

# Singular Value Automata and Approximate Minimization

**Borja Balle**

Amazon Research Cambridge<sup>1</sup>

Weighted Automata: Theory and Applications — May 2018

---

<sup>1</sup>Based on work completed before joining Amazon

## Analytic Automata Theory

More prosaically:

- ▶ The use of tools from mathematical analysis to study questions in automata theory, specifically questions related to approximation and learning
- ▶ Based on joint work with: X. Carreras, P. Gourdeau, M. Mohri, P. Panangaden, D. Precup, G. Rabusseau, A. Quattoni
- ▶ Key references: [\[Bal13, BPP17\]](#)

# Keep It Real!

$\mathbb{R}$

More precisely:

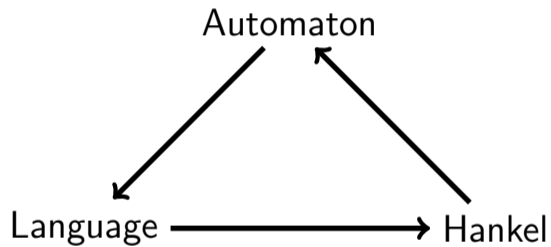
- ▶ Everything works for complex numbers
- ▶ Some things work for arbitrary fields
- ▶ Virtually nothing works for general semi-rings

# Outline

1. Weighted Languages, Weighted Automata, and Hankel Matrices
2. Perturbation Bounds Between Representations
3. Singular Value Automata: Definition
4. Singular Value Automata: Computation
5. Approximate Minimization via SVA Truncation
6. Concluding Remarks

1. Weighted Languages, Weighted Automata, and Hankel Matrices
2. Perturbation Bounds Between Representations
3. Singular Value Automata: Definition
4. Singular Value Automata: Computation
5. Approximate Minimization via SVA Truncation
6. Concluding Remarks

# The Big Picture



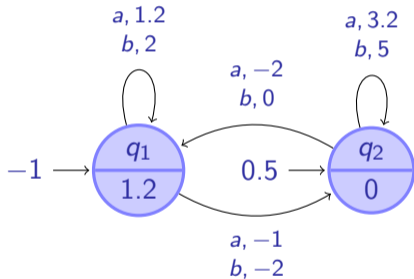
$$f : \Sigma^* \rightarrow \mathbb{R} , \quad f \in \mathbb{R}^{\Sigma^*}$$

## Notation

- ▶ Finite alphabet  $\Sigma$
- ▶ Free monoid  $\Sigma^*$
- ▶ Empty string  $\epsilon$
- ▶ String length  $|x|$
- ▶ String concatenation  $xy = x \cdot y$

# Weighted Finite Automata (WFA)

## Graphical Representation



## Algebraic Representation

$$\alpha = \begin{bmatrix} -1 \\ 0.5 \end{bmatrix} \quad \beta = \begin{bmatrix} 1.2 \\ 0 \end{bmatrix}$$

$$\mathbf{A}_a = \begin{bmatrix} 1.2 & -1 \\ -2 & 3.2 \end{bmatrix}$$

$$\mathbf{A}_b = \begin{bmatrix} 2 & -2 \\ 0 & 5 \end{bmatrix}$$

## Weighted Finite Automaton

A WFA  $A$  with  $n = |A|$  states is a tuple  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\}_{\sigma \in \Sigma} \rangle$  where  $\alpha, \beta \in \mathbb{R}^n$  and  $\mathbf{A}_\sigma \in \mathbb{R}^{n \times n}$



## Language of a WFA

With every WFA  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  with  $n$  states we associate a weighted language  $f_A : \Sigma^* \rightarrow \mathbb{R}$  given by

$$\begin{aligned} f_A(x_1 \cdots x_T) &= \sum_{q_0, q_1, \dots, q_T \in [n]} \alpha(q_0) \left( \prod_{t=1}^T \mathbf{A}_{x_t}(q_{t-1}, q_t) \right) \beta(q_T) \\ &= \alpha^\top \mathbf{A}_{x_1} \cdots \mathbf{A}_{x_T} \beta = \alpha^\top \mathbf{A}_x \beta \end{aligned}$$

### Recognizable/Rational Languages

A weighted language  $f : \Sigma^* \rightarrow \mathbb{R}$  is recognizable/rational if there exists a WFA  $A$  such that  $f = f_A$ . The smallest number of states of such a WFA is  $\text{rank}(f)$ . A WFA  $A$  is minimal if  $|A| = \text{rank}(f_A)$ .

**Observation:** The minimal  $A$  is not unique. Take any invertible matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$ , then

$$\alpha^\top \mathbf{A}_{x_1} \cdots \mathbf{A}_{x_T} \beta = (\alpha^\top \mathbf{Q})(\mathbf{Q}^{-1} \mathbf{A}_{x_1} \mathbf{Q}) \cdots (\mathbf{Q}^{-1} \mathbf{A}_{x_T} \mathbf{Q})(\mathbf{Q}^{-1} \beta)$$

# Hankel Matrices

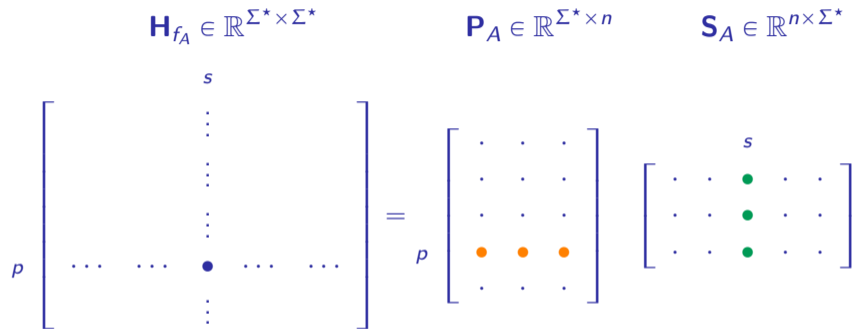
Given a weighted language  $f : \Sigma^* \rightarrow \mathbb{R}$  define its Hankel matrix  $\mathbf{H}_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$  as

$$\mathbf{H}_f = \begin{array}{c} \epsilon \\ a \\ b \\ \vdots \\ p \\ \vdots \end{array} \begin{bmatrix} \epsilon & a & b & \dots & s & \dots \\ f(\epsilon) & f(a) & f(b) & & \vdots & \\ f(a) & f(aa) & f(ab) & & \vdots & \\ f(b) & f(ba) & f(bb) & & \vdots & \\ \dots & \dots & \dots & & f(p \cdot s) & \\ \vdots & & & & & \end{bmatrix}$$

Fliess–Kronecker Theorem [Fli74]

The rank of  $\mathbf{H}_f$  is finite if and only if  $f$  is rational, in which case  $\text{rank}(\mathbf{H}_f) = \text{rank}(f)$

# Structure of Low-Rank Hankel Matrices

$$\mathbf{H}_{f_A} \in \mathbb{R}^{\Sigma^* \times \Sigma^*} \quad \mathbf{P}_A \in \mathbb{R}^{\Sigma^* \times n} \quad \mathbf{S}_A \in \mathbb{R}^{n \times \Sigma^*}$$


$$f_A(p_1 \cdots p_T \cdot s_1 \cdots s_{T'}) = \underbrace{\alpha^\top \mathbf{A}_{p_1} \cdots \mathbf{A}_{p_T}}_{\alpha_A(p)} \underbrace{\mathbf{A}_{s_1} \cdots \mathbf{A}_{s_{T'}}}_{\beta_A(s)} \beta$$

Note: We call  $\mathbf{H}_f = \mathbf{P}_A \mathbf{S}_A$  the *forward-backward factorization* induced by  $A$

# Structure of Shifted Hankel Matrices

$$f(p_1 \cdots p_T s_1 \cdots s_{T'}) = \alpha^\top \mathbf{A}_{p_1} \cdots \mathbf{A}_{p_T} \mathbf{A}_{s_1} \cdots \mathbf{A}_{s_{T'}} \beta$$

$$\mathbf{H} = \begin{matrix} & & & s & & & \\ & & & \cdot & & & \\ & & & \cdot & & & \\ & & & \cdot & & & \\ p & \left[ \begin{array}{cccccc} \cdot & \cdot & f(p_s) & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{array} \right] & = & \left[ \begin{array}{ccc} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{array} \right] & \left[ \begin{array}{cccccc} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{array} \right] \end{matrix}$$

$$f(p_1 \cdots p_T \sigma s_1 \cdots s_{T'}) = \alpha^\top \mathbf{A}_{p_1} \cdots \mathbf{A}_{p_T} \mathbf{A}_a \mathbf{A}_{s_1} \cdots \mathbf{A}_{s_{T'}} \beta$$

$$\mathbf{H}_\sigma = \begin{matrix} & & & s & & & \\ & & & \cdot & & & \\ & & & \cdot & & & \\ & & & \cdot & & & \\ p & \left[ \begin{array}{cccccc} \cdot & \cdot & f(pas) & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{array} \right] & = & \left[ \begin{array}{ccc} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{array} \right] & \left[ \begin{array}{ccc} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{array} \right] & \left[ \begin{array}{cccccc} \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \end{array} \right] \end{matrix}$$

Algebraically: Factorizing  $\mathbf{H}$  lets us solve for  $\mathbf{A}_a$

$$\mathbf{H} = \mathbf{P} \mathbf{S} \quad \implies \quad \mathbf{H}_\sigma = \mathbf{P} \mathbf{A}_\sigma \mathbf{S} \quad \implies \quad \mathbf{A}_\sigma = \mathbf{P}^+ \mathbf{H}_\sigma \mathbf{S}^+$$

## Aside: Moore–Penrose Pseudo-inverse

For any  $\mathbf{M} \in \mathbb{R}^{n \times m}$  there exists a unique *pseudo-inverse*  $\mathbf{M}^+ \in \mathbb{R}^{m \times n}$  satisfying:

- ▶  $\mathbf{M}\mathbf{M}^+\mathbf{M} = \mathbf{M}$ ,  $\mathbf{M}^+\mathbf{M}\mathbf{M}^+ = \mathbf{M}^+$ , and  $\mathbf{M}^+\mathbf{M}$  and  $\mathbf{M}\mathbf{M}^+$  are symmetric
- ▶ If  $\text{rank}(\mathbf{M}) = n$  then  $\mathbf{M}\mathbf{M}^+ = \mathbf{I}$ , and if  $\text{rank}(\mathbf{M}) = m$  then  $\mathbf{M}^+\mathbf{M} = \mathbf{I}$
- ▶ If  $\mathbf{M}$  is square and invertible then  $\mathbf{M}^+ = \mathbf{M}^{-1}$

Given a system of linear equations  $\mathbf{M}\mathbf{u} = \mathbf{v}$ , the following is satisfied:

$$\mathbf{M}^+\mathbf{v} = \underset{\mathbf{u} \in \text{argmin} \|\mathbf{M}\mathbf{u} - \mathbf{v}\|_2}{\text{argmin}} \|\mathbf{u}\|_2 .$$

In particular:

- ▶ If the system is completely determined,  $\mathbf{M}^+\mathbf{v}$  solves the system
- ▶ If the system is underdetermined,  $\mathbf{M}^+\mathbf{v}$  is the solution with smallest norm
- ▶ If the system is overdetermined,  $\mathbf{M}^+\mathbf{v}$  is the minimum norm solution to the least-squares problem  $\min \|\mathbf{M}\mathbf{u} - \mathbf{v}\|_2$

## From Finite Hankel Matrix to WFA

Suppose  $f : \Sigma^* \rightarrow \mathbb{R}$  has rank  $n$  and  $\varepsilon \in \mathcal{P}, \mathcal{S} \subset \Sigma^*$  are such that the sub-block  $\mathbf{H} \in \mathbb{R}^{\mathcal{P} \times \mathcal{S}}$  of  $\mathbf{H}_f$  satisfies  $\text{rank}(\mathbf{H}) = n$ .

Let  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  be obtained as follows:

1. Compute a rank factorization  $\mathbf{H} = \mathbf{P}\mathbf{S}$ ; i.e.  $\text{rank}(\mathbf{P}) = \text{rank}(\mathbf{S}) = \text{rank}(\mathbf{H})$
2. Let  $\alpha^\top$  (resp.  $\beta$ ) be the  $\varepsilon$ -row of  $\mathbf{P}$  (resp.  $\varepsilon$ -column of  $\mathbf{S}$ )
3. Let  $\mathbf{A}_\sigma = \mathbf{P}^+ \mathbf{H}_\sigma \mathbf{S}^+$ , where  $\mathbf{H}_\sigma \in \mathbb{R}^{\mathcal{P} \cdot \sigma \times \mathcal{S}}$  is a sub-block of  $\mathbf{H}_f$

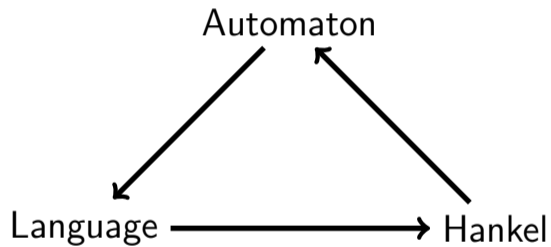
Claim The resulting WFA computes  $f$  and is minimal

### Proof

- ▶ Suppose  $\tilde{A} = \langle \tilde{\alpha}, \tilde{\beta}, \{\tilde{\mathbf{A}}_\sigma\} \rangle$  is a minimal WFA for  $f$ .
- ▶ It suffices to show there exists an invertible  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that  $\alpha^\top = \tilde{\alpha}^\top \mathbf{Q}$ ,  $\mathbf{A}_\sigma = \mathbf{Q}^{-1} \tilde{\mathbf{A}}_\sigma \mathbf{Q}$  and  $\beta = \mathbf{Q}^{-1} \tilde{\beta}$ .
- ▶ By minimality  $\tilde{A}$  induces a rank factorization  $\mathbf{H} = \tilde{\mathbf{P}} \tilde{\mathbf{S}}$  and also  $\mathbf{H}_\sigma = \tilde{\mathbf{P}} \tilde{\mathbf{A}}_\sigma \tilde{\mathbf{S}}$ .
- ▶ Since  $\mathbf{A}_\sigma = \mathbf{P}^+ \mathbf{H}_\sigma \mathbf{S}^+ = \mathbf{P}^+ \tilde{\mathbf{P}} \tilde{\mathbf{A}}_\sigma \tilde{\mathbf{S}} \mathbf{S}^+$ , take  $\mathbf{Q} = \tilde{\mathbf{S}} \mathbf{S}^+$ .
- ▶ Check  $\mathbf{Q}^{-1} = \mathbf{P}^+ \tilde{\mathbf{P}}$  since  $\mathbf{P}^+ \tilde{\mathbf{P}} \tilde{\mathbf{S}} \mathbf{S}^+ = \mathbf{P}^+ \mathbf{H} \mathbf{S}^+ = \mathbf{P}^+ \mathbf{P} \mathbf{S} \mathbf{S}^+ = \mathbf{I}$ .

1. Weighted Languages, Weighted Automata, and Hankel Matrices
2. Perturbation Bounds Between Representations
3. Singular Value Automata: Definition
4. Singular Value Automata: Computation
5. Approximate Minimization via SVA Truncation
6. Concluding Remarks

# The Big Picture





## Weighted Finite Automaton

A WFA with  $n$  states is a tuple  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\}_{\sigma \in \Sigma} \rangle$  where  $\alpha, \beta \in \mathbb{R}^n$  and  $\mathbf{A}_\sigma \in \mathbb{R}^{n \times n}$

Let  $p, q \in [1, \infty]$  be Hölder conjugate  $\frac{1}{p} + \frac{1}{q} = 1$ .

The  $(p, q)$ -norm of a WFA  $A$  is given by

$$\|A\|_{p,q} = \max \left\{ \|\alpha\|_p, \|\beta\|_q, \max_{\sigma \in \Sigma} \|\mathbf{A}_\sigma\|_q \right\},$$

where  $\|\mathbf{A}_\sigma\|_q = \sup_{\|\mathbf{v}\|_q \leq 1} \|\mathbf{A}_\sigma \mathbf{v}\|_q$  is the  $q$ -induced norm.

**Example** For probabilistic automata  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  with  $\alpha$  probability distribution,  $\beta$  acceptance probabilities,  $\mathbf{A}_\sigma$  row (sub-)stochastic matrices we have  $\|A\|_{1,\infty} = 1$

## Perturbation Bounds: Automaton $\rightarrow$ Language [Bal13]

Suppose  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  and  $A' = \langle \alpha', \beta', \{\mathbf{A}'_\sigma\} \rangle$  are WFA with  $n$  states satisfying  $\|A\|_{p,q} \leq \rho$ ,  $\|A'\|_{p,q} \leq \rho$ ,  $\max \{ \|\alpha - \alpha'\|_p, \|\beta - \beta'\|_q, \max_{\sigma \in \Sigma} \|\mathbf{A}_\sigma - \mathbf{A}'_\sigma\|_q \} \leq \Delta$ .

Claim The following holds for any  $x \in \Sigma^*$ :

$$|f_A(x) - f_{A'}(x)| \leq (|x| + 2)\rho^{|x|+1}\Delta .$$

Proof By induction on  $|x|$  we first prove  $\|\mathbf{A}_x - \mathbf{A}'_x\|_q \leq |x|\rho^{|x|-1}\Delta$ :

$$\|\mathbf{A}_{x\sigma} - \mathbf{A}'_{x\sigma}\|_q \leq \|\mathbf{A}_x - \mathbf{A}'_x\|_q \|\mathbf{A}_\sigma\|_q + \|\mathbf{A}'_x\|_q \|\mathbf{A}_\sigma - \mathbf{A}'_\sigma\|_q \leq |x|\rho^{|x|}\Delta + \rho^{|x|}\Delta = (|x| + 1)\rho^{|x|}\Delta .$$

$$\begin{aligned} |f_A(x) - f_{A'}(x)| &= |\alpha^\top \mathbf{A}_x \beta - \alpha'^\top \mathbf{A}'_x \beta'| \leq |\alpha^\top (\mathbf{A}_x \beta - \mathbf{A}'_x \beta')| + |(\alpha - \alpha')^\top \mathbf{A}'_x \beta'| \\ &\leq \|\alpha\|_p \|\mathbf{A}_x \beta - \mathbf{A}'_x \beta'\|_q + \|\alpha - \alpha'\|_p \|\mathbf{A}'_x \beta'\|_q \\ &\leq \|\alpha\|_p \|\mathbf{A}_x\|_q \|\beta - \beta'\|_q + \|\alpha\|_p \|\mathbf{A}_x - \mathbf{A}'_x\|_q \|\beta'\|_q + \|\alpha - \alpha'\|_p \|\mathbf{A}'_x\|_q \|\beta'\|_q \\ &\leq \rho^{|x|+1} \|\beta - \beta'\|_q + \rho^2 \|\mathbf{A}_x - \mathbf{A}'_x\|_q + \rho^{|x|+1} \|\alpha - \alpha'\|_p \\ &\leq \rho^{|x|+1} \Delta + \rho^2 \rho^{|x|-1} |x| \Delta + \rho^{|x|+1} \Delta . \end{aligned}$$

# Norms on Languages

- ▶  $L_p$  norms ( $p \in [1, \infty]$ ),  $\gamma$ -discounted  $L_p$  norms ( $\gamma \in (0, 1)$ )

$$\|f\|_p = \left( \sum_x |f(x)|^p \right)^{1/p} \quad \|f\|_{p,\gamma} = \left( \sum_x \gamma^{p|x|} |f(x)|^p \right)^{1/p}$$

- ▶ Dirichlet norm

$$\|f\|_D = \left( \sum_x (|x| + 1) |f(x)|^2 \right)^{1/2}$$

- ▶ Bisimulation norms **[FZ14, BGP17]**

$$\|f\|_{\infty,\gamma} = \sup_{x \in \Sigma^*} \gamma^{|x|} |f(x)| \quad \|f\|_B = \sup_{x \in \Sigma^\infty} \sum_{k \geq 0} \gamma^k |f(x_{\leq k})|$$

## Aside: Banach and Hilbert Spaces

- ▶ A (possibly infinite-dimensional) vector space  $\mathcal{X}$  equipped with a norm  $\|\bullet\| : \mathcal{X} \rightarrow [0, \infty)$  is a *Banach space* if the pair  $(\mathcal{X}, \|\bullet\|)$  is complete, i.e. Cauchy sequences converge.
  - ▶ Examples:  $\ell_p = \{f : \Sigma^* \rightarrow \mathbb{R} : \|f\|_p < \infty\}$
  - ▶ Exercise: the set of rational  $f \in \ell_p$  is dense in  $\ell_p$  for any  $p \in [1, \infty]$
- ▶ A (real) *Hilbert space* is a Banach space  $(\mathcal{X}, \|\bullet\|)$  equipped with an inner product  $\langle \bullet, \bullet \rangle : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  such that  $\|\mathbf{v}\| = \sqrt{\langle \mathbf{v}, \mathbf{v} \rangle}$ 
  - ▶ Example:  $\ell_2$  with  $\|f\|_2^2 = \langle f, f \rangle = \sum_{x \in \Sigma^*} f(x)^2$
  - ▶ Example  $\ell_D = \{f : \|f\|_D < \infty\}$  with  $\|f\|_D^2 = \langle f, f \rangle_D = \sum_{x \in \Sigma^*} (|x| + 1)f(x)^2$
- ▶ A Hilbert space is *separable* if it admits a countable orthonormal basis.
  - ▶ Examples:  $\ell_2$  and  $\ell_D$  are separable

## Perturbation Bounds: Language $\rightarrow$ Hankel

Consider the Hilbert space  $\ell_D = \{f : \Sigma^* \rightarrow \mathbb{R} : \|f\|_D < \infty\}$  with the Dirichlet inner product

$$\langle f, g \rangle_D = \sum_{x \in \Sigma^*} (|x| + 1) f(x) g(x) .$$

Consider the Frobenius norm on matrices  $\mathbf{T} \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$  given by

$$\|\mathbf{T}\|_F = \sqrt{\sum_{x, y \in \Sigma^*} \mathbf{T}(x, y)^2} .$$

Claim If  $f, f' \in \ell_D$  are two weighted languages such that  $\|f - f'\|_D \leq \Delta$ , then their corresponding Hankel matrices satisfy  $\|\mathbf{H}_f - \mathbf{H}_{f'}\|_F \leq \Delta$ .

Proof

$$\begin{aligned} \|\mathbf{H}_f - \mathbf{H}_{f'}\|_F^2 &= \sum_{x, y \in \Sigma^*} (\mathbf{H}_f(x, y) - \mathbf{H}_{f'}(x, y))^2 = \sum_{x, y \in \Sigma^*} (f(x \cdot y) - f'(x \cdot y))^2 \\ &= \sum_{z \in \Sigma^*} (|z| + 1) (f(z) - f'(z))^2 = \|f - f'\|_D^2 \end{aligned}$$

## Aside: Singular Value Decomposition (SVD)

For any  $\mathbf{M} \in \mathbb{R}^{n \times m}$  with  $\text{rank}(\mathbf{M}) = k$  there exists a *singular value decomposition*

$$\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{V}^T = \sum_{i=1}^k s_i \mathbf{u}_i \mathbf{v}_i^T$$

- ▶  $\mathbf{D} \in \mathbb{R}^{k \times k}$  diagonal contains  $k$  sorted *singular values*  $s_1 \geq s_2 \geq \dots \geq s_k > 0$
- ▶  $\mathbf{U} \in \mathbb{R}^{n \times k}$  contains  $k$  *left singular vectors*, i.e. orthonormal columns  $\mathbf{U}^T \mathbf{U} = \mathbf{I}$
- ▶  $\mathbf{V} \in \mathbb{R}^{m \times k}$  contains  $k$  *right singular vectors*, i.e. orthonormal columns  $\mathbf{V}^T \mathbf{V} = \mathbf{I}$

### Properties of SVD

- ▶  $\mathbf{M} = (\mathbf{U}\mathbf{D}^{1/2})(\mathbf{D}^{1/2}\mathbf{V}^T)$  is a rank factorization
- ▶ Can be used to compute the pseudo-inverse as  $\mathbf{M}^+ = \mathbf{V}\mathbf{D}^{-1}\mathbf{U}^T$
- ▶ Provides optimal low-rank approximations. For  $k' < k$ ,  $\mathbf{M}_{k'} = \mathbf{U}_{k'}\mathbf{D}_{k'}\mathbf{V}_{k'}^T = \sum_{i=1}^{k'} s_i \mathbf{u}_i \mathbf{v}_i^T$  satisfies

$$\mathbf{M}_{k'} \in \underset{\text{rank}(\hat{\mathbf{M}}) \leq k'}{\text{argmin}} \|\mathbf{M} - \hat{\mathbf{M}}\|_2$$

# Perturbation Bounds: Hankel $\rightarrow$ Automaton [Bal13]

- ▶ Suppose  $f : \Sigma^* \rightarrow \mathbb{R}$  has rank  $n$  and  $\epsilon \in \mathcal{P}, \mathcal{S} \subset \Sigma^*$  are such that the sub-block  $\mathbf{H} \in \mathbb{R}^{\mathcal{P} \times \mathcal{S}}$  of  $\mathbf{H}_f$  satisfies  $\text{rank}(\mathbf{H}) = n$
- ▶ Let  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  be obtained as follows:
  1. Compute the **SVD factorization**  $\mathbf{H} = \mathbf{P}\mathbf{S}$ ; i.e.  $\mathbf{P} = \mathbf{U}\mathbf{D}^{1/2}$  and  $\mathbf{S} = \mathbf{D}^{1/2}\mathbf{V}^\top$
  2. Let  $\alpha^\top$  (resp.  $\beta$ ) be the  $\epsilon$ -row of  $\mathbf{P}$  (resp.  $\epsilon$ -column of  $\mathbf{S}$ )
  3. Let  $\mathbf{A}_\sigma = \mathbf{P}^+\mathbf{H}_\sigma\mathbf{S}^+$ , where  $\mathbf{H}_\sigma \in \mathbb{R}^{\mathcal{P} \cdot \sigma \times \mathcal{S}}$  is a sub-block of  $\mathbf{H}_f$
- ▶ Suppose  $\hat{\mathbf{H}} \in \mathbb{R}^{\mathcal{P} \times \mathcal{S}}$  and  $\hat{\mathbf{H}}_\sigma \in \mathbb{R}^{\mathcal{P} \cdot \sigma \times \mathcal{S}}$  satisfy  $\max\{\|\mathbf{H} - \hat{\mathbf{H}}\|_2, \max_\sigma \|\mathbf{H}_\sigma - \hat{\mathbf{H}}_\sigma\|_2\} \leq \Delta$
- ▶ Let  $\hat{A} = \langle \hat{\alpha}, \hat{\beta}, \{\hat{\mathbf{A}}_\sigma\} \rangle$  be obtained as follows:
  1. Compute the **SVD rank- $n$  approximation**  $\hat{\mathbf{H}} \approx \hat{\mathbf{P}}\hat{\mathbf{S}}$ ; i.e.  $\hat{\mathbf{P}} = \hat{\mathbf{U}}_n\hat{\mathbf{D}}_n^{1/2}$  and  $\hat{\mathbf{S}} = \hat{\mathbf{D}}_n^{1/2}\hat{\mathbf{V}}_n^\top$
  2. Let  $\hat{\alpha}^\top$  (resp.  $\hat{\beta}$ ) be the  $\epsilon$ -row of  $\hat{\mathbf{P}}$  (resp.  $\epsilon$ -column of  $\hat{\mathbf{S}}$ )
  3. Let  $\hat{\mathbf{A}}_\sigma = \hat{\mathbf{P}}^+\hat{\mathbf{H}}_\sigma\hat{\mathbf{S}}^+$

Claim For any pair of Hölder conjugate  $(p, q)$  we have

$$\max\{\|\alpha - \hat{\alpha}\|_p, \|\beta - \hat{\beta}\|_q, \max_\sigma \|\mathbf{A}_\sigma - \hat{\mathbf{A}}_\sigma\|_q\} \leq \mathcal{O}(\Delta)$$

## Applications

- ▶ Analysis of machine learning algorithms for WFA [BM12, BCLQ14, BM17]
- ▶ Statistical properties of classes of WFA (e.g. Rademacher complexity) [BM15, BM18]
- ▶ Continuity of operations on WFA and rational languages [BGP17]

## Limitations

- ▶ Automaton  $\rightarrow$  Language: grow with  $|x|$ , depend on representation chosen for  $A$
- ▶ Language  $\rightarrow$  Hankel: only applies to restricted choice of norms (?)
- ▶ Hankel  $\rightarrow$  Automaton: depends on algorithm, cumbersome to prove



1. Weighted Languages, Weighted Automata, and Hankel Matrices
2. Perturbation Bounds Between Representations
- 3. Singular Value Automata: Definition**
4. Singular Value Automata: Computation
5. Approximate Minimization via SVA Truncation
6. Concluding Remarks

# Motivation: Approximate Minimization

- ▶ Suppose  $f$  is a weighted language with  $\text{rank}(f) = n$  and  $\|f\| < \infty$
- ▶ Problem Given  $\hat{n} < n$  find  $\hat{f}$  with  $\text{rank}(\hat{f}) = \hat{n}$  such that

$$\|f - \hat{f}\| \approx \min_{\text{rank}(f') \leq \hat{n}} \|f - f'\|$$

- ▶ Typically,  $f$  is given by a minimal WFA  $A$  and the output is a WFA  $\hat{A}$  with  $|\hat{A}| = \hat{n}$
- ▶ The techniques described so far are too brittle to solve this problem!

## Aside: Operators on Hilbert Spaces

- ▶ Let  $\mathcal{X}_1, \mathcal{X}_2$  be a separable Hilbert spaces. Any linear operator  $\mathbf{T} : \mathcal{X}_1 \rightarrow \mathcal{X}_2$  can be represented as an infinite matrix
- ▶ A linear operator  $\mathbf{T} : \mathcal{X}_1 \rightarrow \mathcal{X}_2$  is *bounded* if  $\|\mathbf{T}\|_{\text{op}} = \sup_{\|\mathbf{v}\|_{\mathcal{X}_1} \leq 1} \|\mathbf{T}\mathbf{v}\|_{\mathcal{X}_2} < \infty$
- ▶ The adjoint  $\mathbf{T}^* : \mathcal{X}_2 \rightarrow \mathcal{X}_1$  of a bounded linear operator  $\mathbf{T}$  is given by  $\langle \mathbf{T}\mathbf{u}, \mathbf{v} \rangle_{\mathcal{X}_2} = \langle \mathbf{u}, \mathbf{T}^*\mathbf{v} \rangle_{\mathcal{X}_1}$
- ▶ A bounded linear operator  $\mathbf{T}$  is *compact* if it is the limit of a sequence of finite-rank operators (w.r.t. the topology induced by  $\|\bullet\|_{\text{op}}$ ).
  - ▶ Example: all finite-rank operators are compact
- ▶ Compact linear operators  $\mathbf{T}$  admit SVD (a.k.a. Hilbert–Schmidt decomposition)

$$\mathbf{T} = \mathbf{U}\mathbf{D}\mathbf{V}^* = \sum_{i=1}^k \mathfrak{s}_i \mathbf{u}_i \langle \mathbf{v}_i, \bullet \rangle_{\mathcal{X}_1} .$$

Here  $k = \text{rank}(\mathbf{T}) \leq \infty$ , and if  $k = \infty$  then  $\lim_{i \rightarrow \infty} \mathfrak{s}_i = 0$ .

- ▶ Finite-rank bounded operators  $\mathbf{T}$  admit a pseudo-inverse  $\mathbf{T}^+$

# Hankel Operators

A Hankel matrix  $\mathbf{H}_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$  can be interpreted as a linear operator  $\mathbf{H}_f : \mathbb{R}^{\Sigma^*} \rightarrow \mathbb{R}^{\Sigma^*}$ :

$$(\mathbf{H}_f g)(x) = \sum_{y \in \Sigma^*} f(x \cdot y) g(y) .$$

- ▶ **Fliess–Kronecker:** Finite rank if and only if  $f$  rational
- ▶ When does it admit an SVD? When it is a compact operator on a Hilbert space!

## Shift Characterization

- ▶ Define the forward/backward left/right shift operators  $\mathbf{L}_\sigma, \mathbf{L}_\sigma^*, \mathbf{R}_\sigma, \mathbf{R}_\sigma^* : \mathbb{R}^{\Sigma^*} \rightarrow \mathbb{R}^{\Sigma^*}$  as:  
 $(\mathbf{L}_\sigma^* f)(x) = f(\sigma x), (\mathbf{R}_\sigma^* f)(x) = f(x\sigma)$

$$(\mathbf{L}_\sigma f)(x) = \begin{cases} f(\sigma^{-1}x) & x_1 = \sigma \\ 0 & \text{otherwise} \end{cases} \quad (\mathbf{R}_\sigma f)(x) = \begin{cases} f(x\sigma^{-1}) & x_{|x|} = \sigma \\ 0 & \text{otherwise} \end{cases}$$

- ▶ Exercise A linear operator  $\mathbf{T} : \mathbb{R}^{\Sigma^*} \rightarrow \mathbb{R}^{\Sigma^*}$  is Hankel if and only if  $\mathbf{R}_\sigma^* \mathbf{T} = \mathbf{T} \mathbf{L}_\sigma, \forall \sigma \in \Sigma$

## Aside: Operator-Theoretic Proof of Fliess' Theorem

Claim Suppose  $\mathbf{H}_f : \ell_2 \rightarrow \ell_2$  is bounded and has finite rank  $n$ . Then there exists a WFA  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  with  $n$  states such that  $f_A = f$

### Proof

Take a rank factorization  $\mathbf{H}_f = \mathbf{P}\mathbf{S}$  and note  $\mathbf{P}$  and  $\mathbf{S}$  are bounded and finite rank. Build the automaton  $A$  by taking:

- ▶  $\alpha^\top$  the  $\epsilon$ -row of  $\mathbf{P}$ ; i.e.  $\alpha^\top = \mathbf{P}(\epsilon, -)$
- ▶  $\beta$  the  $\epsilon$ -column of  $\mathbf{S}$ ; i.e.  $\beta = \mathbf{S}(-, \epsilon)$
- ▶  $\mathbf{A}_\sigma = \mathbf{S}\mathbf{L}_\sigma\mathbf{S}^+$

It suffices to show that for any  $x \in \Sigma^*$  we have  $\alpha^\top \mathbf{A}_x = \mathbf{P}(x, -)$ . By induction on length of  $x$ :

$$\begin{aligned}\alpha^\top \mathbf{A}_x \mathbf{A}_\sigma &= \mathbf{P}(x, -) \mathbf{S} \mathbf{L}_\sigma \mathbf{S}^+ = \Pi_x \mathbf{P} \mathbf{S} \mathbf{L}_\sigma \mathbf{S}^+ = \Pi_x \mathbf{H}_f \mathbf{L}_\sigma \mathbf{S}^+ = \Pi_x \mathbf{R}_\sigma^* \mathbf{H}_f \mathbf{S}^+ \\ &= \Pi_x \mathbf{R}_\sigma^* \mathbf{P} \mathbf{S} \mathbf{S}^+ = \Pi_x \mathbf{R}_\sigma^* \mathbf{P} = \Pi_{x\sigma} \mathbf{P} = \mathbf{P}(x\sigma, -)\end{aligned}$$

# Which Hankel Operators Admit an SVD?

A Hankel matrix  $\mathbf{H}_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$  can be interpreted as a linear operator  $\mathbf{H}_f : \mathbb{R}^{\Sigma^*} \rightarrow \mathbb{R}^{\Sigma^*}$ :

$$(\mathbf{H}_f g)(x) = \sum_{y \in \Sigma^*} f(x \cdot y) g(y) .$$

- ▶ **Fliess–Kronecker:** Finite rank if and only if  $f$  rational
- ▶ When does it admit an SVD? When it is a compact operator on a Hilbert space!
- ▶ Finite rank operators are compact if and only if they are bounded:  
 $\|\mathbf{H}_f\|_{op} = \sup_{\|g\|_2 \leq 1} \|\mathbf{H}_f g\|_2 < \infty$
- ▶ When is a finite rank Hankel operator bounded?

## Boundedness of $\ell_2$ and Dirichlet Norms

Claim Suppose  $f : \Sigma^* \rightarrow \mathbb{R}$  is rational. Then  $\|f\|_2 < \infty$  if and only if  $\|f\|_D < \infty$

Proof One direction is easy:

$$\|f\|_2^2 = \sum_{x \in \Sigma^*} f(x)^2 \leq \sum_{x \in \Sigma^*} (|x| + 1) f(x)^2 = \|f\|_D^2 .$$

The other direction is more technical. Let  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  be a minimal WFA for  $f^2$  with  $n$  states. Then one can show that the spectral radius of  $\mathbf{A} = \sum_\sigma \mathbf{A}_\sigma$  satisfies  $\rho = \rho(\mathbf{A}) < 1$  (see [BPP17]).

$$\begin{aligned} \sum_{x \in \Sigma^t} f(x)^2 &= \sum_{x \in \Sigma^t} \alpha^\top \mathbf{A}_x \beta = \alpha^\top (\mathbf{A}_{\sigma_1} + \cdots + \mathbf{A}_{\sigma_k}) \cdots (\mathbf{A}_{\sigma_1} + \cdots + \mathbf{A}_{\sigma_k}) \beta \\ &= \alpha^\top \mathbf{A}^t \beta \leq \mathcal{O}(t^n \rho^t) . \end{aligned}$$

Therefore, since  $\rho < 1$  we have

$$\|f\|_D^2 = \sum_{x \in \Sigma^*} (|x| + 1) f(x)^2 = \sum_{t \geq 0} \sum_{x \in \Sigma^t} (t + 1) \alpha^\top \mathbf{A}^t \beta \leq \sum_{t \geq 0} \mathcal{O}(t^{n+1} \rho^t) < \infty .$$

# Bounded Hankel Operators of Finite Rank

Let  $\mathbf{H}_f : \ell_2 \rightarrow \ell_2$  be a finite rank Hankel operator.

Theorem The operator  $\mathbf{H}_f$  is bounded if and only if  $f \in \ell_2$ .

Proof Since  $f$  is the first row of  $\mathbf{H}_f$ , from  $\mathbf{H}_f$  bounded to  $\|f\|_2 < \infty$  is easy:

$$\infty > \|\mathbf{H}_f\|_{op} = \sup_{\|g\|_2 \leq 1} \|\mathbf{H}_f g\|_2 \geq \|\mathbf{H}_f \mathbf{e}_\epsilon\|_2 = \|f\|_2 .$$

The other direction uses the boundedness of the Dirichlet norm: let  $\|g\|_2 \leq 1$ , then

$$\begin{aligned} \|H_f g\|_2^2 &= \sum_{x \in \Sigma^*} \left( \sum_{y \in \Sigma^*} f(x \cdot y) g(y) \right)^2 = \sum_{x \in \Sigma^*} \langle \mathbf{L}_x^* f, g \rangle^2 \\ &\leq \|g\|_2^2 \sum_{x \in \Sigma^*} \|\mathbf{L}_x^* f\|_2^2 \leq \sum_{x \in \Sigma^*} \|\mathbf{L}_x^* f\|_2^2 \\ &= \sum_{x \in \Sigma^*} \sum_{y \in \Sigma^*} f(x \cdot y)^2 = \sum_{z \in \Sigma^*} (|z| + 1) f(z)^2 = \|f\|_D^2 < \infty . \end{aligned}$$



# Are We Done Yet?

## Approximate Minimization Strategy

1. Take rational  $f$  with  $\text{rank}(f) = n$  and  $\|f\|_2 < \infty$
2. Since  $\mathbf{H}_f : \ell_2 \rightarrow \ell_2$  is compact, it admits an SVD

$$\mathbf{H}_f = \sum_{i=1}^n s_i \mathbf{u}_i \langle \mathbf{v}_i, \bullet \rangle .$$

3. Given  $\hat{n} < n$  take the corresponding low-rank approximation  $\hat{\mathbf{H}}$

$$\hat{\mathbf{H}} = \sum_{i=1}^{\hat{n}} s_i \mathbf{u}_i \langle \mathbf{v}_i, \bullet \rangle .$$

4. Compute a WFA  $\hat{A}$  from  $\hat{\mathbf{H}}$  ← **NOT NECESSARILY HANKEL!**
5. Bound the error between  $f$  and  $\hat{f} = f_{\hat{A}}$  as

$$\|f - \hat{f}\|_2 \leq \|\mathbf{H}_f - \hat{\mathbf{H}}\|_{op} = s_{\hat{n}+1} .$$

## Duality Between Rank Factorization and Minimal WFA

Well-known fact: If  $\mathbf{M}$  has rank  $n$  and  $\mathbf{M} = \mathbf{P}\mathbf{S} = \mathbf{P}'\mathbf{S}'$  are two rank factorizations, then there exists invertible  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that

$$\mathbf{P}' = \mathbf{P}\mathbf{Q} \quad \mathbf{S}' = \mathbf{Q}^{-1}\mathbf{S}$$

Well-known fact: If  $A = \langle \boldsymbol{\alpha}, \boldsymbol{\beta}, \{\mathbf{A}_\sigma\} \rangle$  and  $A' = \langle \boldsymbol{\alpha}', \boldsymbol{\beta}', \{\mathbf{A}'_\sigma\} \rangle$  are minimal WFA for  $f$  of rank  $n$ , then there exists invertible  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that

$$\boldsymbol{\alpha}'^T = \boldsymbol{\alpha}^T \mathbf{Q} \quad \boldsymbol{\beta}' = \mathbf{Q}^{-1} \boldsymbol{\beta} \quad \mathbf{A}'_\sigma = \mathbf{Q}^{-1} \mathbf{A}_\sigma \mathbf{Q}$$

Less-known fact: From the proof of the Fliess–Kronecker theorem applied to  $f$  of rank  $n$  one obtains a bijection

$$\{(\mathbf{P}, \mathbf{S}) : \mathbf{H}_f = \mathbf{P}\mathbf{S}, \text{rank}(\mathbf{P}) = \text{rank}(\mathbf{S}) = n\} \leftrightarrow \{A = \langle \boldsymbol{\alpha}, \boldsymbol{\beta}, \{\mathbf{A}_\sigma\} \rangle : f_A = f, |A| = n\}$$

# Singular Value Automata

- ▶ Let  $A$  be a minimal WFA with  $n$  states computing  $f$
- ▶ Definition  $A$  is a *singular value automaton* (SVA) if the forward-backward factorization  $\mathbf{H}_f = \mathbf{P}_A \mathbf{S}_A$  comes from a singular value decomposition, i.e.  $\mathbf{P}_A = \mathbf{U} \mathbf{D}^{1/2}$ ,  $\mathbf{S}_A = \mathbf{D}^{1/2} \mathbf{V}^T$ , with  $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}$  and  $\mathbf{D} = \text{diag}(s_1, \dots, s_n)$  with  $s_1 \geq \dots \geq s_n > 0$
- ▶ Theorem Every rational  $f$  with  $\|f\|_2 < \infty$  admits an SVA
- ▶ The SVA of  $f$  is “as unique” as the SVD of  $\mathbf{H}_f$ 
  - ▶ Example: if all inequalities between singular values are strict, SVD is unique up to sign changes in pairs of associated left/right singular vectors  $\Rightarrow$  SVA unique up to sign changes in pairs of associated initial/final weights
- ▶ Given a minimal WFA  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  for  $f$  with  $\|f\|_2 < \infty$  there exists an invertible  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  such that  $A^{\mathbf{Q}} = \langle \mathbf{Q}^T \alpha, \mathbf{Q}^{-1} \beta, \{\mathbf{Q}^{-1} \mathbf{A}_\sigma \mathbf{Q}\} \rangle$  is an SVA for  $f$
- ▶ Definition could be changed to have  $\mathbf{P}_A = \mathbf{U}$  and  $\mathbf{S}_A = \mathbf{D} \mathbf{V}^T$ , or  $\mathbf{P}_A = \mathbf{U} \mathbf{D}$  and  $\mathbf{S}_A = \mathbf{V}^T$ . But the current one makes computation of  $\mathbf{Q}$  above more “symmetric”

## Why Are SVA Special?

- ▶ It *orthogonalizes* the states of a WFA!
- ▶ Suppose  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  is an SVA with  $n$  states for  $f$  inducing the SVD

$$\mathbf{H}_f = \sum_{i=1}^n s_i \mathbf{u}_i \langle \mathbf{v}_i, \bullet \rangle .$$

- ▶ For  $i \in [n]$  let  $A_i = \langle \alpha, \mathbf{e}_i, \{\mathbf{A}_\sigma\} \rangle$  where  $\mathbf{e}_i = (0, \dots, 1, \dots, 0)$  is the  $i$ th coordinate vector
- ▶ The language  $f_i$  of  $A_i$  is given by  $f_i(x) = \alpha^\top \mathbf{A}_x \mathbf{e}_i = \alpha_A(x)^\top [i]$ ; i.e. is the “memory” of state  $i$  after reading  $x$
- ▶ The language  $f_i$  is also the  $i$ th column of the forward matrix  $\mathbf{P}_A = \mathbf{U}\mathbf{D}^{1/2}$ ; i.e.  $f_i = \sqrt{s_i} \mathbf{u}_i$
- ▶ Since the columns of  $\mathbf{U}$  are orthonormal, the languages  $f_i$  and  $f_j$  with  $i \neq j$  are orthogonal

1. Weighted Languages, Weighted Automata, and Hankel Matrices
2. Perturbation Bounds Between Representations
3. Singular Value Automata: Definition
4. Singular Value Automata: Computation
5. Approximate Minimization via SVA Truncation
6. Concluding Remarks

# The Gramians of a WFA

- ▶ Let  $A$  be a minimal WFA for  $f$  with  $n = \text{rank}(f)$  inducing the rank factorization  $\mathbf{H}_f = \mathbf{P}\mathbf{S}$  (i.e.  $\mathbf{P} = \mathbf{P}_A$  and  $\mathbf{S} = \mathbf{S}_A$ )
- ▶ The *reachability Gramian* of  $A$  is the (possibly infinite)  $n \times n$  matrix  $\mathbf{G}_p = \mathbf{P}^\top \mathbf{P}$

$$\mathbf{G}_p = \mathbf{P}^\top \mathbf{P} = \sum_{x \in \Sigma^*} \mathbf{P}(x, -)^\top \mathbf{P}(x, -) = \sum_{x \in \Sigma^*} (\boldsymbol{\alpha}^\top \mathbf{A}_x)^\top (\boldsymbol{\alpha}^\top \mathbf{A}_x)$$

- ▶ The *observability Gramian* of  $A$  is the (possibly infinite)  $n \times n$  matrix  $\mathbf{G}_s = \mathbf{S}\mathbf{S}^\top$  given by

$$\mathbf{G}_s = \mathbf{S}\mathbf{S}^\top = \sum_{x \in \Sigma^*} \mathbf{S}(-, x) \mathbf{S}(-, x)^\top = \sum_{x \in \Sigma^*} (\mathbf{A}_x \boldsymbol{\beta}) (\mathbf{A}_x \boldsymbol{\beta})^\top$$

## Existence of the Gramians

Let  $A$  be a minimal WFA for  $f$  with  $n = \text{rank}(f)$  inducing the rank factorization  $\mathbf{H}_f = \mathbf{P}\mathbf{S}$  (i.e.  $\mathbf{P} = \mathbf{P}_A$  and  $\mathbf{S} = \mathbf{S}_A$ )

Claim The Gramians of  $A$  are finite if and only if  $\|f\|_2 < \infty$

Proof (one direction only)

Suppose  $\|f\|_2 < \infty$  and let  $A' = A^Q = \langle \mathbf{Q}^\top \boldsymbol{\alpha}, \mathbf{Q}^{-1} \boldsymbol{\beta}, \{\mathbf{Q}^{-1} \mathbf{A}_\sigma \mathbf{Q}\} \rangle$  be an SVA for  $f$ . Observe the Gramians  $\mathbf{G}'_p$  and  $\mathbf{G}'_s$  of  $A'$  exist since

$$\mathbf{G}'_p = \mathbf{P}_{A'}^\top \mathbf{P}_{A'} = \mathbf{D}^{1/2} \mathbf{U}^\top \mathbf{U} \mathbf{D}^{1/2} = \mathbf{D}$$

$$\mathbf{G}'_s = \mathbf{S}_{A'} \mathbf{S}_{A'}^\top = \mathbf{D}^{1/2} \mathbf{V}^\top \mathbf{V} \mathbf{D}^{1/2} = \mathbf{D}$$

On the other hand, since  $\mathbf{P}_{A'} = \mathbf{P}_A \mathbf{Q}$  and  $\mathbf{S}_{A'} = \mathbf{Q}^{-1} \mathbf{S}_A$  we have

$$\mathbf{G}'_p = \mathbf{Q}^\top \mathbf{G}_p \mathbf{Q} \quad \mathbf{G}'_s = \mathbf{Q}^{-\top} \mathbf{G}_s \mathbf{Q}^{-1}$$

Therefore  $\mathbf{G}_p$  and  $\mathbf{G}_s$  must be finite

# From Gramians to SVA

- ▶ Let  $A$  be a minimal WFA for  $f$  with  $\|f\|_2 < \infty$
- ▶ Suppose we have the Gramians of  $A$ :  $\mathbf{G}_p$  and  $\mathbf{G}_s$
- ▶ Recall from the previous proof that
  - ▶ If  $A'$  is SVA then  $\mathbf{G}'_p = \mathbf{G}'_s = \mathbf{D} = \text{diag}(s_1, \dots, s_n)$
  - ▶ If  $A' = A^Q$  then  $\mathbf{G}'_p = \mathbf{Q}^\top \mathbf{G}_p \mathbf{Q}$  and  $\mathbf{G}'_s = \mathbf{Q}^{-\top} \mathbf{G}_s \mathbf{Q}^{-1}$
- ▶ Claim The following algorithm returns  $\mathbf{Q}$  such that  $A^Q$  is an SVA
  1. Compute the Cholesky decompositions  $\mathbf{G}_p = \mathbf{L}_p \mathbf{L}_p^\top$  and  $\mathbf{G}_s = \mathbf{L}_s \mathbf{L}_s^\top$
  2. Compute the SVD decomposition  $\mathbf{L}_p^\top \mathbf{L}_s = \mathbf{U} \mathbf{D} \mathbf{V}^\top$
  3. Let  $\mathbf{Q} = \mathbf{L}_p^{-\top} \mathbf{U} \mathbf{D}^{1/2}$
- ▶ In particular, the  $\mathbf{D}$  in this algorithm is the matrix of singular values of  $\mathbf{H}_f$
- ▶ See proof in [BPP17]



# Computing Norms Using Gramians

Suppose  $A$  is a minimal WFA for  $f$  with  $\|f\|_2 < \infty$ .

Let  $\mathbf{G}_p$  and  $\mathbf{G}_s$  be the Gramians of  $A$ .

Then the following hold:

- ▶  $\|f\|_2^2 = \alpha^\top \mathbf{G}_s \alpha = \beta^\top \mathbf{G}_p \beta$
- ▶  $\|f\|_D^2 = \|\mathbf{H}_f\|_F^2 = \text{Tr}(\mathbf{G}_p \mathbf{G}_s)$
- ▶  $\|\mathbf{H}_f\|_{op}^2 = \rho(\mathbf{G}_p \mathbf{G}_s) = \max\{|\lambda| : \det(\mathbf{G}_p \mathbf{G}_s - \lambda \mathbf{I}) = 0\}$

# Computing the Gramians Using Fixed-Points

Let  $A$  be a minimal WFA for  $f$  with  $\|f\|_2 < \infty$ .

Claim  $\mathbf{X} = \mathbf{G}_p$  and  $\mathbf{Y} = \mathbf{G}_s$  are solutions of the fixed-point equations

$$\mathbf{X} = F_p(\mathbf{X}) = \alpha\alpha^\top + \sum_{\sigma} \mathbf{A}_{\sigma}^\top \mathbf{X} \mathbf{A}_{\sigma} \quad \mathbf{Y} = F_s(\mathbf{Y}) = \beta\beta^\top + \sum_{\sigma} \mathbf{A}_{\sigma} \mathbf{Y} \mathbf{A}_{\sigma}^\top$$

Proof Recall  $\mathbf{G}_p = \mathbf{P}_A^\top \mathbf{P}_A = \sum_{x \in \Sigma^*} \mathbf{P}_A(x, -) \mathbf{P}_A(x, -)^\top$  and  $\mathbf{P}_A(x, -) = \alpha^\top \mathbf{A}_x$ . Therefore:

$$\begin{aligned} \mathbf{G}_p &= \sum_{x \in \Sigma^*} (\mathbf{A}_x^\top \alpha)(\alpha^\top \mathbf{A}_x) = \alpha\alpha^\top + \sum_{x \in \Sigma^+} (\mathbf{A}_x^\top \alpha)(\alpha^\top \mathbf{A}_x) \\ &= \alpha\alpha^\top + \sum_{\sigma \in \Sigma} \sum_{x \in \Sigma^*} \mathbf{A}_{\sigma}^\top (\mathbf{A}_x^\top \alpha)(\alpha^\top \mathbf{A}_x) \mathbf{A}_{\sigma} \\ &= \alpha\alpha^\top + \sum_{\sigma \in \Sigma} \mathbf{A}_{\sigma}^\top \left( \sum_{x \in \Sigma^*} (\mathbf{A}_x^\top \alpha)(\alpha^\top \mathbf{A}_x) \right) \mathbf{A}_{\sigma} = \alpha\alpha^\top + \sum_{\sigma \in \Sigma} \mathbf{A}_{\sigma}^\top \mathbf{G}_p \mathbf{A}_{\sigma} \end{aligned}$$

# Solving the Fixed-Point Equations

- ▶ Recall the reachability Gramian  $\mathbf{G}_p$  is a solution of

$$\mathbf{X} = F_p(\mathbf{X}) = \alpha\alpha^\top + \sum_{\sigma} \mathbf{A}_{\sigma}^\top \mathbf{X} \mathbf{A}_{\sigma}$$

- ▶ Let  $\rho$  be the spectral radius of  $\sum_{\sigma} \mathbf{A}_{\sigma} \otimes \mathbf{A}_{\sigma}$ , where  $\otimes$  denotes the Kronecker product (i.e.  $\mathbf{A}_{\sigma} \otimes \mathbf{A}_{\sigma} \in \mathbb{R}^{n^2 \times n^2}$ )
- ▶ We distinguish two cases. If  $\rho < 1$ :
  - ▶  $\mathbf{X} = F_p(\mathbf{X})$  has a *unique* solution
  - ▶ Can be found by solving the linear system with  $n^2$  unknowns obtained through vectorization:  $\text{vec}(\alpha\alpha^\top) = \alpha \otimes \alpha$  and  $\text{vec}(\mathbf{A}_{\sigma}^\top \mathbf{X} \mathbf{A}_{\sigma}) = (\mathbf{A}_{\sigma} \otimes \mathbf{A}_{\sigma})^\top \text{vec}(\mathbf{X})$
- ▶ If  $\rho \geq 1$ :
  - ▶  $\mathbf{X} = F_p(\mathbf{X})$  might have multiple solutions (there is at least one because  $\mathbf{G}_p$  is defined)
  - ▶ In this case rephrase the problem:  $\mathbf{G}_p$  is the least positive semi-definite solution of the linear matrix inequality  $\mathbf{X} \geq F_p(\mathbf{X})$
  - ▶ The solution can be found by semi-definite programming

# Computing SVA: Summary

Suppose  $A$  is a WFA computing a function  $f$ . To compute an SVA for  $f$  do:

1. Test if  $\|f\|_2 < \infty$
2. Minimize  $A$  if necessary
3. Compute Gramians  $\mathbf{G}_p$  and  $\mathbf{G}_s$  (using linear solver or semi-definite solver)
4. Find change of basis  $\mathbf{Q}$  through Cholesky and SVD of finite matrices
5. Return  $A^{\mathbf{Q}}$

## Final remarks

- ▶ Runs in time polynomial in  $|A|$  and  $|\Sigma|$
- ▶ Easy to implement in Python or MATLAB

1. Weighted Languages, Weighted Automata, and Hankel Matrices
2. Perturbation Bounds Between Representations
3. Singular Value Automata: Definition
4. Singular Value Automata: Computation
5. **Approximate Minimization via SVA Truncation**
6. Concluding Remarks

# Approximate Minimization with SVA

- ▶ Suppose  $f$  is a weighted language with  $\text{rank}(f) = n$  and  $\|f\|_2 < \infty$ . Let  $s_i$  be the singular values of  $\mathbf{H}_f$
- ▶ Problem Given  $\hat{n} < n$  find  $\hat{f}$  with  $\text{rank}(\hat{f}) = \hat{n}$  such that

$$\|f - \hat{f}\|_2 \approx \min_{\text{rank}(f') \leq \hat{n}} \|f - f'\|_2$$

- ▶ SVA Solution Compute SVA  $A$  for  $f$  and obtain  $\hat{A}$  by removing the last  $n - \hat{n}$  states

$$\|f - \hat{f}\|_2^2 \leq \sum_{i=\hat{n}+1}^n s_i^2$$

- ▶ Lower Bound Considering approximation in terms of  $\|\bullet\|_D$  instead of  $\|\bullet\|_2$ :

$$\min_{\text{rank}(f') \leq \hat{n}} \|f - f'\|_D^2 \geq \sum_{i=\hat{n}+1}^n s_i^2$$

# Intuition for Removing the Last States from an SVA

- Suppose  $A = \langle \alpha, \beta, \{\mathbf{A}_\sigma\} \rangle$  is an SVA. Since the Gramians satisfy  $\mathbf{G}_p = \mathbf{G}_s = \mathbf{D} = \text{diag}(\varsigma_1, \dots, \varsigma_n)$ , we have

$$\mathbf{D} = \alpha\alpha^\top + \sum_{\sigma} \mathbf{A}_\sigma^\top \mathbf{D} \mathbf{A}_\sigma$$

$$\mathbf{D} = \beta\beta^\top + \sum_{\sigma} \mathbf{A}_\sigma \mathbf{D} \mathbf{A}_\sigma^\top$$

- By looking at the diagonal entries in these equations we can deduce

$$|\mathbf{A}_\sigma(i, j)| \leq \sqrt{\frac{\min\{\varsigma_i, \varsigma_j\}}{\max\{\varsigma_i, \varsigma_j\}}}$$

- For example, connections between the first and last state are weak:

$$|\mathbf{A}_\sigma(1, n)|, |\mathbf{A}_\sigma(n, 1)| \leq \sqrt{\varsigma_n/\varsigma_1}$$

- See [BPP15] for a “pedestrian” bound for  $\|f - \hat{f}\|_2$  based on this idea

# Analysis of SVA Approximate Minimization

SVA

$$\alpha = \begin{bmatrix} \alpha^{(1)} \\ \alpha^{(2)} \end{bmatrix},$$

$$\beta = \begin{bmatrix} \beta^{(1)} \\ \beta^{(2)} \end{bmatrix},$$

$$\mathbf{A}_\sigma = \begin{bmatrix} \mathbf{A}_\sigma^{(11)} & \mathbf{A}_\sigma^{(12)} \\ \mathbf{A}_\sigma^{(21)} & \mathbf{A}_\sigma^{(22)} \end{bmatrix}$$

Truncated SVA

$$\hat{\alpha} = \begin{bmatrix} \alpha^{(1)} \\ \mathbf{0} \end{bmatrix} = \mathbf{\Pi} \alpha,$$

$$\hat{\beta} = \begin{bmatrix} \beta^{(1)} \\ \beta^{(2)} \end{bmatrix} = \beta,$$

$$\hat{\mathbf{A}}_\sigma = \begin{bmatrix} \mathbf{A}_\sigma^{(11)} & \mathbf{0} \\ \mathbf{A}_\sigma^{(21)} & \mathbf{0} \end{bmatrix} = \mathbf{A}_\sigma \mathbf{\Pi}$$

$$\mathbf{\Pi} = \begin{bmatrix} \mathbf{I}_{\hat{n}} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}$$

## Analysis

- ▶ Let  $A$  be SVA for  $f$  and  $\hat{A}$  truncated SVA computing  $\hat{f}$
- ▶ Show  $\|\hat{f}\|_2 \leq \|f\|_2$  (see [BPP17])
- ▶ Show  $\|f - \hat{f}\|_2 \leq s_{\hat{n}+1}^2 + \dots + s_n^2$  (organic free-range proof on the board)



1. Weighted Languages, Weighted Automata, and Hankel Matrices
2. Perturbation Bounds Between Representations
3. Singular Value Automata: Definition
4. Singular Value Automata: Computation
5. Approximate Minimization via SVA Truncation
6. Concluding Remarks

# The Tree Case

- ▶ Take a ranked alphabet  $\Sigma = \Sigma_0 \cup \Sigma_1 \cup \dots$
- ▶ A weighted tree automaton with  $n$  states is a tuple  $A = \langle \alpha, \{\mathbf{T}_\tau\}_{\tau \in \Sigma_{\geq 1}}, \{\beta_\sigma\}_{\sigma \in \Sigma_0} \rangle$  where

$$\alpha, \beta_\sigma \in \mathbb{R}^n \quad \mathbf{T}_\tau \in (\mathbb{R}^n)^{\otimes \text{rk}(\tau)+1}$$

- ▶  $A$  defines a function  $f_A = \text{Trees}_\Sigma \rightarrow \mathbb{R}$  through recursive vector-tensor contractions
- ▶ There exists an analogue of the Hankel matrix for  $f : \text{Trees}_\Sigma \rightarrow \mathbb{R}$  where rows are indexed by contexts and columns by trees
- ▶ The same ideas lead to a notion of *singular value tree automata* [RBC16]
- ▶ In this case the computation of the Gramians is already a highly non-trivial problem

# The One Symbol Case

- ▶ When  $|\Sigma| = 1$ ,  $\Sigma^* = \mathbb{N}$  and one recovers the classical Hankel operators studied in complex analysis and the impulse responses studied in control theory and signal processing
- ▶ A new perspective in terms of functions of one complex variable arises from the power-series point of view: for  $z \in \mathbb{C}$  with small enough modulus

$$f(z) = \sum_{k \geq 0} a_k z^k = \sum_{k \geq 0} \alpha (z\mathbf{A})^k \beta = \alpha^\top (\mathbf{I} - z\mathbf{A})^{-1} \beta = \frac{p(z)}{q(z)}$$

- ▶  $\mathbb{N}$  can be embedded into a locally compact Abelian group  $\mathbb{Z}$ ,  $\ell_2$  gets a new definition in terms of Fourier analysis, Hankel operators get a new definition in terms of Hardy spaces, etc.
- ▶ Example: Nehari's theorem says that  $\|\mathbf{H}_f\|_{op} = \sup_{|z| < 1} |f(z)|$
- ▶ Suggested readings: Peller's "Hankel Operators and Their Applications" [Pel12] and Fuhrmann's "A Polynomial Approach to Linear Algebra" [Fuh11]

- ▶ Complexity of testing  $\|f\|_p < R$ , computing and approximating  $\ell_p$  and other norms on languages
- ▶ Complexity of optimal approximate minimization in terms of  $\|\bullet\|_2$
- ▶ Quality of approximation of SVA truncation in terms of  $\|\bullet\|_2$  or analysis of approximation in terms of  $\|\bullet\|_D$
- ▶ Approximate minimization with other norms

- ▶ **Analytic automata theory** is a vastly understudied area, rich in interesting open problems (for the mathematically adventurous)
- ▶ **Singular value automata** provide a powerful canonical form for WFA over the reals
- ▶ **Approximate minimization** is a generalization of automata minimization with connections to machine learning

Thanks!

# References I



B. Balle.

*Learning Finite-State Machines: Algorithmic and Statistical Aspects.*

PhD thesis, Universitat Politècnica de Catalunya, 2013.



B. Balle, X. Carreras, F.M. Luque, and A. Quattoni.

Spectral learning of weighted automata: A forward-backward perspective.

*Machine Learning*, 2014.



B. Balle, P. Gourdeau, and P. Panangaden.

Bisimulation metrics for weighted automata.

In *ICALP*, 2017.



B. Balle and M. Mohri.

Spectral learning of general weighted automata via constrained matrix completion.

In *NIPS*, 2012.



B. Balle and M. Mohri.

On the rademacher complexity of weighted automata.

In *ALT*, 2015.

## References II



B. Balle and O.-A. Maillard.

Spectral learning from a single trajectory under finite-state policies.

In *ICML*, 2017.



B. Balle and M. Mohri.

Generalization Bounds for Learning Weighted Automata.

*Theoretical Computer Science*, 716:89–106, 2018.



B. Balle, P. Panangaden, and D. Precup.

A canonical form for weighted automata and applications to approximate minimization.

In *LICS*, 2015.



Borja Balle, Prakash Panangaden, and Doina Precup.

Singular value automata and approximate minimization.

*CoRR*, abs/1711.05994, 2017.



M. Fliess.

Matrices de Hankel.

*Journal de Mathématiques Pures et Appliquées*, 1974.



# References III



Paul A Fuhrmann.

*A polynomial approach to linear algebra.*

Springer Science & Business Media, 2011.



Yuan Feng and Lijun Zhang.

When equivalence and bisimulation join forces in probabilistic automata.

In *International Symposium on Formal Methods*, pages 247–262. Springer, 2014.



Vladimir Peller.

*Hankel operators and their applications.*

Springer Science & Business Media, 2012.



G. Rabusseau, B. Balle, and S. B. Cohen.

Low-rank approximation of weighted tree automata.

In *AISTATS*, 2016.

# Singular Value Automata and Approximate Minimization

**Borja Balle**

Amazon Research Cambridge<sup>2</sup>

Weighted Automata: Theory and Applications — May 2018

---

<sup>2</sup>Based on work completed before joining Amazon