334e+10i	0.00126677	1.344	0
9	-2.5567e+08	8.3294e+09i	0.00515879	0.0234752	0
10	4.1068e+06	9.6085e+09i	0.0136536	0.503765	0

(a) Cycle 1

	$\Re(\mu)$	$\Im(\mu)$	<b>rr</b>	<b>wt</b>	<b>keep</b>
1	-1.0463e+07	2.1391e+10i	3.15983e-09	0.151564	0
2	1.3560e-01	2.2714e+10i	7.5411e-09	559182	0
3	-6.6011e+07	2.2864e+10i	4.42028e-08	0.013322	0
4	-2.4243e-02	2.1332e+10i	3.18392e-07	67080.5	0
5	4.2240e+00	2.1293e+10i	2.25896e-06	984.276	0
6	-8.2289e+08	2.0635e+10i	1.3741e-05	0.0239923	0
7	-9.0853e+02	2.0682e+10i	3.62199e-05	243.656	0
8	-2.8524e+07	2.3487e+10i	6.645e-05	0.0287091	0
9	-1.3230e+06	2.0649e+10i	7.29665e-05	0.0201709	0
10	-7.5781e+07	1.9234e+10i	0.0017068	0.0163881	0

(b) Cycle 2

	$\Re(\mu)$	$\Im(\mu)$	<b>rr</b>	<b>wt</b>	<b>keep</b>
1	-4.3476e+06	3.7641e+10i	5.23777e-23	2.51929e-08	0
2	-5.1591e+08	3.7457e+10i	2.16358e-14	3.26875e-07	0
3	-2.2323e+00	3.5908e+10i	3.06443e-08	5.1925e-07	0
4	-5.0044e+00	3.9570e+10i	4.83015e-07	6.49747e-08	0
5	-2.8146e+07	4.0486e+10i	3.1165e-06	1.12279e-07	0
6	-3.6431e+04	3.9670e+10i	3.65723e-06	5.06681e-09	0
7	2.2037e+03	4.0796e+10i	0.000200954	5.49435e-09	0
8	4.6830e+04	3.3983e+10i	0.000950016	1.14312e-07	0
9	-3.9908e+06	4.3477e+10i	0.0023368	1.45873e-07	0
10	-7.1975e+06	3.3975e+10i	0.0024024	5.09175e-09	0

(c) Cycle 3

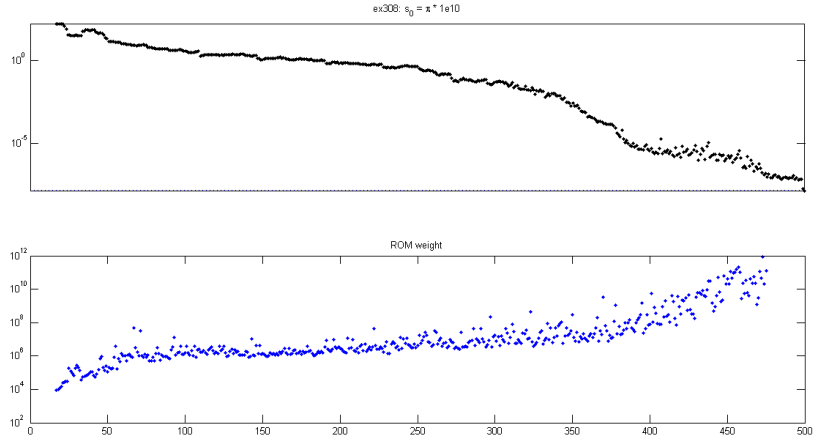
Table 5.3: Here are the most converged eigenvalues of each cycle of the same example, except that we kept no Ritz values from cycle to cycle.

$\sigma$	iterations ( $n$ )	ROM size ( $n'$ )	LI	rel-err	flops	figure
$\pi 10^{10}$	310	310	1	9.1289e-3	1,147,983,165 + M	<b>5.11a</b>
$i\pi 10^{10}$	106	212	1	8.6730e-3	1,490,525,148 + M	<b>5.11b</b>
$(1 + i)\pi 10^{10}$	142	284	1	8.4797e-3	2,015,563,620 + M	<b>5.11c</b>

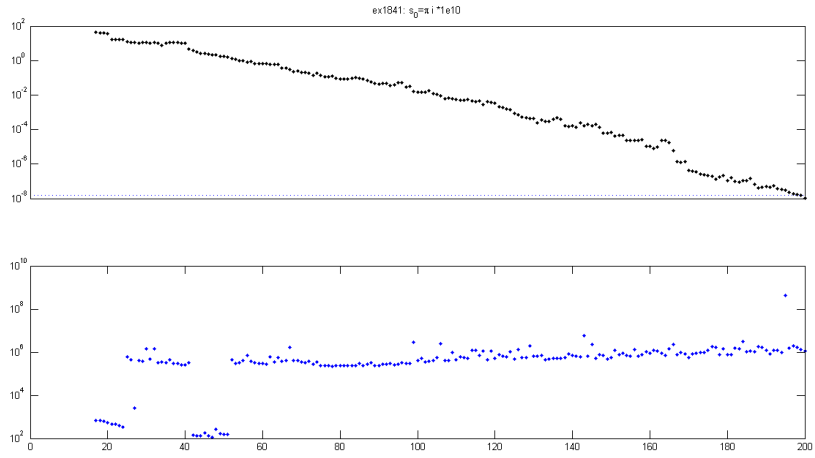
Table 5.4: Benchmark data for **ex1841**. flops is a count of real (in  $\mathbb{R}$ ), non-zero scalar products required for matrix-vector multiplication and inner-products.

### 5.4.2 ex1841

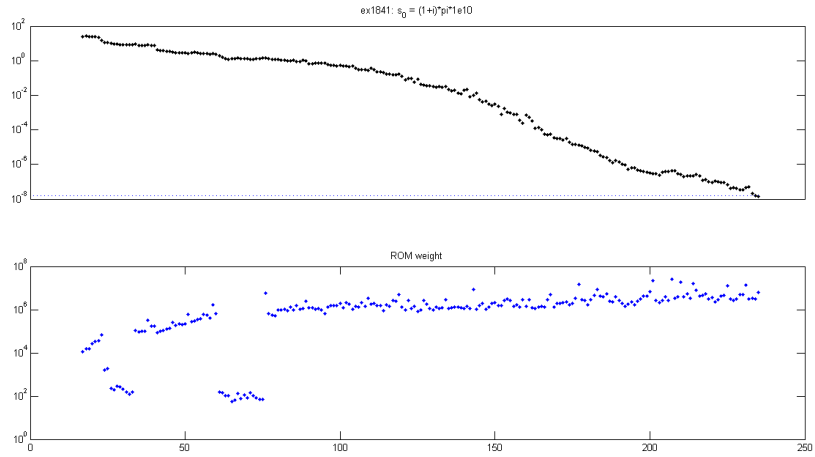
**ex1841** is a  $16 \times 16$  (inputs  $\times$  outputs) MIMO model that comes from a MNA expression of an RCL circuit model of interconnect.



(a)  $\sigma = \pi 10^{10}$



(b)  $\sigma = i\pi 10^{10}$



(c)  $\sigma = (1 + i)\pi 10^{10}$

Figure 5.11: Transfer-function relative-error (5.1) and ROM weight vs.  $n$  for ex1841, at each of the three points shown in figure 5.1.

### ex1841 test1

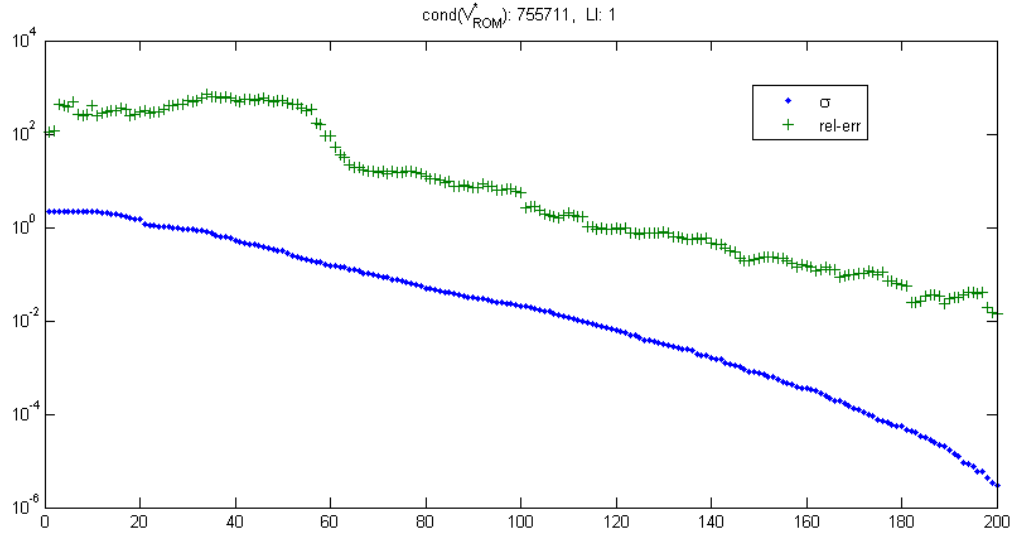
We ran the thick-restarted band-Arnoldi algorithm for 20 iterations at each of 5 interpolation points  $10^3 + 2ip_j \cdot \pi \cdot 10^9$ , for  $p_j = 1, 3, 5, 7, 9$ .

```
test4('ex1841',1e3+[1 3 5 7 9]*2i*pi*1e9,[20 20 20 20 20]);
```

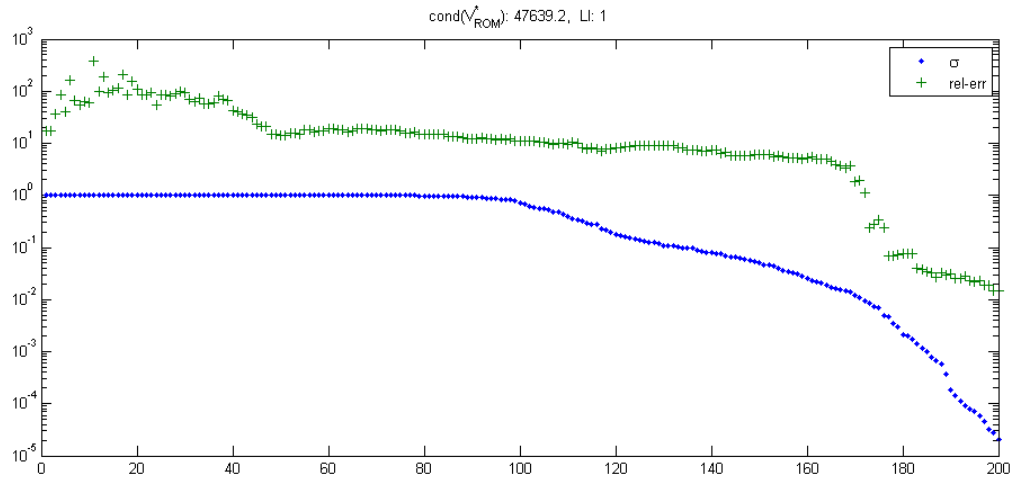
The purpose of this test was to compare ROMs produced by restarted band-Arnoldi with and without “thick-restarting”, i.e. keeping Ritz vectors. For low iteration counts, the difference is almost negligible with this model, even when interpolation points are fairly closely located.

keep_tol	$n$	$n'$	rel-error	flops	cond #
0	100	200	0.0143945	1374490600 + 5M	755711
$\infty$	100	200	0.0143945	4270751800 + 5M	47639

In fact, for this test the only differences in models are the operations required to produce them and the linear-independence of the basis vectors.

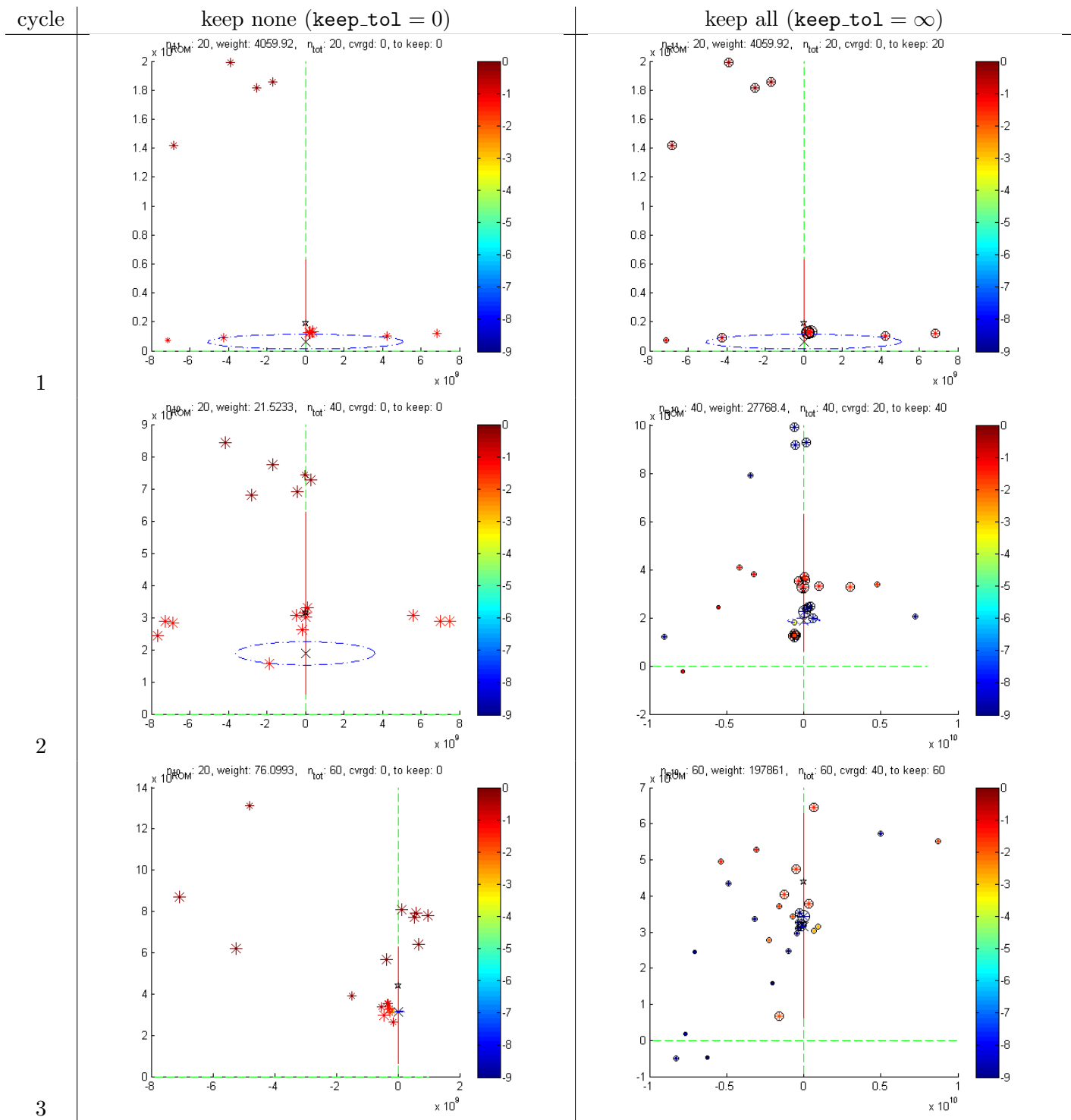


(a) keeping no Ritz-vectors, condition number: 755711

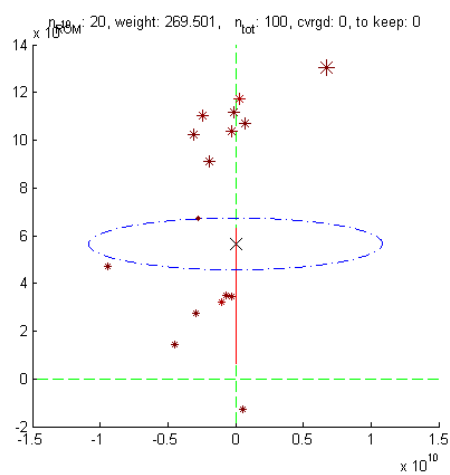
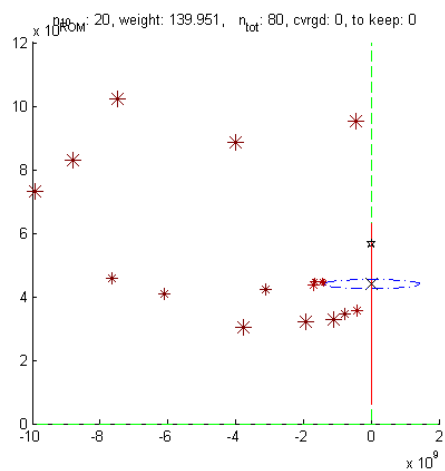


(b) keeping all Ritz-vectors, condition number: 47639

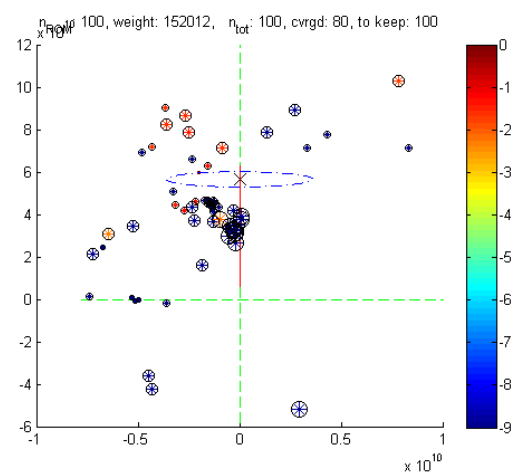
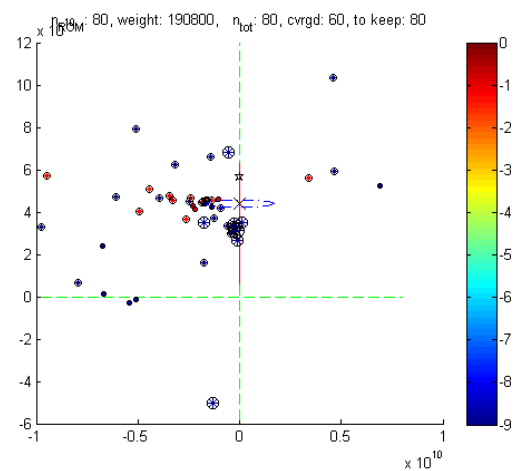
Figure 5.12: singular values and relative transfer-function error for each ROM of test 1.



4



5



# References

- [1] L.A. Aguirre. Quantitative measure of modal dominance for continuous systems. In *Decision and Control, 1993., Proceedings of the 32nd IEEE Conference on*, pages 2405–2410vol.3, 1993. [2.2.5](#)
- [2] Nisar Ahmed and MM Awais. Implicit restart scheme for large scale krylov subspace model reduction method. In *Multi Topic Conference, 2001. IEEE INMIC 2001. Technology for the 21st Century. Proceedings. IEEE International*, pages 131–138. IEEE, 2001. [3.3](#)
- [3] P. Benner, M.E. Hochstenbach, and P. Kurschner. Model order reduction of large-scale dynamical systems with jacobi-davidson style eigensolvers. In *Communications, Computing and Control Applications (CCCA), 2011 International Conference on*, pages 1–6, 2011. [4.1](#)
- [4] W.K. Chen. *The circuits and filters handbook*. The electrical engineering handbook series. CRC Press, 2009.
- [5] David Day and Michael A. Heroux. Solving complex-valued linear systems via equivalent real formulations. *SIAM J. Sci. Comput.*, 23(2):480–498, 2001. [4.3](#)
- [6] V. Druskin and V. Simoncini. Adaptive rational krylov subspaces for large-scale dynamical systems. *Systems & Control Letters*, 60(8):546–560, 2011. [4.1](#)
- [7] P. Feldmann and R.W. Freund. Efficient linear circuit analysis by pade approximation via the lanczos process. *Computer-Aided Design of Integrated Circuits and Systems, IEEE Transactions on*, 14(5):639–649, 1995. [4](#)
- [8] Michalis Frangos and Imad M. Jaimoukha. Adaptive rational interpolation: Arnoldi and lanczos-like equations. *European Journal of Control*, 14(4):342–354, 2008. [4.1](#)
- [9] Roland W. Freund. Krylov-subspace methods for reduced-order modeling in circuit simulation. *J. Comput. Appl. Math.*, 123(1-2):395–421, 2000. [3.4.2](#)
- [10] Roland W. Freund. Model reduction methods based on Krylov subspaces. *Acta Numerica*, 12:267–319, 2003. [5.1](#), [5.2.1](#)
- [11] Kyle Gallivan, G Grimme, and Paul Van Dooren. A rational lanczos algorithm for model reduction. *Numerical Algorithms*, 12(1):33–63, 1996. [4.1](#), [4.1.1](#)
- [12] Juan M Gracia and Francisco E Velasco. Stability of invariant subspaces of regular matrix pencils. *Linear algebra and its applications*, 221:219–226, 1995. [2.2.3](#)



- [13] E Grimme and K Gallivan. Krylov projection methods for rational interpolation. 1997. [4.1](#)
- [14] E Grimme and K Gallivan. A rational lanczos algorithm for model reduction II: Interpolation point selection. In *Numerical Algorithms*, 1998. [4.1](#), [4.2](#)
- [15] Eric James Grimme, Danny C Sorensen, and Paul Van Dooren. Model reduction of state space systems via an implicitly restarted lanczos method. *Numerical Algorithms*, 12(1):1–31, 1996. [2.3](#), [3.3](#)
- [16] Chung-Wen Ho, Albert E. Ruehli, and Pierce A. Brennan. The modified nodal approach to network analysis. *Circuits and Systems, IEEE Transactions on*, 22(6):504–509, 1975. [2.2.4](#)
- [17] Imad M Jaimoukha and Ebrahim M Kasenally. Implicitly restarted krylov subspace methods for stable partial realizations. *SIAM Journal on Matrix Analysis and Applications*, 18(3):633–652, 1997. [2.3](#), [3.3](#)
- [18] Guillaume Lassaux and K Willcox. Model reduction for active control design using multiple-point arnoldi methods. *AIAA Paper*, 616:2003, 2003. [4.1](#)
- [19] Herng-Jer Lee, Chia-Chi Chu, and Wu-Shiung Feng. Multi-point model reductions of VLSI interconnects using the rational arnoldi method with adaptive orders (RAMAO). In *Circuits and Systems, 2004. Proceedings. The 2004 IEEE Asia-Pacific Conference on*, volume 2, pages 1009–1012vol.2, 2004. [4.1](#)
- [20] Herng-Jer Lee, Chia-Chi Chu, and Wu-Shiung Feng. An adaptive-order rational arnoldi method for model-order reductions of linear time-invariant systems. *Linear Algebra and its Applications*, 415:235–261, 2006. [Special Issue on Order Reduction of Large-Scale Systems](#). [4.1](#)
- [21] RB Lehoucq and KJ Maschhoff. Implementation of an implicitly restarted block arnoldi method. *Preprint MCS-P649-0297, Argonne National Lab*, 1997. [5.1](#)
- [22] K Henrik A Olsson and Axel Ruhe. Rational krylov for eigenvalue computation and model order reduction. *BIT Numerical Mathematics*, 46(1):99–111, 2006. [4.1](#)
- [23] Theodore W Palmer. *Banach Algebras and the General Theory of \*-algebras: Volume 1*, volume 2. Cambridge University Press, 2001. [4.3](#)
- [24] Vasilios Papakos and IM Jaimoukha. A deflated implicitly restarted lanczos algorithm for model reduction. In *Decision and Control, 2003. Proceedings. 42nd IEEE Conference on*, volume 3, pages 2902–2907. IEEE, 2003. [3.3](#)
- [25] Beresford N. Parlett and Youcef Saad. Complex shift and invert strategies for real matrices. *Linear Algebra and its Applications*, 8889(0):575–595, 1987. [4.2](#), [4.3](#)
- [26] Beresford N Parlett and David S Scott. The lanczos algorithm with selective orthogonalization. *Mathematics of computation*, 33(145):217–238, 1979. [3](#), [5.2.1](#)
- [27] Axel Ruhe. The rational krylov algorithm for nonsymmetric eigenvalue problems. II. [4.1](#), [4.2.1](#)

- [28] Axel Ruhe. Implementation aspects of band lanczos algorithms for computation of eigenvalues of large sparse symmetric matrices. *Mathematics of Computation*, 33(146):680–687, 1979. 5.1
- [29] Axel Ruhe. Rational krylov sequence methods for eigenvalue computation. *Linear Algebra and its Applications*, 58(0):391–405, 1984. 4.1
- [30] L Miguel Silveira, Mattan Kamon, Ibrahim Elfadel, and Jacob White. A coordinate-transformed Arnoldi algorithm for generating guaranteed stable reduced-order models of RLC circuits. *Computer Methods in Applied Mechanics and Engineering*, 169(3):377–389, 1999. 2.3
- [31] Andreas Stathopoulos, Yousef Saad, and Kesheng Wu. Dynamic thick restarting of the davidson, and the implicitly restarted arnoldi methods. *SIAM J. Sci. Comput*, 19:227–245, 1996. 5.2.1
- [32] Gilbert W Stewart. On the sensitivity of the eigenvalue problem  $ax=\lambda bx$ . *SIAM Journal on Numerical Analysis*, 9(4):669–686, 1972. 2.2.3
- [33] GW Stewart. A krylov–schur algorithm for large eigenproblems. *SIAM Journal on Matrix Analysis and Applications*, 23(3):601–614, 2002. 5.2.1, 13