

Efficient segmental features for speech recognition

Rogier van Dalen
with Mark Gales and Anton Ragni

Google, Mountain View, 23 May 2013

University of Cambridge, Department of Engineering

Segmental features for acoustic modelling

Generative score-spaces

Fast segmental feature extraction

- Log-likelihood score-spaces

- Generative score-spaces

- Experimental results

Conclusion so far

Extensions

- Segmental neural features

- Discriminative training without segmentations

Conclusion

Segmental features for acoustic modelling

Generative score-spaces

Fast segmental feature extraction

Log-likelihood score-spaces

Generative score-spaces

Experimental results

Conclusion so far

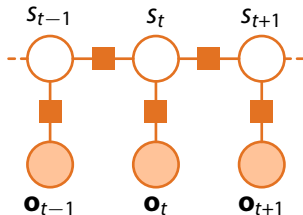
Extensions

Segmental neural features

Discriminative training without segmentations

Conclusion

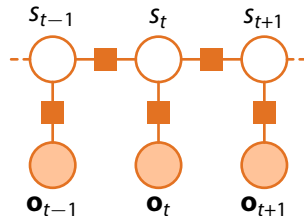
Hidden Markov model:



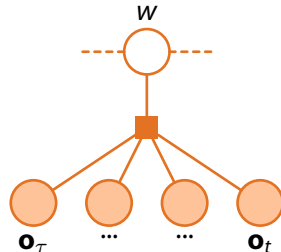
- Markov assumption;
- features (frame-level observations) conditionally independent.
- Training and decoding feasible;
- works in practice, adaptation known.
- Does not model temporal aspects of speech well.

Structured discriminative models

Hidden Markov model



Structured discriminative model

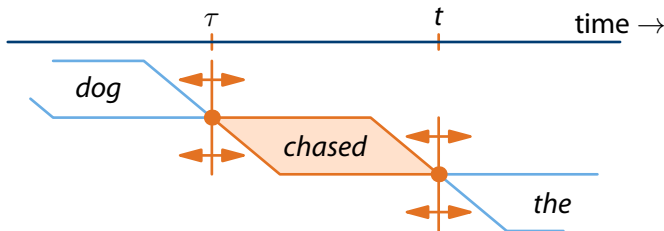


Extract features for a word w (or a phone/triphone).

Segmentations

Extract features for a word w (or a phone/triphone).

⇒ **Unknown start and end times.**



- $1 \leq \tau \leq t \leq T \Rightarrow \Theta(T^2)$ possible segments.
- Average length of segment is $\frac{1}{2}T$.
- Feature depends on all $\frac{1}{2}T$ observations in the segment.

⇒ Extracting features for all segments **normally** $\Omega(T^3)$.

Segmental features for acoustic modelling

Generative score-spaces

Fast segmental feature extraction

Log-likelihood score-spaces

Generative score-spaces

Experimental results

Conclusion so far

Extensions

Segmental neural features

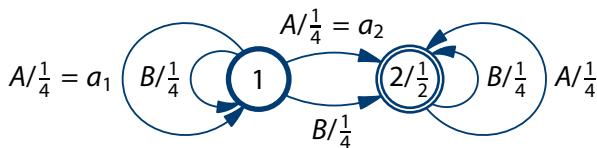
Discriminative training without segmentations

Conclusion

Generative score-spaces: toy example

Two-class, two-symbol problem:

- class 1: $AAAA, BBBB$;
- class 2: $AABB, BBAA$.

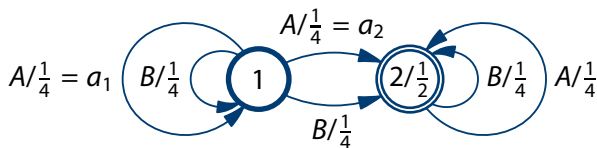


	Class 1		Class 2	
Feature	$AAAA$	$BBBB$	$AABB$	$BBAA$
Log-likelihood	-2.1	-2.1	-2.1	-2.1
∇_{a_1}	6	0	5	1
$\nabla_{a_1} \nabla_{a_2}$	0	0	8	8

Generative score-spaces: toy example

Two-class, two-symbol problem:

- class 1: *AAAA, BBBB*;
- class 2: *AABB, BBAA*.

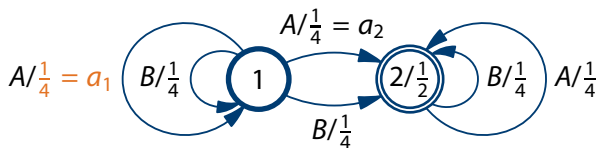


	Class 1		Class 2		
Feature	<i>AAAA</i>	<i>BBBB</i>	<i>AABB</i>	<i>BBAA</i>	
Log-likelihood	-2.1	-2.1	-2.1	-2.1	not separable
∇_{a_1}	6	0	5	1	
$\nabla_{a_1} \nabla_{a_2}$	0	0	8	8	

Generative score-spaces: toy example

Two-class, two-symbol problem:

- class 1: *AAAA, BBBB*;
- class 2: *AABB, BBAA*.

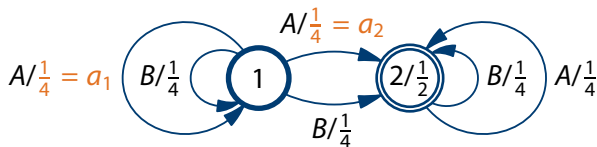


	Class 1		Class 2		
Feature	<i>AAAA</i>	<i>BBBB</i>	<i>AABB</i>	<i>BBAA</i>	
Log-likelihood	-2.1	-2.1	-2.1	-2.1	not separable
∇_{a_1}	6	0	5	1	separable, not linearly
$\nabla_{a_1} \nabla_{a_2}$	0	0	8	8	

Generative score-spaces: toy example

Two-class, two-symbol problem:

- class 1: *AAAA*, *BBBB*;
- class 2: *AABB*, *BBAA*.

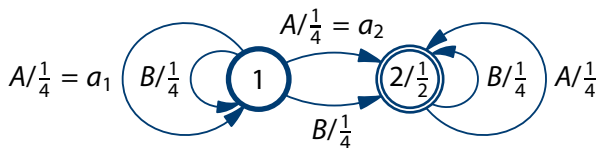


	Class 1		Class 2		
Feature	<i>AAAA</i>	<i>BBBB</i>	<i>AABB</i>	<i>BBAA</i>	
Log-likelihood	-2.1	-2.1	-2.1	-2.1	not separable
∇_{a_1}	6	0	5	1	separable, not linearly
$\nabla_{a_1} \nabla_{a_2}$	0	0	8	8	linearly separable

Generative score-spaces: toy example

Two-class, two-symbol problem:

- class 1: $AAAA, BBBB$;
- class 2: $AABB, BBAA$.



	Class 1		Class 2		
Feature	$AAAA$	$BBBB$	$AABB$	$BBAA$	
Log-likelihood	-2.1	-2.1	-2.1	-2.1	not separable
∇_{a_1}	6	0	5	1	separable, not linearly
$\nabla_{a_1} \nabla_{a_2}$	0	0	8	8	linearly separable

Derivatives **introduce extra dependencies** between observations.

Generative score-spaces:

Likelihood of word HMM $\ell(\mathbf{O}_{\tau:t}; \lambda_w)$ with parameters λ_w .

For every word w in the vocabulary:

$$\phi_w(\mathbf{O}_{\tau:t}) = \begin{bmatrix} \log \ell(\mathbf{O}_{\tau:t}; \lambda_w) \\ \nabla_{\lambda} \log \ell(\mathbf{O}_{\tau:t}; \lambda_w) \end{bmatrix}.$$

Derivatives with respect to HMM parameters:

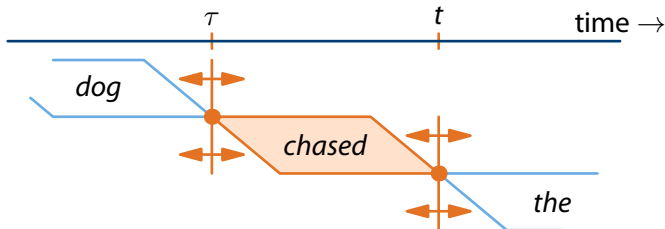
- Introduce dependencies between observations.
 - Feature extraction for all $\Theta(T^2)$ segments **normally** $\Theta(T^3)$.
 - Here: exploit assumptions of HMM.
- ⇒ Feature extraction for $\Theta(T^2)$ segments $\Theta(T^2)$.

Log-linear model

\mathbf{w} is a word sequence, \mathbf{s} is a segmentation of the audio.

$$\arg \max_{\mathbf{w}, \mathbf{s}} P(\mathbf{w}, \mathbf{s} | \mathbf{O}; \alpha) = \arg \max_{\mathbf{w}, \mathbf{s}} \frac{1}{Z(\mathbf{O})} \exp \left(\sum_i \alpha^T \phi(\mathbf{O}_{s_i}, w_i) \right).$$

$\phi(\mathbf{O}_{s_i}, w_i)$ contains HMM log-likelihoods and their derivatives.



Assume the language model fixed.

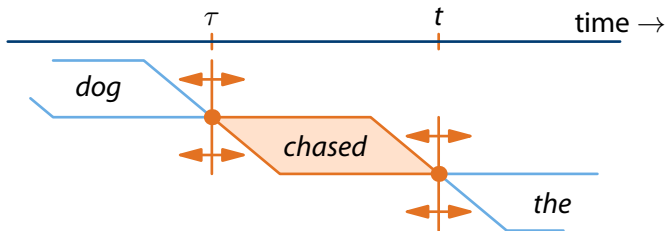
- Decoding is $\Theta(T^2)$ if features for all segments are known.
- ⇒ Feature extraction is the bottleneck.

Log-linear model

\mathbf{w} is a word sequence, \mathbf{s} is a segmentation of the audio.

$$\arg \max_{\mathbf{w}, \mathbf{s}} P(\mathbf{w}, \mathbf{s} | \mathbf{O}; \alpha) = \arg \max_{\mathbf{w}, \mathbf{s}} \sum_i \alpha^T \phi(\mathbf{O}_{s_i}, w_i).$$

$\phi(\mathbf{O}_{s_i}, w_i)$ contains HMM log-likelihoods and their derivatives.



Assume the language model fixed.

- Decoding is $\Theta(T^2)$ if features for all segments are known.
- ⇒ Feature extraction is the bottleneck.

Segmental features for acoustic modelling

Generative score-spaces

Fast segmental feature extraction

Log-likelihood score-spaces

Generative score-spaces

Experimental results

Conclusion so far

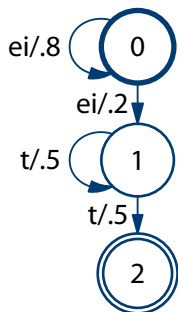
Extensions

Segmental neural features

Discriminative training without segmentations

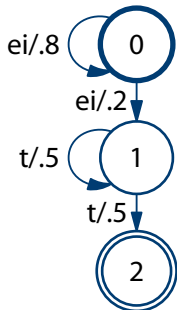
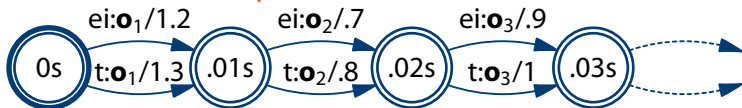
Conclusion

Likelihood of an HMM: marginalise out over discrete sequence.
Produce discrete symbol sequence.



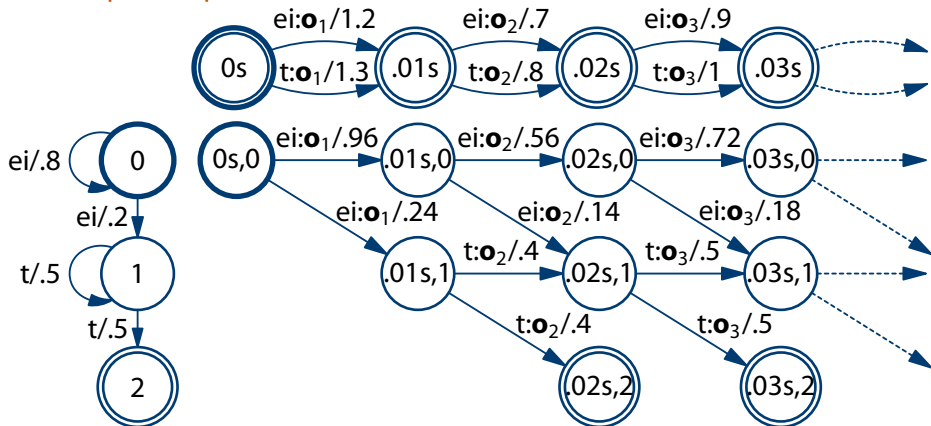
Likelihood of an HMM: marginalise out over discrete sequence.

Produce observations (GMM output densities).



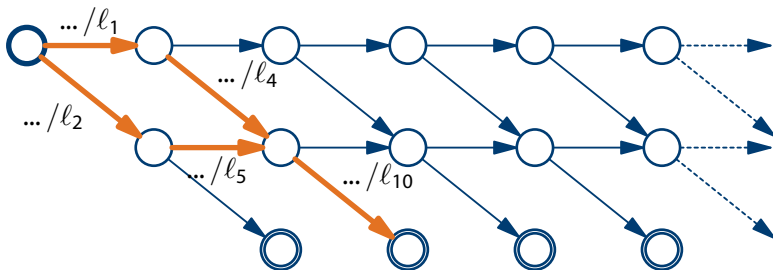
Likelihood of an HMM: marginalise out over discrete sequence.

Compose to produce trellis.



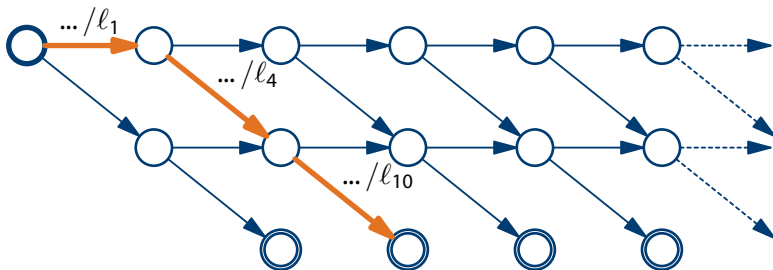
Likelihood

Likelihood of an HMM: sum out over paths in the trellis.



$$\begin{aligned}l(\mathbf{O}_{1:3}; \boldsymbol{\lambda}) &= l_1 \cdot l_4 \cdot l_{10} + l_2 \cdot l_5 \cdot l_{10} \\ &= (l_1 \cdot l_4 + l_2 \cdot l_5) \cdot l_{10}.\end{aligned}$$

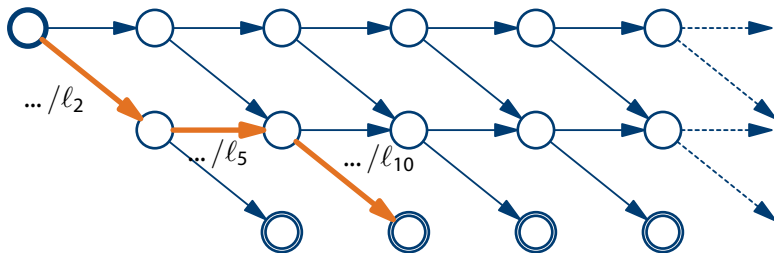
Likelihood of an HMM: sum out over paths in the trellis.



$$\begin{aligned} \ell(\mathbf{O}_{1:3}; \boldsymbol{\lambda}) &= \ell_1 \cdot \ell_4 \cdot \ell_{10} + \ell_2 \cdot \ell_5 \cdot \ell_{10} \\ &= (\ell_1 \cdot \ell_4 + \ell_2 \cdot \ell_5) \cdot \ell_{10}. \end{aligned}$$

Likelihood

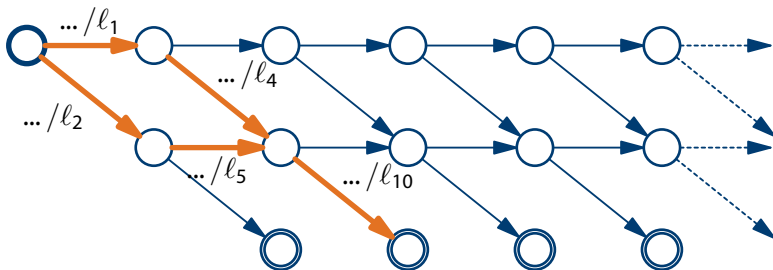
Likelihood of an HMM: sum out over paths in the trellis.



$$\begin{aligned}l(\mathbf{O}_{1:3}; \boldsymbol{\lambda}) &= l_1 \cdot l_4 \cdot l_{10} + l_2 \cdot l_5 \cdot l_{10} \\ &= (l_1 \cdot l_4 + l_2 \cdot l_5) \cdot l_{10}.\end{aligned}$$

Likelihood

Likelihood of an HMM: sum out over paths in the trellis.

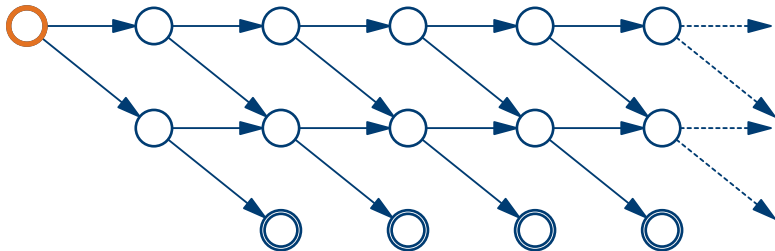


$$\begin{aligned}l(\mathbf{O}_{1:3}; \boldsymbol{\lambda}) &= l_1 \cdot l_4 \cdot l_{10} + l_2 \cdot l_5 \cdot l_{10} \\ &= (l_1 \cdot l_4 + l_2 \cdot l_5) \cdot l_{10}.\end{aligned}$$

⇒ This allows dynamic programming.

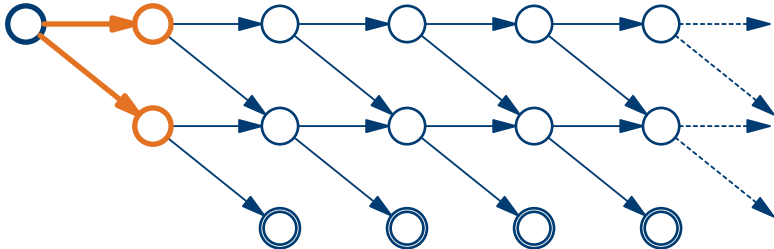
Forward algorithm

Dynamic programming: “frontier” of weights moves forward per time step.



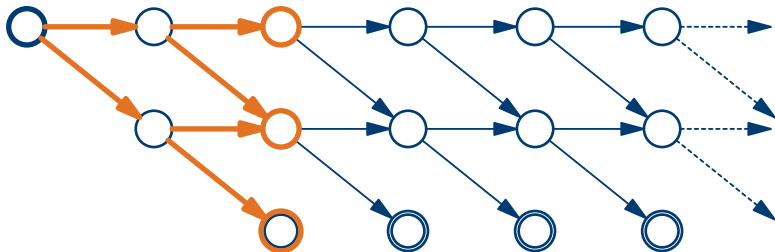
Forward algorithm

Dynamic programming: “frontier” of weights moves forward per time step.



Forward algorithm

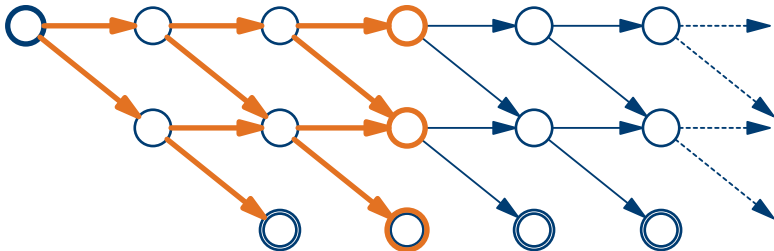
Dynamic programming: “frontier” of weights moves forward per time step.



$$\ell(\mathbf{O}_{1:2}; \lambda)$$

Forward algorithm

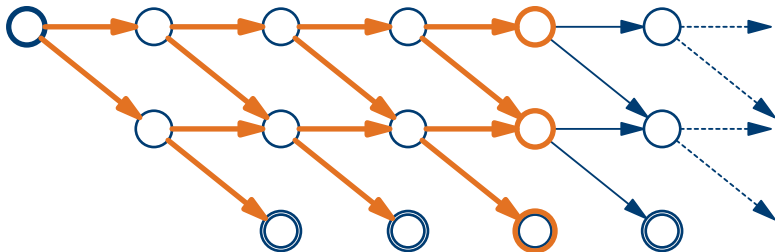
Dynamic programming: “frontier” of weights moves forward per time step.



$$\ell(\mathbf{O}_{1:2}; \boldsymbol{\lambda}); \ell(\mathbf{O}_{1:3}; \boldsymbol{\lambda})$$

Forward algorithm

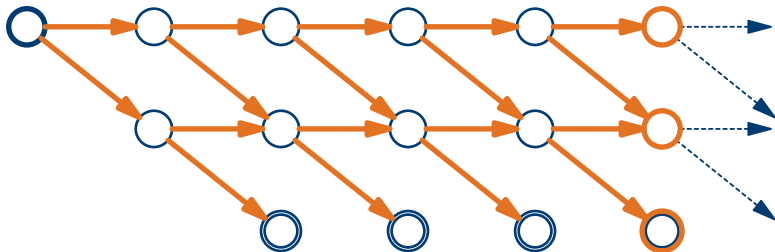
Dynamic programming: “frontier” of weights moves forward per time step.



$$\ell(\mathbf{O}_{1:2}; \lambda); \ell(\mathbf{O}_{1:3}; \lambda); \ell(\mathbf{O}_{1:4}; \lambda)$$

Forward algorithm

Dynamic programming: “frontier” of weights moves forward per time step.



$$\ell(\mathbf{O}_{1:2}; \boldsymbol{\lambda}); \ell(\mathbf{O}_{1:3}; \boldsymbol{\lambda}); \ell(\mathbf{O}_{1:4}; \boldsymbol{\lambda}); \ell(\mathbf{O}_{1:5}; \boldsymbol{\lambda}).$$

- Likelihoods, for all T segments starting at $\tau = 1$, in $\Theta(T)$ time.
- Repeat this for all $\tau = 1 \dots T$:
 - $\Rightarrow \Theta(T^2)$ likelihoods in $\Theta(T^2)$ time.

Add derivatives:

$$\phi_w(\mathbf{O}_{\tau:t}) = \begin{bmatrix} \log \ell(\mathbf{O}_{\tau:t}; \lambda_w) \\ \nabla_{\lambda} \log \ell(\mathbf{O}_{\tau:t}; \lambda_w) \end{bmatrix}, \quad \nabla_{\lambda} \log \ell = \frac{\nabla_{\lambda} \ell}{\ell}.$$

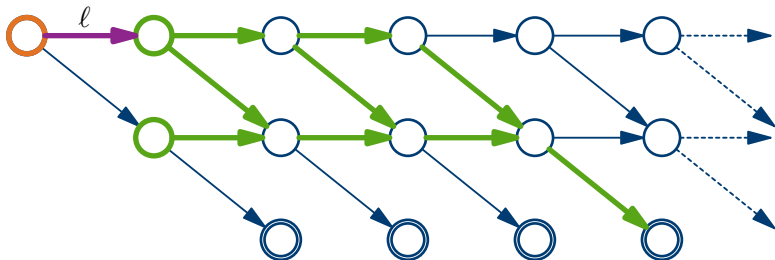
$$\nabla_{\lambda} \ell(\mathbf{O}_{\tau:t}; \lambda_w) = \sum_e \gamma_{\tau:t}(e) \frac{\nabla_{\lambda} \ell[e; \lambda_w]}{\ell[e; \lambda_w]}.$$

$\gamma_{\tau:t}(e)$ is the arc posterior.

This uses the forward-backward algorithm.

Forward-backward algorithm

Forward-backward algorithm to compute the arc posterior $\gamma_{\tau:t}(e)$:

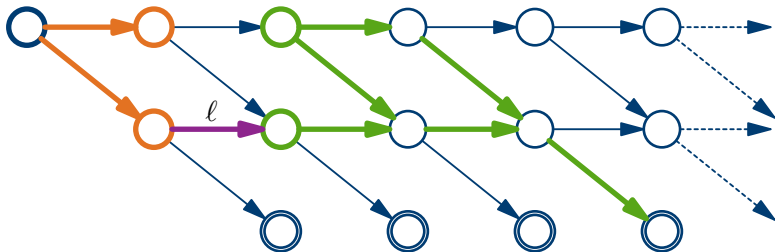


$$\gamma_{\tau:t}(e) \propto \text{forwd}(\dots) \cdot \ell \cdot \text{backwd}(\dots).$$

- Derivatives introduce dependencies between observations.
- This multiplication must be redone for every segment:
 - $\Theta(T^3)$ for all segments.

Forward-backward algorithm

Forward-backward algorithm to compute the arc posterior $\gamma_{\tau:t}(e)$:

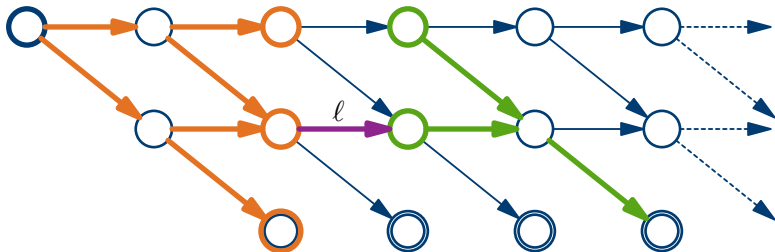


$$\gamma_{\tau:t}(e) \propto \text{forwd}(\dots) \cdot l \cdot \text{backwd}(\dots).$$

- Derivatives introduce dependencies between observations.
- This multiplication must be redone for every segment:
 - $\Theta(T^3)$ for all segments.

Forward-backward algorithm

Forward-backward algorithm to compute the arc posterior $\gamma_{\tau:t}(e)$:

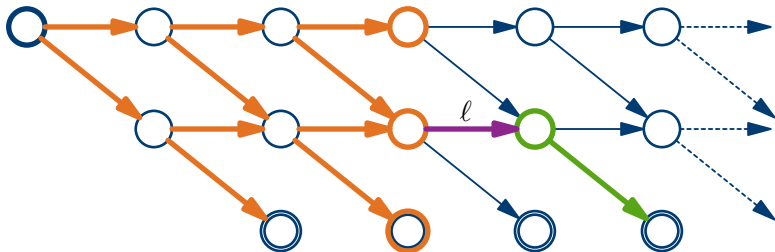


$$\gamma_{\tau:t}(e) \propto \text{forwd}(\dots) \cdot \ell \cdot \text{backwd}(\dots).$$

- Derivatives introduce dependencies between observations.
- This multiplication must be redone for every segment:
 - $\Theta(T^3)$ for all segments.

Forward-backward algorithm

Forward-backward algorithm to compute the arc posterior $\gamma_{\tau:t}(e)$:

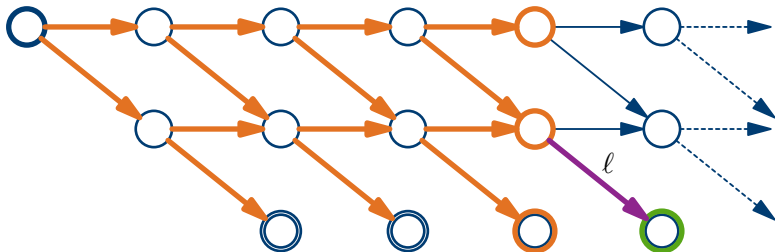


$$\gamma_{\tau:t}(e) \propto \text{forwd}(\dots) \cdot l \cdot \text{backwd}(\dots).$$

- Derivatives introduce dependencies between observations.
- This multiplication must be redone for every segment:
 - $\Theta(T^3)$ for all segments.

Forward-backward algorithm

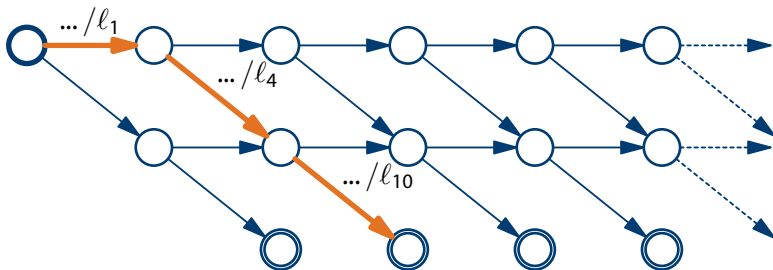
Forward-backward algorithm to compute the arc posterior $\gamma_{\tau:t}(e)$:



$$\gamma_{\tau:t}(e) \propto \text{forwd}(\dots) \cdot \ell \cdot \text{backwd}(\dots).$$

- Derivatives introduce dependencies between observations.
- This multiplication must be redone for every segment:
 - $\Theta(T^3)$ for all segments.

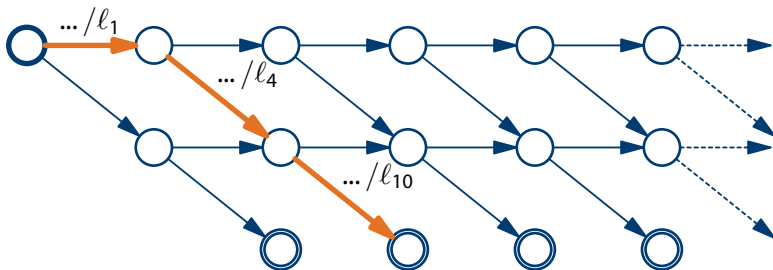
Derivative of the likelihood of an HMM:



$$l(\mathbf{O}_{1:3}; \lambda) = l_1 \cdot l_4 \cdot l_{10} + l_2 \cdot l_5 \cdot l_{10};$$

$$\nabla_{\lambda} l(\mathbf{O}_{1:3}; \lambda) = (\nabla_{\lambda} l_1) \cdot l_4 \cdot l_{10} + l_1 \cdot \nabla_{\lambda} (l_4 \cdot l_{10}) + \dots$$

Derivative of the likelihood of an HMM:



$$l(\mathbf{O}_{1:3}; \lambda) = l_1 \cdot l_4 \cdot l_{10} + l_2 \cdot l_5 \cdot l_{10};$$

$$\nabla_{\lambda} l(\mathbf{O}_{1:3}; \lambda) = (\nabla_{\lambda} l_1) \cdot l_4 \cdot l_{10} + l_1 \cdot \nabla_{\lambda} (l_4 \cdot l_{10}) + \dots$$

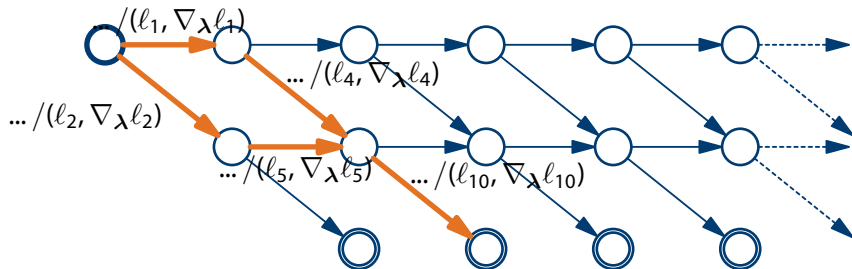
⇒ Generalise this: expectation **semiring**.

What is a **semiring**?

- Set with operations \oplus and \otimes , and values $\bar{0}$ and $\bar{1}$.
- E.g. \mathbb{R} with $+$ and \times , and 0 and 1.
- E.g. $\mathbb{R}^+ \cup \{0\}$ with \max and \times , and 0 and 1
 - turns the forward algorithm into Viterbi.
- E.g. the **expectation semiring**: a tuple (ℓ, c) .

Expectation semiring

Replace weights by tuples $\sigma = (l, \nabla_{\lambda} l)$ in the expectation semiring:



Define operations \otimes and \oplus :

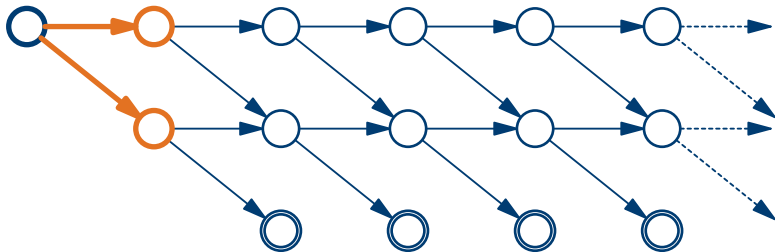
$$(l_1, \nabla_{\lambda} l_1) \otimes (l_4, \nabla_{\lambda} l_4) \triangleq (l_1 \cdot l_4, l_1 \cdot \nabla_{\lambda} l_4 + l_4 \cdot \nabla_{\lambda} l_1)$$

$$(l_1, \nabla_{\lambda} l_1) \oplus (l_2, \nabla_{\lambda} l_2) \triangleq (l_1 + l_2, \nabla_{\lambda} l_1 + \nabla_{\lambda} l_2)$$

$$(\sigma_1 \otimes \sigma_4 \otimes \sigma_{10}) \oplus (\sigma_2 \otimes \sigma_5 \otimes \sigma_{10}) = (l(\mathbf{O}_{1:3}; \lambda), \nabla_{\lambda} l(\mathbf{O}_{1:3}; \lambda)) \triangleq$$

Expectation semiring

The forward algorithm with the expectation semiring:

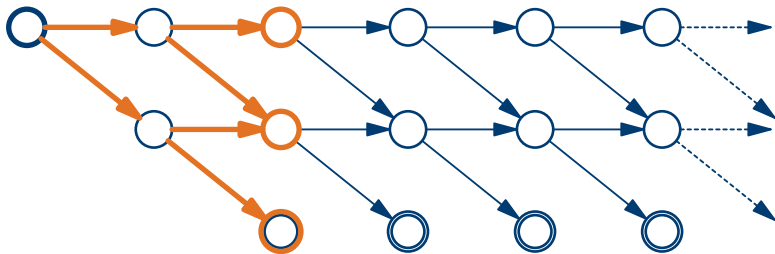


Likelihoods and their derivatives,

- for all T segments starting at $\tau = 1$, in $\Theta(T)$ time.
- for all $\Theta(T^2)$ segments, in $\Theta(T^2)$ time.

Expectation semiring

The forward algorithm with the expectation semiring:



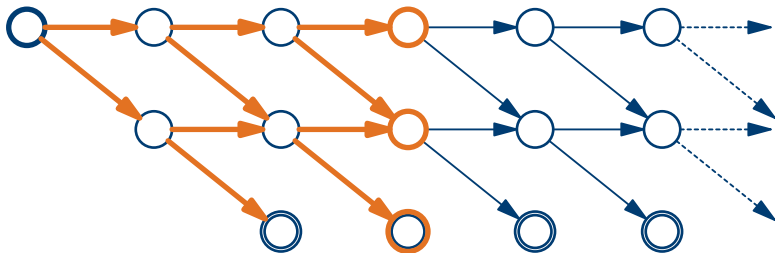
$$(\ell(\mathbf{O}_{1:2}; \boldsymbol{\lambda}), \nabla_{\boldsymbol{\lambda}} \ell(\mathbf{O}_{1:2}; \boldsymbol{\lambda}))$$

Likelihoods and their derivatives,

- for all T segments starting at $\tau = 1$, in $\Theta(T)$ time.
- for all $\Theta(T^2)$ segments, in $\Theta(T^2)$ time.

Expectation semiring

The forward algorithm with the expectation semiring:



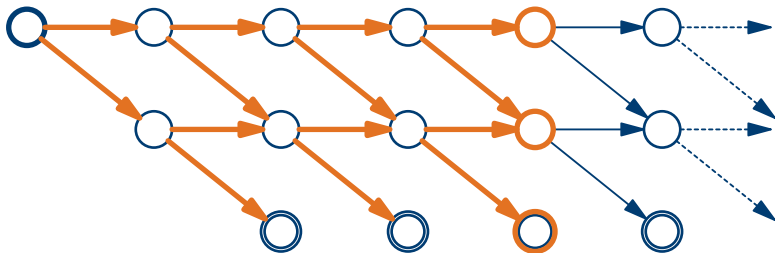
$$(\ell(\mathbf{O}_{1:2}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:2}; \lambda)); \quad (\ell(\mathbf{O}_{1:3}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:3}; \lambda))$$

Likelihoods and their derivatives,

- for all T segments starting at $\tau = 1$, in $\Theta(T)$ time.
- for all $\Theta(T^2)$ segments, in $\Theta(T^2)$ time.

Expectation semiring

The forward algorithm with the expectation semiring:



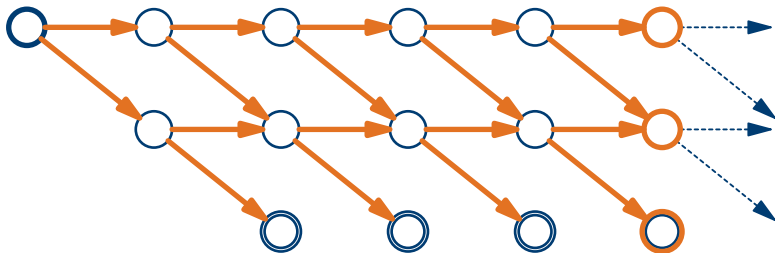
$$\begin{aligned} & (\ell(\mathbf{O}_{1:2}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:2}; \lambda)); & (\ell(\mathbf{O}_{1:3}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:3}; \lambda)); \\ & (\ell(\mathbf{O}_{1:4}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:4}; \lambda)) \end{aligned}$$

Likelihoods and their derivatives,

- for all T segments starting at $\tau = 1$, in $\Theta(T)$ time.
- for all $\Theta(T^2)$ segments, in $\Theta(T^2)$ time.

Expectation semiring

The forward algorithm with the expectation semiring:



$$\begin{aligned} & (\ell(\mathbf{O}_{1:2}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:2}; \lambda)); & (\ell(\mathbf{O}_{1:3}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:3}; \lambda)); \\ & (\ell(\mathbf{O}_{1:4}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:4}; \lambda)); & (\ell(\mathbf{O}_{1:5}; \lambda), \nabla_{\lambda} \ell(\mathbf{O}_{1:5}; \lambda)). \end{aligned}$$

Likelihoods and their derivatives,

- for all T segments starting at $\tau = 1$, in $\Theta(T)$ time.
- for all $\Theta(T^2)$ segments, in $\Theta(T^2)$ time.

AURORA 2: noise-corrupted TIDIGITS

- Minimum Bayes risk training on large lattices.
- Results averaged over 0-20 dB.

System	Segmentation	Test set			Average
		A	B	C	
HMM		9.8	9.1	9.5	9.5
ℓ	HMM	8.1	7.4	8.2	7.8
	optimal	7.8	7.3	8.0	7.6
$\begin{bmatrix} \ell \\ \nabla \ell \end{bmatrix}$	HMM	7.0	6.6	7.6	7.0
	optimal	6.8	6.4	7.3	6.7

Segmental features for acoustic modelling

Generative score-spaces

Fast segmental feature extraction

Log-likelihood score-spaces

Generative score-spaces

Experimental results

Conclusion so far

Extensions

Segmental neural features

Discriminative training without segmentations

Conclusion

- Segmental features: generative score-spaces:
 - derivatives of likelihood;
 - introduce dependencies between observations.
- Extract features in quadratic time.
- Decoding in quadratic time.