

Abstract

- ▶ **Goal:** induce a WCFG from observed strings (e.g. language modeling)
- ▶ **Contribution:** a spectral method for unsupervised WCFG learning
- ▶ **Ingredients:**
 - ▶ Convex optimization of a Hankel matrix for WCFG
 - ▶ Linear constraints to characterize WCFG Hankel matrices
 - ▶ Low-rank objective, derived from the spectral approach

Weighted Context-Free Grammars

- ▶ **Definition:** WCFG $G = \langle \Sigma, n, \alpha_*, \{\beta_\sigma\}, \mathbf{A} \rangle$
 - ▶ Alphabet Σ , n states (states = non-terminals)
 - ▶ Initial vector $\alpha_* \in \mathbb{R}^n$
 - ▶ Terminal vectors $\beta_\sigma \in \mathbb{R}^n$ for $\sigma \in \Sigma$
 - ▶ Bilinear operator $\mathbf{A} \in \mathbb{R}^{n \times n^2}$

- ▶ **Grammar function** $g_G : \Sigma^* \rightarrow \mathbb{R}$

$$g_G(x) = \alpha_*^\top \beta_G(x)$$

- ▶ **Inside function** $\beta_G : \Sigma^+ \rightarrow \mathbb{R}^n$

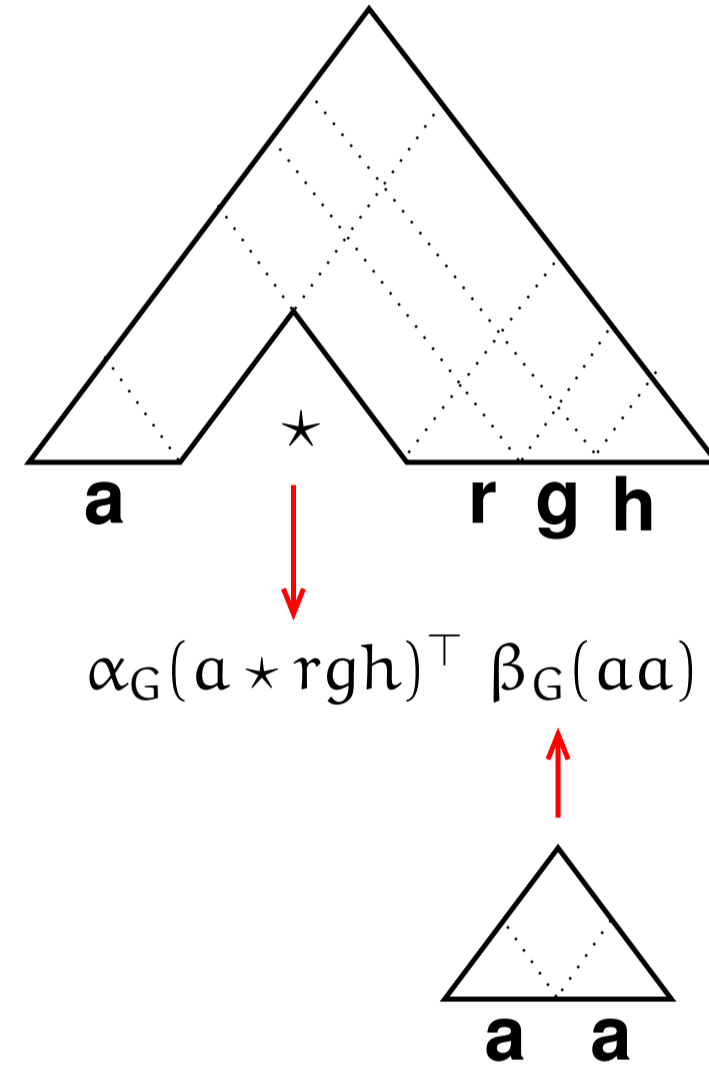
$$\beta_G(\sigma) = \beta_\sigma$$

$$\beta_G(x) = \sum_{\substack{x_1, x_2 \in \Sigma^+ \\ x = x_1 x_2}} \mathbf{A}(\beta_G(x_1) \otimes \beta_G(x_2))$$

- ▶ **Outside function** $\alpha_G : \Sigma^* \times \Sigma^* \rightarrow \mathbb{R}^n$

$$\alpha_G(\langle * \rangle) = \alpha_*$$

$$\alpha_G(\langle x * z \rangle)^\top = \sum_{\substack{x_1 \in \Sigma^+, x_2 \in \Sigma^+ \\ x = x_1 x_2}} \alpha_G(\langle x_1 * z \rangle)^\top \mathbf{A}(\beta_G(x_2) \otimes \mathbf{I}_n) + \sum_{\substack{z_1 \in \Sigma^+, z_2 \in \Sigma^+ \\ z = z_1 z_2}} \alpha_G(\langle x * z_2 \rangle)^\top \mathbf{A}(\mathbf{I}_n \otimes \beta_G(z_1))$$



The Spectral Method

1. Compute Hankel matrix \mathbf{H} using training data
 - ▶ **Supervised:** training data contains derivations, count events
 - ▶ **Unsupervised (this paper):** induce \mathbf{H} from plain strings
2. Compute the SVD of $\mathbf{H}_{\mathcal{O}^1, \mathcal{I}^1} = \mathbf{U} \mathbf{\Lambda} \mathbf{V}^\top$
3. Create a rank- n factorization of $\mathbf{H}_{\mathcal{O}^1, \mathcal{I}^1} \approx \mathbf{F}_n \mathbf{B}_n$

$$\mathbf{F}_n = \mathbf{U}_n \mathbf{\Lambda}_n \quad \mathbf{B}_n = \mathbf{V}_n^\top$$
4. Compute \mathbf{G} :

$$\alpha_*^\top = \mathbf{H}_{*, \mathcal{I}^1} \mathbf{B}_n^+$$

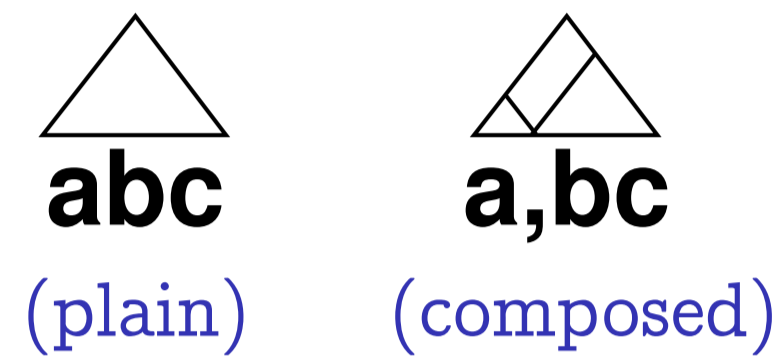
$$\beta_\sigma = \mathbf{F}_n^+ \mathbf{H}_{\mathcal{O}^1, \sigma}$$

$$\mathbf{A} = \mathbf{F}_n^+ \mathbf{H}_{\mathcal{O}^1, \mathcal{I}^2} (\mathbf{B}_n \otimes \mathbf{B}_n)^+$$

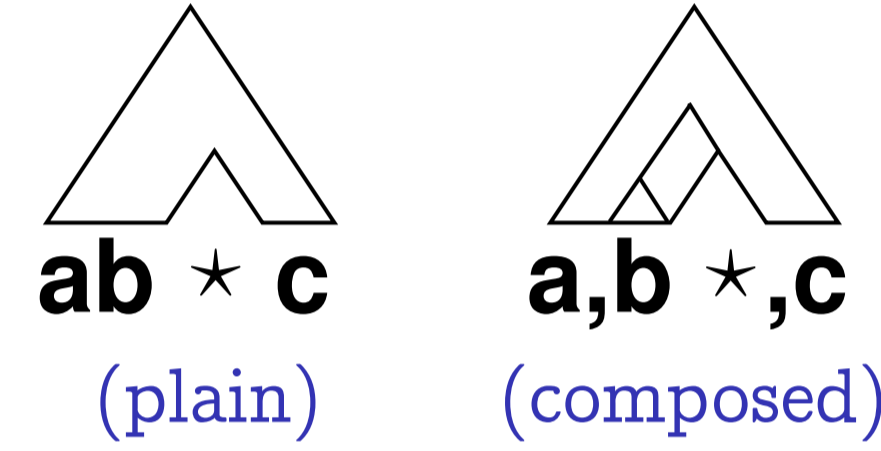
Hankel Matrices for WCFG

- ▶ **Definition:** $\mathbf{H}_g \in \mathbb{R}^{\mathcal{O} \times \mathcal{I}}$ with $\mathbf{H}_g(o, i) = g(o \cdot i)$

- ▶ **Inside strings** \mathcal{I}



- ▶ **Outside contexts** \mathcal{O}



- ▶ **Theorem:** g computed by G with n states

$$\Leftrightarrow \text{rank}(\mathbf{H}_g) = n$$

Intuition:

$$\mathbf{H}_{\mathcal{O}^1, \mathcal{I}^1} = \mathbf{F} \mathbf{B}$$

$$\mathbf{F}(x * z, \cdot) = \alpha(x * z) \in \mathbb{R}^n$$

$$\mathbf{B}(\cdot, x) = \beta(x) \in \mathbb{R}^n$$

$$\mathbf{H}_{\mathcal{O}^1, \mathcal{I}^2} = \mathbf{F} \mathbf{A} (\mathbf{B} \otimes \mathbf{B})$$

$$\mathbf{H}_{*, \mathcal{I}^1} = \alpha_*^\top \mathbf{B}$$

$$\mathbf{H}_{\mathcal{O}^1, \sigma} = \mathbf{F} \beta_\sigma$$

Hankel Induction

- ▶ **Observable statistics (for full strings):**

$$\mathbf{H}(\langle * \rangle, (x)) = p(x)$$

- ▶ **WCFG constraints:**

Hankel constraints:

$$\mathbf{H}(\langle x * z \rangle, (y_1, y_2)) = \mathbf{H}(\langle x, y_1 * , z \rangle, (y_2)) = \mathbf{H}(\langle x, * y_2, z \rangle, (y_1))$$

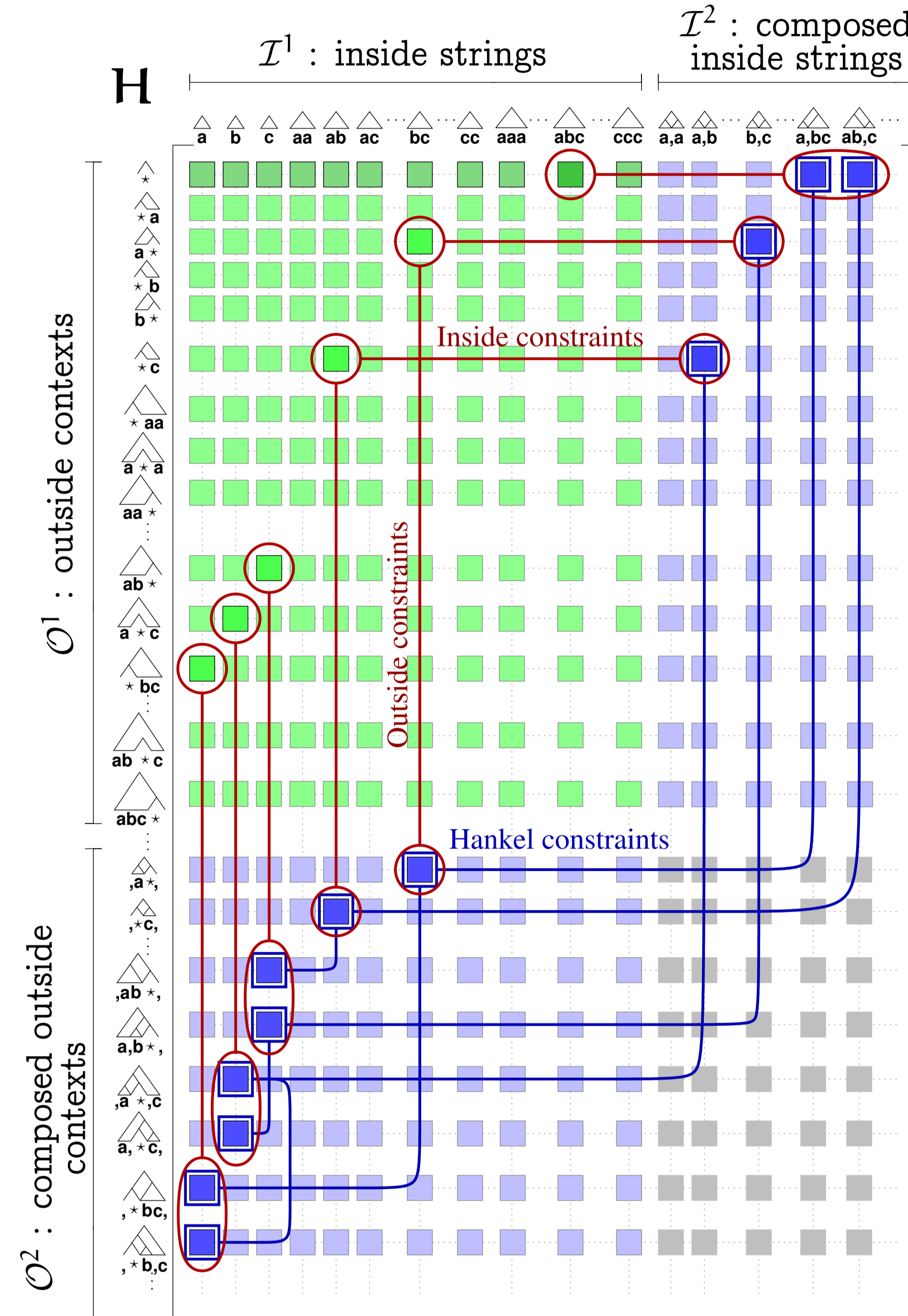
Inside constraints:

$$\mathbf{H}(o, (x)) = \sum_{x = x_1 x_2} \mathbf{H}(o, (x_1, x_2))$$

Outside constraints:

$$\mathbf{H}(\langle x * z \rangle, i) = \sum_{x = x_1 x_2} \mathbf{H}(\langle x_1, x_2 * , z \rangle, i) + \sum_{z = z_1 z_2} \mathbf{H}(\langle x, * z_1, z_2 \rangle, i)$$

- ▶ **Property:** any \mathbf{H} satisfying constraints corresponds to a WCFG

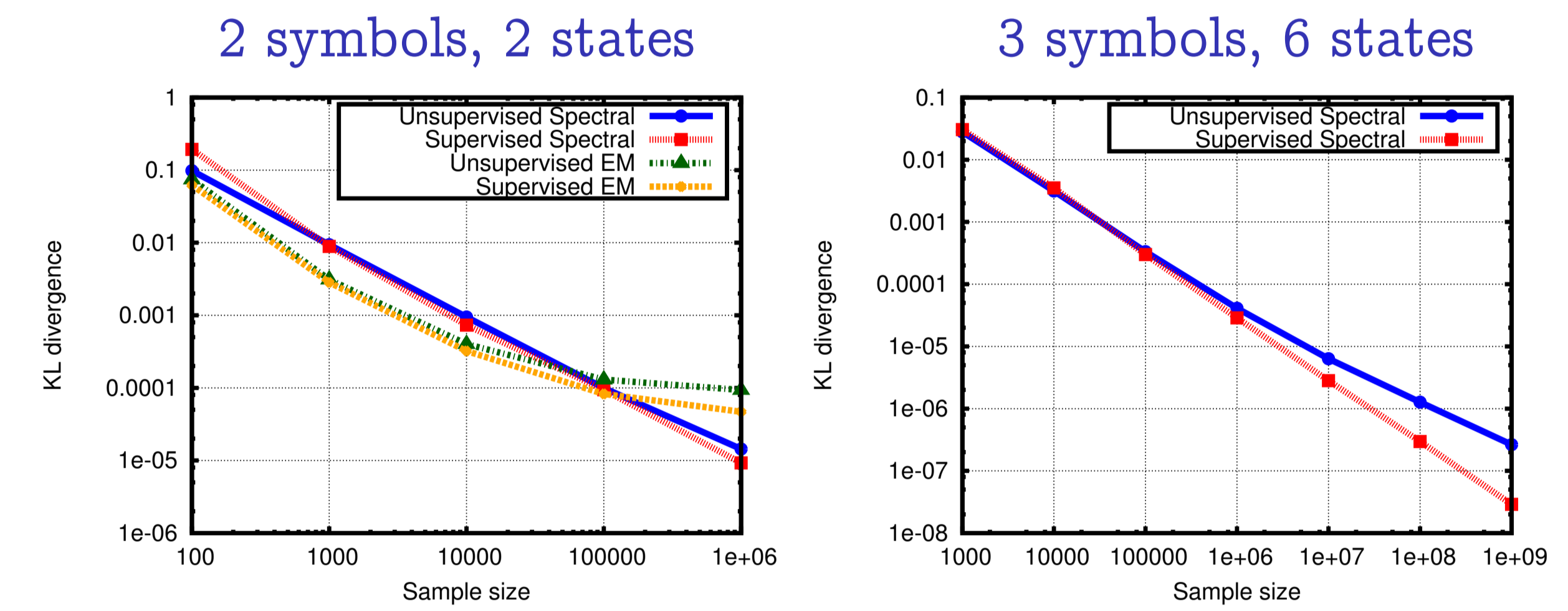


Convex Optimization

$$\begin{aligned} & \underset{\mathbf{H}}{\text{minimize}} \quad \|\mathbf{H}\|_* \\ & \text{subject to} \quad \|\mathbf{O} \cdot \text{vec}(\mathbf{H}) - \mathbf{z}\|_2 \leq \mu \\ & \quad \mathbf{K} \cdot \text{vec}(\mathbf{H}) = 0 \\ & \quad \|\mathbf{H}\|_2 \leq 1. \end{aligned}$$

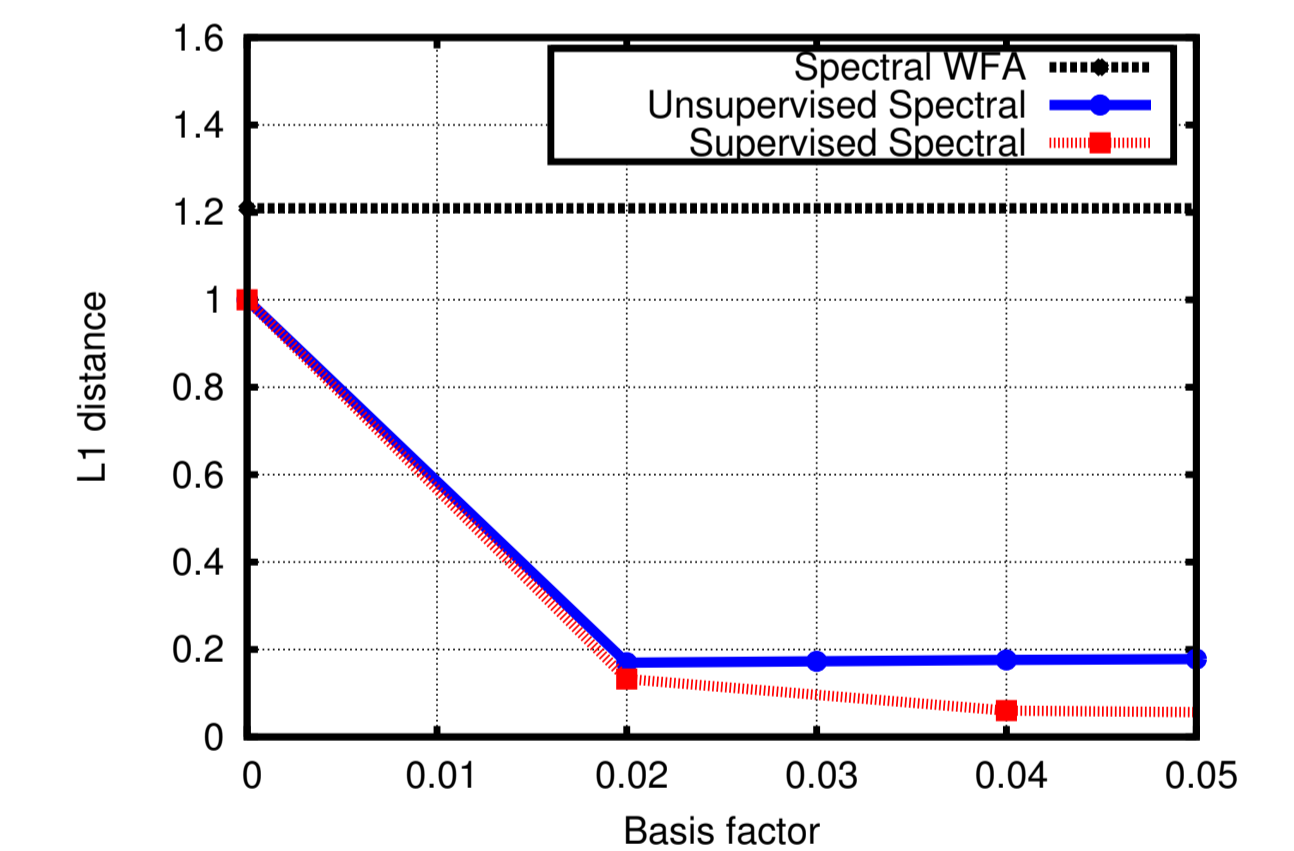
Optimization based on projections for $\|\cdot\|_*$, \mathbf{O} and \mathbf{K} .

Synthetic Experiments



- ▶ Random synthetic targets
- ▶ Fixed inside-outside patterns in Hankel
- ▶ Comparison to EM and a supervised spectral method

Experiments with Dyck Languages



- ▶ The target grammar is:

$$S \rightarrow SS (0.2) \mid a S b (0.4) \mid a b (0.4)$$

- ▶ e.g.: $p(aabb) = 0.16$, $p(aaa) = 0$
- ▶ We vary the size of the Hankel
- ▶ Comparison to spectral algorithm for Weighted Automata

On WSJ-10 data

basis	size of H	obs.	i/o ctr.	basis	size of H	obs.	i/o ctr.
1 × 11	39 × 159	34	162	36 × 34	27,989 × 11,682	916	156,690
6 × 14	1,163 × 764	146	6,360	42 × 37	3,638 × 15,026	1,035	200,346
12 × 18	4,462 × 2,239	322	25,374	48 × 41	45,192 × 18,235	1,157	244,398
18 × 22	9,124 × 4,149	479	52,524	54 × 45	53,741 × 21,196	1,281	284,466
24 × 26	15,755 × 6,858	657	89,718	60 × 48	60,844 × 23,890	1,382	318,354
30 × 29	19,801 × 8,545	769	112,37				