

# Master in Artificial Intelligence

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

## Advanced Human Language Technologies



UNIVERSITAT POLITÈCNICA DE CATALUNYA  
BARCELONATECH

Facultat d'Informàtica de Barcelona



# Outline

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

## 1 Sequence Prediction

- Examples
- Problem Formulation

## 2 Approaches

- Local Classifiers
- HMMs
- Global Predictors

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)
- Conditional Random Fields (CRF)

# Outline

## 1 Sequence Prediction

### ■ Examples

### ■ Problem Formulation

## 2 Approaches

### ■ Local Classifiers

### ■ HMMs

### ■ Global Predictors

## 3 Log-linear Models for Sequence Prediction

### ■ Maximum Entropy Markov Models (MEMMs)

### ■ Conditional Random Fields (CRF)

Sequence  
Prediction

Examples

Approaches

Log-linear  
Models for  
Sequence  
Prediction

# Examples - Named Entity Recognition (NER)

Sequence  
Prediction

Examples

Approaches

Log-linear  
Models for  
Sequence  
Prediction

<b>y</b>	PER	-	QNT	-	-	ORG	ORG	-	TIME
<b>x</b>	Jim	bought	300	shares	of	Acme	Corp.	in	2006

# Examples - Named Entity Recognition (NER)

Sequence  
Prediction

Examples

Approaches

Log-linear  
Models for  
Sequence  
Prediction

y	PER	-	QNT	-	-	ORG	ORG	-	TIME
x	Jim	bought	300	shares	of	Acme	Corp.	in	2006

y	PER	PER	-	-	LOC
x	Jack	London	went	to	Paris

y	PER	PER	-	-	LOC
x	Paris	Hilton	went	to	London

# Examples - Part-of-Speech (PoS) Tagging

Sequence  
Prediction

Examples

Approaches

Log-linear  
Models for  
Sequence  
Prediction

<b>y</b>	DT	NN	VBZ	IN	DT	JJ	NN	.
<b>x</b>	The	fox	jumps	over	the	lazy	dog	.

# Examples - Part-of-Speech (PoS) Tagging

Sequence  
Prediction

Examples

Approaches

Log-linear  
Models for  
Sequence  
Prediction

<b>y</b>	DT	NN	VBZ	IN	DT	JJ	NN	.
<b>x</b>	The	fox	jumps	over	the	lazy	dog	.

<b>y</b>	DT	NN	NN	VBD	DT	JJ	NN	.
<b>x</b>	The	fox	jumps	scared	the	lazy	dog	.

# Outline

## 1 Sequence Prediction

- Examples

- Problem Formulation

## 2 Approaches

- Local Classifiers

- HMMs

- Global Predictors

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)

- Conditional Random Fields (CRF)

Sequence  
Prediction

Problem Formulation

Approaches

Log-linear  
Models for  
Sequence  
Prediction



# Problem Formulation

Sequence  
Prediction  
Problem Formulation

Approaches

Log-linear  
Models for  
Sequence  
Prediction

- $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \dots \mathbf{x}_n$  are input sequences,  $\mathbf{x}_i \in \mathcal{X}$
- $\mathbf{y} = \mathbf{y}_1\mathbf{y}_2 \dots \mathbf{y}_n$  are output sequences,  $\mathbf{y}_i \in \{1, \dots, L\}$

- **Goal:** given training data

$$\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$$

learn a predictor  $\mathbf{x} \rightarrow \mathbf{y}$  that **works well** on unseen inputs  $\mathbf{x}$

- What is the form of our prediction model?

# Outline

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

- 1 Sequence Prediction
  - Examples
  - Problem Formulation
- 2 Approaches
  - Local Classifiers
  - HMMs
  - Global Predictors
- 3 Log-linear Models for Sequence Prediction
  - Maximum Entropy Markov Models (MEMMs)
  - Conditional Random Fields (CRF)

# Outline

## 1 Sequence Prediction

- Examples
- Problem Formulation

## 2 Approaches

- Local Classifiers
- HMMs
- Global Predictors

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)
- Conditional Random Fields (CRF)

Sequence  
Prediction

Approaches

Local Classifiers

Log-linear  
Models for  
Sequence  
Prediction

# Approach 1: Local Classifiers

Jack ? London went to Paris

Decompose the sequence into  $n$  classification problems:

- A classifier predicts individual labels at each position

$$\hat{y}_i = \underset{y \in \{\text{LOC}, \text{PER}, -\}}{\operatorname{argmax}} \quad \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, y)$$

- $\mathbf{f}(\mathbf{x}, i, y)$  represents an assignment of label  $y$  for  $x_i$
- $\mathbf{w}$  is a vector of parameters, has a weight for each feature of  $\mathbf{f}$ 
  - Use standard classification methods to learn  $\mathbf{w}$

# Approach 1: Local Classifiers

Jack ? London went to Paris

Decompose the sequence into  $n$  classification problems:

- A classifier predicts individual labels at each position

$$\hat{y}_i = \operatorname{argmax}_{y \in \{\text{LOC}, \text{PER}, -\}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, y)$$

- $\mathbf{f}(\mathbf{x}, i, y)$  represents an assignment of label  $y$  for  $x_i$
- $\mathbf{w}$  is a vector of parameters, has a weight for each feature of  $\mathbf{f}$ 
  - Use standard classification methods to learn  $\mathbf{w}$
- At test time, predict the best sequence by
  - a simple concatenation of the best label for each position

# Indicator Features

- $\mathbf{f}(\mathbf{x}, i, y)$  is a vector of  $d$  features representing label  $y$  for  $\mathbf{x}_i$

$$\mathbf{f}(\mathbf{x}, i, y) = ( f_1(\mathbf{x}, i, y), \dots, f_j(\mathbf{x}, i, y), \dots, f_d(\mathbf{x}, i, y) )$$

- What's in a feature  $f_j(\mathbf{x}, i, y)$ ?
  - Anything we can compute using  $\mathbf{x}$  and  $i$  and  $y$
  - Anything that indicates whether  $y$  is a good (or bad) label for  $\mathbf{x}_i$
  - **Indicator features:** binary-valued features looking at a single simple property

$$f_j(\mathbf{x}, i, y) = \begin{cases} 1 & \text{if } \mathbf{x}_i = \text{London and } y = \text{LOC} \\ 0 & \text{otherwise} \end{cases}$$

$$f_k(\mathbf{x}, i, y) = \begin{cases} 1 & \text{if } \mathbf{x}_{i+1} = \text{went and } y = \text{LOC} \\ 0 & \text{otherwise} \end{cases}$$

# More Features for NE Recognition

Jack      <sup>PER</sup>  
London    went   to   Paris

In practice, construct  $f(\mathbf{x}, i, y)$  by ...

- Define a number of simple patterns of  $\mathbf{x}$  and  $i$ 
  - current word  $x_i$
  - is  $x_i$  capitalized?
  - $x_i$  has digits?
  - prefixes/suffixes of size 1, 2, 3,  
...
  - is  $x_i$  a known location?
  - is  $x_i$  a known person?
- next word
- previous word
- current and next words together
- other combinations
- Generate features by combining patterns with possible labels  $y$

# More Features for NE Recognition

PER      PER      -  
Jack   London   went   to   Paris

In practice, construct  $f(\mathbf{x}, i, y)$  by ...

- Define a number of simple patterns of  $\mathbf{x}$  and  $i$ 
  - current word  $x_i$
  - is  $x_i$  capitalized?
  - $x_i$  has digits?
  - prefixes/suffixes of size 1, 2, 3,  
...
  - is  $x_i$  a known location?
  - is  $x_i$  a known person?
- next word
- previous word
- current and next words together
- other combinations
- Generate features by combining patterns with possible labels  $y$

**Main limitation:** features can't capture interactions between labels!



# Outline

## 1 Sequence Prediction

- Examples
- Problem Formulation

## 2 Approaches

- Local Classifiers
- **HMMs**
- Global Predictors

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)
- Conditional Random Fields (CRF)

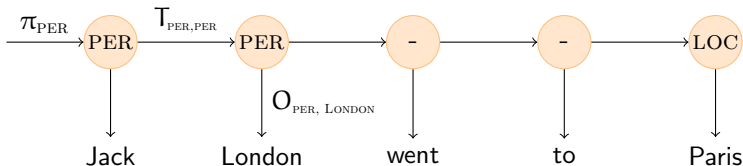
Sequence  
Prediction

Approaches

HMMs

Log-linear  
Models for  
Sequence  
Prediction

## Approach 2: HMM for Sequence Prediction

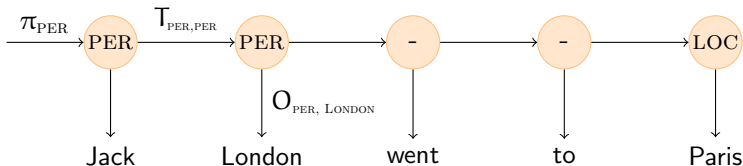


- Define an HMM where each label is a state
- Model parameters:
  - $\pi_y$  : probability of starting with label  $y$
  - $T_{yy'}$  : probability of transitioning from label  $y$  to  $y'$
  - $O_{yx}$  : probability of generating symbol  $x$  given label  $y$
- Predictions:

$$p(\mathbf{x}, \mathbf{y}) = \pi_{y_1} O_{y_1 x_1} \prod_{i>1} T_{y_{i-1} y_i} O_{y_i x_i}$$

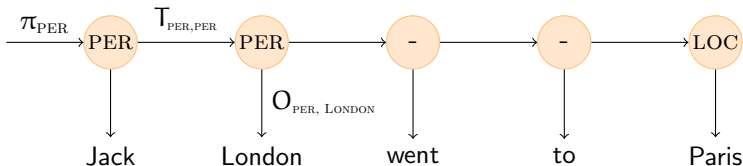
- Learning: relative counts + smoothing
- Prediction: Viterbi algorithm

## Approach 2: Representation in HMM



- Label interactions are captured in the transition parameters
- But interactions between symbols and labels are quite limited!
  - Only  $O_{y_i x_i} = p(x_i | y_i)$
  - Not clear how to exploit patterns such as:
    - Capitalization, digits
    - Prefixes and suffixes
    - Next word, previous word
    - Combinations of these with label transitions

## Approach 2: Representation in HMM



Sequence  
Prediction

Approaches

HMMs

Log-linear  
Models for  
Sequence  
Prediction

- Label interactions are captured in the transition parameters
- But interactions between symbols and labels are quite limited!
  - Only  $O_{y_i x_i} = p(x_i | y_i)$
  - Not clear how to exploit patterns such as:
    - Capitalization, digits
    - Prefixes and suffixes
    - Next word, previous word
    - Combinations of these with label transitions
- Why? HMM independence assumptions:  
given label  $y_i$ , token  $x_i$  is independent of anything else

# Local Classifiers vs. HMM

## LOCAL CLASSIFIERS

- Form:

$$\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, y)$$

- Learning: standard classifiers
- Prediction: independent for each  $\mathbf{x}_i$
- Advantage: feature-rich
- Drawback: no label interactions

## HMM

- Form:

$$\pi_{y_1} O_{y_1, x_1} \prod_{i>1} T_{y_{i-1}, y_i} O_{y_i, x_i}$$

- Learning: relative counts
- Prediction: Viterbi
- Advantage: label interactions
- Drawback: no fine-grained features

# Outline

## 1 Sequence Prediction

- Examples
- Problem Formulation

## 2 Approaches

- Local Classifiers
- HMMs
- **Gobal Predictors**

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)
- Conditional Random Fields (CRF)

Sequence  
Prediction

Approaches  
Gobal Predictors

Log-linear  
Models for  
Sequence  
Prediction

## Approach 3: Global Sequence Predictors

<b>y:</b>	PER	PER	-	-	LOC
<b>x:</b>	Jack	London	went	to	Paris

Learn a single classifier from  $\mathbf{x} \rightarrow \mathbf{y}$

$$\text{predict}(\mathbf{x}_{1:n}) = \underset{\mathbf{y} \in \mathcal{Y}^n}{\operatorname{argmax}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y})$$

## Approach 3: Global Sequence Predictors

<b>y:</b>	PER	PER	-	-	LOC
<b>x:</b>	Jack	London	went	to	Paris

Learn a single classifier from  $\mathbf{x} \rightarrow \mathbf{y}$

$$\text{predict}(\mathbf{x}_{1:n}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y})$$

But ...

- How do we represent entire sequences in  $\mathbf{f}(\mathbf{x}, \mathbf{y})$ ?
- There are **exponentially-many** sequences  $\mathbf{y}$  for a given  $\mathbf{x}$ , how do we solve the **argmax** problem?



# Factored Representations

y:	PER	PER	-	-	LOC
x:	Jack	London	went	to	Paris

- How do we represent entire sequences in  $f(\mathbf{x}, \mathbf{y})$ ?

Sequence  
Prediction

Approaches

Global Predictors

Log-linear  
Models for  
Sequence  
Prediction

# Factored Representations

y:	PER	PER	-	-	LOC
x:	Jack	London	went	to	Paris

- How do we represent entire sequences in  $f(\mathbf{x}, \mathbf{y})$ ?
  - Look at the full label sequence  $\mathbf{y}$  (intractable)

# Factored Representations

y:	PER	PER	-	-	LOC
x:	Jack	London	went	to	Paris

- How do we represent entire sequences in  $f(\mathbf{x}, \mathbf{y})$ ?
  - Look at the full label sequence  $\mathbf{y}$  (intractable)
  - Look at  $n$ -grams of output labels  $\langle \mathbf{y}_{i-n+1}, \dots, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (too expensive)

# Factored Representations

<b>y:</b>	PER	PER	-	-	LOC
<b>x:</b>	Jack	London	went	to	Paris

- How do we represent entire sequences in  $\mathbf{f}(\mathbf{x}, \mathbf{y})$ ?
  - Look at the full label sequence  $\mathbf{y}$  (intractable)
  - Look at **n-grams** of output labels  $\langle \mathbf{y}_{i-n+1}, \dots, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (too expensive)
  - Look at **trigrams** of output labels  $\langle \mathbf{y}_{i-2}, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (possible for small  $|\mathcal{Y}|$ )

# Factored Representations

y:	PER	PER	-	-	LOC
x:	Jack	London	went	to	Paris

- How do we represent entire sequences in  $f(\mathbf{x}, \mathbf{y})$ ?
  - Look at the full label sequence  $\mathbf{y}$  (intractable)
  - Look at **n-grams** of output labels  $\langle \mathbf{y}_{i-n+1}, \dots, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (too expensive)
  - Look at **trigrams** of output labels  $\langle \mathbf{y}_{i-2}, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (possible for small  $|\mathcal{Y}|$ )
  - Look at **bigrams** of output labels  $\langle \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (definitely tractable)

# Factored Representations

y:	PER	PER	-	-	LOC
x:	Jack	London	went	to	Paris

- How do we represent entire sequences in  $f(\mathbf{x}, \mathbf{y})$ ?
  - Look at the full label sequence  $\mathbf{y}$  (intractable)
  - Look at **n-grams** of output labels  $\langle \mathbf{y}_{i-n+1}, \dots, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (too expensive)
  - Look at **trigrams** of output labels  $\langle \mathbf{y}_{i-2}, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (possible for small  $|\mathcal{Y}|$ )
  - Look at **bigrams** of output labels  $\langle \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (definitely tractable)
  - Look at individual assignments  $\mathbf{y}_i$  (standard classification)

# Factored Representations

y:	PER	PER	-	-	LOC
x:	Jack	London	went	to	Paris

- How do we represent entire sequences in  $f(\mathbf{x}, \mathbf{y})$ ?
  - Look at the full label sequence  $\mathbf{y}$  (intractable)
  - Look at **n-grams** of output labels  $\langle \mathbf{y}_{i-n+1}, \dots, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (too expensive)
  - Look at **trigrams** of output labels  $\langle \mathbf{y}_{i-2}, \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (possible for small  $|\mathcal{Y}|$ )
  - Look at **bigrams** of output labels  $\langle \mathbf{y}_{i-1}, \mathbf{y}_i \rangle$  (definitely tractable)
  - Look at individual assignments  $\mathbf{y}_i$  (standard classification)
- A factored representation will lead to a tractable model

# Bigram Indicator Features

	1	2	3	4	5
y	PER	PER	-	-	LOC
x	Jack	London	went	to	Paris

- Indicator features:

$$f_j(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) = \begin{cases} 1 & \text{if } \mathbf{x}_i = \text{"London"} \text{ and} \\ & \mathbf{y}_{i-1} = \text{PER and } \mathbf{y}_i = \text{PER} \\ 0 & \text{otherwise} \end{cases}$$

e.g.,  $f_j(\mathbf{x}, 2, \text{PER}, \text{PER}) = 1$ ,  $f_j(\mathbf{x}, 3, \text{PER}, -) = 0$



# More Bigram Indicator Features

	1	2	3	4	5
x	Jack	London	went	to	Paris
y	PER	PER	-	-	LOC
y'	PER	LOC	-	-	LOC
y''	-	-	-	LOC	-
x'	My	trip	to	London	...

$f_1(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{PER}$

$f_2(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{LOC}$

$f_3(x, i, y_{i-1}, y_i) = 1$  iff  $x_{i-1} \sim /(\text{in}|\text{to}|\text{at})/$  &  $x_i \sim /^{[A-Z]}/$  &  $y_i = \text{LOC}$

$f_4(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{LOC}$  &  $\text{WORLD-CITIES}(x_i) = 1$

$f_5(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{PER}$  &  $\text{FIRST-NAMES}(x_i) = 1$

# More Bigram Indicator Features

	1	2	3	4	5
x	Jack	London	went	to	Paris
y	PER	PER	-	-	LOC
y'	PER	LOC	-	-	LOC
y''	-	-	-	LOC	-
x'	My	trip	to	London	...

$f_1(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{PER}$

$f_2(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{LOC}$

$f_3(x, i, y_{i-1}, y_i) = 1$  iff  $x_{i-1} \sim /(\text{in}|\text{to}|\text{at})/$  &  $x_i \sim /^{[A-Z]}/$  &  $y_i = \text{LOC}$

$f_4(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{LOC}$  &  $\text{WORLD-CITIES}(x_i) = 1$

$f_5(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{PER}$  &  $\text{FIRST-NAMES}(x_i) = 1$

# More Bigram Indicator Features

	1	2	3	4	5
x	Jack	London	went	to	Paris
y	PER	PER	-	-	LOC
y'	PER	LOC	-	-	LOC
y''	-	-	-	LOC	-
x'	My	trip	to	London	...

$f_1(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{PER}$

$f_2(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{LOC}$

$f_3(x, i, y_{i-1}, y_i) = 1$  iff  $x_{i-1} \sim /(\text{in|to|at})/$  &  $x_i \sim /^{[A-Z]}/$  &  $y_i = \text{LOC}$

$f_4(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{LOC}$  &  $\text{WORLD-CITIES}(x_i) = 1$

$f_5(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{PER}$  &  $\text{FIRST-NAMES}(x_i) = 1$

# More Bigram Indicator Features

	1	2	3	4	5
x	Jack	London	went	to	Paris
y	PER	PER	-	-	LOC
y'	PER	LOC	-	-	LOC
y''	-	-	-	LOC	-
x'	My	trip	to	London	...

$f_1(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{PER}$

$f_2(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{LOC}$

$f_3(x, i, y_{i-1}, y_i) = 1$  iff  $x_{i-1} \sim /(\text{in}|\text{to}|\text{at})/$  &  $x_i \sim /^{[A-Z]}/$  &  $y_i = \text{LOC}$

$f_4(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{LOC}$  &  $\text{WORLD-CITIES}(x_i) = 1$

$f_5(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{PER}$  &  $\text{FIRST-NAMES}(x_i) = 1$

# More Bigram Indicator Features

	1	2	3	4	5
x	Jack	London	went	to	Paris
y	PER	PER	-	-	LOC
y'	PER	LOC	-	-	LOC
y''	-	-	-	LOC	-
x'	My	trip	to	London	...

$f_1(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{PER}$

$f_2(x, i, y_{i-1}, y_i) = 1$  iff  $x_i = \text{"London"}$  &  $y_{i-1} = \text{PER}$  &  $y_i = \text{LOC}$

$f_3(x, i, y_{i-1}, y_i) = 1$  iff  $x_{i-1} \sim /(\text{in|to|at})/$  &  $x_i \sim /^{[A-Z]}/$  &  $y_i = \text{LOC}$

$f_4(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{LOC}$  &  $\text{WORLD-CITIES}(x_i) = 1$

$f_5(x, i, y_{i-1}, y_i) = 1$  iff  $y_i = \text{PER}$  &  $\text{FIRST-NAMES}(x_i) = 1$

# Bigram-Factored Representations

y:	PER	PER	-	-	LOC
x:	Jack	London	went	to	Paris

- $\mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) = (\mathbf{f}_1(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i), \dots, \mathbf{f}_d(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i))$ 
  - A d-dimensional feature vector of a label bigram at i
  - Each dimension is typically a boolean indicator (0 or 1)
- $\mathbf{f}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$ 
  - A d-dimensional feature vector of the entire y
  - Aggregated representation by summing bigram feature vectors
  - Each dimension is now a **count** of a feature pattern

# Linear Sequence Prediction

$$\text{best}(\mathbf{x}_{1:n}) = \underset{\mathbf{y} \in \mathcal{Y}^n}{\operatorname{argmax}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) = \underset{\mathbf{y} \in \mathcal{Y}^n}{\operatorname{argmax}} \mathbf{w} \cdot \sum_{i=1}^n \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$

Sequence  
Prediction

Approaches

Global Predictors

Log-linear  
Models for  
Sequence  
Prediction

# Linear Sequence Prediction

Sequence  
Prediction

Approaches

Global Predictors

Log-linear  
Models for  
Sequence  
Prediction

$$\text{best}(\mathbf{x}_{1:n}) = \underset{\mathbf{y} \in \mathcal{Y}^n}{\operatorname{argmax}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) = \underset{\mathbf{y} \in \mathcal{Y}^n}{\operatorname{argmax}} \mathbf{w} \cdot \sum_{i=1}^n \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$

■ Note the linearity of the expression:

$$\begin{aligned} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) &= \mathbf{w} \cdot \sum_{i=1}^n \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) = \sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) \\ &= \sum_{i=1}^n \sum_{j=1}^d \mathbf{w}_j \mathbf{f}_j(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) \end{aligned}$$



# Linear Sequence Prediction

Sequence  
Prediction

Approaches

Global Predictors

Log-linear  
Models for  
Sequence  
Prediction

$$\text{best}(\mathbf{x}_{1:n}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \mathbf{w} \cdot \sum_{i=1}^n \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$

- Note the linearity of the expression:

$$\begin{aligned} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}) &= \mathbf{w} \cdot \sum_{i=1}^n \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) = \sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) \\ &= \sum_{i=1}^n \sum_{j=1}^d \mathbf{w}_j \mathbf{f}_j(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) \end{aligned}$$

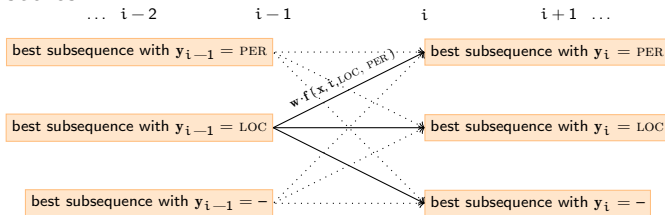
- Next questions:
  - How do we solve the **argmax** problem?
  - How do we learn  $\mathbf{w}$ ?

# Predicting with Factored Sequence Models

- Consider a fixed  $\mathbf{w}$ . Given  $\mathbf{x}_{1:n}$  find:

$$\operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$

- We can use the Viterbi algorithm, takes  $O(n|\mathcal{Y}|^2)$
- Intuition: output sequences that share bigrams will share scores



# Viterbi for Linear Factored Predictors

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}^n} \sum_{i=1}^n w \cdot f(x, i, y_{i-1}, y_i)$$

- **Definition:** score of optimal sequence for  $x_{1:i}$  ending with  $a \in \mathcal{Y}$

$$\delta_i(a) = \max_{y \in \mathcal{Y}^i : y_i = a} \sum_{j=1}^i w \cdot f(x, j, y_{j-1}, y_j)$$

- Use the following recursions, for all  $a \in \mathcal{Y}$ :

$$\delta_1(a) = w \cdot f(x, 1, y_0 = \text{NULL}, a)$$

$$\delta_i(a) = \max_{b \in \mathcal{Y}} (\delta_{i-1}(b) + w \cdot f(x, i, b, a))$$

- The optimal score for  $x$  is  $\max_{a \in \mathcal{Y}} \delta_n(a)$
- The optimal sequence  $\hat{y}$  can be recovered through *pointers*

# Linear Factored Sequence Prediction

Sequence  
Prediction

Approaches

Global Predictors

Log-linear  
Models for  
Sequence  
Prediction

$$\text{predict}(\mathbf{x}_{1:n}) = \underset{\mathbf{y} \in \mathcal{Y}^n}{\operatorname{argmax}} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y})$$

- Factored representation, e.g. based on bigrams
- Flexible, arbitrary features of full  $\mathbf{x}$  and the factors
- Efficient prediction using Viterbi
- Next topic: learning  $\mathbf{w}$ :
  - Maximum-Entropy Markov Models (local)
  - Conditional Random Fields (global)

# Outline

## 1 Sequence Prediction

- Examples
- Problem Formulation

## 2 Approaches

- Local Classifiers
- HMMs
- Global Predictors

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)
- Conditional Random Fields (CRF)

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

# Sequence Tagging with Log-Linear Models

- $\mathbf{x}$  are input sequences (e.g. sentences of words)
- $\mathbf{y}$  are output sequences (e.g. sequences of NE tags)
- **Goal:** given training data  
 $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$

learn a model  $\mathbf{x} \rightarrow \mathbf{y}$

- Log-linear models:

$$\operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}{Z(\mathbf{x}; \mathbf{w})}$$

# Sequence Tagging with Log-Linear Models

- $\mathbf{x}$  are input sequences (e.g. sentences of words)
- $\mathbf{y}$  are output sequences (e.g. sequences of NE tags)
- **Goal:** given training data
$$\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$$
learn a model  $\mathbf{x} \rightarrow \mathbf{y}$

- Log-linear models:

$$\operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}{Z(\mathbf{x}; \mathbf{w})}$$

- Exponentially many  $\mathbf{y}$ 's for a given input  $\mathbf{x}$

# Sequence Tagging with Log-Linear Models

- $\mathbf{x}$  are input sequences (e.g. sentences of words)
- $\mathbf{y}$  are output sequences (e.g. sequences of NE tags)
- **Goal:** given training data
$$\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$$
learn a model  $\mathbf{x} \rightarrow \mathbf{y}$

- Log-linear models:

$$\operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} P(\mathbf{y} | \mathbf{x}; \mathbf{w}) = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^n} \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}{Z(\mathbf{x}; \mathbf{w})}$$

- Exponentially many  $\mathbf{y}$ 's for a given input  $\mathbf{x}$ 
  - **Solution 1:** decompose  $P(\mathbf{y} | \mathbf{x})$  (MEMMs)
  - **Solution 2:** decompose  $\mathbf{f}(\mathbf{x}, \mathbf{y})$  (CRFs)



# Outline

## 1 Sequence Prediction

- Examples
- Problem Formulation

## 2 Approaches

- Local Classifiers
- HMMs
- Global Predictors

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)
- Conditional Random Fields (CRF)

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

Maximum Entropy  
Markov Models  
(MEMMs)

# Maximum Entropy Markov Models (MEMMs)

(McCallum, Freitag, Pereira 2000)

- Notation:  $\mathbf{x}_{1:n} = \mathbf{x}_1 \dots \mathbf{x}_n$
- Similarly to HMMs:

$$\begin{aligned} P(\mathbf{y}_{1:n} \mid \mathbf{x}_{1:n}) &= P(\mathbf{y}_1 \mid \mathbf{x}_{1:n}) \times P(\mathbf{y}_{2:n} \mid \mathbf{x}_{1:n}, \mathbf{y}_1) \\ &= P(\mathbf{y}_1 \mid \mathbf{x}_{1:n}) \times \prod_{i=2}^n P(\mathbf{y}_i \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:i-1}) \\ &= P(\mathbf{y}_1 \mid \mathbf{x}_{1:n}) \times \prod_{i=2}^n P(\mathbf{y}_i \mid \mathbf{x}_{1:n}, \mathbf{y}_{i-1}) \end{aligned}$$

- Assumption under MEMMs:

$$P(\mathbf{y}_i \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:i-1}) = P(\mathbf{y}_i \mid \mathbf{x}_{1:n}, \mathbf{y}_{i-1})$$

# Decoding with MEMMs

- Decompose tagging problem:

$$P(\mathbf{y}_{1:n} \mid \mathbf{x}_{1:n}) = P(\mathbf{y}_1 \mid \mathbf{x}_{1:n}) \times \prod_{i=2}^n P(\mathbf{y}_i \mid \mathbf{x}_{1:n}, i, \mathbf{y}_{i-1})$$

- Given  $\mathbf{w}$ , given  $\mathbf{x}$ , find:

$$\begin{aligned} \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^*} P(\mathbf{y} \mid \mathbf{x}; \mathbf{w}) &= \operatorname{argmax}_{\mathbf{y}} \prod_{i=1}^n P(\mathbf{y}_i \mid \mathbf{x}, \mathbf{y}_{i-1}) \\ &= \operatorname{argmax}_{\mathbf{y}} \frac{\prod_{i=1}^n \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i))}{\prod_{i=1}^n Z(\mathbf{x}, i; \mathbf{w})} \\ &= \operatorname{argmax}_{\mathbf{y}} \prod_{i=1}^n \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)) \\ &= \operatorname{argmax}_{\mathbf{y}} \sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) \end{aligned}$$

- We can use the Viterbi algorithm

# Parameter Estimation with MEMMs

- Learn *local* log-linear distributions (i.e. MaxEnt)

$$p(y_i | \mathbf{x}, i, y_{i-1}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, y_{i-1}, y_i))}{Z(\mathbf{x}, i, y_{i-1})}$$

where

- $\mathbf{x}$  is an input sequence
- $y_i$  and  $y_{i-1}$  are tags
- $\mathbf{f}(\mathbf{x}, i, y_{i-1}, y_i)$  is a feature vector of  $\mathbf{x}$ , the position to be tagged, the previous tag and the current tag

# Outline

## 1 Sequence Prediction

- Examples
- Problem Formulation

## 2 Approaches

- Local Classifiers
- HMMs
- Global Predictors

## 3 Log-linear Models for Sequence Prediction

- Maximum Entropy Markov Models (MEMMs)
- Conditional Random Fields (CRF)

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

Conditional Random  
Fields (CRF)

# Conditional Random Fields

(Lafferty, McCallum, Pereira 2001)

- Log-linear model of the conditional distribution:

$$P(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}{Z(\mathbf{x})}$$

where

- $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \dots \mathbf{x}_n \in \mathcal{X}^*$
- $\mathbf{y} = \mathbf{y}_1\mathbf{y}_2 \dots \mathbf{y}_n \in \mathcal{Y}^*$  and  $\mathcal{Y} = \{1, \dots, L\}$
- $\mathbf{f}(\mathbf{x}, \mathbf{y})$  is a feature vector of  $\mathbf{x}$  and  $\mathbf{y}$
- $\mathbf{w}$  are model parameters

- To predict the best sequence

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^*} P(\mathbf{y}|\mathbf{x})$$

# Conditional Random Fields

(Lafferty, McCallum, Pereira 2001)

- Log-linear model of the conditional distribution:

$$P(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}{Z(\mathbf{x})}$$

where

- $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \dots \mathbf{x}_n \in \mathcal{X}^*$
- $\mathbf{y} = \mathbf{y}_1\mathbf{y}_2 \dots \mathbf{y}_n \in \mathcal{Y}^*$  and  $\mathcal{Y} = \{1, \dots, L\}$
- $\mathbf{f}(\mathbf{x}, \mathbf{y})$  is a feature vector of  $\mathbf{x}$  and  $\mathbf{y}$
- $\mathbf{w}$  are model parameters

- To predict the best sequence

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^*} P(\mathbf{y}|\mathbf{x})$$

- Exponentially many  $\mathbf{y}$ 's for a given input  $\mathbf{x}$

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

Conditional Random  
Fields (CRF)

# Conditional Random Fields

(Lafferty, McCallum, Pereira 2001)

- Log-linear model of the conditional distribution:

$$P(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, \mathbf{y}))}{Z(\mathbf{x})}$$

where

- $\mathbf{x} = \mathbf{x}_1\mathbf{x}_2 \dots \mathbf{x}_n \in \mathcal{X}^*$
- $\mathbf{y} = \mathbf{y}_1\mathbf{y}_2 \dots \mathbf{y}_n \in \mathcal{Y}^*$  and  $\mathcal{Y} = \{1, \dots, L\}$
- $\mathbf{f}(\mathbf{x}, \mathbf{y})$  is a feature vector of  $\mathbf{x}$  and  $\mathbf{y}$
- $\mathbf{w}$  are model parameters

- To predict the best sequence

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^*} P(\mathbf{y}|\mathbf{x})$$

- Exponentially many  $\mathbf{y}$ 's for a given input  $\mathbf{x}$
- Choose  $\mathbf{f}(\mathbf{x}, \mathbf{y})$  so that  $\hat{\mathbf{y}}$  can be computed efficiently



# Conditional Random Fields (CRFs)

- The model form is:

$$P(\mathbf{y}|\mathbf{x}; \mathbf{w}) = \frac{\exp(\sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i))}{Z(\mathbf{x}, \mathbf{w})}$$

where

$$Z(\mathbf{x}, \mathbf{w}) = \sum_{\mathbf{z} \in \mathcal{Y}^*} \exp(\sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{z}_{i-1}, \mathbf{z}_i))$$

- Features  $\mathbf{f}(\dots)$  are given (they are problem-dependent)
- $\mathbf{w} \in \mathbb{R}^D$  are the parameters of the model
- CRFs are **log-linear models** on the feature functions

# Decoding with CRFs

- Given  $\mathbf{w}$ , given  $\mathbf{x}$ , find:

$$\begin{aligned}\operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^*} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) &= \operatorname{argmax}_{\mathbf{y}} \frac{\exp(\sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i))}{Z(\mathbf{x}; \mathbf{w})} \\ &= \operatorname{argmax}_{\mathbf{y}} \exp(\sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)) \\ &= \operatorname{argmax}_{\mathbf{y}} \sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)\end{aligned}$$

- We can use the Viterbi algorithm

# Parameter Estimation in CRFs

- How to estimate model parameters  $\mathbf{w}$  given a training set:

$$\left\{ (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)}) \right\}$$

- Let's define the conditional log-likelihood of the data:

$$L(\mathbf{w}) = \frac{1}{m} \sum_{k=1}^m \log P(\mathbf{y}^{(k)} | \mathbf{x}^{(k)}; \mathbf{w})$$

- $L(\mathbf{w})$  measures how well  $\mathbf{w}$  explains the data. A good value for  $\mathbf{w}$  will give a high value for  $P(\mathbf{y}^{(k)} | \mathbf{x}^{(k)}; \mathbf{w})$  for all  $k = 1 \dots m$ .
- We want  $\mathbf{w}$  that **maximizes**  $L(\mathbf{w})$

# Learning the Parameters of a CRF

- Recall previous lecture on log-linear / maximum-entropy models
- Find:

$$\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^D} L(\mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|^2$$

where

- The first term is the log-likelihood of the data
- The second term is a regularization term, it penalizes solutions with large norm
- $\lambda$  is a parameter to control the trade-off between fitting the data and model complexity

# Learning the Parameters of a CRF

- So we want to find:

$$\begin{aligned}\mathbf{w}^* &= \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d} L'(\mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w} \in \mathbb{R}^d} \left( \frac{1}{m} \sum_{k=1}^m \log P(\mathbf{y}^{(k)} | \mathbf{x}^{(k)}; \mathbf{w}) - \frac{\lambda}{2} \|\mathbf{w}\|^2 \right)\end{aligned}$$

- In general there is no analytical solution to this optimization
- ... but it is a **convex** function  $\Rightarrow$  We use iterative techniques, i.e. gradient-based optimization
- Very fast algorithms exist (e.g. LBFGS)

# Learning the Parameters of a CRF: Gradient step

- Initialize  $\mathbf{w} = \mathbf{0}$
- Repeat
  - Compute gradient  $\delta = (\delta_1, \dots, \delta_d)$ , where:

$$\delta_j = \frac{\partial L'(\mathbf{w})}{\partial \mathbf{w}_j} \quad \forall j = 1 \dots d$$

- Compute step size

$$\beta^* = \operatorname{argmax}_{\beta \in \mathbb{R}} L'(\mathbf{w} + \beta \delta)$$

- Move  $\mathbf{w}$  in the direction of the gradient

$$\mathbf{w} \leftarrow \mathbf{w} + \beta^* \delta$$

- until convergence ( $\|\delta\| < \epsilon$ )

# Computing the gradient

$$\begin{aligned}\frac{\partial L'(\mathbf{w})}{\partial \mathbf{w}_j} &= \frac{1}{m} \sum_{k=1}^m \mathbf{f}_j(\mathbf{x}^{(k)}, \mathbf{y}^{(k)}) \\ &\quad - \sum_{k=1}^m \sum_{\mathbf{y} \in \mathcal{Y}^{n_k}} P(\mathbf{y} | \mathbf{x}^{(k)}; \mathbf{w}) \mathbf{f}_j(\mathbf{x}^{(k)}, \mathbf{y}) \\ &\quad - \lambda \mathbf{w}_j\end{aligned}$$

where

$$\mathbf{f}_j(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \mathbf{f}_j(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$

- First term: observed mean feature value
- Second term: expected feature value under current  $\mathbf{w}$

# Computing the gradient

- The first term is easy to compute, by counting explicitly over all sequence elements:

$$\frac{1}{m} \sum_{k=1}^m \sum_{i=1}^{n_k} \mathbf{f}_j(\mathbf{x}^{(k)}, i, \mathbf{y}_{i-1}^{(k)}, \mathbf{y}_i^{(k)})$$

- The second term is more involved, because it sums over all sequences  $\mathbf{y} \in \mathcal{Y}^{n_k}$

$$\sum_{k=1}^m \sum_{\mathbf{y} \in \mathcal{Y}^{n_k}} P(\mathbf{y} | \mathbf{x}^{(k)}; \mathbf{w}) \sum_{i=1}^{n_k} \mathbf{f}_j(\mathbf{x}^{(k)}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$



# Computing the gradient

- For a given training example  $(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$ :

$$\sum_{\mathbf{y} \in \mathcal{Y}^{n_k}} P(\mathbf{y} | \mathbf{x}^{(k)}; \mathbf{w}) \sum_{i=1}^{n_k} \mathbf{f}_j(\mathbf{x}^{(k)}, i, \mathbf{y}_{i-1}, \mathbf{y}_i) =$$
$$\sum_{i=1}^{n_k} \sum_{a, b \in \mathcal{Y}} \mu_i^k(a, b) \mathbf{f}_j(\mathbf{x}^{(k)}, i, a, b)$$

where

$$\mu_i^k(a, b) = \sum_{\mathbf{y} \in \mathcal{Y}^{n_k} : \mathbf{y}_{i-1}=a, \mathbf{y}_i=b} P(\mathbf{y} | \mathbf{x}^{(k)}; \mathbf{w})$$

- The quantities  $\mu_i^k$  can be computed efficiently in  $O(n|\mathcal{Y}|^2)$  using the forward-backward algorithm

# Forward-Backward for CRFs

- Assume fixed  $\mathbf{x}$ . Calculate in  $O(n|Y|^2)$

$$\mu_i(a, b) = \sum_{\mathbf{y} \in Y^n: \mathbf{y}_{i-1}=a, \mathbf{y}_i=b} P(\mathbf{y}|\mathbf{x}; \mathbf{w}) \quad , \quad 1 \leq i \leq n; \quad a, b \in Y$$

- Define (forward and backward quantities):

$$\alpha_i(a) = \sum_{\mathbf{y} \in Y^i: \mathbf{y}_i=a} \exp(\sum_{j=1}^i \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, j, \mathbf{y}_{j-1}, \mathbf{y}_j))$$

$$\beta_i(b) = \sum_{\mathbf{y} \in Y^{(n-i+1)}: \mathbf{y}_1=b} \exp(\sum_{j=2}^{n-i+1} \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i+j-1, \mathbf{y}_{j-1}, \mathbf{y}_j))$$

- Compute recursively  $\alpha_i(a)$  and  $\beta_i(b)$  (similar to Viterbi)
- $Z = \sum_a \alpha_n(a)$
- $\mu_i(a, b) = \alpha_{i-1}(a) \cdot \exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, a, b)) \cdot \beta_i(b) / Z$

# Compute the probability of a label sequence

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

Conditional Random  
Fields (CRF)

$$P(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \frac{\exp(\sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i))}{Z(\mathbf{x}; \mathbf{w})}$$

where

$$Z(\mathbf{x}; \mathbf{w}) = \sum_{\mathbf{z} \in \mathcal{Y}^n} \exp(\sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{z}_{i-1}, \mathbf{z}_i))$$

- Compute  $Z(\mathbf{x}; \mathbf{w})$  efficiently, using the forward algorithm

# CRFs: summary so far

- Log-linear models for sequence prediction,  $P(\mathbf{y}|\mathbf{x}; \mathbf{w})$
- Computations factorize on label bigrams
- Model form:

$$\operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^*} \sum_i \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$

- Decoding: uses Viterbi (from HMMs)
- Parameter estimation:
  - Gradient-based methods, in practice L-BFGS
  - Computation of gradient uses forward-backward (from HMMs)

# CRFs: summary so far

- Log-linear models for sequence prediction,  $P(\mathbf{y}|\mathbf{x}; \mathbf{w})$
- Computations factorize on label bigrams
- Model form:

$$\operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}^*} \sum_i \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i)$$

- Decoding: uses Viterbi (from HMMs)
- Parameter estimation:
  - Gradient-based methods, in practice L-BFGS
  - Computation of gradient uses forward-backward (from HMMs)
- **Next Questions:** MEMMs or CRFs? HMMs or CRFs?

# MEMMs and CRFs

$$\text{MEMMs: } P(\mathbf{y} \mid \mathbf{x}) = \prod_{i=1}^n \frac{\exp(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i))}{Z(\mathbf{x}, i, \mathbf{y}_{i-1}; \mathbf{w})}$$

$$\text{CRFs: } P(\mathbf{y} \mid \mathbf{x}) = \frac{\exp(\sum_{i=1}^n \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, \mathbf{y}_{i-1}, \mathbf{y}_i))}{Z(\mathbf{x})}$$

- MEMMs locally normalized; CRFs globally normalized
- MEMM assume that
$$P(\mathbf{y}_i \mid \mathbf{x}_{1:n}, \mathbf{y}_{1:i-1}) = P(\mathbf{y}_i \mid \mathbf{x}_{1:n}, \mathbf{y}_{i-1})$$
- Both exploit the same factorization, i.e. same features
- Same computations to compute  $\arg\max_{\mathbf{y}} P(\mathbf{y} \mid \mathbf{x})$
- MEMMs are cheaper to train
- CRFs are easier to extend to other structures (e.g. parsing trees)

# HMMs for sequence prediction

- $\mathbf{x}$  are the observations,  $\mathbf{y}$  are the (un)hidden states
- HMMs model the joint distribution  $P(\mathbf{x}, \mathbf{y})$
- Parameters: (assume  $\mathcal{X} = \{1, \dots, k\}$  and  $\mathcal{Y} = \{1, \dots, l\}$ )
  - $\pi \in \mathbb{R}^l$ ,  $\pi_a = P(\mathbf{y}_1 = a)$
  - $T \in \mathbb{R}^{l \times l}$ ,  $T_{a,b} = P(\mathbf{y}_i = b | \mathbf{y}_{i-1} = a)$
  - $O \in \mathbb{R}^{l \times k}$ ,  $O_{a,c} = P(\mathbf{x}_i = c | \mathbf{y}_i = a)$
- Model form

$$P(\mathbf{x}, \mathbf{y}) = \pi_{y_1} O_{y_1, x_1} \prod_{i=2}^n T_{y_{i-1}, y_i} O_{y_i, x_i}$$

- Parameter Estimation: maximum likelihood by counting events and normalizing

# HMMs and CRFs

- In CRFs:  $\hat{y} = \text{amax}_y \sum_i \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, y_{i-1}, y_i)$

- In HMMs:

$$\begin{aligned}\hat{y} &= \text{amax}_y \pi_{y_1} O_{y_1, x_1} \prod_{i=2}^n T_{y_{i-1}, y_i} O_{y_i, x_i} \\ &= \text{amax}_y \log(\pi_{y_1} O_{y_1, x_1}) + \sum_{i=2}^n \log(T_{y_{i-1}, y_i} O_{y_i, x_i})\end{aligned}$$

- An HMM can be ported into a CRF by setting:

$\mathbf{f}_j(\mathbf{x}, i, y, y')$	$\mathbf{w}_j$



# HMMs and CRFs

- In CRFs:  $\hat{y} = \text{amax}_y \sum_i w \cdot f(x, i, y_{i-1}, y_i)$

- In HMMs:

$$\begin{aligned}\hat{y} &= \text{amax}_y \pi_{y_1} O_{y_1, x_1} \prod_{i=2}^n T_{y_{i-1}, y_i} O_{y_i, x_i} \\ &= \text{amax}_y \log(\pi_{y_1} O_{y_1, x_1}) + \sum_{i=2}^n \log(T_{y_{i-1}, y_i} O_{y_i, x_i})\end{aligned}$$

- An HMM can be ported into a CRF by setting:

$f_j(x, i, y, y')$	$w_j$
$i = 1 \ \& \ y' = a$	$\log(\pi_a)$

# HMMs and CRFs

- In CRFs:  $\hat{y} = \text{amax}_y \sum_i \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, y_{i-1}, y_i)$

- In HMMs:

$$\begin{aligned}\hat{y} &= \text{amax}_y \pi_{y_1} O_{y_1, x_1} \prod_{i=2}^n T_{y_{i-1}, y_i} O_{y_i, x_i} \\ &= \text{amax}_y \log(\pi_{y_1} O_{y_1, x_1}) + \sum_{i=2}^n \log(T_{y_{i-1}, y_i} O_{y_i, x_i})\end{aligned}$$

- An HMM can be ported into a CRF by setting:

$\mathbf{f}_j(\mathbf{x}, i, y, y')$	$\mathbf{w}_j$
$i = 1 \ \& \ y' = a$	$\log(\pi_a)$
$i > 1 \ \& \ y = a \ \& \ y' = b$	$\log(T_{a,b})$

# HMMs and CRFs

- In CRFs:  $\hat{y} = \text{amax}_y \sum_i \mathbf{w} \cdot \mathbf{f}(\mathbf{x}, i, y_{i-1}, y_i)$

- In HMMs:

$$\begin{aligned}\hat{y} &= \text{amax}_y \pi_{y_1} O_{y_1, x_1} \prod_{i=2}^n T_{y_{i-1}, y_i} O_{y_i, x_i} \\ &= \text{amax}_y \log(\pi_{y_1} O_{y_1, x_1}) + \sum_{i=2}^n \log(T_{y_{i-1}, y_i} O_{y_i, x_i})\end{aligned}$$

- An HMM can be ported into a CRF by setting:

$\mathbf{f}_j(\mathbf{x}, i, y, y')$	$\mathbf{w}_j$
$i = 1 \ \& \ y' = a$	$\log(\pi_a)$
$i > 1 \ \& \ y = a \ \& \ y' = b$	$\log(T_{a,b})$
$y' = a \ \& \ x_i = c$	$\log(O_{a,b})$

- Hence, HMM parameters  $\subset$  CRF parameters

# HMMs and CRFs: main differences

Sequence  
Prediction

Approaches

Log-linear  
Models for  
Sequence  
Prediction

Conditional Random  
Fields (CRF)

- Representation:
  - HMM “features” are tied to the generative process.
  - CRF features are **very** flexible. They can look at the whole input  $x$  paired with a label bigram  $(y, y')$ .
  - In practice, for prediction tasks, “good” discriminative features can improve accuracy **a lot**.
- Parameter estimation:
  - HMMs focus on explaining the data, both  $x$  and  $y$ .
  - CRFs focus on the mapping from  $x$  to  $y$ .
  - A priori, it is hard to say which paradigm is better.
  - Same dilemma as Naive Bayes vs. Maximum Entropy.