

PAPER

Singular value automata and approximate minimization

Borja Balle^{1,*†}, Prakash Panangaden² and Doina Precup²

¹Amazon Research, Cambridge, UK and ²School of Computer Science, McGill University, Montreal, Canada

*Corresponding author. Email: borja.balle@gmail.com

(Received 19 January 2018; revised 07 February 2019; accepted 18 April 2019; first published online 27 May 2019)

Abstract

The present paper uses spectral theory of linear operators to construct approximately minimal realizations of weighted languages. Our new contributions are: (i) a new algorithm for the singular value decomposition (SVD) decomposition of finite-rank infinite Hankel matrices based on their representation in terms of weighted automata, (ii) a new canonical form for weighted automata arising from the SVD of its corresponding Hankel matrix, and (iii) an algorithm to construct approximate minimizations of given weighted automata by truncating the canonical form. We give bounds on the quality of our approximation.

1. Introduction

When one considers *quantitative systems*, it becomes meaningful to ask about the *approximate* minimization of transition systems or automata. This concept, meaningless for ordinary automata, is appropriate for many types of systems: weighted automata, probabilistic automata of various kinds, and timed automata. The present paper focuses on weighted automata where we are able to exploit spectral theory of linear operators to construct approximately minimal realizations of weighted languages. Our main contributions are:

- A new algorithm for the singular value decomposition (SVD) of finite-rank infinite Hankel matrices based on their representation in terms of weighted automata (Sections 5 and 6).
- A new canonical form for weighted automata arising from the SVD of its corresponding Hankel matrix (Section 4).
- An algorithm to construct approximate minimizations of given weighted automata by truncating the canonical form (Section 7).

Minimization of automata has been a major subject since the 1950s, starting with the now classical work of the pioneers of automata theory. Recently there has been activity on novel algorithms for minimization based on duality (Bezhaniashvili et al. 2012; Bonchi et al. 2014) which are ultimately based on a remarkable algorithm due to Brzozowski from the 1960s (Brzozowski 1962). The general co-algebraic framework permits one to generalize Brzozowski's algorithm to other classes of automata like weighted automata.

Weighted automata are also used in a variety of practical settings, such as machine learning where they are used to represent predictive models for time series data and text. For example, weighted automata are commonly used for pattern recognition in sequences occurring in speech

[†]This work was completed while the authors were at Lancaster University.

recognition (Mohri et al. 2008), image compression (Albert and Kari 2009), natural language processing (Knight and May 2009), model checking (Baier et al. 2009), and machine translation (de Gispert et al. 2010). The machine learning motivations of our work are discussed at greater length in Section 8, as they are the main impetus for the present work. There has also been interest in this type of representations in the general theory of quantitative systems, including concurrency theory (Boreale 2009) and semantics (Bonchi et al. 2012).

While the detailed discussion of the machine learning motivations appears in the related work section, it is appropriate to make a few points at the outset. First, the formalism of weighted finite automata (WFA) serves as a unifying formalism; examples of models that are subsumed include hidden Markov models (HMM), predictive state representations (PSR), and probabilistic automata of various kinds. Second, in many learning scenarios one has to make a guess of the number of states in advance of the learning process; the resulting algorithm is then trying to construct as best it can a minimal realization within the given constraint. Thus our work gives a general framework for the analysis of these types of learning scenarios.

The present paper extends and improves the results of our previous work (Balle et al. 2015), where the singular value automaton was defined for the first time. The contents of this paper are organized as follows. Section 2 defines the notation that will be used throughout the paper and reviews a series of well-known results that will be needed. Section 3 develops some basic results on analytic properties of rational series computed by weighted automata. Section 4 establishes the existence of the singular value automaton, a canonical form for weighted automata computing square-summable rational series. Section 5 proves some fundamental equations satisfied by singular value automata and provides an algorithm for computing the canonical form. Section 6 shows how to implement the algorithms from the previous section using two different methods for computing the Gramian matrices associated with a factorization of the Hankel matrix. Section 7 describes the main application of singular value automata to approximate minimization and proves an important approximation result. Section 8 discusses related work in approximate minimization, spectral learning of weighted automata, and the theory of linear dynamical systems. We conclude with Section 9, where we point out interesting future research directions.

2. Notation and preliminaries

Given a positive integer d , we denote $[d] = \{1, \dots, d\}$. We use \mathbb{R} to denote the field of real numbers, and $\mathbb{N} = \{0, 1, \dots\}$ for the commutative monoid of natural numbers. In this section we present notation and preliminary results about linear algebra, functional analysis, and weighted automata that will be used throughout the paper. We state all our results in terms of real numbers because this is the most common choice in the literature on weighted automata, but all our results remain true (and the proofs are virtually the same) if one considers automata with weights in the field of complex numbers \mathbb{C} .

2.1 Linear algebra and functional analysis

We use bold letters to denote vectors $\mathbf{v} \in \mathbb{R}^d$ and matrices $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$. Unless explicitly stated, all vectors are column vectors. We write \mathbf{I} for the identity matrix, $\text{diag}(a_1, \dots, a_n)$ for a diagonal matrix with a_1, \dots, a_n in the diagonal, and $\text{diag}(\mathbf{M}_1, \dots, \mathbf{M}_n)$ for the block-diagonal matrix containing the square matrices \mathbf{M}_i along the diagonal. The i th coordinate vector $(0, \dots, 0, 1, 0, \dots, 0)^\top$ is denoted by \mathbf{e}_i and the all ones vector $(1, \dots, 1)^\top$ is denoted by $\mathbf{1}$. For a matrix $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$, $i \in [d_1]$, and $j \in [d_2]$, we use $\mathbf{M}(i, \cdot)$ and $\mathbf{M}(\cdot, j)$ to denote the i th row and the j th column of \mathbf{M} respectively. Given a matrix $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$ we denote by $\text{vec}(\mathbf{M}) \in \mathbb{R}^{d_1 \cdot d_2}$ the vector

obtained by concatenating the columns of \mathbf{M} so that $\text{vec}(\mathbf{M})((i-1)d_2 + j) = \mathbf{M}(i, j)$. Given two matrices $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$ and $\mathbf{M}' \in \mathbb{R}^{d'_1 \times d'_2}$ we denote their tensor (or Kronecker) product by $\mathbf{M} \otimes \mathbf{M}' \in \mathbb{R}^{d_1 d'_1 \times d_2 d'_2}$, with entries given by $(\mathbf{M} \otimes \mathbf{M}')((i-1)d'_1 + i', (j-1)d'_2 + j') = \mathbf{M}(i, j)\mathbf{M}'(i', j')$, where $i \in [d_1]$, $j \in [d_2]$, $i' \in [d'_1]$, and $j' \in [d'_2]$. For simplicity, we will sometimes write $\mathbf{M}^{\otimes 2} = \mathbf{M} \otimes \mathbf{M}$, and similarly for vectors. A *rank factorization* of a rank n matrix $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$ is an expression of the form $\mathbf{M} = \mathbf{Q}\mathbf{R}$ where $\mathbf{Q} \in \mathbb{R}^{d_1 \times n}$ and $\mathbf{R} \in \mathbb{R}^{n \times d_2}$ are full-rank matrices; i.e., $\text{rank}(\mathbf{Q}) = \text{rank}(\mathbf{R}) = \text{rank}(\mathbf{M}) = n$. When \mathbf{Q} is a square invertible matrix, we use the shorthand notation $\mathbf{Q}^{-\top}$ to denote the transpose of its inverse $(\mathbf{Q}^{-1})^\top$.

Given a matrix $\mathbf{M} \in \mathbb{R}^{d_1 \times d_2}$ of rank n , its *SVD*¹ is a decomposition of the form $\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ where $\mathbf{U} \in \mathbb{R}^{d_1 \times n}$, $\mathbf{D} \in \mathbb{R}^{n \times n}$, and $\mathbf{V} \in \mathbb{R}^{d_2 \times n}$ are such that: $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}$, and $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_n)$ with $\sigma_1 \geq \dots \geq \sigma_n > 0$. The columns of \mathbf{U} and \mathbf{V} are thus orthonormal and are called left and right *singular vectors* respectively, and the σ_i are its *singular values*. The SVD is unique (up to sign changes in associate singular vectors) whenever all inequalities between singular values are strict. The Moore–Penrose pseudo-inverse of \mathbf{M} is denoted by \mathbf{M}^\dagger and is the *unique* matrix (if it exists) such that $\mathbf{M}\mathbf{M}^\dagger\mathbf{M} = \mathbf{M}$ and $\mathbf{M}^\dagger\mathbf{M}\mathbf{M}^\dagger = \mathbf{M}^\dagger$. It can be computed from the SVD $\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ as $\mathbf{M}^\dagger = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^\top$.

For $1 \leq p \leq \infty$ we will write $\|\mathbf{v}\|_p$ for the ℓ^p norm of vector \mathbf{v} . The corresponding *induced norm* on matrices is $\|\mathbf{M}\|_p = \sup_{\|\mathbf{v}\|_p=1} \|\mathbf{M}\mathbf{v}\|_p$. We recall the following characterizations for induced norms with $p \in \{1, \infty\}$: $\|\mathbf{M}\|_1 = \max_j \sum_i |\mathbf{M}(i, j)|$ and $\|\mathbf{M}\|_\infty = \max_i \sum_j |\mathbf{M}(i, j)|$. In addition to the induced norms, we will also use Schatten norms. If \mathbf{M} is a rank- n matrix with singular values $\mathbf{s} = (\sigma_1, \dots, \sigma_n)$, the *Schatten p -norm* of \mathbf{M} is given by $\|\mathbf{M}\|_{S,p} = \|\mathbf{s}\|_p$. Most of these norms have given names: $\|\cdot\|_2 = \|\cdot\|_{S,\infty} = \|\cdot\|_{\text{op}}$ is the *operator (or spectral) norm*; $\|\cdot\|_{S,2} = \|\cdot\|_F$ is the *Frobenius norm*; and $\|\cdot\|_{S,1} = \|\cdot\|_{\text{tr}}$ is the *trace (or nuclear) norm*. For a matrix \mathbf{M} the *spectral radius* is the largest modulus $\rho(\mathbf{M}) = \max_i |\lambda_i(\mathbf{M})|$ among the eigenvalues $\lambda_i(\mathbf{M})$ of \mathbf{M} . For a square matrix \mathbf{M} , the series $\sum_{k \geq 0} \mathbf{M}^k$ converges if and only if $\rho(\mathbf{M}) < 1$, in which case the sum yields $(\mathbf{I} - \mathbf{M})^{-1}$.

Recall that if a square matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ is symmetric, then all its eigenvalues are real. A symmetric matrix \mathbf{M} is *positive semi-definite* when all its eigenvalues are nonnegative; we denote this fact by writing $\mathbf{M} \geq \mathbf{0}$, where $\mathbf{0}$ is a zero $d \times d$ matrix. The Loewner partial ordering on the set of all $d \times d$ matrices is obtained by defining $\mathbf{M}_1 \geq \mathbf{M}_2$ to mean $\mathbf{M}_1 - \mathbf{M}_2 \geq \mathbf{0}$. The fact that this gives a partial order follows from the fact that the positive semi-definite operators form a cone. In particular, $\mathbf{M}_1 \geq \mathbf{M}_2$ implies the trace inequality $\text{Tr}(\mathbf{M}_1) \geq \text{Tr}(\mathbf{M}_2)$.

Sometimes we will name the columns and rows of a matrix using ordered index sets \mathcal{I} and \mathcal{J} . In this case we will write $\mathbf{M} \in \mathbb{R}^{\mathcal{I} \times \mathcal{J}}$ to denote a matrix of size $|\mathcal{I}| \times |\mathcal{J}|$ with rows indexed by \mathcal{I} and columns indexed by \mathcal{J} .

Recall that a *Banach space* is a complete normed vector space $(X, \|\cdot\|_X)$. A *Hilbert space* is a Banach space $(X, \|\cdot\|_X)$ where the norm arises from an inner product: $\|x\|_X^2 = \langle x, x \rangle_X$. A Hilbert space is separable if it admits a countable orthonormal basis. The *operator norm* of a linear operator $T : X \rightarrow Y$ between two Banach spaces is given by $\|T\|_{\text{op}} = \sup_{\|x\|_X \leq 1} \|Tx\|_Y$. The operator is *bounded* (and continuous) if $\|T\|_{\text{op}}$ is finite. An operator $T : X \rightarrow Y$ is *compact* if the closure in the topology of Y of the image under T of the unit ball in X is a compact set in Y . A sufficient condition for compactness is to be a bounded finite-rank operator.

Our main interest in compact operators is motivated by the existence of a decomposition equivalent to SVD for compact operators in Hilbert spaces. Note that for a bounded operator $T : X \rightarrow Y$ between separable Hilbert spaces, it is possible to choose countable orthonormal basis $F = (f_j)_{j \in \mathcal{J}}$ and $E = (e_i)_{i \in \mathcal{I}}$ for X and Y , respectively, and write down an infinite matrix $\mathbf{T} \in \mathbb{R}^{\mathcal{I} \times \mathcal{J}}$ for T with entries given by $\mathbf{T}(i, j) = \langle e_i, T f_j \rangle_Y$. In the case of finite-rank bounded operators, the *Hilbert–Schmidt decomposition* (Zhu 1990) provides a decomposition for the infinite matrix associated

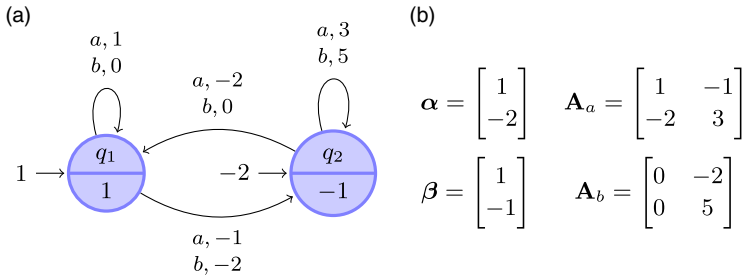


Figure 1. (a) Example of WFA A with two states. Within each circle we denote the state number q_i and the corresponding final weight. The initial weights are denoted using arrows pointing to each state, and the transition weights are given by arrows between states. For example, $f_A(ba) = 1 \times (-2) \times 3 \times (-1) + 1 \times (-2) \times (-2) \times 5 \times 3 \times (-1) + (-2) \times 5 \times (-2) \times 1 = 60$. (b) Corresponding initial vector α , final vector β , and transition matrices A_a and A_b .

with an operator analogous to the compact SVD decomposition for finite matrices. In particular, if T has rank n , then the decomposition theorem yields singular values $\sigma_1 \geq \dots \geq \sigma_n > 0$ and singular vectors $v_i \in X$ and $u_i \in Y$ for $i \in [n]$ such that for all $x \in X$ we have

$$Tx = \sum_{i=1}^n \sigma_i \langle v_i, x \rangle_X u_i. \tag{1}$$

By writing the singular vectors u_i and v_i in terms of the bases E and F , we can write this decomposition as $T = UDV^T$ with $U \in \mathbb{R}^{\mathcal{L} \times n}$ and $V \in \mathbb{R}^{\mathcal{J} \times n}$ satisfying the same properties as the SVD for finite matrices.

2.2 Weighted automata and rational functions

Let Σ be a fixed finite alphabet with $|\Sigma| < \infty$ symbols, and Σ^* the set of all finite strings with symbols in Σ . We use ε to denote the empty string. Given two strings $p, s \in \Sigma^*$, we write $w = ps$ for their concatenation, in which case we say that p is a prefix of w and s is a suffix of w . We denote by $|w|$ the length (number of symbols) in a string $w \in \Sigma^*$. Given a set of strings $X \subseteq \Sigma^*$ and a function $f : \Sigma^* \rightarrow \mathbb{R}$, we denote by $f(X)$ the summation $\sum_{x \in X} f(x)$ if defined. For example, we will write $f(\Sigma^t) = \sum_{|x|=t} f(x)$ for any $t \geq 0$. The notation $\Sigma^{<t}$ (resp. $\Sigma^{\leq t}$) denotes all string of length less than (resp. at most) t . As customary, we use Σ^+ to denote the set of nonempty strings.

Now we introduce our notation for weighted automata. We want to note that we will not work with weights in arbitrary semi-rings; this paper only considers automata with real weights and the usual addition and multiplication operations. In addition, instead of resorting to the usual description of automata as directed graphs with labeled nodes and edges, we will use a linear-algebraic representation which is more convenient for our purposes. Thus, a WFA of dimension n over Σ is a tuple $A = \langle \alpha, \beta, \{A_a\}_{a \in \Sigma} \rangle$ where $\alpha \in \mathbb{R}^n$ is the vector of *initial weights*, $\beta \in \mathbb{R}^n$ is the vector of *final weights*, and for each symbol $a \in \Sigma$ the matrix $A_a \in \mathbb{R}^{n \times n}$ contains the *transition weights* associated with a . An example is provided in Figure 1. Note that in this representation a fixed initial state is given by α (as opposed to formalisms that only specify a transition structure), and the transition endomorphisms A_a and the final linear form β are given in a fixed basis on \mathbb{R}^n (as opposed to abstract descriptions where these objects are represented as basis-independent elements objects on an abstract n -dimensional vector space).

We will use $\dim(A)$ to denote the dimension of a WFA, to which we sometimes also refer to as the number of states in the WFA. The state-space of a WFA of dimension n is identified

with the integer set $[n]$. Every WFA A realizes a function $f_A : \Sigma^* \rightarrow \mathbb{R}$ which, given a string $x = x_1 \cdots x_t \in \Sigma^*$, produces

$$f_A(x) = \alpha^\top \mathbf{A}_{x_1} \cdots \mathbf{A}_{x_t} \beta = \alpha^\top \mathbf{A}_x \beta,$$

where we defined the shorthand notation $\mathbf{A}_x = \mathbf{A}_{x_1} \cdots \mathbf{A}_{x_t}$ that will be used throughout the paper. In terms of the graphical description of A , the value $f_A(x)$ can be interpreted as the sum of the weights of all paths labeled by x from an initial to a final state, where the weight of a path is the product of the initial weight, the corresponding transition weights, and the final weight:

$$f_A(x) = \sum_{(q_0, \dots, q_t) \in [n]^{t+1}} \alpha(q_0) \left(\prod_{i=1}^t \mathbf{A}_{x_i}(q_{i-1}, q_i) \right) \beta(q_t),$$

where $t = |x|$. A function $f : \Sigma^* \rightarrow \mathbb{R}$ is called *rational*² if there exists a WFA A such that $f = f_A$. The *rank* of a rational function f is the dimension of the smallest WFA realizing f . We say that a WFA A is *minimal* if $\dim(A) = \text{rank}(f_A)$.

Hankel matrices provide a powerful characterization of rationality that will be heavily used in the sequel. Let $\mathbf{H} \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ be an infinite matrix whose rows and columns are indexed by strings. We say that \mathbf{H} is *Hankel*³ if for all strings $p, p', s, s' \in \Sigma^*$ such that $ps = p's'$ we have $\mathbf{H}(p, s) = \mathbf{H}(p', s')$. Given a function $f : \Sigma^* \rightarrow \mathbb{R}$ we can associate with it a Hankel matrix $\mathbf{H}_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ with entries $\mathbf{H}_f(p, s) = f(ps)$. Conversely, given a matrix $\mathbf{H} \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ with the Hankel property, there exists a unique function $f : \Sigma^* \rightarrow \mathbb{R}$ such that $\mathbf{H}_f = \mathbf{H}$. The following well-known theorem characterizes all Hankel matrices of finite rank.

Theorem 2.1. (Berstel and Reutenauer 2011). *For any function $f : \Sigma^* \rightarrow \mathbb{R}$, the Hankel matrix \mathbf{H}_f has finite rank n if and only if f is rational with $\text{rank}(f) = n$. In other words, $\text{rank}(f) = \text{rank}(\mathbf{H}_f)$ for any function $f : \Sigma^* \rightarrow \mathbb{R}$.*

2.3 Probabilistic automata

Probabilistic automata will be used as a recurring example throughout the paper. Here we introduce the main definitions and stress some key differences arising from subtle changes in the definition that can make a difference in terms of the analytic properties of this kind of automata. Generally speaking, a probabilistic automaton is a WFA A whose weights have a probabilistic interpretation, and such that the values $f_A(x)$ of the function computed by A represent the likelihood of an event associated with string x .

A *generative probabilistic automaton* (GPA) is a WFA A such that the function f_A computes a probability distribution on Σ^* . That is, we have $f_A(x) \geq 0$ and $\sum_{x \in \Sigma^*} f_A(x) = 1$. In addition, we say a GPA $A = \langle \alpha, \beta, \{\mathbf{A}_a\} \rangle$ is *proper* (pGPA) if its weights have a probabilistic interpretation, i.e.

- (1) Initial weights represent a probability distribution over states: $\alpha \geq 0$ and $\alpha^\top \mathbf{1} = 1$.
- (2) Transition weights and final weights represent probabilities of emitting a symbol and transitioning to a next state or terminating: $\mathbf{A}_\sigma \geq 0^4$, $\beta \geq 0$, and $\sum_{\sigma \in \Sigma} \mathbf{A}_\sigma \mathbf{1} + \beta = \mathbf{1}$.

An example is provided in Figure 2. It is shown in Denis and Esposito (2008) that not all GPA are pGPA, and that there exists probability distributions on Σ^* that cannot be computed by any pGPA.

A *dynamic probabilistic automaton* (DPA) is a WFA A defining a probability distribution D_A over streams in Σ^ω and such that the function f_A on finite strings computes the probability under D_A of cones of the form $x\Sigma^\omega$ for $x \in \Sigma^*$. That is, we have the semantics $f_A(x) = \mathbb{P}_{D_A}[x\Sigma^\omega]$, which

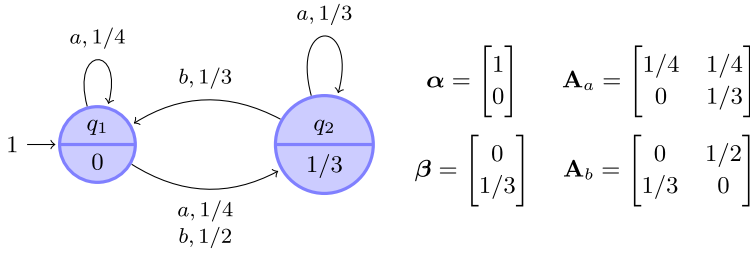


Figure 2. Example of pGPA with two states.

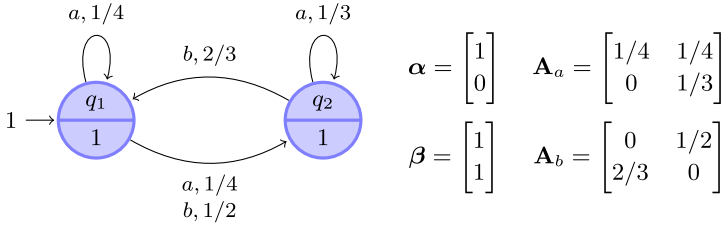


Figure 3. Example of det-free pDPA with two states.

implies that $f_A(\Sigma^t) = 1$ for all $t \geq 0$. Again, we say that a DPA $A = \langle \alpha, \beta, \{A_\sigma\} \rangle$ is *proper* (pDPA) if its weights have a probabilistic interpretation as follows:

- (1) Initial weights represent a probability distribution over states: $\alpha \geq 0$ and $\alpha^T \mathbf{1} = 1$.
- (2) Final weights are all equal to one: $\beta = \mathbf{1}$.
- (3) Transition weights represent probabilities of emitting a symbol and transitioning to a next state: $A_\sigma \geq 0$ and $\sum_{\sigma \in \Sigma} A_\sigma \mathbf{1} = \mathbf{1}$.

An example is provided in Figure 3. As with GPA, there exist improper DPA, and distributions D_A on Σ^ω that cannot be computed by any pDPA (Denis and Esposito 2008). An important subclass of pDPA is those for which there is no state with deterministic emissions. A pDPA $A = \langle \alpha, \beta, \{A_\sigma\} \rangle$ is *det-free* if we have $\|A_\sigma\|_\infty < 1$ for each $\sigma \in \Sigma$. Note that if A has n states and there exists σ such that $\|A_\sigma\|_\infty = 1$, then there exists $i \in [n]$ such that $A_\sigma \mathbf{1}(i) = 1$ and therefore from state i the automaton A always emits symbol σ .

3. Banach and Hilbert spaces of rational functions

In the literature on formal language theory, functions $f : \Sigma^* \rightarrow \mathbb{R}$ are sometimes regarded as weighted languages and weighted automata computing them as linear representations. From an algebraic point of view, one can identify a weighted language f with an element of the vector space \mathbb{R}^{Σ^*} . This vector space contains several subspaces that play an important role in the theory developed in this paper. Furthermore, some of these spaces can be endowed with additional operations and norms, yielding a wide variety of algebraic and analytic structures. To the best of our knowledge, analytic properties of these spaces have never been systematically studied before in the automata theory literature. This section introduces the basic facts and definitions that will be needed in the rest of the paper. We also take this as an opportunity to prove basic facts about these spaces and pinpoint ideas that need to be developed further. Overall, we hope this provides the foundations for a much needed *analytic theory of rational functions*.

A fundamental linear subspace of \mathbb{R}^{Σ^*} is the space of all rational functions, which we denote by $\mathcal{R}(\Sigma)$. That $\mathcal{R}(\Sigma)$ is a linear subspace follows from the simple observations that if $f, g \in \mathcal{R}(\Sigma)$

and $c \in \mathbb{R}$, then cf and $f + g$ are both rational (Berstel and Reutenauer 2011). An important subspace of $\mathcal{R}(\Sigma)$ is the space of all $f \in \mathbb{R}^{\Sigma^*}$ with finite support, which we denote by $\mathcal{C}_{00}(\Sigma)$. That is, $f \in \mathcal{C}_{00}(\Sigma)$ if and only if $|\text{supp}(f)| < \infty$, where $\text{supp}(f) = \{x : f(x) \neq 0\}$ is the support of f . It is immediate from this definition that $\mathcal{C}_{00}(\Sigma)$ is a linear subspace of \mathbb{R}^{Σ^*} . The containment $\mathcal{C}_{00}(\Sigma) \subset \mathcal{R}(\Sigma)$ follows from observing that every function with finite support is rational (Berstel and Reutenauer 2011).

Another important family of subspaces of \mathbb{R}^{Σ^*} are the ones containing all functions with finite p -norm for some $1 \leq p \leq \infty$, which is given by $\|f\|_p^p = \sum_{x \in \Sigma^*} |f(x)|^p$ for finite p , and $\|f\|_\infty = \sup_{x \in \Sigma^*} |f(x)|$; we denote this space by $\ell^p(\Sigma)$. Note that these are Banach spaces, and as with the usual theory of Banach spaces over sequences we have $\ell^p(\Sigma) \subset \ell^q(\Sigma)$ for $p < q$. Of these, $\ell^2(\Sigma)$ can be endowed with the structure of a separable Hilbert space with the inner product $\langle f, g \rangle = \sum_{x \in \Sigma^*} f(x)g(x)$. Recall that in this case we have the *Cauchy-Schwarz inequality* $\langle f, g \rangle^2 \leq \|f\|_2^2 \|g\|_2^2$. In addition, we have its generalization, *Hölder's inequality*: given $f \in \ell^p(\Sigma)$ and $g \in \ell^q(\Sigma)$ with $p^{-1} + q^{-1} \leq 1$, then $\|f \cdot g\|_1 \leq \|f\|_p \|g\|_q$, where $(f \cdot g)(x) = f(x)g(x)$ is the *Hadamard product* between two languages.

By intersecting any of the previous subspaces with $\mathcal{R}(\Sigma)$ one obtains $\ell^p_{\mathcal{R}}(\Sigma) = \mathcal{R}(\Sigma) \cap \ell^p(\Sigma)$, the normed vector space containing all rational functions with finite p -norm. In most cases the alphabet Σ will be clear from the context, and we will just write $\mathcal{R}, \mathcal{C}_{00}, \ell^p$, and $\ell^p_{\mathcal{R}}$. It is important to note that although the ℓ^p spaces can be endowed with the structure of a Banach or Hilbert space, the $\ell^p_{\mathcal{R}}$ spaces cannot, because they are not complete; i.e., it is possible to find sequences of functions in $\ell^p_{\mathcal{R}}$ whose limit in the topology induced by the corresponding norm is not rational. For example, consider the function given by $f(x) = (k + 1)^{-|x|}$ if x is a palindrome and $f(x) = 0$ otherwise. Since $\text{supp}(f)$ is the set of all palindromes, then f is not rational (Berstel and Reutenauer 2011), and in addition $\|f\|_1 < \infty$ by construction. Thus, we have $f \in \ell^p \setminus \mathcal{R}$ for any $1 \leq p \leq \infty$. Now, for any $l \geq 0$ let $f_l(x) = f(x)$ if $|x| \leq l$ and $f_l(x) = 0$ otherwise. Since f_l has finite support for every $l \geq 0$, we have $f_l \in \ell^p_{\mathcal{R}}$. Finally, it is easy to check that $\lim_{l \rightarrow \infty} \|f - f_l\|_p = 0$, implying that we have a sequence of functions in $\ell^p_{\mathcal{R}}$ converging to a nonrational function. Therefore, none of the $\ell^p_{\mathcal{R}}$ spaces is complete. Nonetheless, the following result shows that all ℓ^p spaces with $1 \leq p < \infty$ can be obtained as the completion of their corresponding $\ell^p_{\mathcal{R}}$ space.

Theorem 3.1. *For any $1 \leq p < \infty$, the Banach space ℓ^p coincides with the completion of $\ell^p_{\mathcal{R}}$ with respect to $\|\cdot\|_p$.*

Proof. Fix $1 \leq p < \infty$. Since $\mathcal{C}_{00} \subset \ell^p_{\mathcal{R}}$, it suffices to show that \mathcal{C}_{00} is dense in ℓ^p with respect to the topology induced by $\|\cdot\|_p$. Let $f \in \ell^p$ and for any $l \geq 0$ define $f_l(x) = f(x)$ if $|x| \leq l$ and $f_l(x) = 0$ otherwise. Clearly we have $f_l \in \mathcal{C}_{00}$ by construction. To see that $f_l \rightarrow f$ in the topology of ℓ^p as $l \rightarrow \infty$ we write $s_l = \|f_l - f\|_p^p = \sum_{|x| > l} |f(x)|^p$ and observe that we must have $s_l \rightarrow 0$. Otherwise we would have $\lim_{l \rightarrow \infty} \sum_{|x| > l} |f(x)|^p > 0$, which is a contradiction with $\|f\|_p^p = \sum_{x \in \Sigma^*} |f(x)|^p < \infty$. □

3.1 Bounded Hankel operators

Recall that Theorem 2.1 gives a characterization of the functions $f : \Sigma^* \rightarrow \mathbb{R}$ which have a Hankel matrix with finite rank. Using the concepts introduced above, we can interpret the Hankel matrix as an operator on Hilbert spaces and ask when this operator satisfies some nice properties. The main result of this section is a characterization of the rational functions whose Hankel matrix induces a bounded operator on ℓ^2 .

Recall that a matrix $\mathbf{T} \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ can be interpreted as the expression of a (possibly unbounded) linear operator $T : \ell^2 \rightarrow \ell^2$ in terms of the canonical basis $(\mathbf{e}_x)_{x \in \Sigma^*}$. In the case of a Hankel matrix \mathbf{H}_f , we can see it is associated with an operator H_f corresponding to the operation $g \mapsto H_f g$ with $(H_f g)(x) = \sum_y f(xy)g(y)$ (assuming the series converges). An operator $T : \ell^2 \rightarrow \ell^2$ is bounded if $\|T\|_{\text{op}} = \sup_{\|g\|_2 \leq 1} \|Tg\|_2 < \infty$. Not all Hankel operators H_f are bounded, but we shall give a necessary and sufficient condition for H_f to be bounded when f is rational. We start with the following technical lemma.

Lemma 3.2. *Let $A = \langle \alpha, \beta, \{\mathbf{A}_a\} \rangle$ be a WFA such that $f_A(x) \geq 0$ for all $x \in \Sigma^*$. Define $\mathbf{A} = \sum_{\sigma} \mathbf{A}_{\sigma}$ and let $\rho = \rho(\mathbf{A})$ be its spectral radius. Then the following hold:*

- (1) *If A is minimal and $f_A \in \ell^1_{\mathcal{R}}$, then $\rho < 1$.*
- (2) *If $\rho < 1$, then $f_A \in \ell^1_{\mathcal{R}}$.*

Proof. We start by recalling that if $A = \langle \alpha, \beta, \{\mathbf{A}_a\} \rangle$ is a minimal WFA with n states, then there exist sets of prefixes $\mathcal{P} = \{p_1, \dots, p_n\}$ and suffixes $\mathcal{S} = \{s_1, \dots, s_n\}$ such that the sets of vectors $\{\alpha^{\top} \mathbf{A}_{p_1}, \dots, \alpha^{\top} \mathbf{A}_{p_n}\}$ and $\{\mathbf{A}_{s_1} \beta, \dots, \mathbf{A}_{s_n} \beta\}$ define two bases for \mathbb{R}^n (Berstel and Reutenauer 2011). For convenience we will write $\alpha^{\top}_{p_i} = \alpha^{\top} \mathbf{A}_{p_i}$ and $\beta_{s_j} = \mathbf{A}_{s_j} \beta$.

Now assume $f_A \in \ell^1_{\mathcal{R}}$ and suppose λ is an arbitrary eigenvalue of \mathbf{A} . We need to show that $|\lambda| < 1$. Let \mathbf{v} be any eigenvector with eigenvalue λ and suppose $\|\mathbf{v}\|_2 = 1$. Using the basis given by \mathcal{P} and \mathcal{S} , we can find coefficients such that $\mathbf{v} = \sum_{i \in [n]} \gamma_i \alpha_{p_i} = \sum_{j \in [n]} \delta_j \beta_{s_j}$. For any $k \geq 0$, we can now write the following:

$$\begin{aligned} |\lambda|^k &= |\lambda^k \mathbf{v}^{\top} \mathbf{v}| = |\mathbf{v}^{\top} (\lambda^k \mathbf{v})| = |\mathbf{v}^{\top} \mathbf{A}^k \mathbf{v}| \\ &= \left| \left(\sum_i \gamma_i \alpha_{p_i}^{\top} \right) \left(\sum_{\sigma} \mathbf{A}_{\sigma} \right)^k \left(\sum_j \delta_j \beta_{s_j} \right) \right| \\ &\leq \sum_{ij} |\gamma_i| |\delta_j| \sum_{x \in p_i \Sigma^k s_j} |f_A(x)|. \end{aligned}$$

Since we have $f_A \in \ell^1_{\mathcal{R}}$ by hypothesis, for fixed i and j we have $\sum_{k \geq 0} \sum_{x \in p_i \Sigma^k s_j} |f_A(x)| \leq \sum_{x \in \Sigma^*} |f_A(x)| < \infty$. Therefore we can conclude that $\sum_{k \geq 0} |\lambda|^k < \infty$, which necessarily implies $|\lambda| < 1$.

To obtain the converse suppose $\rho(\mathbf{A}) < 1$ and note that because f_A is nonnegative, we have

$$\|f_A\|_1 = \sum_{x \in \Sigma^*} |f(x)| = \sum_{x \in \Sigma^*} f(x) = \sum_{k \geq 0} \alpha^{\top} \mathbf{A}^k \beta < \infty. \tag{2}$$

Note that this implication does not require the minimality of A . □

The following theorem is the main result of this section.

Theorem 3.3. *Let $f : \Sigma^* \rightarrow \mathbb{R}$ be a rational function. The Hankel operator H_f is bounded if and only if $f \in \ell^2_{\mathcal{R}}$.*

Proof. It is easy to see that the membership $f \in \ell^2$ is a necessary condition for the boundedness of \mathbf{H}_f . Indeed, by noting that f appears as the first column of \mathbf{H}_f we have $f = H_f e_{\varepsilon}$, and since $\|e_{\varepsilon}\|_2 = 1$ we have $\|f\|_2 = \|H_f e_{\varepsilon}\|_2 \leq \|H_f\|_{\text{op}}$.

Next we prove sufficiency. Let $g \in \ell^2$ with $\|g\|_2 = 1$ and for any $x \in \Sigma^*$ define the function $f_x(y) = f(xy)$. With this notation we can write

$$\begin{aligned} \|H_f g\|_2^2 &= \sum_{x \in \Sigma^*} \left(\sum_{y \in \Sigma^*} f(xy)g(y) \right)^2 = \sum_{x \in \Sigma^*} \langle f_x, g \rangle^2 \\ &\leq \|g\|_2^2 \sum_{x \in \Sigma^*} \|f_x\|_2^2 = \sum_{x \in \Sigma^*} \sum_{y \in \Sigma^*} f(xy)^2 \\ &= \sum_{z \in \Sigma^*} (1 + |z|)f(z)^2, \end{aligned} \tag{3}$$

where we used Cauchy–Schwarz’s inequality, and the fact that a string z can be written as $z = xy$ in $1 + |z|$ different ways.

Recall that $f \in \ell^2_{\mathcal{R}}$ implies $f^2 \in \ell^1_{\mathcal{R}}$. Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be a minimal WFA for f^2 and write $\mathbf{A} = \sum_{\sigma} A_{\sigma}$. Note we have $\rho = \rho(\mathbf{A}) < 1$ by Lemma 3.2. Suppose $\mathbf{A} = \mathbf{W}\mathbf{J}\mathbf{W}^{-1}$ is the Jordan canonical form of \mathbf{A} , and let m denote the maximum algebraic multiplicity of any eigenvalue of \mathbf{A} . By computing the k th power of the largest Jordan block

$$\mathbf{J}_{\lambda} = \begin{bmatrix} \lambda & 1 & 0 & \cdots & 0 \\ 0 & \lambda & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \lambda & 1 \\ 0 & 0 & 0 & 0 & \lambda \end{bmatrix} \in \mathbb{R}^{m' \times m'}$$

associated with the maximal eigenvalue $|\lambda| = \rho$ (with $m' \leq m$), one can see there exists a constant $c > 0$ such that the following holds for all $k \geq 0$:

$$\sum_{x \in \Sigma^k} f(x)^2 = \alpha^{\top} \mathbf{A}^k \beta = \alpha^{\top} \mathbf{W}\mathbf{J}^k \mathbf{W}^{-1} \beta \leq ck^{m-1} \rho^k.$$

This is a standard calculation in the analysis of nonreversible Markov chains; see Fact 3 in Rosenthal (1995) for more details. Now we use that $\rho < 1$, in which case this bound yields

$$\sum_{z \in \Sigma^*} |z|f(z)^2 = \sum_{k \geq 0} k \sum_{z \in \Sigma^k} f(z)^2 \leq c \sum_{k \geq 0} k^m \rho^k < \infty.$$

Plugging this into (3) we can conclude that $\|H_f g\|_2$ is finite and therefore H_f is bounded. □

4. The singular value automaton

The central object of study in this paper is the singular value automaton (SVA). Essentially, this is a canonical form for weighted automata which is tightly connected to the SVD of the corresponding Hankel matrix. We will start this section by establishing some fundamental preliminary results on the relation between minimal WFAs and rank factorizations of Hankel matrices. By assuming that one such Hankel matrix admits an SVD, the relation above will lead us directly to the definition of singular value automaton. We then proceed to explore necessary conditions for the existence of SVA. These will essentially say that only rational functions in $\ell^2_{\mathcal{R}}$ admit a singular value automaton, provide some easily testable conditions, and guarantee the existence of an SVA for a large class of probabilistic automata.

4.1 Correspondence between minimal WFA and rank factorizations

An important operation on WFA is conjugation by an invertible matrix. Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be a WFA of dimension n and suppose $Q \in \mathbb{R}^{n \times n}$ is invertible. Then we can define the *conjugate* of A by Q as:

$$A' = A^Q = Q^{-1}AQ = \langle Q^T \alpha, Q^{-1} \beta, \{Q^{-1}A_a Q\} \rangle. \tag{4}$$

It follows immediately that $f_A = f_{A'}$ since, at every step in the computation of $f_{A'}(x)$, the products QQ^{-1} vanish. This means that the function computed by a WFA is invariant under conjugation, and that given a rational function f , there exist infinitely many WFA realizing f . The following result offers a full characterization of all minimal WFA realizing a particular rational function.

Theorem 4.1. (Berstel and Reutenauer 2011). *If A and B are minimal WFA realizing the same function, then $B = A^Q$ for some invertible Q .*

The goal of this section is to provide a “lifted” version of this result establishing a connection between every pair of rank factorizations of the Hankel matrix H_f , and then show that these rank factorizations are in bijection with all minimal WFA for f . We start by recalling how every minimal WFA realizing f induces a rank factorization for H_f .

Suppose f is a rational function and $A = \langle \alpha, \beta, \{A_a\} \rangle$ is a WFA realizing f . The *forward matrix* of A is defined as the infinite matrix $P_A \in \mathbb{R}^{\Sigma^* \times n}$ with entries given by $P_A(p, :) = \alpha^T A_p$ for any string $p \in \Sigma^*$; sometimes we will refer to the strings indexing rows in a forward matrix as *prefixes*. Similarly, let $S_A \in \mathbb{R}^{\Sigma^* \times n}$ be the *backward matrix* of A given by $S_A(s, :) = (A_s \beta)^T$ for any string $s \in \Sigma^*$; strings indexing rows in a backward matrix are commonly called *suffixes*. Now note that for every $p, s \in \Sigma^*$ we have

$$H_f(p, s) = f(ps) = (\alpha^T A_p) (A_s \beta) = \sum_{i \in [n]} P_A(p, i) S_A(s, i) = P_A(p, :) S_A^T(:, s). \tag{5}$$

Therefore, we see that the forward and backward matrix of A yield the factorization $H_f = P_A S_A^T$. This is known as the *forward-backward* (FB) factorization of H_f induced by A (Balle et al. 2014a).

Recall that a WFA A with n states is called *reachable* when the space spanned by all the forward vectors has dimension n ; that is:

$$\dim \text{span}\{\alpha^T A_x \mid x \in \Sigma^*\} = \text{rank}(P_A) = n. \tag{6}$$

Similarly, A is called *observable* if the dimension of the space spanned by the backward vectors equals n ; that is:

$$\dim \text{span}\{A_x \beta \mid x \in \Sigma^*\} = \text{rank}(S_A) = n. \tag{7}$$

Note that when A is minimal, the number of columns of the forward and backward matrices equals the rank of H_f , and therefore the FB factorization is a rank factorization. Therefore, it follows from Theorem 2.1 the useful characterization of minimality saying that a WFA A is minimal if and only if it is both reachable and observable.

The following result shows that every rank factorization of H_f is actually an FB factorization. We can understand this result as a refinement of Theorem 2.1 in the sense that given a finite-rank Hankel matrix, it provides a characterization of all its possible rank factorizations.

Proposition 4.2. *Let f be rational and suppose $H_f = PS^T$ is a rank factorization. Then there exists a minimal WFA A realizing f which induces this factorization.*

Proof. Let B be any minimal WFA realizing f and denote $n = \text{rank}(f)$. Then we have two rank factorizations $PS^T = P_B S_B^T$ for the Hankel matrix H_f . Therefore, the columns of P and P_B both

span the same n -dimensional sub-space of \mathbb{R}^{Σ^*} , and there exists a change of basis $\mathbf{Q} \in \mathbb{R}^{n \times n}$ such that $\mathbf{P}_B \mathbf{Q} = \mathbf{P}$. This implies we must also have $\mathbf{S}^\top = \mathbf{Q}^{-1} \mathbf{S}_B^\top$. It follows that $A = B^{\mathbf{Q}}$ is a minimal WFA for f inducing the desired rank factorization. \square

4.2 Definition of singular value automaton

It is well known that the compact SVD of a matrix is a rank-revealing decomposition in the sense that the intermediate dimensions of the decomposition correspond to the rank of the matrix. This decomposition can be used to construct rank factorizations for said matrix. The singular value automaton links this idea with the minimal WFA identified in Proposition 4.2.

Recall that if \mathbf{H}_f is a Hankel matrix of rank n admitting an SVD, then there exists a square matrix $\mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_n) \in \mathbb{R}^{n \times n}$ and two infinite matrices $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{\Sigma^* \times n}$ with orthonormal columns (i.e., $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}$) such that $\mathbf{H}_f = \mathbf{U} \mathbf{D} \mathbf{V}^\top$ with $\mathbf{U}, \mathbf{V} \in \mathbb{R}^{\Sigma^* \times n}$. By splitting this decomposition into two parts, we obtain the rank factorization $\mathbf{H}_f = (\mathbf{U} \mathbf{D}^{1/2})(\mathbf{V} \mathbf{D}^{1/2})^\top$. Thus, whenever \mathbf{H}_f admits an SVD, we can invoke Proposition 4.2 to conclude there exists a minimal WFA realizing f whose induced FB rank factorization coincides with the one we obtained above from SVD. Putting this into a formal statement we get the following theorem.

Theorem 4.3. *Let f be a rational function and suppose \mathbf{H}_f admits a compact SVD $\mathbf{H}_f = \mathbf{U} \mathbf{D} \mathbf{V}^\top$. Then there exists a minimal WFA A for f inducing the rank factorization $\mathbf{H}_f = (\mathbf{U} \mathbf{D}^{1/2})(\mathbf{V} \mathbf{D}^{1/2})^\top$. That is, A is a WFA for f with FB rank factorization given by $\mathbf{P}_A = \mathbf{U} \mathbf{D}^{1/2}$ and $\mathbf{S}_A = \mathbf{V} \mathbf{D}^{1/2}$.*

The WFA given by the above theorem can be considered as a canonical form for a rational function whose Hankel matrix admits an SVD. This is made formal in the following definition. Next section will provide conditions for the existence of this automaton.

Definition 4.4. *Let $f \in \ell^2_{\mathcal{R}}$. A singular value automaton (SVA) for f is a minimal WFA A realizing f such that the FB rank factorization of \mathbf{H}_f induced by A has the form given in Theorem 4.3.*

Note the SVA provided by Theorem 4.3 is unique up to the same conditions in which SVD is unique. In particular, it is easy to verify that if the Hankel singular values of $f \in \ell^2_{\mathcal{R}}$ satisfy the strict inequalities $\sigma_1 > \dots > \sigma_n$, then the transition weights of the SVA A of f are uniquely defined, and the initial and final weights are uniquely defined up to sign changes.

4.3 Rational functions admitting an SVA

By leveraging the fact that every compact operator on a Hilbert space admits an SVD and our Theorem 3.3 characterizing rational functions with bounded Hankel operator, we immediately get a characterization of rational functions admitting an SVA.

Theorem 4.5. *A rational function $f : \Sigma^* \rightarrow \mathbb{R}$ admits an SVA if and only if $f \in \ell^2_{\mathcal{R}}$.*

Proof. Since a finite-rank bounded operator is compact and therefore admits a compact SVD, Theorems 3.3 and 4.3 imply that every $f \in \ell^2_{\mathcal{R}}$ admits an SVA. On the other hand, if a rational function admits an SVA, then its Hankel \mathbf{H}_f matrix admits a compact SVD and therefore H_f is bounded. Applying Theorem 3.3 we see that this implies $f \in \ell^2_{\mathcal{R}}$. \square

In view of this result, when given a rational function as a WFA, one just has to check that the function has finite ℓ^2 norm to ensure the existence of an SVA for that function. A direct way to test this based on Lemma 3.2 is given below.

Theorem 4.6. *Let A be a WFA and let B be a minimization of the automaton $A \otimes A$ computing f_A^2 . Then we have $f_A \in \ell_{\mathcal{R}}^2$ if and only if $\rho(\sum_{\sigma \in \Sigma} \mathbf{B}_{\sigma}) < 1$.*

Proof. Let $\mathbf{B} = \sum_{\sigma \in \Sigma} \mathbf{B}_{\sigma}$. The if part follows from observing that $\rho(\sum_{\sigma \in \Sigma} \mathbf{B}_{\sigma}) < 1$ implies that $\sum_{x \in \Sigma^*} \mathbf{B}_x = \sum_{t \geq 0} \mathbf{B}^t$ converges, and therefore $\|f_A\|_2^2 = \sum_{x \in \Sigma^*} \beta_0^{\top} \mathbf{B}_x \beta_{\infty}$ is finite. The only if part is a direct application of Lemma 3.2. □

The above theorem gives a direct way to check if for a given A we have $f_A \in \ell_{\mathcal{R}}^2$ by using a WFA minimization algorithm and computing the spectral radius of a given matrix. If A has n states, then B can be obtained by minimizing an automaton with n^2 states, which takes time $O(n^6)$ (Berstel and Reutenauer 2011) and yields a WFA B with $n' \leq n^2$ states. Computing the spectral radius of \mathbf{B} takes time $O(n'^3)$ (Trefethen and Bau III 1997), so the overall complexity of testing $f_A \in \ell_{\mathcal{R}}^2$ based on the above theorem is $O(n^6)$. The following theorem gives sufficient conditions for $f_A \in \ell_{\mathcal{R}}^2$, some of which can be checked without the need to run a WFA minimization algorithm.

Theorem 4.7. *Let A be a WFA computing a function f_A . Any of the following conditions implies $f_A \in \ell_{\mathcal{R}}^2$:*

- (1) $f_A \in \ell_{\mathcal{R}}^1$,
- (2) $\rho(\sum_{\sigma} \mathbf{A}_{\sigma} \otimes \mathbf{A}_{\sigma}) < 1$,
- (3) $\|\sum_{\sigma} \mathbf{A}_{\sigma} \otimes \mathbf{A}_{\sigma}\|_p < 1$ for some $1 \leq p \leq \infty$,
- (4) $\|\sum_{\sigma} \mathbf{A}_{\sigma} \mathbf{A}_{\sigma}^{\top}\|_2 < 1$.

Proof. The first item follows from the inclusion $\ell_{\mathcal{R}}^1 \subset \ell_{\mathcal{R}}^2$. To get (2) note that by Lemma 3.2 the condition implies $f_A^2 \in \ell_{\mathcal{R}}^1$ and therefore $f_A \in \ell_{\mathcal{R}}^2$. Condition (3) follows from the property of the spectral radius $\rho(\mathbf{M}) \leq \|\mathbf{M}\|_p$. The last condition follows from the main result in Lototsky (2015) showing that $\rho(\sum_{\sigma} \mathbf{A}_{\sigma} \otimes \mathbf{A}_{\sigma}) \leq \|\sum_{\sigma} \mathbf{A}_{\sigma} \mathbf{A}_{\sigma}^{\top}\|_2$. □

We can use these conditions to identify classes of probabilistic automata that compute functions in $\ell_{\mathcal{R}}^2$ and therefore have an SVA. We will need the following technical lemma.

Lemma 4.8. *The following inequality holds for any set of square matrices $\{\mathbf{A}_1, \dots, \mathbf{A}_m\}$:*

$$\begin{aligned} \left\| \sum_{k \in [m]} \mathbf{A}_k \otimes \mathbf{A}_k \right\|_{\infty} &\leq \|[\mathbf{A}_1 \dots \mathbf{A}_m]\|_{\infty} \left\| \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_m \end{bmatrix} \right\|_{\infty} \\ &= \|[\mathbf{A}_1 \dots \mathbf{A}_m]\|_{\infty} \|[\mathbf{A}_1^{\top} \dots \mathbf{A}_m^{\top}]\|_1. \end{aligned}$$

Proof. Recall that the induced matrix norm with $p = \infty$ is given by $\|\mathbf{M}\|_{\infty} = \max_i \sum_j |\mathbf{M}(i, j)|$. Then the desired inequality can be obtained as follows:

$$\begin{aligned} \left\| \sum_{k \in [m]} \mathbf{A}_k \otimes \mathbf{A}_k \right\|_{\infty} &= \max_{i_1, i_2 \in [n]} \sum_{j_1, j_2=1}^n \left| \sum_k \mathbf{A}_k(i_1, j_1) \mathbf{A}_k(i_2, j_2) \right| \\ &\leq \max_{i_1, i_2 \in [n]} \sum_k \sum_{j_1, j_2=1}^n |\mathbf{A}_k(i_1, j_1)| |\mathbf{A}_k(i_2, j_2)| \end{aligned}$$

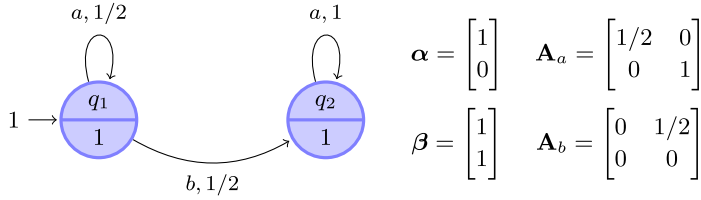


Figure 4. Example of pDPA with two states which is not det-free. Note that $f_A(bc^k) = 1/2$ for all $k \geq 0$ and therefore $f_A \notin \ell^2_{\mathcal{R}}$.

$$\begin{aligned}
 &= \max_{i_1, i_2 \in [n]} \sum_k \left(\sum_{j_1=1}^n |\mathbf{A}_k(i_1, j_1)| \right) \left(\sum_{j_2=1}^n |\mathbf{A}_k(i_2, j_2)| \right) \\
 &\leq \max_{i_1} \sum_k \left(\sum_{j_1=1}^n |\mathbf{A}_k(i_1, j_1)| \right) \left(\max_{i_2} \sum_{j_2=1}^n |\mathbf{A}_k(i_2, j_2)| \right) \\
 &= \max_i \sum_k \|\mathbf{A}_k\|_{\infty} \left(\sum_{j=1}^n |\mathbf{A}_k(i, j)| \right) \\
 &\leq \left(\max_k \|\mathbf{A}_k\|_{\infty} \right) \left(\max_i \sum_k \sum_{j=1}^n |\mathbf{A}_k(i, j)| \right) \\
 &= \left\| \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_m \end{bmatrix} \right\|_{\infty} \|\mathbf{A}_1 \dots \mathbf{A}_m\|_{\infty}.
 \end{aligned}$$

The second equality follows from the duality between the norms $\|\cdot\|_1$ and $\|\cdot\|_{\infty}$. □

Corollary 4.9. *If A is a GPA or a det-free pDPA, then $f_A \in \ell^2_{\mathcal{R}}$.*

Proof. For A GPA it follows directly from Theorem 4.7 by noting that we have $\|f_A\|_1 = 1$. Now suppose $A = \langle \alpha, \beta, \{\mathbf{A}_a\} \rangle$ be a det-free pDPA, so by construction we have $\sum_{a \in \Sigma} \mathbf{A}_a \mathbf{1} = \mathbf{1}$ and $\|\mathbf{A}_a\|_{\infty} < 1$ for all $a \in \Sigma$. Note that the first property implies $\|\mathbf{A}_{a_1} \dots \mathbf{A}_{a_k}\|_{\infty} = 1$ and the second property implies

$$\left\| \begin{bmatrix} \mathbf{A}_{a_1} \\ \vdots \\ \mathbf{A}_{a_k} \end{bmatrix} \right\|_{\infty} < 1. \tag{8}$$

Therefore, using Lemma 4.8 we see that $\|\sum_{a \in \Sigma} \mathbf{A}_a \otimes \mathbf{A}_a\|_{\infty} < 1$ and therefore by (3) in Theorem 4.7 we get $f_A \in \ell^2_{\mathcal{R}}$. □

Note that the det-free condition on pDPA is necessary to ensure $f_A \in \ell^2_{\mathcal{R}}$ as witnessed by the example in Figure 4.

5. Fundamental equations of SVA

In this section we establish two fundamental facts about SVA that follow from a systematic study of the properties of its observability and reachability Gramian matrices (cf. definitions in Section 5.1). These matrices, which can be defined for any WFA realizing a function in $\ell^2_{\mathcal{R}}$, bear a strong relation with the change of basis needed to transform an arbitrary minimal WFA into its SVA form. By studying this relation we will derive an efficient algorithm for the computation of SVA canonical forms provided that we know how to compute the Gramians associated with a WFA. Two algorithms for computing such Gramians are developed in Section 6. The second of these algorithms is based on fixed-point equations for the Gramians that are derived in Section 5.3, which also play a key role on the analysis of an approximate minimization approach given in Section 7.

5.1 Observability and reachability Gramians

Let f be rational function and $\mathbf{H}_f = \mathbf{P}\mathbf{S}^\top$ be a FB factorization for the Hankel matrix of f induced by a (non-necessarily minimal) WFA A with n states. Suppose that \mathbf{P} is such that the inner products of its columns $\langle \mathbf{P}(\cdot, i), \mathbf{P}(\cdot, j) \rangle = \sum_{x \in \Sigma^*} \mathbf{P}(x, i)\mathbf{P}(x, j)$ are finite for every $i, j \in [n]$. Then the positive semi-definite matrix $\mathbf{G}_p = \mathbf{P}^\top \mathbf{P} \in \mathbb{R}^{n \times n}$ is well defined. We call \mathbf{G}_p the *reachability Gramian* of A . Similarly, suppose the same condition on the inner products holds for the columns of \mathbf{S} . Then the matrix $\mathbf{G}_s = \mathbf{S}^\top \mathbf{S} \in \mathbb{R}^{n \times n}$ is well defined and we will call it the *observability Gramian* of A . These definitions are motivated by the following result.

Proposition 5.1. *Let A be a WFA with n states and suppose that its reachability and observability Gramians are well defined. Then the following hold:*

- (1) *A is reachable if and only if $\text{rank}(\mathbf{G}_p) = n$;*
- (2) *A is observable if and only if $\text{rank}(\mathbf{G}_s) = n$;*
- (3) *A is minimal if and only if $\text{rank}(\mathbf{G}_p) = \text{rank}(\mathbf{G}_s) = n$.*

Proof. Recall that A is reachable whenever $\text{rank}(\mathbf{P}) = n$, which implies that \mathbf{G}_p is the Gramian of n linearly independent vectors and therefore $\text{rank}(\mathbf{G}_p) = n$. On the other hand, if $\text{rank}(\mathbf{G}_p) = n$, then by the bound on the rank of a product of matrices we have

$$n = \text{rank}(\mathbf{G}_p) = \text{rank}(\mathbf{P}^\top \mathbf{P}) \leq \max\{\text{rank}(\mathbf{P}^\top), \text{rank}(\mathbf{P})\} = \text{rank}(\mathbf{P}) \leq n, \tag{9}$$

from where we conclude that $\text{rank}(\mathbf{P}) = n$ and therefore A is reachable.

The observable case follows exactly the same reasoning, and the claim about minimality is just a consequence of recalling that A is minimal if and only if it is both reachable and observable. \square

Note that the above result assumed the Gramians are well defined in the first place. Nonetheless, a similar result can be obtained without such assumptions if one is willing to work with finite versions of these matrices obtained by summing only strings up to some fixed (large enough) length. In particular, defining for any $t \geq 0$ the matrices

$$\mathbf{G}_p^{(t)} = \sum_{x \in \Sigma^{\leq t}} \mathbf{P}(x, \cdot)^\top \mathbf{P}(x, \cdot), \tag{10}$$

$$\mathbf{G}_s^{(t)} = \sum_{x \in \Sigma^{\leq t}} \mathbf{S}(x, \cdot)^\top \mathbf{S}(x, \cdot), \tag{11}$$

it is possible to see that when $t \geq n$ we have $\text{rank}(\mathbf{G}_p^{(t)}) = \text{rank}(\mathbf{P})$ and $\text{rank}(\mathbf{G}_s^{(t)}) = \text{rank}(\mathbf{S})$. However, we shall not pursue this direction here. Instead we look for necessary and sufficient conditions guaranteeing the finiteness of the Gramian matrices.

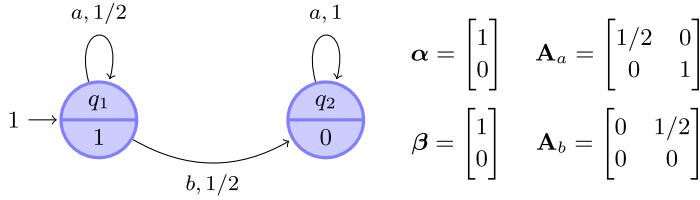


Figure 5. Example of non-minimal WFA computing a function in $\ell^2_{\mathcal{R}}$ for which the forward Gramian is not defined. To see that A is not minimal, note that f_A can be computed by the one state WFA obtained by removing q_2 from A . Note that $\mathbf{G}_p(2, 2)$ is not defined since $(\alpha^\top \mathbf{A}_b \mathbf{A}_a^k \mathbf{e}_2)^2 = 1/4$ for all $k \geq 0$.

Proposition 5.2. *Let A be a minimal WFA realizing a rational function f . The reachability and observability Gramians of A are well defined if and only if $f \in \ell^2_{\mathcal{R}}$.*

Proof. Suppose A is a minimal WFA with n states realizing a function $f \in \ell^2_{\mathcal{R}}$. It follows from Theorems 4.3 and 3.3 that there exists an invertible matrix \mathbf{Q} such that $B = A^{\mathbf{Q}}$ is an SVA. Since B induces the FB factorization given by $\mathbf{P}_B = \mathbf{U}\mathbf{D}^{1/2}$ and $\mathbf{S}_B = \mathbf{V}\mathbf{D}^{1/2}$, we see that the corresponding Gramian matrices are well defined and since $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbf{I}$ we have $\mathbf{G}_{B,p} = \mathbf{G}_{B,s} = \mathbf{D}$. Now recall that the FB factorization induced by A has $\mathbf{P}_A \mathbf{Q} = \mathbf{P}_B$ and $\mathbf{Q}^{-1} \mathbf{S}_A^\top = \mathbf{S}_B^\top$. Therefore, the Gramian matrices associated with A are also well defined since they can be obtained as $\mathbf{G}_{A,p} = \mathbf{Q}^{-\top} \mathbf{G}_{B,p} \mathbf{Q}^{-1}$ and $\mathbf{G}_{A,s} = \mathbf{Q} \mathbf{G}_{B,s} \mathbf{Q}^\top$.

Now suppose A has well-defined Gramian matrices $\mathbf{G}_p = \mathbf{P}^\top \mathbf{P}$ and $\mathbf{G}_s = \mathbf{S}^\top \mathbf{S}$. This implies that the trace $\text{Tr}(\mathbf{G}_p \mathbf{G}_s)$ is finite, which can be used to show that $f \in \ell^2_{\mathcal{R}}$ as follows:

$$\|f\|_2^2 = \sum_{x \in \Sigma^*} f(x)^2 \leq \sum_{x \in \Sigma^*} (|x| + 1) f(x)^2 = \text{Tr}(\mathbf{H}_f \mathbf{H}_f^\top) \tag{12}$$

$$= \text{Tr}(\mathbf{P} \mathbf{S}^\top \mathbf{S} \mathbf{P}^\top) = \text{Tr}(\mathbf{P}^\top \mathbf{P} \mathbf{S} \mathbf{S}^\top) = \text{Tr}(\mathbf{G}_p \mathbf{G}_s) < \infty. \tag{13}$$

□

Note that the minimality assumption is not needed when showing that A having well-defined Gramians implies $f_A \in \ell^2_{\mathcal{R}}$. On the other hand, the minimality of A is essential to show that $f_A \in \ell^2_{\mathcal{R}}$ implies that both Gramians are well defined, as witnessed by the example in Figure 5.

5.2 Gramians and SVA

The reason for introducing the reachability and observability Gramians in the previous section is because these matrices can be used to reduce any given (minimal) WFA to its SVA form. The details of this construction are presented in this section, and they draw upon some ideas already present in the proof of Proposition 5.2. Essentially, this section provides a reduction from the computation of the SVA to the computation of the Gramians. The later problem is studied in detail in Section 6.

Let A be a minimal WFA with n states realizing a function $f \in \ell^2_{\mathcal{R}}$. By Proposition 5.2 we know that the Gramians $\mathbf{G}_{A,p}$ and $\mathbf{G}_{A,s}$ are defined. Furthermore, Theorems 4.3 and 3.3 guarantee the existence of an invertible matrix \mathbf{Q} such that $B = A^{\mathbf{Q}}$ is an SVA for f . Let \mathbf{D} be the diagonal matrix containing the singular values of the Hankel matrix of f . By inspecting the proof of Proposition 5.2, we see that these facts imply the following important equations:

$$\mathbf{G}_{B,p} = \mathbf{D} = \mathbf{Q}^\top \mathbf{G}_{A,p} \mathbf{Q}, \tag{14}$$

$$\mathbf{G}_{B,s} = \mathbf{D} = \mathbf{Q}^{-1} \mathbf{G}_{A,s} \mathbf{Q}^{-\top}. \tag{15}$$

These equations say that given A we can obtain its corresponding SVA by finding an invertible matrix \mathbf{Q} simultaneously transforming the Gramians of A into two equal diagonal matrices. The

following results provide a way to do this by taking the Cholesky decompositions of the Gramian matrices and computing an additional SVD.

Lemma 5.3. *Let A be a minimal WFA with n states realizing a function $f \in \ell_{\mathcal{R}}^2$. Suppose the Gramians \mathbf{G}_p and \mathbf{G}_s satisfy $\mathbf{G}_p = \mathbf{G}_s = \mathbf{D} = \text{diag}(\sigma_1, \dots, \sigma_n)$ with $\sigma_1 \geq \dots \geq \sigma_n > 0$. Then A is an SVA, and \mathbf{D} is the matrix of singular values of \mathbf{H}_f .*

Proof. Let $\mathbf{H}_f = \mathbf{P}\mathbf{S}^\top$ be the FB factorization induced by A . Since $\mathbf{G}_p = \mathbf{P}^\top\mathbf{P}$ and $\mathbf{G}_s = \mathbf{S}^\top\mathbf{S}$ are diagonal and full rank, we see that the columns of \mathbf{P} (resp. \mathbf{S}) are orthogonal. Now take $\mathbf{U} = \mathbf{P}\mathbf{D}^{-1/2}$ and $\mathbf{V} = \mathbf{S}\mathbf{D}^{-1/2}$ and note that these two matrices have orthonormal columns since $\mathbf{U}^\top\mathbf{U} = \mathbf{V}^\top\mathbf{V} = \mathbf{I}$. Noting that $\mathbf{H}_f = \mathbf{P}\mathbf{S}^\top = \mathbf{U}\mathbf{D}\mathbf{V}^\top$ is a decomposition satisfying the constraints of an SVD, we conclude that A is an SVA. \square

Theorem 5.4. *Let A be a minimal WFA with n states realizing a function $f \in \ell_{\mathcal{R}}^2$ with Gramians \mathbf{G}_s and \mathbf{G}_p . Let $\mathbf{G}_s = \mathbf{L}_s\mathbf{L}_s^\top$ and $\mathbf{G}_p = \mathbf{L}_p\mathbf{L}_p^\top$ be their Cholesky decompositions. Suppose $\mathbf{L}_p^\top\mathbf{L}_s$ has SVD $\mathbf{U}\mathbf{D}\mathbf{V}^\top$. Then the WFA $B = A^\mathbf{Q}$ with $\mathbf{Q} = \mathbf{L}_p^{-\top}\mathbf{U}\mathbf{D}^{1/2}$ is an SVA for A . Furthermore, we have $\mathbf{Q}^{-1} = \mathbf{D}^{1/2}\mathbf{V}^\top\mathbf{L}_s^{-1}$.*

Proof. In the first place, note that minimality of A implies that \mathbf{G}_p and \mathbf{G}_s are full rank. Thus the factors \mathbf{L}_p and \mathbf{L}_s are invertible, $\mathbf{L}_p^\top\mathbf{L}_s$ has full rank, and both \mathbf{Q} and \mathbf{Q}^{-1} are well defined. To check the equality $\mathbf{Q}^{-1} = \mathbf{D}^{1/2}\mathbf{V}^\top\mathbf{L}_s^{-1}$, we just write

$$\left(\mathbf{L}_p^{-\top}\mathbf{U}\mathbf{D}^{1/2}\right)\left(\mathbf{D}^{1/2}\mathbf{V}^\top\mathbf{L}_s^{-1}\right) = \mathbf{L}_p^{-\top}\left(\mathbf{U}\mathbf{D}\mathbf{V}^\top\right)\mathbf{L}_s^{-1} = \mathbf{I}. \tag{16}$$

Next we check that \mathbf{Q} is such that $\mathbf{G}_{B,p} = \mathbf{G}_{B,s} = \mathbf{D}$:

$$\begin{aligned} \mathbf{G}_{B,p} &= \mathbf{Q}^\top\mathbf{G}_p\mathbf{Q} = \left(\mathbf{D}^{1/2}\mathbf{U}^\top\mathbf{L}_p^{-1}\right)\left(\mathbf{L}_p\mathbf{L}_p^\top\right)\left(\mathbf{L}_p^{-\top}\mathbf{U}\mathbf{D}^{1/2}\right) = \mathbf{D}, \\ \mathbf{G}_{B,s} &= \mathbf{Q}^{-1}\mathbf{G}_s\mathbf{Q}^{-\top} = \left(\mathbf{D}^{1/2}\mathbf{V}^\top\mathbf{L}_s^{-1}\right)\left(\mathbf{L}_s\mathbf{L}_s^\top\right)\left(\mathbf{L}_s^{-\top}\mathbf{V}\mathbf{D}^{1/2}\right) = \mathbf{D}, \end{aligned}$$

where we used that $\mathbf{U}^\top\mathbf{U} = \mathbf{V}^\top\mathbf{V} = \mathbf{I}$. Therefore, we can apply Lemma 5.3 to conclude that B is an SVA. \square

The previous theorem motivates the following simple algorithm for computing the SVA of a function $f \in \ell_{\mathcal{R}}^2$ provided that a minimal WFA A and its corresponding Gramian matrices are given. We shall address the computation of the Gramian matrices in the next section. For now we note that the constraint of A being minimal is not essential, since it is possible to minimize a WFA with n states in time $O(n^3)$ (Berstel and Reutenauer 2011). Furthermore, given a minimal WFA A , it is possible to check the membership $f_A \in \ell_{\mathcal{R}}^2$ using any of the tests discussed in Section 4.3, which provides a way to verify the pre-condition necessary to ensure the existence of the Gramian matrices.

Algorithm 1: ComputeSVA

Input: A minimal WFA A realizing $f \in \ell_{\mathcal{R}}^2$, and the Gramians $\mathbf{G}_{A,p}$ and $\mathbf{G}_{A,s}$

Output: An SVA B for f

- 1 Compute the Cholesky decompositions $\mathbf{G}_s = \mathbf{L}_s\mathbf{L}_s^\top$ and $\mathbf{G}_p = \mathbf{L}_p\mathbf{L}_p^\top$
 - 2 Compute the SVD $\mathbf{U}\mathbf{D}\mathbf{V}^\top$ of $\mathbf{L}_p^\top\mathbf{L}_s$
 - 3 Let $B = A^\mathbf{Q}$ with $\mathbf{Q} = \mathbf{L}_p^{-\top}\mathbf{U}\mathbf{D}^{1/2}$
 - 4 **return** B
-

The running time of $\text{ComputeSVA}(A)$ in terms of *floating point operations* (flops) can be bounded using the following well-known facts about numerical linear algebra (see, e.g., Trefethen and Bau III 1997):

- Computing the product of two matrices $d \times d$ matrices takes time $O(d^3)$ if implemented naively, and can be done in time $O(d^\omega)$ for some constant $\omega < 2.4$ using sophisticated algorithms that only yield practical improvements on very large matrices.
- The SVD of a matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$ can be computed in time $O(d^3)$, and the Cholesy decomposition of a positive definite matrix $\mathbf{G} \in \mathbb{R}^{d \times d}$ can also be computed in time $O(d^3)$.
- The inverse of an invertible lower triangular matrix $\mathbf{L} \in \mathbb{R}^{d \times d}$ can be computed in time $O(d^3)$ using Gaussian elimination.

Therefore, if the input A to Algorithm 1 has n states, its total running time is $O(n^3 + |\Sigma|n^\omega)$.

The following important observation about the product of two Gramians follows from the results showing how to compute an SVA from the Gramian matrices of a minimal WFA.

Corollary 5.5. *Let A be a minimal WFA with n states realizing a function $f \in \ell^2_{\mathcal{R}}$. Then the product of the Gramians $\mathbf{W} = \mathbf{G}_{A,s}\mathbf{G}_{A,p}$ is a diagonalizable matrix with eigenvalues given by $\lambda_i(\mathbf{W}) = \sigma_i(\mathbf{H}_f)^2$ for $i \in [n]$. Furthermore, if \mathbf{Q} is an invertible matrix such that $A^{\mathbf{Q}}$ is an SVA, then \mathbf{Q} diagonalizes \mathbf{W} ; that is $\mathbf{W} = \mathbf{Q}\mathbf{D}^2\mathbf{Q}^{-1}$.*

Proof. Let $B = A^{\mathbf{Q}}$ be an SVA for f as above. By multiplying (14) and (15) together, we see that

$$\mathbf{G}_{B,s}\mathbf{G}_{B,p} = \mathbf{D}^2 = \mathbf{Q}^{-1}\mathbf{G}_{A,s}\mathbf{G}_{A,p}\mathbf{Q} = \mathbf{Q}^{-1}\mathbf{W}\mathbf{Q}. \tag{17}$$

Therefore, \mathbf{W} is diagonalizable and its eigenvalues are the squares of the Hankel singular values of f . Additionally, the above expression shows that \mathbf{Q} necessarily is a matrix of eigenvectors for \mathbf{W} . □

5.3 Gramian fixed-point equations

In addition to their definitions in terms of a FB factorization, the Gramian matrices of a WFA can be characterized in terms of fixed-point equations. This point of view will prove useful later both for theoretical arguments and for developing algorithms for computing them.

Theorem 5.6. *Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be a WFA with n states such that the corresponding Gramians \mathbf{G}_p and \mathbf{G}_s are well defined. Then $\mathbf{X} = \mathbf{G}_p$ and $\mathbf{Y} = \mathbf{G}_s$ are solutions to the following fixed-point equations:*

$$\mathbf{X} = \alpha\alpha^\top + \sum_{a \in \Sigma} A_a^\top \mathbf{X} A_a, \tag{18}$$

$$\mathbf{Y} = \beta\beta^\top + \sum_{a \in \Sigma} A_a \mathbf{Y} A_a^\top. \tag{19}$$

Proof. Recall that $\mathbf{G}_p = \mathbf{P}^\top \mathbf{P}$ with $\mathbf{P} \in \mathbb{R}^{\Sigma^* \times n}$, and the row of \mathbf{P} corresponding to $x \in \Sigma^*$ is given by $\alpha^\top A_x$. Expanding these definitions we get

$$\begin{aligned} \mathbf{G}_p &= \sum_{x \in \Sigma^*} (A_x^\top \alpha)(\alpha^\top A_x) \\ &= \alpha\alpha^\top + \sum_{x \in \Sigma^+} (A_x^\top \alpha)(\alpha^\top A_x) \end{aligned}$$

$$\begin{aligned} &= \alpha\alpha^\top + \sum_{a \in \Sigma} \sum_{x \in \Sigma^*} \mathbf{A}_a^\top (\mathbf{A}_x^\top \alpha) (\alpha^\top \mathbf{A}_x) \mathbf{A}_a \\ &= \alpha\alpha^\top + \sum_{a \in \Sigma} \mathbf{A}_a^\top \left(\sum_{x \in \Sigma^*} (\mathbf{A}_x^\top \alpha) (\alpha^\top \mathbf{A}_x) \right) \mathbf{A}_a, \end{aligned}$$

where we just used that $\mathbf{A}_x^\top = (\mathbf{A}_{x_1} \cdots \mathbf{A}_{x_t})^\top = \mathbf{A}_{x_t}^\top \cdots \mathbf{A}_{x_1}^\top$ and that any string $y \in \Sigma^+$ satisfies $y = xa$ for some $x \in \Sigma^*$ and $a \in \Sigma$. The derivation for \mathbf{G}_s follows exactly the same pattern. \square

We note here that in the simple case where $|\Sigma| = 1$ equations (18) and (19) are special cases of the well-known *discrete Lyapunov equation*.

Another important remark about this result is that the same argument used in the proof can be used to show that the matrices $\mathbf{G}_p^{(t)}$ and $\mathbf{G}_s^{(t)}$ defined in equations (10) and (11) satisfy the following recurrence relations for any $t \geq 0$:

$$\mathbf{G}_p^{(t+1)} = \alpha\alpha^\top + \sum_{a \in \Sigma} \mathbf{A}_a^\top \mathbf{G}_p^{(t)} \mathbf{A}_a, \tag{20}$$

$$\mathbf{G}_s^{(t+1)} = \beta\beta^\top + \sum_{a \in \Sigma} \mathbf{A}_a \mathbf{G}_s^{(t)} \mathbf{A}_a^\top. \tag{21}$$

Thus, for any WFA A with n states, it will be convenient to define the mappings $F_p, F_s : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$ given by

$$F_p(\mathbf{X}) = \alpha\alpha^\top + \sum_{a \in \Sigma} \mathbf{A}_a^\top \mathbf{X} \mathbf{A}_a, \tag{22}$$

$$F_s(\mathbf{Y}) = \beta\beta^\top + \sum_{a \in \Sigma} \mathbf{A}_a \mathbf{Y} \mathbf{A}_a^\top. \tag{23}$$

With this notation, the results from this section can be summarized by saying that for any $t \geq 0$ we have $\mathbf{G}_p^{(t)} = F_p^{t+1}(\mathbf{0})$, and when the Gramian \mathbf{G}_p is defined, then it is a fixed point of the form $F_p(\mathbf{X}) = \mathbf{X}$ which can be obtained as $\lim_{t \rightarrow \infty} F_p^t(\mathbf{0})$. The same results apply to \mathbf{G}_s by replacing F_p with F_s .

These maps satisfy an important property when applied to positive semi-definite matrices.

Lemma 5.7. *The maps F_p and F_s defined in (22) and (23) are monotonically increasing with respect to the Loewner order.*

Proof. Let \mathbf{X} and \mathbf{Y} be positive semi-definite matrices satisfying $\mathbf{X} \geq \mathbf{Y}$. We need to show $F_p(\mathbf{X}) \geq F_p(\mathbf{Y})$. Recalling that for any matrices $\mathbf{M} \geq \mathbf{0}$ and \mathbf{Q} one has $\mathbf{Q}^\top \mathbf{M} \mathbf{Q} \geq \mathbf{0}$, we see that

$$F_p(\mathbf{X}) - F_p(\mathbf{Y}) = \sum_a \mathbf{A}_a^\top (\mathbf{X} - \mathbf{Y}) \mathbf{A}_a \geq \mathbf{0}, \tag{24}$$

since positive semi-definite matrices are closed under addition. The claim for F_s follows from a similar argument. \square

Finally, we conclude this section by stating a simple observation about the sequences $\mathbf{G}_p^{(t)}$ and $\mathbf{G}_s^{(t)}$ that will prove useful in the sequel.

Lemma 5.8. *One has $\mathbf{G}_p^{(t+1)} \geq \mathbf{G}_p^{(t)}$ and $\mathbf{G}_s^{(t+1)} \geq \mathbf{G}_s^{(t)}$ for any t .*

Proof. These just follow from (10) and (11) by observing that the differences

$$\begin{aligned} \mathbf{G}_p^{(t+1)} - \mathbf{G}_p^{(t)} &= \sum_{|x|=t+1} \mathbf{P}(x, :)^T \mathbf{P}(x, :), \\ \mathbf{G}_s^{(t+1)} - \mathbf{G}_s^{(t)} &= \sum_{|x|=t+1} \mathbf{S}(x, :)^T \mathbf{S}(x, :), \end{aligned}$$

are positive semi-definite matrices. □

5.4 Applications of Gramians

We have seen so far that having the Gramians of a minimal WFA A computing a rational function $f_A \in \ell^2_{\mathcal{R}}$ is enough to efficiently find the SVA of A . We now show how having the Gramians of A is also useful to compute several other quantities associated with f_A , including its ℓ^2 norm.

Theorem 5.9. *Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be a WFA computing a rational function f . Then the following hold:*

- (1) *If the Gramian \mathbf{G}_s is defined, then $\|f\|_2^2 = \alpha^T \mathbf{G}_s \alpha$.*
- (2) *If the Gramian \mathbf{G}_p is defined, then $\|f\|_2^2 = \beta^T \mathbf{G}_p \beta$.*
- (3) *If both Gramians are defined, then $\|\mathbf{H}_f\|_{\text{op}}^2 = \rho(\mathbf{G}_p \mathbf{G}_s)$ and $\|\mathbf{H}_f\|_{S,2}^2 = \text{Tr}(\mathbf{G}_p \mathbf{G}_s)$.*

Proof. Suppose \mathbf{G}_s is defined. Letting \bar{x} denote the reverse of a string x , the first equation follows from

$$\begin{aligned} \|f\|_2^2 &= \sum_{x \in \Sigma^*} f(x)^2 = \sum_{x \in \Sigma^*} (\alpha^T \mathbf{A}_x \beta) (\beta^T \mathbf{A}_{\bar{x}}^T \alpha) = \alpha^T \left(\sum_{x \in \Sigma^*} \mathbf{A}_x \beta \beta^T \mathbf{A}_{\bar{x}}^T \right) \alpha \\ &= \alpha^T \left(\sum_{x \in \Sigma^*} \mathbf{S}(x, :)^T \mathbf{S}(x, :)^T \right) \alpha = \alpha^T \mathbf{G}_s \alpha. \end{aligned}$$

By writing $f(x)^2 = (\beta^T \mathbf{A}_{\bar{x}}^T \alpha) (\alpha^T \mathbf{A}_x \beta)$, the proof of $\|f\|_2^2 = \beta^T \mathbf{G}_p \beta$ follows from the same argument.

Now suppose both Gramians are defined and recall from Proposition 5.2 that this implies $f \in \ell^2_{\mathcal{R}}$. Therefore, $\|\mathbf{H}_f\|_{\text{op}}$ and $\|\mathbf{H}_f\|_{S,2}$ are both finite. The desired equations follow directly from Corollary 5.5 by noting that $\rho(\mathbf{G}_p \mathbf{G}_s) = \lambda_1(\mathbf{G}_p \mathbf{G}_s) = \sigma_1(\mathbf{H}_f)^2 = \|\mathbf{H}_f\|_{\text{op}}^2$ and $\text{Tr}(\mathbf{G}_p \mathbf{G}_s) = \sum_{i=1}^n \lambda_i(\mathbf{G}_p \mathbf{G}_s) = \sum_{i=1}^n \sigma_i(\mathbf{H}_f)^2 = \|\mathbf{H}_f\|_{S,2}^2$. □

Note that this last result shows that if either the reachability or observability Gramian of a possibly non-minimal WFA A are defined, then we have $f_A \in \ell^2_{\mathcal{R}}$. This gives a criterion for testing a WFA for finite ℓ^2 norm in addition to those provided by Theorem 4.7. It is also interesting to contrast this results with Proposition 5.2, in which we showed that if A is minimal and $f_A \in \ell^2_{\mathcal{R}}$, then both Gramians are necessarily defined.

6. Computing the gramians

In this section, we present several algorithmic approaches for computing the SVA of a rational function in $\ell^2_{\mathcal{R}}$ given in the form of an arbitrary minimal WFA. By Algorithm 1 this problem reduces to that of computing the Gramian matrices associated with the given WFA. The first

approach works in the particular case where the fixed-point Gramian equations have a unique solution, in which case the Gramians can be efficiently computed by solving a system of linear equations. The second, more general algorithm is based on the computation of the least solution to a semi-definite system of matrix inequalities.

6.1 The unique solution case

The main idea behind our first algorithm for computing the Gramian matrices of a WFA is based on directly exploiting the definitions of these matrices. In particular, since $G_p = P^T P$, we have that $G_p(i, j)$ is the inner product between the i th and the j th columns of P . By noting that each of these columns is in fact a rational function, we see that computing G_p can be reduced to the problem of computing the inner product of two rational functions. Since it is possible to find a closed-form solution to this inner product computation, this observation can be exploited to compute G_p directly by obtaining these inner products one at a time. However, we will observe that a significant amount of these calculations can actually be reused from entry to entry. This motivates the development of an improved procedure that efficiently exploits this structure by amortizing the shared computations among all entries in G_p . Of course, by symmetry the very same arguments can be applied to the Gramian G_s .

We start with the following simple observation about solutions to the Gramian fixed-point equations.

Lemma 6.1. *Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be a WFA with n states and $X \in \mathbb{R}^{n \times n}$ an arbitrary matrix. Recall that $x = \text{vec}(X) \in \mathbb{R}^{n^2}$ is the vector obtained by concatenating the columns of X . Then the following hold:*

- (1) X is a solution of $X = \alpha\alpha^T + \sum_a A_a^T X A_a$ if and only if x is a solution of $(\alpha \otimes \alpha)^T = x^T (I - \sum_a A_a \otimes A_a)$,
- (2) X is a solution of $X = \beta\beta^T + \sum_a A_a X A_a^T$ if and only if x is a solution of $(\beta \otimes \beta) = (I - \sum_a A_a \otimes A_a)x$.

Proof. The result follows immediately from the well-known relations $\text{vec}(v v^T) = v \otimes v$ and $\text{vec}(A X B^T) = (B \otimes A) \text{vec}(X)$, and the linearity of the $\text{vec}(\bullet)$ operation. □

Now we can show that the fixed-point equations have a unique solution when a simple condition is satisfied. This yields an efficient algorithm for computing G_p and G_s when an easily testable condition holds.

Theorem 6.2. *Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be a WFA with n states and denote by ρ the spectral radius of the matrix $\sum_a A_a \otimes A_a$. If $\rho < 1$, then the following are satisfied:*

- (1) $x = \text{vec}(G_p)$ is the unique solution to $(\alpha \otimes \alpha)^T = x^T (I - \sum_a A_a \otimes A_a)$
- (2) $y = \text{vec}(G_s)$ is the unique solution to $(\beta \otimes \beta) = (I - \sum_a A_a \otimes A_a)y$

Proof. Recall that the WFA $B = \langle \alpha \otimes \alpha, \beta \otimes \beta, \{A_a \otimes A_a\} \rangle$ satisfies $f_B = f_A^2$. Therefore, we have $f_B(x) \geq 0$ for all $x \in \Sigma^*$. Using the assumption on ρ and Lemma 3.2 we see that $f_B \in \ell_{\mathcal{R}}^1$ and therefore $f_A \in \ell_{\mathcal{R}}^2$, which by Proposition 5.2 implies that $G_{A,p}$ and $G_{A,s}$ are well defined. Therefore,

Theorem 5.6 and Lemma 6.1 tell us that both $(\alpha \otimes \alpha)^\top = \mathbf{x}^\top(\mathbf{I} - \sum_a \mathbf{A}_a \otimes \mathbf{A}_a)$ and $(\beta \otimes \beta) = (\mathbf{I} - \sum_a \mathbf{A}_a \otimes \mathbf{A}_a)\mathbf{y}$ have at least one solution.

Suppose $\mathbf{y}, \mathbf{y}' \in \mathbb{R}^{n^2}$ are two solutions to equation $(\beta \otimes \beta) = (\mathbf{I} - \sum_a \mathbf{A}_a \otimes \mathbf{A}_a)\mathbf{y}$. This implies that $(\mathbf{I} - \sum_a \mathbf{A}_a \otimes \mathbf{A}_a)\mathbf{y} = (\mathbf{I} - \sum_a \mathbf{A}_a \otimes \mathbf{A}_a)\mathbf{y}'$, from where we deduce that $\mathbf{y} - \mathbf{y}' = (\sum_a \mathbf{A}_a \otimes \mathbf{A}_a)(\mathbf{y} - \mathbf{y}')$. Thus, either $\mathbf{y} = \mathbf{y}'$ or $\mathbf{y} - \mathbf{y}'$ is an eigenvector of $\sum_a \mathbf{A}_a \otimes \mathbf{A}_a$ with eigenvalue 1. Since the latter is not possible because we assumed $\rho < 1$, we conclude that the solution is unique. The same argument applies to $(\alpha \otimes \alpha)^\top = \mathbf{x}^\top(\mathbf{I} - \sum_a \mathbf{A}_a \otimes \mathbf{A}_a)$. \square

6.2 The general case

In the case where the automaton $A = \langle \alpha, \beta, \{\mathbf{A}_a\} \rangle$ is such that $\lambda = 1$ is an eigenvalue of $\sum_a \mathbf{A}_a \otimes \mathbf{A}_a$, then the linear system considered in the previous section will not have a unique solution. For example, this might occur when A is minimal but $A \otimes A$ is not. Therefore, in general we will need some extra information about the Gramian matrices in order to find them among the subset of possible solutions of the linear systems given by Lemma 6.1. This information is provided by our next lemma, which states that the Gramian matrices are the least positive-semi-definite solutions of some linear matrix inequalities. Throughout this section we assume that $A = \langle \alpha, \beta, \{\mathbf{A}_a\} \rangle$ is a WFA with n states such that the corresponding Gramians \mathbf{G}_p and \mathbf{G}_s are well defined, and therefore the linear systems in Lemma 6.1 admit at least one solution.

Lemma 6.3. *The following hold:*

(1) *The Gramian \mathbf{G}_p is the least positive semi-definite solution to the linear matrix inequality*

$$\mathbf{X} \geq \alpha\alpha^\top + \sum_{a \in \Sigma} \mathbf{A}_a^\top \mathbf{X} \mathbf{A}_a. \tag{25}$$

(2) *The Gramian \mathbf{G}_s is the least positive semi-definite solution to the linear matrix inequality*

$$\mathbf{Y} \geq \beta\beta^\top + \sum_{a \in \Sigma} \mathbf{A}_a \mathbf{Y} \mathbf{A}_a^\top. \tag{26}$$

Proof. Since the proofs of both statements follow exactly the same structure, we give only the proof for \mathbf{G}_p . From Theorem 5.6 it follows that \mathbf{G}_p satisfies (25). Now let \mathbf{X} be another positive semi-definite matrix satisfying (25). We will show by induction that for every $t \geq 0$ we have

$$\mathbf{X} \geq \sum_{x \in \Sigma^{\leq t}} \mathbf{A}_x^\top \alpha\alpha^\top \mathbf{A}_x + \sum_{x \in \Sigma^{t+1}} \mathbf{A}_x^\top \mathbf{X} \mathbf{A}_x. \tag{27}$$

First note that the case $t = 0$ follows immediately from (25). Now assume that the inequality is true for some t and consider the case $t + 1$. We have

$$\begin{aligned} \mathbf{X} &\geq \sum_{x \in \Sigma^{\leq t}} \mathbf{A}_x^\top \alpha\alpha^\top \mathbf{A}_x + \sum_{x \in \Sigma^{t+1}} \mathbf{A}_x^\top \mathbf{X} \mathbf{A}_x \\ &\geq \sum_{x \in \Sigma^{\leq t}} \mathbf{A}_x^\top \alpha\alpha^\top \mathbf{A}_x + \sum_{x \in \Sigma^{t+1}} \mathbf{A}_x^\top \alpha\alpha^\top \mathbf{A}_x + \sum_{x \in \Sigma^{t+2}} \mathbf{A}_x^\top \mathbf{X} \mathbf{A}_x \\ &= \sum_{x \in \Sigma^{\leq t+1}} \mathbf{A}_x^\top \alpha\alpha^\top \mathbf{A}_x + \sum_{x \in \Sigma^{t+2}} \mathbf{A}_x^\top \mathbf{X} \mathbf{A}_x, \end{aligned}$$

where the second inequality uses (25) and the fact that $\mathbf{Y} \geq \mathbf{Z}$ implies $\mathbf{M}^\top \mathbf{Y} \mathbf{M} \geq \mathbf{M}^\top \mathbf{Z} \mathbf{M}$ for any matrix \mathbf{M} . By rewriting (27) and noting that $\sum_{x \in \Sigma^{t+1}} \mathbf{A}_x^\top \mathbf{X} \mathbf{A}_x \geq \mathbf{0}$ for any $t \geq 0$, we see that

$$\sum_{x \in \Sigma^{\leq t}} \mathbf{A}_x^\top \boldsymbol{\alpha} \boldsymbol{\alpha}^\top \mathbf{A}_x \leq \mathbf{X} - \sum_{x \in \Sigma^{t+1}} \mathbf{A}_x^\top \mathbf{X} \mathbf{A}_x \leq \mathbf{X}. \tag{28}$$

Since \mathbf{G}_p is defined, we must have $\mathbf{G}_p = \lim_{t \rightarrow \infty} \sum_{x \in \Sigma^{\leq t}} \mathbf{A}_x^\top \boldsymbol{\alpha} \boldsymbol{\alpha}^\top \mathbf{A}_x$, and therefore $\mathbf{G}_p \leq \mathbf{X}$. \square

As a direct consequence of the above lemma, we get the following characterization for the Gramian matrices of any WFA A with $f_A \in \ell_{\mathcal{R}}^2$.

Theorem 6.4. *The Gramian \mathbf{G}_p (resp. \mathbf{G}_s) is the least positive semi-definite fixed point of (18) (resp. (19)).*

Proof. For \mathbf{G}_p , the result follows from Lemma 6.3 since any fixed-point of (18) satisfies (25); the same holds for \mathbf{G}_s . \square

Using this characterization we can derive an efficient algorithm for finding the Gramian matrices even when the linear systems given by Lemma 6.1 have more than one solution. The solution is based on solving a semi-definite optimization program. For simplicity we only present the optimization problem for finding the Gramian \mathbf{G}_s and note that a completely symmetric argument also works for \mathbf{G}_p . We start by introducing some notation. Let $\mathbf{M} = \mathbf{I} + \sum_a \mathbf{A}_a \otimes \mathbf{A}_a \in \mathbb{R}^{n^2 \times n^2}$ and $\mathbf{y}_0 = \mathbf{M}^\dagger (\boldsymbol{\beta} \otimes \boldsymbol{\beta})$. Also, let $\mathbf{y}_1, \dots, \mathbf{y}_d \in \mathbb{R}^{n^2}$ be a basis of linearly independent vectors for the column-space of the matrix $\mathbf{I} - \mathbf{M}^\dagger \mathbf{M}$. For $0 \leq i \leq d$, we write $\mathbf{Y}_i \in \mathbb{R}^{n \times n}$ to denote the matrix such that $\mathbf{y}_i = \text{vec}(\mathbf{Y}_i)$. Finally, we let π denote the linear map representing the orthogonal projection onto the space of $n \times n$ symmetric matrices, which is given by $\pi(\mathbf{Y}) = (\mathbf{Y} + \mathbf{Y}^\top)/2$. With this notation we define the following semi-definite optimization problem:

$$\text{minimize}_{t_1, \dots, t_d \in \mathbb{R}} \sum_{i=1}^d t_i \text{Tr}(\mathbf{Y}_i) \tag{29}$$

$$\text{subject to } \pi(\mathbf{Y}_0) + \sum_{i=1}^d t_i \pi(\mathbf{Y}_i) \geq \mathbf{0}, \tag{30}$$

$$\mathbf{Y}_0 - \mathbf{Y}_0^\top + \sum_{i=1}^d t_i (\mathbf{Y}_i - \mathbf{Y}_i^\top) = \mathbf{0}. \tag{31}$$

Theorem 6.5. *Let t_1^*, \dots, t_d^* be the optimal solution to (29). Then the matrix $\mathbf{Y}^* = \mathbf{Y}_0 + \sum_{i=1}^d t_i^* \mathbf{Y}_i$ is the least positive semi-definite solution of $\mathbf{Y} = \boldsymbol{\beta} \boldsymbol{\beta}^\top + \sum_a \mathbf{A}_a \mathbf{Y} \mathbf{A}_a^\top$.*

Proof. We start by observing that all solutions to $\mathbf{Y} = \boldsymbol{\beta} \boldsymbol{\beta}^\top + \sum_a \mathbf{A}_a \mathbf{Y} \mathbf{A}_a^\top$ are of the form $\mathbf{Y} = \mathbf{Y}_0 + \sum_{i=1}^d t_i \mathbf{Y}_i$ for some t_1, \dots, t_d . This follows from the fact that the Moore–Penrose pseudo-inverse can be used to show that every solution of the linear system $\boldsymbol{\beta} \otimes \boldsymbol{\beta} = (\mathbf{I} - \sum_a \mathbf{A}_a \otimes \mathbf{A}_a) \mathbf{y}$ can be written in the form $\mathbf{M}^\dagger (\boldsymbol{\beta} \otimes \boldsymbol{\beta}) + (\mathbf{I} - \mathbf{M}^\dagger \mathbf{M}) \mathbf{z}$ for some $\mathbf{z} \in \mathbb{R}^{n^2}$. Since any solution of this form can be rewritten as $\mathbf{y}_0 + \sum_{i=1}^d t_i \mathbf{y}_i$, the claim follows directly by the linearity of the $\text{vec}(\cdot)$ operation.

Next we show that any matrix of the form $\mathbf{Y} = \mathbf{Y}_0 + \sum_{i=1}^d t_i \mathbf{Y}_i$ satisfying (30) and (31) is symmetric and positive semi-definite. Indeed, if (31) is satisfied, then $\pi(\mathbf{Y}) = \mathbf{Y}$ since

$$\begin{aligned} \mathbf{Y} - \pi(\mathbf{Y}) &= \left(\mathbf{Y}_0 - \frac{\mathbf{Y}_0 + \mathbf{Y}_0^\top}{2} \right) + \sum_{i=1}^d t_i \left(\mathbf{Y}_i - \frac{\mathbf{Y}_i + \mathbf{Y}_i^\top}{2} \right) \\ &= \frac{\mathbf{Y}_0 - \mathbf{Y}_0^\top}{2} + \sum_{i=1}^d t_i \frac{\mathbf{Y}_i - \mathbf{Y}_i^\top}{2} = \mathbf{0}. \end{aligned}$$

Therefore, \mathbf{Y} is symmetric and (30) implies $\mathbf{Y} = \pi(\mathbf{Y}) \geq \mathbf{0}$, so \mathbf{Y} is positive semi-definite.

Finally suppose \mathbf{Y} and \mathbf{Y}' are two positive semi-definite solutions of (19) with $\mathbf{Y} \leq \mathbf{Y}'$. Then by linearity of the trace we have $\text{Tr}(\mathbf{Y}) \leq \text{Tr}(\mathbf{Y}')$. Therefore, the least positive semi-definite solution to (19) is also the positive semi-definite solution with minimum trace \mathbf{Y}^* obtained by solving (29). □

7. Application: Approximate minimization of WFA

The fact that given a (minimal) WFA realizing a function in $\ell^2_{\mathcal{R}}$ we can efficiently compute its corresponding SVA opens the door to multiple applications. In this section, we focus on the application of SVA to the design and analysis of algorithms for model reduction. To motivate the need for such algorithms, consider the situation where one has a WFA modeling a system of interest and the need arises for testing whether the system satisfies a given property. If testing this property requires multiple evaluations of the function computed by the WFA, the cost of this computation will grow with the number of states, and if the system is large the repeated evaluation of millions of queries might take a very long time. But if the decision about the property being satisfied does not depend too much on the individual answers of each query, it might be acceptable to provide *approximate* answers for each of these queries. If in addition these approximate queries can be performed much faster than exact queries, then the whole testing process can be sped up by trading-off accuracy and query processing time. In the rest of this section, we formalize the problem of approximate evaluations of WFA and provide a solution based on the truncation of SVA canonical forms.

7.1 Problem formulation

We now proceed to give a formal definition of the approximate minimization problem for WFA. Roughly speaking, this corresponds to finding a small WFA computing a good approximation to the function realized by a large minimal WFA. A solution to this problem will yield a way to speed up approximate evaluation of WFA.

Let A be a minimal WFA with n states computing a rational function $f \in \ell^2_{\mathcal{R}}$. Given a target number of states $\hat{n} < n$, we want to find a WFA \hat{A} with \hat{n} states computing a function \hat{f} which minimizes $\|f - \hat{f}\|_2$ among all rational function of rank at most \hat{n} . This problem can be formulated as an optimization problem as follows:

$$\inf_{\text{rank}(\hat{f}) \leq \hat{n}} \|f - \hat{f}\|_2. \tag{32}$$

The first observation we make about this problem is that although it is not explicitly encoded in (32), any solution \hat{f} will have finite ℓ^2 norm. Indeed, it is easy to see that

$$\|\hat{f}\|_2 \leq \|f\|_2 + \|f - \hat{f}\|_2 = \|f\|_2 + \inf_{\text{rank}(\hat{f}) \leq \hat{n}} \|f - \hat{f}\|_2 \leq 2\|f\|_2, \tag{33}$$

where the last inequality uses that the rational function 0 has rank 1 and therefore it is a feasible point of the optimization (32).

The second important observation is that, like rank constrained optimizations over finite matrices, the optimization in (32) is not convex. To see this, let us write $A = \langle \alpha, \beta, \{A_a\} \rangle$ for the original automaton and $\hat{A} = \langle \hat{\alpha}, \hat{\beta}, \{\hat{A}_a\} \rangle$ for the automaton we are looking for, noting that any automaton with at most \hat{n} states can be written as a (non-minimal) WFA with \hat{n} . Then, by the monotonicity of $z \mapsto z^2$ we can replace the objective $\|f - \hat{f}\|_2$ with $\|f - \hat{f}\|_2^2$, and see that, using the WFA representation for $f - \hat{f}$ and the closed-form expression for $\|f - \hat{f}\|_2^2$ in terms of this WFA representation, (32) can be rewritten as the minimization over \hat{A} of the quantity

$$[\alpha^\top \ \hat{\alpha}^\top] \otimes [\alpha^\top \ \hat{\alpha}^\top] \left(\mathbf{I} - \sum_{a \in \Sigma} \begin{bmatrix} A_a & \mathbf{0} \\ \mathbf{0} & \hat{A}_a \end{bmatrix} \otimes \begin{bmatrix} A_a & \mathbf{0} \\ \mathbf{0} & \hat{A}_a \end{bmatrix} \right)^{-1} \begin{bmatrix} \beta \\ -\hat{\beta} \end{bmatrix} \otimes \begin{bmatrix} \beta \\ -\hat{\beta} \end{bmatrix}. \tag{34}$$

Since this equivalent objective function is not convex, we have little hope of being able to efficiently solve (32) exactly. Instead, we will take a different approach and see how truncating the SVA for A to have \hat{n} states yields an approximate solution which can be efficiently computed.

7.2 SVA truncation

In this section, we describe an approximate minimization algorithm for WFA realizing a function in $\ell^2_{\mathcal{R}}$. The algorithm takes as input a minimal WFA A with n states and a target number of states \hat{n} , and outputs a new WFA \hat{A} with \hat{n} states approximating the original WFA A . To obtain \hat{A} , we first compute the SVA A' associated with A , and then remove the $n - \hat{n}$ states associated with the smallest singular values of \mathbf{H}_{f_A} . More formally, by writing the block decomposition of the operators associated with the SVA A' shown below, we get the operators for \hat{A} by taking the sub-block in the top left containing the first \hat{n} rows and \hat{n} columns:

$$A'_a = \begin{bmatrix} A_a^{(11)} & A_a^{(12)} \\ A_a^{(21)} & A_a^{(22)} \end{bmatrix}, \quad \hat{A}_a = [A_a^{(11)}]. \tag{35}$$

Note that if we define the matrix $\Gamma = [\mathbf{I}_{\hat{n}} \ \mathbf{0}] \in \mathbb{R}^{\hat{n} \times n}$, then we have $\hat{A}_a = \Gamma A'_a \Gamma^\top$. To reflect this fact, we shall sometimes write $\hat{A} = \Gamma A' \Gamma^\top$. Algorithm 2 provides a description of the full procedure, which we call `SVATruncation`. Since the algorithm only involves a call to `ComputeSVA` and a simple algebraic manipulation of the resulting WFA, the running time of `SVATruncation` is dominated by the complexity of `ComputeSVA`, and hence is polynomial in $|\Sigma|$, $\dim(A)$, and \hat{n} .

Algorithm 2: SVATruncation

Input: A minimal WFA A with n states, a target number of states $\hat{n} < n$

Output: A WFA \hat{A} with \hat{n} states

- 1 Let $A' \leftarrow \text{ComputeSVA}(A)$
 - 2 Let $\Gamma = [\mathbf{I}_{\hat{n}} \ \mathbf{0}] \in \mathbb{R}^{\hat{n} \times n}$
 - 3 Let $\hat{A}_a = \Gamma A'_a \Gamma^\top$ for all $a \in \Sigma$
 - 4 Let $\hat{\alpha} = \Gamma \alpha'$
 - 5 Let $\hat{\beta} = \Gamma \beta'$
 - 6 Let $\hat{A} = \langle \hat{\alpha}, \hat{\beta}, \{\hat{A}_a\} \rangle$
 - 7 **return** \hat{A}
-

Roughly speaking, the rationale behind SVATruncation is that given an SVA, the states corresponding to the smallest singular values are the ones with less influence on the Hankel matrix, and therefore should also be the ones with less influence on the associated rational function. However, the details are more tricky than this simple intuition. The reason being that a low rank approximation to \mathbf{H}_f obtained by truncating its SVD is not in general a Hankel matrix, and therefore does not correspond to any rational function. In particular, the Hankel matrix of the function \hat{f} computed by \hat{A} is not obtained by truncating the SVD of \mathbf{H}_f . This makes our analysis more involved than just applying the well-known bounds for low-rank approximation based on SVD. Nonetheless, we are able to obtain a bound of the same form that one would expect by measuring the error of a low-rank approximation using the Frobenius norm. Along these lines, our main result is the following theorem, which bounds the ℓ^2 -distance between the rational function f realized by the original WFA A , and the rational function \hat{f} realized by the output WFA \hat{A} .

Theorem 7.1. *Let A be a minimal WFA with n states computing a function $f \in \mathcal{L}_{\mathcal{R}}^2$ with Hankel singular values $\sigma_1 \geq \dots \geq \sigma_n > 0$. Let \hat{f} denote the function computed by the truncated SVA \hat{A} with $1 \leq \hat{n} < n$ states. Then the following holds:*

$$\|f - \hat{f}\|_2^2 \leq \sum_{i=\hat{n}+1}^n \sigma_i^2. \tag{36}$$

The proof will be given in Section 7.4. First, a few remarks about this result are in order. The first is to observe that because $\sigma_1 \geq \dots \geq \sigma_n$, the error decreases when \hat{n} increases, which is the desired behavior: the more states \hat{A} has, the closer it is to A . The second is that (36) does not depend on which representation A of f is given as input to SVATruncation. This is a consequence of first obtaining the corresponding SVA A' before truncating. Obviously, one could obtain another approximate minimization by truncating A directly. However, in that case the final error would depend on the initial A and in general it does not seem possible to use this approach for providing representation independent bounds on the quality of approximation. To see the importance of starting the truncation procedure from the SVA canonical form, let us consider the following result which follows directly from the Gramian fixed-point equations for SVA.

Lemma 7.2. *Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be an SVA with n states realizing a function $f \in \mathcal{L}_{\mathcal{R}}^2$ with Hankel singular values $\sigma_1 \geq \dots \geq \sigma_n$. Then the following are satisfied:*

- (1) For all $j \in [n]$, $\sum_i \sigma_i \sum_a A_a(i, j)^2 = \sigma_j - \alpha(j)^2$,
- (2) For all $i \in [n]$, $\sum_j \sigma_j \sum_a A_a(i, j)^2 = \sigma_i - \beta(i)^2$.

Proof. These equations correspond to the diagonal entries of the Gramian fixed-point equations for SVA

$$\mathbf{D} = \alpha\alpha^\top + \sum_a A_a^\top \mathbf{D} A_a, \tag{37}$$

$$\mathbf{D} = \beta\beta^\top + \sum_a A_a \mathbf{D} A_a^\top. \tag{38}$$

□

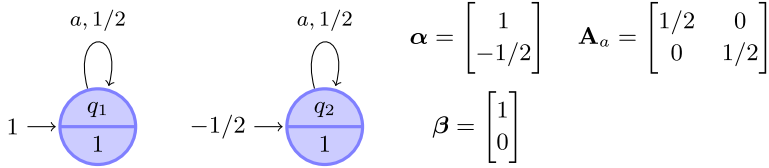


Figure 6. Example of WFA A such that $\|f_{\hat{A}}\|_2 \leq \|f_A\|_2$, where \hat{A} is the automaton obtained by removing the state q_2 . In particular, $\|f_A\|_2^2 = 1/3$ and $\|f_{\hat{A}}\|_2^2 = 4/3$.

To see why this lemma justifies the truncation of an SVA, we consider the following simple consequence. By fixing $i, j \in [n]$ and $a \in \Sigma$, we can use the first equation to get

$$\sigma_i \mathbf{A}_a(i, j)^2 = \sigma_j - \alpha(j)^2 - \left(\sum_i \sigma_i \sum_a \mathbf{A}_a(i, j)^2 - \sigma_i \mathbf{A}_a(i, j)^2 \right) \leq \sigma_j.$$

Applying a similar argument to the second equation, we conclude that

$$|\mathbf{A}_a(i, j)| \leq \min \left\{ \sqrt{\frac{\sigma_i}{\sigma_j}}, \sqrt{\frac{\sigma_j}{\sigma_i}} \right\} = \sqrt{\frac{\min\{\sigma_i, \sigma_j\}}{\max\{\sigma_i, \sigma_j\}}}.$$

This bound is telling us that in an SVA, transition weights further away from the diagonals of the \mathbf{A}_a are going to be small whenever there is a wide spread between the largest and smallest singular values; for example, $|\mathbf{A}_a(1, n)| \leq \sqrt{\sigma_n/\sigma_1}$. Intuitively, this means that in an SVA the last states are very weakly connected to the first states, and therefore removing these connections should not affect the output of the WFA too much. The proof of Theorem 7.1 exploits this intuition, while at the same time leverages the full power of the fixed-point SVA Gramian equations.

We finish this section by stating another result about \hat{A} : SVA truncation always reduces the norm of the original function. Logically speaking, this is a preliminary to Theorem 7.1, since it shows that the function computed by \hat{A} has finite ℓ^2 norm and already implies the finiteness of $\|f - \hat{f}\|_2$. From an approximation point of view, this result basically says that SVA truncation can be interpreted as an algorithm for approximate minimization “from below,” which might be a desirable property in some applications.

Theorem 7.3. *Let A be a WFA computing a function $f \in \ell^2_{\mathcal{R}}$ of rank n and \hat{A} a truncation of the SVA of A with $\hat{n} < n$ states. The function \hat{f} computed by \hat{A} satisfies $\hat{f} \in \ell^2_{\mathcal{R}}$ and $\|\hat{f}\|_2 \leq \|f\|_2$.*

It is important to note that in general truncating an arbitrary WFA does not always reduce its norm as shown by the example in Figure 6.

7.3 SVA truncation: Bounding the norm

The proof of the Theorem 7.3 illustrates how having different ways to represent the function \hat{f} computed by the SVA truncation can be useful; this fact will also be essential in the proof of Theorem 7.1. In particular, given an SVA A , we note that the automaton \tilde{A} obtained by padding with zeros all the coefficients in the initial and transition weights of A that are removed when taking its truncation in SVATruncation computes the same function as \hat{A} .

More concretely, let us recall the notation from (35) splitting of the weights conforming A into a block corresponding to states 1 to \hat{n} , and another block containing states $\hat{n} + 1$ to n . We can define a similar partition for the initial and final weights of A . In particular, we write the following:

$$\begin{aligned} \alpha &= \begin{bmatrix} \alpha^{(1)} \\ \alpha^{(2)} \end{bmatrix}, \\ \beta &= \begin{bmatrix} \beta^{(1)} \\ \beta^{(2)} \end{bmatrix}, \\ A_a &= \begin{bmatrix} A_a^{(11)} & A_a^{(12)} \\ A_a^{(21)} & A_a^{(22)} \end{bmatrix}. \end{aligned}$$

Now the SVA truncation $\hat{A} = \Gamma A \Gamma^\top = \langle \hat{\alpha}, \hat{\beta}, \{\hat{A}_a\} \rangle$ with $\Gamma = [I_{\hat{n}} \mathbf{0}]$ can be written in terms of this block decomposition as $\hat{\alpha} = \alpha^{(1)}$, $\hat{\beta} = \beta^{(1)}$, and $\hat{A}_a = A_a^{(11)}$.

The important observation here is that starting from A we can write other WFA computing the same function as \hat{A} . The following construction yields a WFA \tilde{A} of size n with this property. Define the $n \times n$ matrix

$$\Pi = \begin{bmatrix} I_{\hat{n}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \tag{39}$$

and let $\tilde{A} = \langle \tilde{\alpha}, \tilde{\beta}, \{\tilde{A}_a\} \rangle$ with $\tilde{\alpha} = \Pi \alpha$, $\tilde{\beta} = \beta$, and $\tilde{A}_a = A_a \Pi^\top = A_a \Pi$. For convenience we shall sometimes write $\tilde{A} = A \Pi$. Note that the weights of \tilde{A} are given by:

$$\begin{aligned} \tilde{\alpha} &= \begin{bmatrix} \alpha^{(1)} \\ \mathbf{0} \end{bmatrix}, \\ \tilde{\beta} &= \begin{bmatrix} \beta^{(1)} \\ \beta^{(2)} \end{bmatrix}, \\ \tilde{A}_a &= \begin{bmatrix} A_a^{(11)} & \mathbf{0} \\ A_a^{(21)} & \mathbf{0} \end{bmatrix}. \end{aligned}$$

Lemma 7.4. *Let $A = \langle \alpha, \beta, \{A_a\} \rangle$ be an SVA with n states. If $\hat{A} = \Gamma A \Gamma^\top = \langle \hat{\alpha}, \hat{\beta}, \{\hat{A}_a\} \rangle$ is the truncation of A with \hat{n} states, and $\tilde{A} = A \Pi = \langle \tilde{\alpha}, \tilde{\beta}, \{\tilde{A}_a\} \rangle$ is the WFA with n states defined above, then \hat{A} and \tilde{A} compute the same function \hat{f} .*

Proof. Given $x \in \Sigma^*$ define $\tilde{\alpha}_x^\top = \tilde{\alpha}^\top \tilde{A}_x$ and $\hat{\alpha}_x^\top = \hat{\alpha}^\top \hat{A}_x$. By using the pattern of zeros in \tilde{A}_x , a simple induction argument on the length of x shows that the following is always satisfied:

$$\tilde{\alpha}_x^\top = [\hat{\alpha}_x^\top \mathbf{0}]. \tag{40}$$

Therefore for any $x \in \Sigma^*$ we have $f_{\tilde{A}}(x) = \tilde{\alpha}_x^\top \tilde{\beta} = \hat{\alpha}_x^\top \tilde{\beta}^{(1)} = f_{\hat{A}}(x)$. □

The advantage of having a WFA with n states computing the same function as the SVA truncation is that now both A and \tilde{A} have Gramians of the same dimensions which can be compared. The following lemma provides such comparison.

Lemma 7.5. *Let A be an SVA with n states and reachability Gramian $\mathbf{G}_p = \mathbf{D}$. Let $\tilde{A} = A\Pi$ the WFA with n states computing the same function as the truncation of A with \hat{n} states. Then the reachability Gramian $\tilde{\mathbf{G}}_p$ of \tilde{A} is defined and satisfies $\mathbf{G}_p \geq \tilde{\mathbf{G}}_p$.*

Proof. Recall the definition of the map F_p for A from Section 5.3, and note that the corresponding map for \tilde{A} satisfies

$$\tilde{F}_p(\mathbf{X}) = \tilde{\alpha}\tilde{\alpha}^\top + \sum_a \tilde{\mathbf{A}}_a^\top \mathbf{X} \tilde{\mathbf{A}}_a = \Pi F_p(\mathbf{X}) \Pi. \tag{41}$$

Taking $\tilde{\mathbf{G}}_p^{(0)} = \mathbf{0}$ and $\tilde{\mathbf{G}}_p^{(t+1)} = \tilde{F}_p(\tilde{\mathbf{G}}_p^{(t)})$, we have $\tilde{\mathbf{G}}_p^{(t+1)} \geq \tilde{\mathbf{G}}_p^{(t)}$ for all $t \geq 0$ (Lemma 5.8). Furthermore, $\tilde{\mathbf{G}}_p = \lim_{t \rightarrow \infty} \tilde{\mathbf{G}}_p^{(t)}$ if the limit is defined.

We will simultaneously show that the limit above is defined and satisfies $\tilde{\mathbf{G}}_p \leq \mathbf{D}$. Define the sequence $\mathbf{X}_0 = \mathbf{D}$ and $\mathbf{X}_{t+1} = \tilde{F}_p(\mathbf{X}_t)$ for $t \geq 0$. Clearly all the matrices in the sequence are positive semi-definite, and furthermore we claim that they satisfy $\mathbf{X}_t \geq \mathbf{X}_{t+1}$ for all t . The case $t = 0$ is immediate since $\mathbf{X}_1 = \tilde{F}_p(\mathbf{D}) = \Pi F_p(\mathbf{D}) \Pi = \Pi \mathbf{D} \Pi = \Pi \mathbf{D} \leq \mathbf{D} = \mathbf{X}_0$. For $t > 0$ we use induction and the fact that \tilde{F}_p is monotonous (Lemma 5.7): if $\mathbf{X}_t \geq \mathbf{X}_{t+1}$, then $\mathbf{X}_{t+1} = \tilde{F}_p(\mathbf{X}_t) \geq \tilde{F}_p(\mathbf{X}_{t+1}) = \mathbf{X}_{t+2}$. Thus, since $\tilde{\mathbf{G}}_p^{(0)} = \mathbf{0} \leq \mathbf{D} = \mathbf{X}_0$, for all $t \geq 0$ we have $\tilde{\mathbf{G}}_p^{(t)} \leq \mathbf{X}_t \leq \mathbf{D}$. This implies that the monotonously increasing sequence $\tilde{\mathbf{G}}_p^{(t)}$ is bounded by \mathbf{D} , and therefore its limit exists and is upper bounded by \mathbf{D} . \square

The above lemma will be enough to prove the desired upper bound on the norm of \hat{f} . On the other hand, we note that because \tilde{A} is not a minimal WFA, the boundedness of \hat{f} or the existence of the reachability Gramian $\tilde{\mathbf{G}}_p$ do not immediately imply the existence of the observability Gramian $\tilde{\mathbf{G}}_s$ for \tilde{A} ; we will see in the next section that in fact this Gramian is also defined.

Proof of Theorem 7.3. By Lemma 7.4 we can work with \tilde{A} instead of \hat{A} . Now, Lemma 7.5 shows that the Gramian $\tilde{\mathbf{G}}_p$ of \tilde{A} is defined, so by Theorem 5.9 the function \hat{f} computed by \tilde{A} has finite ℓ^2 norm. Furthermore, since $\mathbf{G}_{A,p} \geq \mathbf{G}_{\tilde{A},p}$, the expressions for the norm $\|f\|_2$ in Theorem 5.9 imply that $\|f\|_2^2 = \beta^\top \mathbf{G}_{A,p} \beta \geq \beta^\top \mathbf{G}_{\tilde{A},p} \beta = \|\hat{f}\|_2^2$. \square

7.4 SVA truncation: Error analysis

In this section, we prove the bound on $\|f - \hat{f}\|_2$ given in Theorem 7.1, where \hat{f} is the function computed by the SVA truncation of f with \hat{n} states. In fact, the bound will follow from an exact closed-form expression for the error $\|f - \hat{f}\|_2$ given in terms of the Gramians of a WFA computing $\tilde{f} = f - \hat{f}$.

We recall from last section the automaton $\tilde{A} = A\Pi$ with n states computing \hat{f} , where $\Pi = \text{diag}(\mathbf{I}_{\hat{n}}, \mathbf{0})$. Now we proceed to combine A and \tilde{A} to obtain a WFA computing the difference $\tilde{f} = f - \hat{f}$. The construction follows the same argument used to show that the difference of two rational functions is rational and yields the WFA $\bar{A} = \langle \bar{\alpha}, \bar{\beta}, \{\bar{\mathbf{A}}_a\} \rangle$ with $2n$ states given by

$$\begin{aligned} \bar{\alpha} &= \begin{bmatrix} \alpha \\ \tilde{\alpha} \end{bmatrix}, \\ \bar{\beta} &= \begin{bmatrix} \beta \\ -\tilde{\beta} \end{bmatrix}, \\ \bar{A}_a &= \begin{bmatrix} A_a & \mathbf{0} \\ \mathbf{0} & \tilde{A}_a \end{bmatrix} = \text{diag}(A_a, \tilde{A}_a). \end{aligned}$$

It is immediate to check from these constructions that \bar{A} satisfies $f_{\bar{A}} = \bar{f}$. The following lemmas establish a few preliminary facts about \bar{A} .

Lemma 7.6. *The observability Gramian \tilde{G}_s of \tilde{A} is defined.*

Proof. Let $H_{\tilde{f}} = \tilde{P}\tilde{S}^\top$ be the factorization induced by \tilde{A} and recall that $\tilde{G}_s = \tilde{S}^\top\tilde{S}$ if the corresponding inner products between the columns of \tilde{S} are defined. Thus, to prove that \tilde{G}_s is defined it suffices to show that all the columns of \tilde{S} have finite ℓ^2 norm, which is equivalent to showing that $\|\tilde{S}\|_F < \infty$. Expanding this Frobenius norm we have

$$\begin{aligned} \|\tilde{S}\|_F^2 &= \sum_{x \in \Sigma^*} \|\tilde{A}_x\beta\|_2^2 \\ &= \sum_{x \in \Sigma^*} \beta^\top \tilde{A}_x^\top \tilde{A}_x \beta \\ &= \sum_{x \in \Sigma^*} \text{Tr}(\tilde{A}_x\beta\beta^\top \tilde{A}_x^\top), \end{aligned} \tag{42}$$

where the last equality uses the cyclic property of the trace. Now note that using the SVA fixed-point equation $D = \beta\beta^\top + \sum_a A_a D A_a^\top$ we can rewrite any term in the infinite sum as

$$\begin{aligned} \text{Tr}(\tilde{A}_x\beta\beta^\top \tilde{A}_x^\top) &= \text{Tr}\left(\tilde{A}_x\left(D - \sum_a A_a D A_a^\top\right)\tilde{A}_x^\top\right) \\ &= \text{Tr}\left(\tilde{A}_x\left(D - \sum_a A_a \Pi D A_a^\top - \sum_a A_a (\mathbf{I} - \Pi) D A_a^\top\right)\tilde{A}_x^\top\right). \end{aligned}$$

Since Π is idempotent and commutes with diagonal matrices, we have $A_a \Pi D A_a^\top = A_a \Pi D \Pi A_a^\top = \tilde{A}_a D \tilde{A}_a^\top$. Therefore, by linearity of the matrix trace, we can plug the last two observations into (42) and get

$$\begin{aligned} \|\tilde{S}\|_F^2 &= \text{Tr}(D) - \sum_a \text{Tr}(\tilde{A}_a D \tilde{A}_a^\top) - \sum_a \text{Tr}(A_a (\mathbf{I} - \Pi) D A_a^\top) \\ &\quad + \sum_{x \in \Sigma^+} \text{Tr}(\tilde{A}_x D \tilde{A}_x^\top) - \sum_{x \in \Sigma^+} \sum_a \text{Tr}(\tilde{A}_x \tilde{A}_a D \tilde{A}_a^\top \tilde{A}_x^\top) \\ &\quad - \sum_{x \in \Sigma^+} \sum_a \text{Tr}(\tilde{A}_x A_a (\mathbf{I} - \Pi) D A_a^\top \tilde{A}_x^\top). \end{aligned}$$

By aggregating terms we see that all terms of the form $\text{Tr}(\tilde{\mathbf{A}}_x \mathbf{D} \tilde{\mathbf{A}}_x^\top)$ for $x \in \Sigma^+$ cancel and finally get

$$\|\tilde{\mathbf{S}}\|_{\mathbb{F}}^2 = \text{Tr}(\mathbf{D}) - \sum_{x \in \Sigma^*} \sum_a \text{Tr}(\tilde{\mathbf{A}}_x \mathbf{A}_a (\mathbf{I} - \Pi) \mathbf{D} \mathbf{A}_a^\top \tilde{\mathbf{A}}_x^\top) \leq \text{Tr}(\mathbf{D}), \tag{43}$$

where we used that $\tilde{\mathbf{A}}_x \mathbf{A}_a (\mathbf{I} - \Pi) \mathbf{D} \mathbf{A}_a^\top \tilde{\mathbf{A}}_x^\top \geq \mathbf{0}$ and the trace of a positive semi-definite matrix is always nonnegative. \square

Lemma 7.7. *Let $\mathbf{H}_f = \mathbf{P}\mathbf{S}^\top$ be the factorization induced by the SVA A , and $\mathbf{H}_{\tilde{f}} = \tilde{\mathbf{P}}\tilde{\mathbf{S}}^\top$ be the factorization induced by \tilde{A} . Then the WFA \tilde{A} computing $\tilde{f} = f - \hat{f}$ induces the factorization $\mathbf{H}_{\tilde{f}} = \tilde{\mathbf{P}}\tilde{\mathbf{S}}^\top$ with $\tilde{\mathbf{P}} = [\mathbf{P} \ \tilde{\mathbf{P}}]$ and $\tilde{\mathbf{S}} = [\mathbf{S} \ -\tilde{\mathbf{S}}]$. Furthermore, the Gramians $\tilde{\mathbf{G}}_p$ and $\tilde{\mathbf{G}}_s$ are defined and can be written as*

$$\begin{aligned} \tilde{\mathbf{G}}_p &= \tilde{\mathbf{P}}^\top \tilde{\mathbf{P}} = \begin{bmatrix} \mathbf{G}_p & \mathbf{P}^\top \tilde{\mathbf{P}} \\ \tilde{\mathbf{P}}^\top \mathbf{P} & \tilde{\mathbf{G}}_p \end{bmatrix}, \\ \tilde{\mathbf{G}}_s &= \tilde{\mathbf{S}}^\top \tilde{\mathbf{S}} = \begin{bmatrix} \mathbf{G}_s & -\mathbf{S}^\top \tilde{\mathbf{S}} \\ -\tilde{\mathbf{S}}^\top \mathbf{S} & \tilde{\mathbf{G}}_s \end{bmatrix}. \end{aligned}$$

Proof. The structure of $\tilde{\mathbf{P}}$, $\tilde{\mathbf{S}}$, $\tilde{\mathbf{G}}_p$, and $\tilde{\mathbf{G}}_s$ follow from a straightforward computation. That these Gramians are defined follows from noting that because all the Gramians of A and \tilde{A} are defined, then all the columns of \mathbf{P} , \mathbf{S} , $\tilde{\mathbf{P}}$, and $\tilde{\mathbf{S}}$ have finite ℓ^2 norm. \square

Now we are ready to prove the main result of this section giving an exact closed-form expression for the ℓ^2 distance between f and \hat{f} .

Theorem 7.8. *For any truncation size $1 \leq \hat{n} < n$, we have*

$$\|f - \hat{f}\|_2^2 = \text{Tr} \left(\mathbf{D}^{1/2} (\mathbf{I} - \Pi) (\tilde{\mathbf{S}}\mathbf{S}^\top + \mathbf{S}\tilde{\mathbf{S}}^\top - \tilde{\mathbf{S}}\tilde{\mathbf{S}}^\top) (\mathbf{I} - \Pi) \mathbf{D}^{1/2} \right). \tag{44}$$

Proof. Recall from Theorem 5.9 that $\|\tilde{f}\|_2^2 = \tilde{\boldsymbol{\alpha}}^\top \tilde{\mathbf{G}}_s \tilde{\boldsymbol{\alpha}} = \text{Tr}(\tilde{\boldsymbol{\alpha}}^\top \tilde{\mathbf{S}}^\top \tilde{\mathbf{S}} \tilde{\boldsymbol{\alpha}}) = \text{Tr}(\tilde{\mathbf{S}} \tilde{\boldsymbol{\alpha}} \tilde{\boldsymbol{\alpha}}^\top \tilde{\mathbf{S}}^\top)$, where the last equality follows from a standard property of the trace. Note that by construction of \tilde{A} we have $\tilde{\boldsymbol{\alpha}}^\top = \boldsymbol{\alpha}^\top [\mathbf{I} \ \Pi]$, which when plugged in the previous equation yields

$$\|\tilde{f}\|_2^2 = \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \Pi \end{bmatrix} \boldsymbol{\alpha} \boldsymbol{\alpha}^\top [\mathbf{I} \ \Pi] \tilde{\mathbf{S}}^\top \right). \tag{45}$$

Recall that A is an SVA, and therefore the fixed-point equation (18) applied to A yields $\boldsymbol{\alpha} \boldsymbol{\alpha}^\top = \mathbf{D} - \sum_a \mathbf{A}_a^\top \mathbf{D} \mathbf{A}_a$. When combined with (45) we obtain, by linearity of the trace:

$$\begin{aligned} \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \Pi \end{bmatrix} \boldsymbol{\alpha} \boldsymbol{\alpha}^\top [\mathbf{I} \ \Pi] \tilde{\mathbf{S}}^\top \right) &= \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \Pi \end{bmatrix} \mathbf{D} [\mathbf{I} \ \Pi] \tilde{\mathbf{S}}^\top \right) \\ &\quad - \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \Pi \end{bmatrix} \left(\sum_a \mathbf{A}_a^\top \mathbf{D} \mathbf{A}_a \right) [\mathbf{I} \ \Pi] \tilde{\mathbf{S}}^\top \right). \end{aligned} \tag{46}$$

Using that $[\mathbf{I} \ \Pi] = [\mathbf{I} \ \mathbf{I}] - [\mathbf{0} \ \mathbf{I} - \Pi]$, we decompose the first term as:

$$\begin{aligned} & \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} \mathbf{D}[\mathbf{I} \ \mathbf{I}] \tilde{\mathbf{S}}^\top \right) + \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{0} \\ \mathbf{I} - \Pi \end{bmatrix} \mathbf{D}[\mathbf{0} \ \mathbf{I} - \Pi] \tilde{\mathbf{S}}^\top \right) \\ & - \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} \mathbf{D}[\mathbf{0} \ \mathbf{I} - \Pi] \tilde{\mathbf{S}}^\top \right) - \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{0} \\ \mathbf{I} - \Pi \end{bmatrix} \mathbf{D}[\mathbf{I} \ \mathbf{I}] \tilde{\mathbf{S}}^\top \right). \end{aligned} \tag{47}$$

We now proceed to bound the sum of the last three terms in this expression. Note in the first place that each of these terms is of the form $\text{Tr}(\mathbf{M}\mathbf{D}\mathbf{N}) = \text{Tr}(\mathbf{D}^{1/2}\mathbf{N}\mathbf{M}\mathbf{D}^{1/2}) = \text{Tr}^{\mathbf{D}}(\mathbf{N}\mathbf{M})$, where in the last step we just introduced a bit of convenient notation. Furthermore, recall that by definition of $\bar{\mathbf{A}}$ we have $\tilde{\mathbf{S}} = [\mathbf{S} - \tilde{\mathbf{S}}]$. With these observations we obtain the following three equations:

$$\text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{0} \\ \mathbf{I} - \Pi \end{bmatrix} \mathbf{D}[\mathbf{0} \ \mathbf{I} - \Pi] \tilde{\mathbf{S}}^\top \right) = \text{Tr}^{\mathbf{D}} \left((\mathbf{I} - \Pi) \tilde{\mathbf{S}}^\top \tilde{\mathbf{S}} (\mathbf{I} - \Pi) \right), \tag{48}$$

$$- \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} \mathbf{D}[\mathbf{0} \ \mathbf{I} - \Pi] \tilde{\mathbf{S}}^\top \right) = \text{Tr}^{\mathbf{D}} \left((\mathbf{I} - \Pi) \tilde{\mathbf{S}}^\top (\mathbf{S} - \tilde{\mathbf{S}}) \right), \tag{49}$$

$$- \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{0} \\ \mathbf{I} - \Pi \end{bmatrix} \mathbf{D}[\mathbf{I} \ \mathbf{I}] \tilde{\mathbf{S}}^\top \right) = \text{Tr}^{\mathbf{D}} \left((\mathbf{S}^\top - \tilde{\mathbf{S}}^\top) \tilde{\mathbf{S}} (\mathbf{I} - \Pi) \right). \tag{50}$$

By observing that we have $\text{Tr}^{\mathbf{D}}((\mathbf{I} - \Pi)\mathbf{M}) = \text{Tr}^{\mathbf{D}}(\mathbf{M}(\mathbf{I} - \Pi)) = \text{Tr}^{\mathbf{D}}((\mathbf{I} - \Pi)\mathbf{M}(\mathbf{I} - \Pi))$ for any square matrix \mathbf{M} , we conclude that the sum of the last three terms in (47) equals

$$\text{Tr}^{\mathbf{D}} \left((\mathbf{I} - \Pi) (\tilde{\mathbf{S}}^\top \mathbf{S} + \mathbf{S}^\top \tilde{\mathbf{S}} - \tilde{\mathbf{S}}^\top \tilde{\mathbf{S}}) (\mathbf{I} - \Pi) \right). \tag{51}$$

To complete the proof of the equation, we will now show that the sum of the remaining terms in (46) vanishes; that is:

$$\text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} \mathbf{D}[\mathbf{I} \ \mathbf{I}] \tilde{\mathbf{S}}^\top \right) - \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \Pi \end{bmatrix} \left(\sum_a \mathbf{A}_a^\top \mathbf{D}\mathbf{A}_a \right) [\mathbf{I} \ \Pi] \tilde{\mathbf{S}}^\top \right) = 0. \tag{52}$$

We start by noting the following identity:

$$\mathbf{A}_a [\mathbf{I} \ \Pi] = [\mathbf{A}_a \ \hat{\mathbf{A}}_a] = [\mathbf{I} \ \mathbf{I}] \bar{\mathbf{A}}_a. \tag{53}$$

Therefore, using the fixed-point equation (19) we see that

$$\begin{aligned} \text{Tr} \left(\tilde{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \Pi \end{bmatrix} \left(\sum_a \mathbf{A}_a^\top \mathbf{D}\mathbf{A}_a \right) [\mathbf{I} \ \Pi] \tilde{\mathbf{S}}^\top \right) &= \text{Tr} \left(\tilde{\mathbf{S}} \left(\sum_a \bar{\mathbf{A}}_a^\top \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} \mathbf{D}[\mathbf{I} \ \mathbf{I}] \bar{\mathbf{A}}_a \right) \tilde{\mathbf{S}}^\top \right) \\ &= \text{Tr}^{\mathbf{D}} \left([\mathbf{I} \ \mathbf{I}] \left(\sum_a \bar{\mathbf{A}}_a \tilde{\mathbf{S}}^\top \tilde{\mathbf{S}} \bar{\mathbf{A}}_a^\top \right) \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} \right) \\ &= \text{Tr}^{\mathbf{D}} \left([\mathbf{I} \ \mathbf{I}] \left(\tilde{\mathbf{S}}^\top \tilde{\mathbf{S}} - \bar{\beta} \bar{\beta}^\top \right) \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} \right). \end{aligned}$$

Now (52) follows from simply observing that by the construction of $\bar{\mathbf{A}}$ we have $[\mathbf{I} \ \mathbf{I}] \bar{\beta} = 0$. □

Finally we can show how the bound in Theorem 7.1 follows directly from the exact expression for the error obtained in Theorem 7.8.

Proof of Theorem 7.1. We start noting that (44) can be rewritten as

$$\text{Tr}^{\mathbf{D}} \left((\mathbf{I} - \Pi) \mathbf{S}^{\top} \mathbf{S} (\mathbf{I} - \Pi) \right) - \text{Tr}^{\mathbf{D}} \left((\mathbf{I} - \Pi) [\mathbf{I} \mathbf{I}] \bar{\mathbf{S}}^{\top} \bar{\mathbf{S}} \begin{bmatrix} \mathbf{I} \\ \mathbf{I} \end{bmatrix} (\mathbf{I} - \Pi) \right). \tag{54}$$

Note that the second term has the form $\text{Tr}(\mathbf{M}\mathbf{M}^{\top})$ and therefore is nonnegative. Using that A is an SVA, we see that the first term is

$$\text{Tr}^{\mathbf{D}} \left((\mathbf{I} - \Pi) \mathbf{S}^{\top} \mathbf{S} (\mathbf{I} - \Pi) \right) = \text{Tr} \left(\mathbf{D}^{1/2} (\mathbf{I} - \Pi) \mathbf{D} (\mathbf{I} - \Pi) \mathbf{D}^{1/2} \right) = \sum_{i=\hat{n}+1}^n \sigma_i^2. \tag{55}$$

Thus, it follows from the last two observations that (44) is at most $\sum_{i=\hat{n}+1}^n \sigma_i^2$. □

8. Related work

In this section, we provide wider context for our work by relating it to recent developments in machine learning and to well-established results in the theory of linear dynamical systems.

Spectral techniques for learning weighted automata and other latent variable models have recently drawn a lot of attention in the machine learning community. Following the significant milestone papers (Hsu et al. 2012; Bailly et al. 2009), in which an efficient spectral algorithm for learning HMM and stochastic rational languages was given, the field has grown very rapidly. The original algorithm, which is based on SVDs of finite sub-blocks of Hankel matrices, has been extended to reduced-rank HMMs (Siddiqi et al. 2010), PSR (Boots et al. 2009), finite-state transducers (Balle et al. 2011; Bailly et al. 2013), and many other classes of functions on strings (Bailly 2011; Balle and Mohri 2012; Recasens and Quattoni 2013). Although each of these papers works with slightly different problems and analysis techniques, the key ingredient turns out to be always the same: parametrize the target model as a WFA and learn this WFA from the SVD of a finite sub-block of its Hankel matrix (Balle et al. 2014a). Therefore, it is possible (and desirable) to study all these learning algorithms from the point of view of rational series, which are exactly the class of real-valued functions on strings that can be computed by WFA.

The appeal of spectral learning techniques comes from their computational superiority when compared to iterative algorithms like Expectation–Maximization (EM) (Dempster et al. 1977). Another very attractive property of spectral methods is the possibility of proving rigorous statistical guarantees about the learned automaton. For example, under a realizability assumption, these methods are known to be consistent and amenable to finite-sample analysis in the PAC sense (Hsu et al. 2012). An important detail is that, in addition to realizability, these results work under the assumption that the user correctly guesses the number of latent states of the target distribution. Though this is not a real caveat when it comes to using these algorithms in practice – the optimal number of states can be identified using a model selection procedure (Balle et al. 2014b) – it is one of the barriers in extending the statistical analysis of spectral methods to the non-realizable setting.

Tackling the non-realizability question requires, as a special case, dealing with the situation in which data is generated from a WFA with n states and the learning algorithm is asked to produce a WFA with $\hat{n} < n$ states. This case is already a nontrivial problem which – barring the noisiness introduced by estimating the Hankel matrix from observed data – can in fact be interpreted as an approximate minimization of WFA. From this point of view, we believe our results provide the fundamental tools necessary for addressing important problems in the theory of learning weighted automata, including the robust statistical analysis of spectral learning algorithms.

A connection between spectral learning algorithms and approximate minimization for a small class of HMM was considered in Kulesza et al. (2014). This paper also presents a theoretical result bounding the error between the original and minimized HMM in terms of the total variation distance. The bounds in this paper are incomparable to ours. However, in a follow-up work (Kulesza

et al. 2015), published concurrently with our original paper on SVA (Balle et al. 2015), a problem similar to the one considered here is addressed, albeit different methods are used and the results are less general than our approximate minimization method. Another paper on which the issue of approximate minimization of weighted automata is considered in a tangential manner is Kiefer and Wachter (2014). In this case, the authors again focus on an ℓ^1 -like accuracy measure to compare two automata: an original one, and another one obtained by removing transitions with small weights occurring during an exact minimization procedure. Though the removal operation is introduced as a means of obtaining a numerically stable minimization algorithm, the paper also presents some experiments exploring the effect of removing transitions with larger weights. With the exception of these timid results, the problem of approximate minimization for general WFA remained largely unstudied before our paper.

However, the case of an alphabet with one symbol, $|\Sigma| = 1$, has been thoroughly studied from multiple points of view. In the control theory literature, several methods have been proposed for approximate minimization of time-invariant linear dynamical systems under the names of model reduction, truncation, and approximation; see Antoulas (2005) for a comprehensive presentation. One possible approach to the model reduction problem is to consider the so-called balanced realizations of a linear dynamical system and apply a convenient truncation method to the balanced realization to obtain a smaller system (Enns 1984). In the one symbol case, the connection with weighted automata arises from observing that the impulse response of a time-invariant linear dynamical system can be parametrized as a weighted automata with one letter (and possibly vector-valued outputs for multiple-input multiple-output systems) (Antoulas 2005). From this point of view, the canonical form for weighted automata given by our SVA can be interpreted as a generalization of balanced realizations to the case where the alphabet has two or more letters.

The study of model reduction techniques in the one symbol case can also be connected to sophisticated ideas in the study of approximations for Hankel operators in the functional and complex analysis literatures; see, e.g., Peller (2012) for a comprehensive treatment of the theory of Hankel operators. In the same way we do in Section 3, when the alphabet Σ has only one symbol the Hankel matrix of a rational function yields a linear operator between Hilbert or Banach spaces of sequences. The spectral properties of these Hankel operators have been thoroughly studied. For example, deep connections to the theory of complex function on the unit disk and Fourier analysis have been uncovered (Nikol'skii 2012). Along these lines one finds the celebrated AAK theorem characterizing optimal approximations of Hankel operators by Hankel operators of bounded rank (Adamyan et al. 1971). This theorem has been widely exploited in control theory to provide alternative approaches to balanced realizations for model reduction, thus providing a link between the abstract setting of Hankel operators and the concrete problem of approximating of linear dynamical systems (Glover 1984) (see also Fuhrmann 2011). One of the fundamental ideas in this line of work is realizing that for $|\Sigma| = 1$ the free monoid Σ^* can be identified with the natural numbers \mathbb{N} , which can be canonically embedded in the *abelian* group \mathbb{Z} . Unfortunately for us, this approach cannot be directly generalized to the case $|\Sigma| > 1$ because in this case the corresponding embedding yields a free non-abelian group, and standard Fourier analysis on those groups is not available. Although some recent attempts have been made to extend some of the results about Hankel operators to the noncommutative case using methods from functional analysis (Popescu 2003), this theory is still largely underdeveloped, and the few existing results can only be obtained via nonconstructive arguments.

9. Conclusion and future work

In the present paper, we have given a new approximate minimization technique based on spectral theory ideas. The essential point was to use the SVD to decide how to truncate the original automaton without losing too much accuracy. We have given quantitative bounds on how close the approximate machine is to the original.

One crucial aspect that we have not addressed is the question of constructing the *best possible* approximation given a bound on the size of the state space or, equivalently, the dimension of the vector space on which the machine is defined. In the one-letter case, sophisticated results from the theory of Hankel operators (Adamyán et al. 1971; Peller 2012) provide a satisfactory answer to this problem. However, extending this to the multiple-letter case means extending an already deep and difficult theory to the noncommutative case. Nevertheless, it remains an exciting challenge.

A different approach is to change the approximation measure from ℓ^2 to a more natural metric between WFA. In recent work (Balle et al. 2017), we developed a metric to measure the distance between WFAs based on bisimulation. This metric has interesting properties, but unfortunately it is hard to compute for the present type of approximation. Nevertheless, it might be fruitful to explore approximation schemes based on approximate bisimulation as has been done for some types of Markov processes (Desharnais et al. 2003). It would be interesting to compare the quality of such approximation schemes with the present one.

Acknowledgements. We thank François Denis for sharing the proof of Lemma 3.2 with us, Guillaume Rabusseau for useful discussions, and an anonymous referee for suggesting improvements to the presentation of this paper.

Notes

- 1 To be more precise, this is a *compact* SVD, since the inner dimensions of the decomposition are all equal to the rank. In this paper, we shall always use the term SVD to mean compact SVD.
- 2 Some authors call these functions *recognizable* and use a notion of rationality associated with belonging to a set of functions closed under certain operations. Since both notions are equivalent for the computation model of WFA we consider in this paper, we purposefully avoid the distinction between rationality and recognizability.
- 3 In real analysis a matrix \mathbf{M} is Hankel if $\mathbf{M}(i, j) = \mathbf{M}(k, l)$ whenever $i + j = k + l$, which implies that \mathbf{M} is symmetric. In our case we have $\mathbf{H}(p, s) = \mathbf{H}(p', s')$ whenever $ps = p's'$, but \mathbf{H} is not symmetric because string concatenation is not commutative whenever $|\Sigma| > 1$.
- 4 Note that these inequalities have scalars in their RHS and should be interpreted as entry-wise inequalities, and not as claims about positive semi-definite matrices.

References

- Adamyán, V. M., Arov, D. Z. and Krein, M. G. (1971). Analytic properties of Schmidt pairs for a Hankel operator and the generalized Schur–Takagi problem. *Matematicheskii Sbornik* **128** (1) 34–75.
- Albert, J. and Kari, J. (2009). Digital image compression. In: Manfred D., Werner K., and Heiko V. (eds.) *Handbook of Weighted Automata*, Springer, 453–479. https://doi.org/10.1007/978-3-642-01492-5_11.
- Antoulas, A. C. (2005). *Approximation of large-scale dynamical systems*. SIAM, Philadelphia, PA. <https://doi.org/10.1137/1.9780898718713>.
- Baier, C., Größer, M. and Ciesinski, F. (2009). Model checking linear-time properties of probabilistic systems. In: Manfred D., Werner K., and Heiko V. (eds.) *Handbook of Weighted Automata*, Springer, 519–570. https://doi.org/10.1007/978-3-642-01492-5_13.
- Bailly, R. (2011). Quadratic weighted automata: Spectral algorithm and likelihood maximization. In: *Asian Conference on Machine Learning*, 147–163.
- Bailly, R., Denis, F. and Ralaivola, L. (2009). Grammatical inference as a principal component analysis problem. In: *Proceedings of the 26th Annual International Conference on Machine Learning*, ACM, 33–40.
- Bailly, R., Carreras, X. and Quattoni, A. (2013). Unsupervised spectral learning of finite state transducers. In: *Advances in Neural Information Processing Systems*, 800–808.
- Balle, B. and Mohri, M. (2012). Spectral learning of general weighted automata via constrained matrix completion. In: *Advances in neural information processing systems*, 2159–2167.
- Balle, B., Quattoni, A. and Carreras, X. (2011). A spectral learning algorithm for finite state transducers. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, Springer, 156–171.
- Balle, B., Carreras, X., Luque, F. and Quattoni, A. (2014a). Spectral learning of weighted automata: A forward-backward perspective. *Machine Learning* **96** (1–2) 33–63. https://doi.org/10.1007/978-3-642-01492-5_13
- Balle, B., Hamilton, W. and Pineau, J. (2014b). Methods of moments for learning stochastic languages: Unified presentation and empirical comparison. In: *International Conference on Machine Learning*, 1386–1394.

- Balle, B., Panangaden, P. and Precup, D. (2015). A canonical form for weighted automata and applications to approximate minimization. In: *2015 30th Annual ACM/IEEE Symposium on Logic in Computer Science (LICS)*, IEEE, 701–712.
- Balle, B., Gourdeau, P. and Panangaden, P. (2017). Bisimulation metrics for weighted finite automata. In: *Proceedings of the 44th International Colloquium On Automata Languages and Programming Warsaw*, vol. 103, 1–14.
- Berstel, J. and Reutenauer, C. (2011). *Noncommutative Rational Series with Applications*. Cambridge University Press.
- Bezhanišvili, N., Kupke, C. and Panangaden, P. (2012). Minimization via duality. In: *Logic, Language, Information and Computation—19th International Workshop, WoLLIC 2012, Buenos Aires, Argentina, September 3–6, 2012. Proceedings*, volume 7456 of *Lecture Notes in Computer Science*, Springer, 191–205.
- Bonchi, F., Bonsangue, M., Boreale, M., Rutten, J. and Silva, A. (2012). A coalgebraic perspective in linear weighted automata. *Information and Computation* 211 77–105.
- Bonchi, F., Bonsangue, M. M., Hansen, H. H., Panangaden, P., Rutten, J. and Silva, A. (2014). Algebra-coalgebra duality in Brzozowski's minimization algorithm. *ACM Transactions on Computational Logic* 15 (1) 3:1–3:29.
- Boots, B., Siddiqi, S. and Gordon, G. (2009). Closing the learning-planning loop with predictive state representations. In: *Proceedings of Robotics: Science and Systems VI*.
- Boreale, M. (2009). Weighted bisimulation in linear algebraic form. In: *CONCUR 2009-Concurrency Theory*, Springer, 163–177.
- Brzozowski, J. A. (1962). Canonical regular expressions and minimal state graphs for definite events. In: Fox, J. (ed.) *Proceedings of the Symposium on Mathematical Theory of Automata*, number 12 in MRI Symposia Series, Polytechnic Press of the Polytechnic Institute of Brooklyn, April 1962, 529–561. Book appeared in 1963.
- de Gispert, A., Iglesias, G., Blackwood, G., Banga, E. and Byrne, W. (2010). Hierarchical phrase-based translation with weighted finite-state transducers and shallow-n grammars. *Computational Linguistics* 36 (3), 505–533.
- Dempster, A. P., Laird, N. M. and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)* 39 (1), 1–22.
- Denis, F. and Esposito, Y. (2008). On rational stochastic languages. *Fundamenta Informaticae* 86 (1,2) 41–77.
- Desharnais, J., Gupta, V., Jagadeesan, R. and Panangaden, P. (2003). Approximating labeled Markov processes. *Information and Computation* 184 (1) 160–200.
- Enns, D. F. (1984). Model reduction with balanced realizations: An error bound and a frequency weighted generalization. In: *The 23rd IEEE Conference on Decision and Control, 1984*, vol. 23, IEEE, 127–132.
- Fuhrmann, P. A. (2011). *A Polynomial Approach to Linear Algebra*. Springer Science & Business Media.
- Glover, K. (1984). All optimal hankel-norm approximations of linear multivariable systems and their l₈-error bounds. *International journal of control* 39 (6) 1115–1193.
- Hsu, D., Kakade, S. M. and Zhang, T. (2012). A spectral algorithm for learning hidden Markov models. *Journal of Computer and System Sciences* 78 (5) 1460–1480.
- Kiefer, S. and Wachter, B. (2014) Stability and complexity of minimising probabilistic automata. In: Javier E., Pierre F., Thore H., and Elias K. (eds.) *Proceedings of the 41st International Colloquium on Automata, Languages and Programming (ICALP), part II*, vol 8573, LNCS, Copenhagen, Denmark, Springer, 268–279.
- Knight, K. and May, J. (2009). Applications of weighted automata in natural language processing. In: Manfred D., Werner K., and Heiko V. (eds.) *Handbook of Weighted Automata*, Springer, 571–596.
- Kulesza, A., Rao, N. R. and Singh, S. (2014). Low-rank spectral learning. In: *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Statistics*, 522–530.
- Kulesza, A., Jiang, N. and Singh, S. (2015). Low-rank spectral learning with weighted loss functions. In: *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics*.
- Lototsky, S. V. (2015). Simple spectral bounds for sums of certain Kronecker products. *Linear Algebra and its Applications* 469 114–129.
- Mohri, M., Pereira, F. C. N. and Riley, M. (2008). Speech recognition with weighted finite-state transducers. In: *Handbook on Speech Processing and Speech Communication*.
- Nikol'skii, N. K. (2012). *Treatise on the Shift Operator: Spectral Function Theory*, vol. 273, Springer Science & Business Media.
- Peller, V. (2012). *Hankel Operators and Their Applications*. Springer Science & Business Media.
- Popescu, G. (2003). Multivariable Nehari problem and interpolation. *Journal of Functional Analysis* 200 (2) 536–581.
- Recasens, A. and Quattoni, A. (2013). Spectral learning of sequence taggers over continuous sequences. In: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, Springer, 289–304.
- Rosenthal, J. S. (1995). Convergence rates for Markov chains. *Siam Review* 37 (3) 387–405.
- Siddiqi, S., Boots, B. and Gordon, G. (2010). Reduced-rank hidden Markov models. In: *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, 741–748.
- Trefethen, L. N. and Bau III, D. (1997). *Numerical Linear Algebra*. Siam.
- Zhu, K. (1990). *Operator Theory in Function Spaces*, vol. 138. American Mathematical Society.