

# Spectral Learning of General Weighted Automata via Constrained Matrix Completion

*Borja Balle*<sup>1</sup>    Mehryar Mohri<sup>2</sup>

<sup>1</sup>Universitat Politècnica de Catalunya

<sup>2</sup>Courant Institute of Mathematical Sciences and Google Research

NIPS 2012, Lake Tahoe

## The Problem – Regression over Strings

Data: i.i.d. sample  $S$  with strings + real labels from distribution  $\mathcal{D}$

$$S = (\text{abbca}, 3.4) (\text{baa}, 0.6) (\text{ccaaaabba}, -2.9) (\text{abba}, 1.1) \dots$$

Goal: Learn a regressor  $\hat{f} : \Sigma^* \rightarrow \mathbb{R}$  with small generalization error

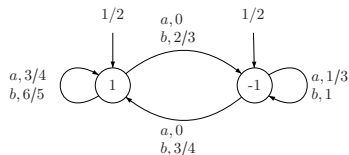
$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \ell(\hat{f}(x), y) \right]$$

Examples:

- ▶ *Reward modeling* in reinforcement learning
- ▶ *Biological measurement* as a function of DNA/AA sequence
- ▶ Learn from *expert labeling* in natural language processing

# Hypothesis Class: Weighted Finite Automata (WFA)

## Graphical representation



## Algebraic representation

$$\alpha^\top = \begin{bmatrix} 1/2 & 1/2 \end{bmatrix} \quad \beta^\top = \begin{bmatrix} 1 & -1 \end{bmatrix}$$
$$\mathbf{A}_a = \begin{bmatrix} 3/4 & 0 \\ 0 & 1/3 \end{bmatrix} \quad \mathbf{A}_b = \begin{bmatrix} 6/5 & 2/3 \\ 3/4 & 1 \end{bmatrix}$$

Compute function  $f : \Sigma^* \rightarrow \mathbb{R}$ :

$$f(x_1 \cdots x_t) = \alpha^\top \mathbf{A}_{x_1} \cdots \mathbf{A}_{x_t} \beta \quad (n^2|\Sigma| + 2n) \text{ parameters}$$

## Why WFA?

- ▶ Expressive well-studied class [DKV09]
- ▶ Rich family of algorithms [Mohri09]  
(weighted minimization, determinization,  $\epsilon$ -removal)
- ▶ Widely used in applications [MPR08, AK09, BGC09, KM09]  
(speech recognition, image processing, OCR, system testing)

# Overview

**Our Result:** A supervised learning algorithm for WFA that combines spectral learning and matrix completion

In the rest of the talk I will...

1. Recall the **spectral method** in a nutshell
2. Describe a family of **learning algorithms**
3. Give a **generalization bound**

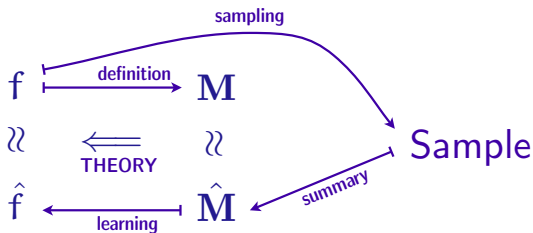
# Spectral Learning in a Nutshell

(Workshop on Friday!)

Key Ideas:

- ▶ Matrix of observables  $\mathbf{M}$  contains sufficient information
- ▶ SVD decomposition of  $\mathbf{M}$  is used to recover model
- ▶ Computation is noise tolerant

General Scheme:



## Spectral Learning by...

Has been applied to many models:

- ▶ Sequential models  
[HKZ09,BDR09,SBG10,BSG10,BQC11,Bailly11,BQC12]
- ▶ Tree-like structures  
[BHD10,PSX11,ACHKSZ11,LQBC12,CSCFU12,DRCFU12]
- ▶ Other graphical models  
[SBSGS10,AFHKL12,AHK12,PSITX12]

In the particular case of WFA  $\mathbf{M}$  is a well-known matrix...

# The Hankel Matrix

**Definition:** Hankel matrix  $\mathbf{H}_f$  of a function  $f : \Sigma^* \rightarrow \mathbb{R}$  is such that

- ▶ rows are indexed by *prefixes*  $\mathbf{u} \in \mathcal{P}$
- ▶ columns are indexed by *suffixes*  $\mathbf{v} \in \mathcal{S}$
- ▶ entries are evaluations  $\mathbf{H}_f(\mathbf{u}, \mathbf{v}) = f(\mathbf{u}\mathbf{v})$

**Example:**  $\Sigma = \{a, b\}$  and  $f(x) = \#$  of  $a$ 's in  $x$

$$\begin{aligned}\mathcal{P} &= \{a, b, aa\} \text{ (rows)} \\ \mathcal{S} &= \{\epsilon, a, b\} \text{ (columns)}\end{aligned}$$

$$\mathbf{H}_f = \begin{matrix} & \epsilon & a & b \\ \begin{matrix} a \\ b \\ aa \end{matrix} & \begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 0 \\ 2 & 3 & 2 \end{bmatrix} \end{matrix}$$

**Note:** Entrywise redundancies

$$w = u_1v_1 = u_2v_2 \Rightarrow \mathbf{H}_f(u_1, v_1) = \mathbf{H}_f(u_2, v_2) = f(w)$$

# Spectral Learning for WFA

In the particular case of WFA

- ▶ Hankel matrices play the role of  $\mathbf{M}$
- ▶ Redundancies (entrywise & rank) are relevant
- ▶  $f$  computed by WFA with  $\text{rank}(\mathbf{H}_f)$  states  
(when  $\mathcal{P}$ ,  $\mathcal{S}$  big enough)

Example:  $\Sigma = \{a, b\}$  and  $f(x) = \#$  of  $a$ 's in  $x$

$\mathcal{P} = \{a, b, aa\}$  (rows)

$\mathcal{S} = \{\epsilon, a, b\}$  (columns)

$\text{rank}(\mathbf{H}_f) = 2$

$$\mathbf{H}_f = \begin{array}{c} \epsilon \quad a \quad b \\ a \\ b \\ aa \end{array} \begin{bmatrix} 1 & 2 & 1 \\ 0 & 1 & 0 \\ 2 & 3 & 2 \end{bmatrix}$$

But: how do you obtain a Hankel matrix in the regression setting?



# Missing Entries

In usual applications...

entries in  $\mathbf{H}_f$  are empirical counts, e.g.  $f(x) = \Pr[x]$

$$\left\{ \begin{array}{l} aa, b, bab, a, \\ b, a, ab, aa, \\ ba, b, aa, a, \\ aa, bab, b, aa \end{array} \right\} \longrightarrow \begin{array}{l} \epsilon \\ a \\ b \\ ba \end{array} \begin{array}{cc} a & b \\ \left[ \begin{array}{cc} .19 & .25 \\ .31 & .06 \\ .06 & .00 \\ .00 & .13 \end{array} \right] \end{array}$$

But in this case...

entries in  $\mathbf{H}_f$  are labels observed in the sample

$$\left\{ \begin{array}{l} (bab,1) \\ (bbb,0) \\ (aaa,3) \\ (a,1) \\ (ab,1) \\ (aa,2) \\ (aba,2) \\ (bb,0) \end{array} \right\} \longrightarrow \begin{array}{l} a \\ b \\ aa \\ ab \\ ba \\ bb \end{array} \begin{array}{ccc} \epsilon & a & b \\ \left[ \begin{array}{ccc} 1 & 2 & 1 \\ * & * & 0 \\ 2 & 3 & * \\ 1 & 2 & * \\ * & * & 1 \\ 0 & * & 0 \end{array} \right] \end{array}$$

# Constrained Matrix Completion

Solution: Apply *matrix completion* to  $\hat{\mathbf{H}}_f$   
[CR09,CP10,CT10,FSSS11,Recht11,FS11,NW12]

But: Constrain completed matrix to be *Hankel*

Algorithm: Use convex optimization

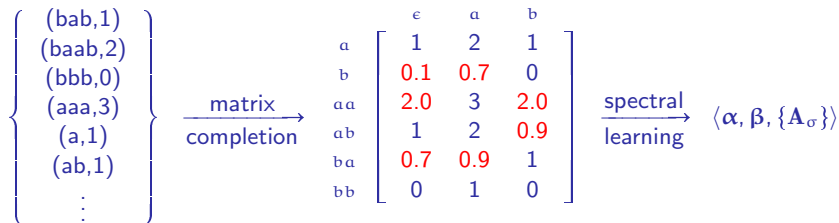
$$\hat{\mathbf{H}} = \underset{\mathbf{H} \in \mathbb{H}}{\operatorname{argmin}} \ell(\mathbf{H}; \mathbf{S}) + \lambda \cdot \mathbf{R}(\mathbf{H})$$

- ▶ Loss  $\ell$  controls agreement of  $\mathbf{H}$  with sample
- ▶ Regularizer  $\mathbf{R}$  controls complexity of  $\mathbf{H}$  (e.g. Schatten norm)
- ▶  $\mathbf{H} \in \mathbb{H}$  imposes *convex constraints* (equalities between entries)

(Double) Role of Regularization:

- ▶ Solve ill-posedness of matrix completion problem
- ▶ Less complex  $\hat{\mathbf{H}}$  will lead to simpler WFA

# A Family of Algorithms



Family of algorithms parametrized by:

- ▶ Choice of rows and columns in  $\mathbf{H}$
- ▶ A constrained matrix completion algorithm
- ▶ Regularization parameters

**Question:** Can these algorithms provably succeed?

# Generalization Bound

Hypotheses:

- ▶ Reasonable assumptions on distribution  $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$
- ▶ Completion loss  $\ell(\mathbf{H}; S) = \sum_{(\mathbf{x}, \mathbf{y}) \in S} |\mathbf{f}_{\mathbf{H}}(\mathbf{x}) - \mathbf{y}|$
- ▶ Completion regularizer  $R(\mathbf{H}) = \|\mathbf{H}\|_F^2$

**Theorem:** with high probability over  $S \sim \mathcal{D}^m$ , the output  $\mathbf{f}_S$  of the algorithm satisfies:

$$\mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}} [|\mathbf{f}_S(\mathbf{x}) - \mathbf{y}|] \leq \hat{\mathbb{E}}_S [|\mathbf{f}_S(\mathbf{x}) - \mathbf{y}|] + \mathcal{O}\left(\frac{\ln m}{m^{1/3}}\right)$$

**Proof:** joint *stability analysis* of matrix completion and spectral learning

Want to Know More?

Poster T47

# Spectral Learning of General Weighted Automata via Constrained Matrix Completion

*Borja Balle*<sup>1</sup>    Mehryar Mohri<sup>2</sup>

<sup>1</sup>Universitat Politècnica de Catalunya

<sup>2</sup>Courant Institute of Mathematical Sciences and Google Research

NIPS 2012, Lake Tahoe