Hankel-Based Methods for Latent-State Structured Prediction

Ariadna Quattoni Xerox Research Centre Europe Collaborators: R. Bailly, B. Balle, X.Carreras, A. Globerson, F. Luque, A. Mayo

LxMLS 2014

Structured Prediction



Latent State Models

Consider a sequence abaccb

$$f(abaccb) = \alpha_f(ab) \cdot \beta_f(accb)$$
$$= \alpha_f(ab) \cdot A_a \cdot \beta_f(ccb)$$
$$= \alpha_f(aba) \cdot A_c \cdot \beta_f(cb)$$

where:

- n is the dimension of the model
- α_f maps prefixes to \mathbb{R}^n
- β_f maps suffixes to \mathbb{R}^n
- A_a is a bilinear operator in $\mathbb{R}^{n \times n}$

F: functions computed by WFAs (HMMs is a subclass)

Latent State Models

Consider a sequence abaccb

$$f(abaccb) = \alpha_f(ab) \cdot \beta_f(accb)$$
$$= \alpha_f(ab) \cdot A_a \cdot \beta_f(ccb)$$
$$= \alpha_f(aba) \cdot A_c \cdot \beta_f(cb)$$

where:

- n is the dimension of the model
- α_f maps prefixes to \mathbb{R}^n
- β_f maps suffixes to \mathbb{R}^n
- A_a is a bilinear operator in $\mathbb{R}^{n \times n}$

F: functions computed by WFAs (HMMs is a subclass)

- Finding f ∈ F: Estimate α_f, β_f and A_σ from S
 EM
- \blacktriangleright Many ways of expressing the same $f \in F$
 - Algebraic Methods

- Spectral Learning of WFA
- Hankel-Based Optimization for Structured Prediction

- Spectral Learning of WFA
- Hankel-Based Optimization for Structured Prediction



'Classic spectral learning algorithms'

 Hankel-Based Optimization for Structured Prediction





- Spectral Learning of WFA
- Hankel-Based Optimization for Structured Prediction





$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$



• Class of WA parametrized by alphabet Σ and number of states n



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\mathbf{A} = \langle \alpha_1, \alpha_\infty, \{A_\sigma\}_{\sigma \in \Sigma} \rangle$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\ 0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\ 0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\ 0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\ 0.1 & 0.1 \end{bmatrix}$$

$$\mathbf{A} = \langle \alpha_1, \alpha_\infty, \{A_\sigma\}_{\sigma \in \Sigma} \rangle$$

$$\alpha_1 \in \mathbb{R}^n$$



$$\begin{aligned} &\alpha_1 = \begin{bmatrix} 1.0\\ 0.0 \end{bmatrix} \quad A_\alpha = \begin{bmatrix} 0.4 & 0.2\\ 0.1 & 0.1 \end{bmatrix} \\ &\alpha_\infty = \begin{bmatrix} 0.0\\ 0.6 \end{bmatrix} \quad A_b = \begin{bmatrix} 0.1 & 0.3\\ 0.1 & 0.1 \end{bmatrix} \end{aligned}$$

$$\mathbf{A} = \langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \rangle$$

$$\alpha_1 \in \mathbb{R}^n \Longrightarrow$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\ 0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\ 0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\ 0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\ 0.1 & 0.1 \end{bmatrix}$$

$$\mathbf{A} = \langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \rangle$$

$$\alpha_1 \in \mathbb{R}^n \Longrightarrow \text{(initial weights)}$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\mathbf{A} = \left\langle \alpha_1, \alpha_\infty, \{A_\sigma\}_{\sigma \in \Sigma} \right\rangle$$
$$\alpha_1 \in \mathbb{R}^n \Longrightarrow \text{(initial weights)} \quad \alpha_\infty \in \mathbb{R}^n$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\mathbf{A} = \langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \rangle$$

$$\alpha_1 \in \mathbb{R}^n \Longrightarrow \text{(initial weights)} \quad \alpha_{\infty} \in \mathbb{R}^n \Longrightarrow$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\mathbf{A} = \langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \rangle$$

$$\alpha_1 \in \mathbb{R}^n \Longrightarrow \text{(initial weights)} \quad \alpha_{\infty} \in \mathbb{R}^n \Longrightarrow \text{(terminal weights)}$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\begin{split} \mathbf{A} &= \left\langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \right\rangle \\ \alpha_1 \in \mathbb{R}^n \Longrightarrow (\text{initial weights}) \quad \alpha_{\infty} \in \mathbb{R}^n \Longrightarrow (\text{terminal weights}) \\ A_{\sigma} \in \mathbb{R}^{n \times n} \end{split}$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\begin{split} \mathbf{A} &= \langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \rangle \\ \alpha_1 \in \mathbb{R}^n \Longrightarrow (\text{initial weights}) \quad \alpha_{\infty} \in \mathbb{R}^n \Longrightarrow (\text{terminal weights}) \\ A_{\sigma} \in \mathbb{R}^{n \times n} \Longrightarrow \end{split}$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\begin{split} \mathbf{A} &= \langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \rangle \\ \alpha_1 \in \mathbb{R}^n \Longrightarrow (\text{initial weights}) \quad \alpha_{\infty} \in \mathbb{R}^n \Longrightarrow (\text{terminal weights}) \\ A_{\sigma} \in \mathbb{R}^{n \times n} \Longrightarrow (\text{transition weights}) \end{split}$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\begin{split} \mathbf{A} &= \langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \rangle \\ \alpha_1 \in \mathbb{R}^n \Longrightarrow (\text{initial weights}) \quad \alpha_{\infty} \in \mathbb{R}^n \Longrightarrow (\text{terminal weights}) \\ A_{\sigma} \in \mathbb{R}^{n \times n} \Longrightarrow (\text{transition weights}) \end{split}$$

• Computes a function
$$f_{\mathbf{A}}: \Sigma^{\star} \to \mathbb{R}$$



$$\alpha_{1} = \begin{bmatrix} 1.0\\0.0 \end{bmatrix} \quad A_{\alpha} = \begin{bmatrix} 0.4 & 0.2\\0.1 & 0.1 \end{bmatrix}$$
$$\alpha_{\infty} = \begin{bmatrix} 0.0\\0.6 \end{bmatrix} \quad A_{b} = \begin{bmatrix} 0.1 & 0.3\\0.1 & 0.1 \end{bmatrix}$$

$$\begin{split} \mathbf{A} &= \left\langle \alpha_1, \alpha_{\infty}, \{A_{\sigma}\}_{\sigma \in \Sigma} \right\rangle \\ \alpha_1 \in \mathbb{R}^n \Longrightarrow (\text{initial weights}) \quad \alpha_{\infty} \in \mathbb{R}^n \Longrightarrow (\text{terminal weights}) \\ A_{\sigma} \in \mathbb{R}^{n \times n} \Longrightarrow (\text{transition weights}) \end{split}$$

• Computes a function
$$f_{\mathbf{A}}: \Sigma^{\star} \to \mathbb{R}$$

$$f_{\mathbf{A}}(\mathbf{x}) = f_{\mathbf{A}}(\mathbf{x}_1 \cdots \mathbf{x}_t) = \alpha_1^{\top} A_{\mathbf{x}_1} \cdots A_{\mathbf{x}_t} \alpha_{\infty}$$

Examples of functions computed by WFAs

Generates infinite strings, computes probabilities of prefixes $\mathbb{P}[x\Sigma^*]$ Emission and transition are conditionally independent given state



Examples of functions computed by WFAs

Compute conditional probabilities $\mathbb{P}[y|x] = \alpha_1^\top A_x^y \alpha_\infty$ for pairs $(x, y) \in (\Sigma \times \Delta)^*$, must have |x| = |y|Can also assume models factorized like in HMM





Any WA ${\bf A}$ defines forward and backward maps

Any WA A defines *forward* and *backward* maps

Any WA A defines forward and backward maps $\implies \alpha_A$, $\beta_A : \Sigma^* \to \mathbb{R}^n$

$$\alpha_{\mathbf{A}}(\mathbf{p})^{\mathsf{T}} = \alpha_{1}^{\mathsf{T}} \mathcal{A}_{p_{1}} \cdots \mathcal{A}_{p_{t}}$$

$$\alpha_{\mathbf{A}}(\mathbf{p})^{\mathsf{T}} = \alpha_{1}^{\mathsf{T}} A_{\mathbf{p}_{1}} \cdots A_{\mathbf{p}_{t}}$$
$$\beta_{\mathbf{A}}(\mathbf{s})^{\mathsf{T}} = A_{\mathbf{s}_{1}} \cdots A_{\mathbf{s}_{t'}} \alpha_{\infty}$$
Any WA **A** defines *forward* and *backward* maps $\implies \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(s)^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t'}} \boldsymbol{\alpha}_{\infty} \\ \boldsymbol{f}_{\mathbf{A}}(x) &= \boldsymbol{\alpha}_{\mathbf{A}}(p)^{\top} \boldsymbol{\beta}_{\mathbf{A}}(s) \end{aligned}$$

Any WA **A** defines *forward* and *backward* maps $\implies \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\alpha_{\mathbf{A}}(\mathbf{p})^{\top} = \alpha_{1}^{\top} A_{p_{1}} \cdots A_{p_{t}}$$
$$\beta_{\mathbf{A}}(\mathbf{s})^{\top} = A_{s_{1}} \cdots A_{s_{t}}, \alpha_{\infty}$$
$$f_{\mathbf{A}}(\mathbf{x}) = \alpha_{\mathbf{A}}(\mathbf{p})^{\top} \beta_{\mathbf{A}}(\mathbf{s})$$

Any WA **A** defines *forward* and *backward* maps $\implies \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(\mathbf{s})^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t'}} \boldsymbol{\alpha}_{\infty} \\ \boldsymbol{f}_{\mathbf{A}}(\mathbf{x}) &= \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} \boldsymbol{\beta}_{\mathbf{A}}(\mathbf{s}) \end{aligned}$$



Any WA **A** defines *forward* and *backward* maps $\implies \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(s)^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t'}} \boldsymbol{\alpha}_{\infty} \\ \mathbf{f}_{\mathbf{A}}(x) &= \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} \boldsymbol{\beta}_{\mathbf{A}}(s) \end{aligned}$$

$$[\alpha_{A}(p)]_{\mathfrak{i}}=\mathbb{P}[p\ ,\ h_{+1}=\mathfrak{i}]$$

Any WA **A** defines *forward* and *backward* maps $\implies \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(s)^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t'}} \boldsymbol{\alpha}_{\infty} \\ \mathbf{f}_{\mathbf{A}}(x) &= \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} \boldsymbol{\beta}_{\mathbf{A}}(s) \end{aligned}$$

$$\begin{split} & [\alpha_{A}(p)]_{\mathfrak{i}} = \mathbb{P}[p \text{ , } h_{+1} = \mathfrak{i}] \\ & [\beta_{A}(s)]_{\mathfrak{i}} = \mathbb{P}[s \mid h = \mathfrak{i}] \end{split}$$

Any WA **A** defines *forward* and *backward* maps $\implies \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(s)^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t}}, \boldsymbol{\alpha}_{\infty} \\ \boldsymbol{f}_{\mathbf{A}}(x) &= \boldsymbol{\alpha}_{\mathbf{A}}(p)^{\top} \boldsymbol{\beta}_{\mathbf{A}}(s) \end{aligned}$$

Example, for HMMs:

Consequences:

$$\begin{split} & [\alpha_{\mathbf{A}}(p)]_{\mathfrak{i}} = \mathbb{P}[p , h_{+1} = \mathfrak{i}] \\ & [\beta_{\mathbf{A}}(s)]_{\mathfrak{i}} = \mathbb{P}[s \mid h = \mathfrak{i}] \end{split}$$

Any WA **A** defines *forward* and *backward* maps $\Longrightarrow \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(s)^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t'}} \boldsymbol{\alpha}_{\infty} \\ \mathbf{f}_{\mathbf{A}}(\mathbf{x}) &= \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} \boldsymbol{\beta}_{\mathbf{A}}(s) \end{aligned}$$

$$\begin{split} & [\alpha_{\mathbf{A}}(p)]_{\mathfrak{i}} = \mathbb{P}[p , h_{+1} = \mathfrak{i}] \\ & [\beta_{\mathbf{A}}(s)]_{\mathfrak{i}} = \mathbb{P}[s \mid h = \mathfrak{i}] \end{split}$$



Any WA **A** defines *forward* and *backward* maps $\Longrightarrow \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(\mathbf{s})^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t'}} \boldsymbol{\alpha}_{\infty} \\ \boldsymbol{f}_{\mathbf{A}}(\mathbf{x}) &= \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} \boldsymbol{\beta}_{\mathbf{A}}(\mathbf{s}) \end{aligned}$$

Example, for HMMs:Consequences: $[\alpha_A(p)]_i = \mathbb{P}[p, h_{+1} = i]$ A_σ can be recovered from
 $f_A(p\sigma s), \alpha_A(p), and \beta_A(s)$:

Any WA **A** defines *forward* and *backward* maps $\implies \alpha_{\mathbf{A}}, \beta_{\mathbf{A}} : \Sigma^* \to \mathbb{R}^n$ such that for any splitting $x = p \cdot s$ one has

$$\begin{aligned} \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} &= \boldsymbol{\alpha}_{1}^{\top} \boldsymbol{A}_{p_{1}} \cdots \boldsymbol{A}_{p_{t}} \\ \boldsymbol{\beta}_{\mathbf{A}}(\mathbf{s})^{\top} &= \boldsymbol{A}_{s_{1}} \cdots \boldsymbol{A}_{s_{t'}} \boldsymbol{\alpha}_{\infty} \\ \boldsymbol{f}_{\mathbf{A}}(\mathbf{x}) &= \boldsymbol{\alpha}_{\mathbf{A}}(\mathbf{p})^{\top} \boldsymbol{\beta}_{\mathbf{A}}(\mathbf{s}) \end{aligned}$$

Example, for HMMs:Consequences: $[\alpha_A(p)]_i = \mathbb{P}[p, h_{+1} = i]$ A_σ can be recovered from
 $f_A(p\sigma s), \alpha_A(p), and \beta_A(s)$: $[\beta_A(s)]_i = \mathbb{P}[s \mid h = i]$ $f_A(p\sigma s) = \alpha_A(p)^\top A_\sigma \beta_A(s)$

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is \square

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$

For $p, s \in \Sigma^{\star}$, entries are defined by

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$

For $p, s \in \Sigma^*$, entries are defined by

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ For $p, s \in \Sigma^*$, entries are defined by $\longrightarrow H_f(p, s) = f(p \cdot s)$

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ For $p, s \in \Sigma^*$, entries are defined by $\longrightarrow H_f(p, s) = f(p \cdot s)$

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ For $p, s \in \Sigma^*$, entries are defined by $\Longrightarrow H_f(p, s) = f(p \cdot s)$ b ab e a aa \cdots f(aab) e a Ъ aa f(aab) ab

The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ For $p, s \in \Sigma^*$, entries are defined by $\longrightarrow H_f(p, s) = f(p \cdot s)$

Very *redundant* representation



The Hankel matrix of $f: \Sigma^* \to \mathbb{R}$ is $\longrightarrow H_f \in \mathbb{R}^{\Sigma^* \times \Sigma^*}$ For $p, s \in \Sigma^*$, entries are defined by $\longrightarrow H_f(p, s) = f(p \cdot s)$

Very *redundant* representation

we will work with finite sub-blocks of H



Schützenberger's Theorem

Schützenberger's Theorem

Schützenberger's Theorem

Theorem: $rank(H_f) = n \iff f = f_A$ with minimal A of size n

Schützenberger's Theorem

Theorem: $\mathsf{rank}(H_f) = \mathfrak{n} \Longleftrightarrow f = f_\mathbf{A}$ with minimal \mathbf{A} of size \mathfrak{n}

Remarks (\Leftarrow)

Schützenberger's Theorem

Theorem: $\mathsf{rank}(H_f) = \mathfrak{n} \Longleftrightarrow f = f_\mathbf{A}$ with minimal \mathbf{A} of size \mathfrak{n}

Remarks (\Leftarrow) Write $F = \alpha_{\mathbf{A}}(\Sigma^{\star})^{\top} \in \mathbb{R}^{\Sigma^{\star} \times n}$ and $B = \beta_{\mathbf{A}}(\Sigma^{\star}) \in \mathbb{R}^{n \times \Sigma^{\star}}$ Note $H_{f} = F \cdot B$ Then, $\operatorname{rank}(H_{f}) = n$

Schützenberger's Theorem

Theorem: $\mathsf{rank}(H_f) = \mathfrak{n} \Longleftrightarrow f = f_\mathbf{A}$ with minimal \mathbf{A} of size \mathfrak{n}

Remarks (
$$\Leftarrow$$
)
Write $F = \alpha_{\mathbf{A}}(\Sigma^{\star})^{\top} \in \mathbb{R}^{\Sigma^{\star} \times n}$ and $B = \beta_{\mathbf{A}}(\Sigma^{\star}) \in \mathbb{R}^{n \times \Sigma^{\star}}$
Note $H_{f} = F \cdot B$
Then, rank $(H_{f}) = n$

Remarks (\Rightarrow)

Schützenberger's Theorem

Theorem: $\mathsf{rank}(H_f) = \mathfrak{n} \Longleftrightarrow f = f_\mathbf{A}$ with minimal \mathbf{A} of size \mathfrak{n}

Remarks (\Leftarrow) Write $F = \alpha_{\mathbf{A}}(\Sigma^{\star})^{\top} \in \mathbb{R}^{\Sigma^{\star} \times n}$ and $B = \beta_{\mathbf{A}}(\Sigma^{\star}) \in \mathbb{R}^{n \times \Sigma^{\star}}$ Note $H_{f} = F \cdot B$ Then, rank $(H_{f}) = n$

Remarks (\Rightarrow)

Assume $rank(H_f) = n$

Take rank factorization $H_f = F \cdot B$ with $F \in \mathbb{R}^{\Sigma^* \times n}$ and $B \in \mathbb{R}^{n \times \Sigma^*}$ A WA that computes f can be constructed using F and B

Low-rank Hankel and Operators

Low-rank Hankel and Operators



Low-rank Hankel and Operators



$H_{\sigma} = F A_{\sigma} B \implies A_{\sigma} = F^+ H_{\sigma} B^+$

Idea: Use SVD decomposition to obtain a factorization of H • Choose a basis: sets of prefixes \mathcal{P} and suffixes \mathcal{S} • Obtain H and H_{σ} over basis $(\mathcal{P}, \mathcal{S})$ • Compute *compact* SVD as $H = U\Lambda V^{\top}$ with $\mathbf{U} \in \mathbb{R}^{\mathcal{P} \times \mathbf{n}} \qquad \boldsymbol{\Lambda} \in \mathbb{R}^{\mathbf{n} \times \mathbf{n}} \qquad \mathbf{V} \in \mathbb{R}^{\mathcal{S} \times \mathbf{n}}$ • Construct the WA using $F = U\Lambda$ and $B = V^{\top}$: • $A_{\sigma} = F^+ H_{\sigma} V^+$ for all $\sigma \in \Sigma$ • $\alpha_1^\top = H(\lambda, \cdot)^\top B^+$, $\alpha_{\infty} = F^+ H(\cdot, \lambda)$

Idea: Use SVD decomposition to obtain a factorization of H Choose a basis: sets of prefixes P and suffixes S • Obtain H and H_{σ} over basis $(\mathcal{P}, \mathcal{S})$ • Compute *compact* SVD as $H = U\Lambda V^{\top}$ with $\mathbf{U} \in \mathbb{R}^{\mathcal{P} \times \mathbf{n}} \qquad \boldsymbol{\Lambda} \in \mathbb{R}^{\mathbf{n} \times \mathbf{n}} \qquad \mathbf{V} \in \mathbb{R}^{\mathcal{S} \times \mathbf{n}}$ • Construct the WA using $F = U\Lambda$ and $B = V^{\top}$: • $A_{\sigma} = F^+ H_{\sigma} V^+$ for all $\sigma \in \Sigma$ • $\alpha_1^\top = H(\lambda, \cdot)^\top B^+$, $\alpha_{\infty} = F^+ H(\cdot, \lambda)$

Properties

- Easy to implement: just linear algebra
- Fast to compute: $O(\max\{|\mathcal{P}|, |\mathcal{S}|\}^3)$
- ▶ Noise tolerant: $\hat{H} \approx H$ and $\hat{H}_{\sigma} \approx H_{\sigma}$ implies $\hat{A}_{\sigma} \approx A_{\sigma}$



► A dimension n: controls complexity of the WA family

Experiment: PoS-Tag Sequence Models



PTB sequences of simplified PoS tags (Petrov et al 2010) Configuration: expectations on frequent substrings Metric: error rate on predicting next symbol in test sequences
Experiment: PoS-Tag Sequence Models



Comparison with a bigram baseline and EM Metric: error rate on predicting next symbol in test sequences At training, the Spectral Method is > 100 faster than EM

- Considers stochastic functions
- Assumes there is a model that generates sequences
- We can sample from this model

- Considers stochastic functions
- Assumes there is a model that generates sequences
- We can sample from this model



- Considers stochastic functions
- Assumes there is a model that generates sequences
- We can sample from this model



- Considers stochastic functions
- Assumes there is a model that generates sequences
- We can sample from this model



- Considers stochastic functions
- Assumes there is a model that generates sequences
- We can sample from this model



- Considers stochastic functions
- Assumes there is a model that generates sequences
- We can sample from this model



- Necessary statistics are hard to estimate
- In some cases we can not directly estimate the Hankel from samples

eg. Learning a function that is not a distribution Learning a PCFG when we don't have derivations Learning transducers when we don't have alignments

- Necessary statistics are hard to estimate
- In some cases we can not directly estimate the Hankel from samples:

eg. Learning a function that is not a distribution Learning a PCFG when we don't have derivations Learning transducers when we don't have alignments

It would be better not to have to assume that samples are generated by some model !!

- Necessary statistics are hard to estimate
- In some cases we can not directly estimate the Hankel from samples:

eg. Learning a function that is not a distribution Learning a PCFG when we don't have derivations Learning transducers when we don't have alignments

It would be better not to have to assume that samples are generated by some model !!



Outline

- Spectral Learning of WFA
- Hankel-Based Optimization for Structured Prediction

Outline



Outline



Sequence Tagging

output: hlp-xpatxmxs input: hippopotamus

Fully Observable Models

Latent-variable Models





- + Making predictions is tractable
- + Learning is convex
- Performance crucially depends on features
- Hidden layer provides more expressivity
- Making predicitons is not tractable
- Learning is non-convex (this paper)

Learning Structured Predictors with Latent Variables

Desiderata:

- Expressive scoring functions
- Tractable prediction function
- Effective regularizer
- Convex training procedure

Main Idea: Change of Representation + Relaxation

Problem Formulation

- Scoring functions computed by Input-Output OOMs
- Piecewise Prediction and Loss Function

Solving the Learning Problem

Spectral trick:

optimize over parameters of $f \rightarrow$ optimize low-rank matrix H

- Relax low-rank constraint using nuclear norm of H
- Recover parameters of f from H using the spectral method

Scoring Functions Computed by IO-OOM

Latent Score $\theta(x, y, h)$:

$$\alpha(h_0) \; \prod_{t=1}^T A_{y_t}^{x_t}(h_{t-1},h_t) \; \beta(h_T)$$

Scoring Function
$$F_A(x, y)$$
:

$$\sum_{h} \theta(x, y, h) = \alpha^{T} A_{y_{1}}^{x_{1}} \dots A_{y_{T}}^{x_{T}} \beta$$

- Model: A : $\langle \alpha, \beta, \{A_b^a\} \rangle$
- Number of states: n
- Initial Weights: $\alpha \in \mathbb{R}^n$
- Final Weights: $\beta \in \mathbb{R}^n$
- Observable Operators $A_b^a \in \mathbb{R}^{n \times n}$

- ► Expressive Function Family → e.g. it includes HMM
- Making Predictions (i.e. maximizing $F_A(x, y)$) \rightarrow NP-hard

Piecewise Prediction and Loss for IO-OOM

Approximation: $F_A^k(x, y)$:

$$\sum_{t=1}^{T-(k-1)} F_{A}(x_{t:t+k-1}, y_{t:t+k-1})$$

Sum k–grams

• Task loss: l(y, z)

Loss $L_k(x, y, F_A)$:

e.g. hamming distance

$$\max_{z} [F_A^k(x,z) - F_A^k(x,y) + l(y,z))$$

 Prediction and Loss Function → computed in O(T|Y|^k) using the Viterbi Algorithm

Discrete Regularizer for IO-OOM

Learning Problem:

$$\operatorname{argmin}_{A \in \mathcal{F}} \sum_{i=1}^{m} L_k(x^i, y^i, F_A) + \tau |A|)$$

- Function class (IO-OOM): F
- Training Examples: $\langle x^i, y^i \rangle$
- Loss Function: L_k
- Regularizer \rightarrow number of states: |A|
- Trade-off constant: τ
- $k \ge 2 \rightarrow \text{Non-convex}$ dependence of L_k on parameters of A
- L_k involves polynomials of order k+3

Optimization Strategy

• L_k convex on values of $A \rightarrow$ optimization over $(X \times Y)^k$ values

Three challenges

- 1. Table of values \rightarrow must correspond to valid IO-OOM
- 2. Regularizer over table \rightarrow must correspond to #states of IO-OOM
- 3. Recover parameters of A from this table

Optimization Strategy

• L_k convex on values of $A \rightarrow$ optimization over $(X \times Y)^k$ values

Three challenges

- 1. Table of values \rightarrow must correspond to valid IO-OOM
- 2. Regularizer over table \rightarrow must correspond to #states of IO-OOM
- 3. Recover parameters of A from this table

Solution: the Spectral Trick









IO-OOM and Hankel Matrices



 $X = \{a, b, c\} \quad Y = \{\diamondsuit, \heartsuit\}$

IO-OOM and Hankel Matrices

$$X = \{a, b, c\} \quad Y = \{\diamondsuit, \heartsuit\}$$



Hankel Structure:

Fundamental Theorem:

- Equality constraints
- Low-rank constraints

F is realized by an n-state IO-OOM

H has rank at most n for every basis

Max-Margin Completion of Hankel Matrices

Optimization with rank regularization:

 $\underset{\mathsf{H} \in \mathbb{H}(\mathsf{P},\mathsf{S})}{\operatorname{argmin}} \sum_{\mathtt{i}=1}^m L_k(x^{\mathtt{i}}, \mathtt{y}^{\mathtt{i}}, \mathsf{H}) + \tau \; \mathsf{rank}(\mathsf{H})$

Convex relaxation:

 $\underset{H \in \mathbb{H}(\mathsf{P}, \mathsf{S})}{\operatorname{argmin}} \sum_{i=1}^m L_k(x^i, y^i, H) + \tau \, ||H||_*$

- Set of Hankel Matrices over some basis: III(P, S)
- Rank regularizer: rank(H)
- Nuclear norm relaxation: ||H||*

- Optimization almost equivalent \rightarrow we search over IO-OOM that can be recovered from $H \in \mathbb{H}(P, S)$
- Once we solve for H we can recover parameters using the spectral technique

Estimation of Hankel Matrices via Convex Optimization

FOBOS Algorithm: Minimization of $L(H) + \tau ||H||_*$

- Initialize: $H_0 = 0$
- while $t \leq MaxIter do$
 - + Set G_t to a subgradient of L(H) at H_t
 - Set $H_{t+0.5} = H_t \frac{c}{\sqrt{t}}G_t$
 - Calculate the SVD of $H_{t+0.5} = U\Sigma V^{T}$
 - Define a diagonal matrix Σ' such that $\sigma'_i = \max[\sigma_i \nu_t \tau, 0]$
 - set $H_{t+1} = U\Sigma'V^{\top}$

end while

Spectral Recovery using the method by (Hsu et al. 2009)

- Spectral Algorithm for IO-OOM
 - Assume F is realized by a minimal n-state IO-OOM A
 - We are given a basis (P, S) such that H has rank n
 - We are given corresponding H^a_b
 - To recover parameters of A:
 - Perform SVD to get $H = U\Sigma V^{\top}$
 - Define $A_b^a = (HV)^+ H_b^a V$
- Typical spectral algorithms assume that we can estimate H
- In contrast, we regard H as an optimization variable in a loss minimization procedure

Experiments

Task: Phonetic Transcription (UCI "Nettalk" Dataset)





Feature Functionfeature vector $\phi_{abc};ABC(x,y,t) \Longrightarrow$ $1 \text{ if } \begin{array}{c} y_{t-1}y_ty_{t+1} = ABC \\ x_{t-1}x_tx_{t+1} = abc \\ 0 \text{ otherwise} \end{array}$ $\phi(x,y,t) \in \mathbb{R}^{|X|^3 \times |Y|^3}$ $\phi(x,y,t) \in \mathbb{R}^{|X|^3 \times |Y|^3}$ $\psi(x,y,t) \in \mathbb{R}^{|X|^3 \times |Y|^3}$



$$\lambda(x, y) = \sum_{t} w^{t} \phi(x, y, t) \implies \sum_{t} w(x_{t-1}, x_{t}, x_{t+1}; y_{t-1}, y_{t}, y_{t+1})$$








New Regularization for Factorized Linear Models



Specifically designed for structured prediction

Learning parameters of latent-state models









Take-home message: Fundamental ideas behind spectral learning have a wide range of applicability for structured prediction



Thanks!