

Lecture 5: Linear Sequential Models

André Martins



Deep Structured Learning Course, Fall 2020

Announcements

Deadline for Homework 1 is today!

Homework 2 is out. Deadline is November 18. Start early!!!

Next week: guest lecture on Reinforcement Learning by Francisco Melo!

That lecture will be 30 min earlier (it will start at 13:30)!

Today's Roadmap

Today we're starting to talk about **structure**, more specifically **sequences**:

- Generative sequence models: (hidden) Markov models
- Dynamic programming: the Viterbi and Forward-Backward algorithms
- Viterbi decoding and minimum risk decoding
- Unsupervised learning with the Baum-Welch (EM) algorithm

① Structured Prediction

② Sequence Models

Markov Models

Hidden Markov Models

Conditional Random Fields

Structured Perceptron

Structured Support Vector Machines

Structured Prediction

Structured prediction: a machine learning framework for predicting structured, constrained, and interdependent outputs

NLP deals with *structured* and *ambiguous* textual data:

- machine translation, speech recognition, parsing, ...

Computer vision, computational biology, etc. deal with structured outputs too:

- image segmentation, image parsing, ...
- protein folding, protein design, ...
- time series prediction, ...

Sequence Labeling

Some of these tasks are **sequential** in nature...

Example: POS Tagging

Map **sentences** to sequences of **part-of-speech tags**.

Time	flies	like	an	arrow	.
noun	verb	prep	det	noun	.

- Need to predict a morphological tag for each word of the sentence
- High correlation between adjacent words!

(Ratnaparkhi, 1999; Brants, 2000; Toutanova et al., 2003)

Example: Named Entity Recognition

From **sentences** extract **named entities**.

Louis	Elsevier	was	born	in	Leuven	.
B-PER	I-PER	O	O	O	B-LOC	.

- Identify word segments that refer to entities (person, organization, location)
- Typically done with sequence models and B-I-O tagging

(Zhang and Johnson, 2003; Ratnov and Roth, 2009)

Notation

- Input set \mathcal{X} (e.g., sentences)

Notation

- Input set \mathcal{X} (e.g., sentences)
- Output set \mathcal{Y} (e.g. POS sequences)—**large and structured!**
 - How many possible outputs for a sentence with length L , assuming K possible tags?

Notation

- Input set \mathcal{X} (e.g., sentences)
- Output set \mathcal{Y} (e.g. POS sequences)—**large and structured!**
 - How many possible outputs for a sentence with length L , assuming K possible tags? $O(K^L)$

Notation

- Input set \mathcal{X} (e.g., sentences)
- Output set \mathcal{Y} (e.g. POS sequences)—**large and structured!**
 - How many possible outputs for a sentence with length L , assuming K possible tags? $O(K^L)$
- Weights \mathbf{W} (with lots of parameter sharing), features $\phi(x)$
- Linear compatibility function $F_{\mathbf{W}}(x, y) = w_y \cdot \phi(x)$

Notation

- Input set \mathcal{X} (e.g., sentences)
- Output set \mathcal{Y} (e.g. POS sequences)—**large and structured!**
 - How many possible outputs for a sentence with length L , assuming K possible tags? $O(K^L)$
- Weights \mathbf{W} (with lots of parameter sharing), features $\phi(\mathbf{x})$
- Linear compatibility function $F_{\mathbf{W}}(\mathbf{x}, \mathbf{y}) = \mathbf{w}_{\mathbf{y}} \cdot \phi(\mathbf{x})$
- **Training problem:** learn \mathbf{W} from data $\mathcal{T} = \{\langle \mathbf{x}_t, \mathbf{y}_t \rangle\}_{t=1}^{|\mathcal{T}|}$

Notation

- Input set \mathcal{X} (e.g., sentences)
- Output set \mathcal{Y} (e.g. POS sequences)—**large and structured!**
 - How many possible outputs for a sentence with length L , assuming K possible tags? $O(K^L)$
- Weights \mathbf{W} (with lots of parameter sharing), features $\phi(x)$
- Linear compatibility function $F_{\mathbf{W}}(x, y) = w_y \cdot \phi(x)$
- **Training problem:** learn \mathbf{W} from data $\mathcal{T} = \{\langle x_t, y_t \rangle\}_{t=1}^{|\mathcal{T}|}$
- **Decoding problem:**

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} F_{\mathbf{W}}(x, y)$$

This is what makes structured prediction special!

The Decoding Problem

Also called “the inference problem”:

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} F_W(x, y)$$

- In multi-class classification, we just enumerate the scores and pick the argmax
- **But in structure prediction, \mathcal{Y} is too large—it is intractable to enumerate!!**

Key assumption: F_W decomposes into (overlapping) *parts*

- For example: for sequences, scores may factor as a sum over label n -grams

Roadmap: Models for Structured Prediction

Binary/Multi-class	Structured Prediction
--------------------	-----------------------

Naive Bayes	?
-------------	---

Logistic Regression	?
---------------------	---

Perceptron	?
------------	---

SVMs	?
------	---

① Structured Prediction

② Sequence Models

Markov Models

Hidden Markov Models

Conditional Random Fields

Structured Perceptron

Structured Support Vector Machines

① Structured Prediction

② Sequence Models

Markov Models

Hidden Markov Models

Conditional Random Fields

Structured Perceptron

Structured Support Vector Machines

Key Problem

How to define a probability distribution over a string (a sequence of discrete symbols)?

- If the length of the string is unbounded, there are infinitely many strings
- How to make sure the distribution normalizes to 1?
- So far, we only talked about distributions on finite random variables...

Probability Distributions Over Strings

- Alphabet Σ (e.g. words in the vocabulary)
- Set of strings: $\Sigma^* := \{\epsilon\} \cup \Sigma \cup \Sigma^2 \dots$
- Σ^* is a countably infinite set
- We want to define a probability distribution over Σ^*
- This distribution must normalize properly, $\sum_{x \in \Sigma^*} \mathbb{P}(x) = 1$

Trivial Choices

- Probability 0 to sequences greater than L , uniform otherwise
- Given a sample of N sequences, assign a uniform distribution of $1/N$ to each observed sequence, 0 if the sentence is not in the sample

What if we want *every* sequence to have some probability?

Lower Extreme: Bag-of-Words Model

$$\mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L, \text{STOP}) = \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{x}_i)$$

- Also called “unigram” model
- Assumes every word is generated independently of other words
- Probability of a string is insensitive to word order
- How many parameters?
- **Not structured prediction!**

Upper Extreme: Full History Model

$$\mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L, \text{STOP}) = \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{x}_i | \mathbf{x}_0, \dots, \mathbf{x}_{i-1})$$

- Assumes the generation of each word depends on the entire history (*all* the previous words)
- Huge expressive power!
- But: too many parameters to estimate! (How many?)
- Cannot generalize well

In-Between: First Order Markov Models

$$\mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L, \text{STOP}) = \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1})$$

- Each word only depends on the previous word
- Which parameters need to be estimated?
- Transition probabilities $\mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1}) \dots$
- ... including initial and final probabilities $\mathbb{P}(\mathbf{x}_1 | \text{START})$ and $\mathbb{P}(\text{STOP} | \mathbf{x}_L)$
- Total number of parameters: $O(|\Sigma|^2)$

In-Between: K th Order Markov Models

$$\mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_L, \text{STOP}) = \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1}, \dots, \mathbf{x}_{i-K})$$

- Each word depends on the K th previous words
- Transition probabilities $\mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1}, \dots, \mathbf{x}_{i-K})$
- Total number of parameters: $O(|\Sigma|^{K+1})$
- **Widely used in language modeling**
- **Applications: machine translation, speech recognition, OCR...**

Markov Models: How to Estimate the Parameters

How to estimate these transition probabilities $\mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1}, \dots, \mathbf{x}_{i-K})$ from data?

- Maximum likelihood estimation
- Boils down to **counting and normalizing**
- What if we never observe a particular K -gram at training time? What happens at test time?

Problem: Data Sparsity

- As K grows, we suffer more from insufficient data: there are many K -grams with zero counts, therefore the maximum likelihood estimate will be $\mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1}, \dots, \mathbf{x}_{i-K}) = 0$
- This is undesirable!
- Example: a language model which gives probability 0 to unseen words
- Just because an event has never been observed in training data does not mean it cannot occur in test data
- So, if the count is zero, what should $\mathbb{P}(\mathbf{x}_i | \mathbf{x}_{i-1}, \dots, \mathbf{x}_{i-K}) = 0$ be?
- We need **smoothing**!

Smoothing

Many smoothing techniques have been proposed:

- Additive smoothing (e.g. add-one, a.k.a. Laplace smoothing)
- Good-Turing estimate
- Jelinek-Mercer smoothing (interpolation)
- Katz smoothing (backoff)
- Witten-Bell smoothing
- Absolute discounting
- Kneser-Ney smoothing
- Modified Kneser-Ney smoothing (considered the best)

This is a black art...

Markov models are extremely simple.

How well do they perform?

Example

- Unigram model estimated on 2.8M words of political blog text.

*this trying our putting and funny and among it herring it obama
but certainly foreign my c on byron again but from i i so and i
chuck yeah the as but but republicans if this stay oh so or it mccain
bush npr this with what and they right i while because obama*

(Adapted from Noah Smith's slides.)

Example

- Bigram model estimated on 2.8M words of political blog text.

the lack of the senator mccain hadn t keep this story backwards while showering praise of the kind of gop weakness it was mistaken for american economist anywhere in the white house press hounded the absence of those he s as a wide variety of this election day after the candidate b richardson was polled ri in hempstead moderated by the convention that he had zero wall street journal argues sounds like you may be the primary but even close the bill told c e to take the obama on the public schools and romney fred flinstone s see how a lick skillet road it s little sexist remarks

(Adapted from Noah Smith's slides.)

Example

- Trigram model estimated on 2.8M words of political blog text.

as i can pin them all none of them want to bet that any of the might be conservatism unleashed into the privacy rule book and when told about what paul fans organized another massive fundraising initiative yesterday and i don t know what the rams supposedly want ooh but she did but still victory dinner alone among republicans there are probably best not all of the fundamentalist community asked for an independent maverick now for crystallizing in one especially embarrassing

(Adapted from Noah Smith's slides.)

Example

- 5-gram model estimated on 2.8M words of political blog text.

he realizes fully how shallow and insincere conservative behavior has been he realizes that there is little way to change the situation this recent arianna huffington item about mccain issuing heartfelt denials of things that were actually true or for that matter about the shia sunni split and which side iran was on would get confused about this any more than someone with any knowledge of us politics would get confused about whether neo confederates were likely to be supporting the socialist workers party at the end of the world and i m not especially discouraged now that newsweek shows obama leading by three now

(Adapted from Noah Smith's slides.)

① Structured Prediction

② Sequence Models

Markov Models

Hidden Markov Models

Conditional Random Fields

Structured Perceptron

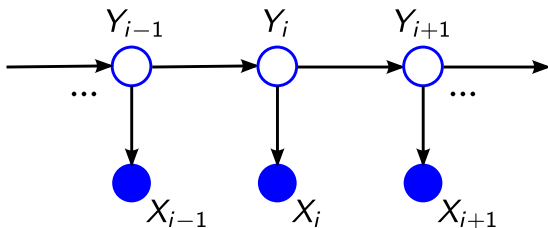
Structured Support Vector Machines

Hidden Markov Models

- Assume each word x_i is associated with a **hidden state** y_i
- Two alphabets: word alphabet Σ **and state alphabet** Λ
- States are generated according to a (first-order) Markov model
- Word emissions are conditioned on each state
- **Applications: POS tagging, named entity recognition, shallow parsing, ...**

Hidden Markov Models

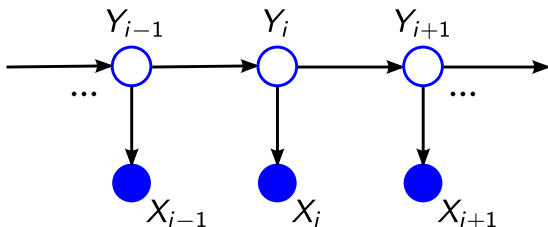
$$\mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{y}_1, \mathbf{x}_2, \mathbf{y}_2, \dots, \mathbf{x}_L, \mathbf{y}_L, \text{STOP}) = \prod_{i=1}^{L+1} \underbrace{\mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1})}_{\text{transitions}} \times \prod_{i=1}^L \underbrace{\mathbb{P}(\mathbf{x}_i | \mathbf{y}_i)}_{\text{emissions}}$$



What are the parameters?

Hidden Markov Models

$$\mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{y}_1, \mathbf{x}_2, \mathbf{y}_2, \dots, \mathbf{x}_L, \mathbf{y}_L, \text{STOP}) = \prod_{i=1}^{L+1} \underbrace{\mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1})}_{\text{transitions}} \times \prod_{i=1}^L \underbrace{\mathbb{P}(\mathbf{x}_i | \mathbf{y}_i)}_{\text{emissions}}$$

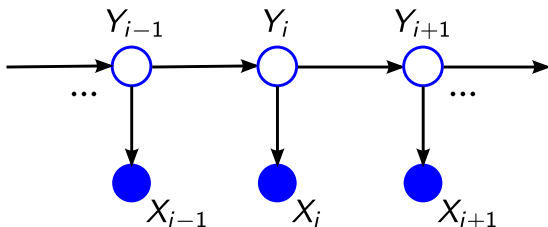


What are the parameters?

- **Transition probabilities** $\mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1})$
- **Emission probabilities** $\mathbb{P}(\mathbf{x}_i | \mathbf{y}_i)$

Hidden Markov Models

$$\mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{y}_1, \mathbf{x}_2, \mathbf{y}_2, \dots, \mathbf{x}_L, \mathbf{y}_L, \text{STOP}) = \prod_{i=1}^{L+1} \underbrace{\mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1})}_{\text{transitions}} \times \prod_{i=1}^L \underbrace{\mathbb{P}(\mathbf{x}_i | \mathbf{y}_i)}_{\text{emissions}}$$



What are the parameters?

- **Transition probabilities** $\mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1})$
- **Emission probabilities** $\mathbb{P}(\mathbf{x}_i | \mathbf{y}_i)$
- Total number of parameters: $O(|\Sigma||\Lambda| + |\Lambda|^2)$

Sampling from an HMM

Two equivalent ways of sampling:

- 1 Sample the sequence of hidden states using a Markov model $\mathbb{P}(\mathbf{y}_i \mid \mathbf{y}_{i-1})$; then for each state \mathbf{y}_i , sample a word from $\mathbb{P}(\mathbf{x}_i \mid \mathbf{y}_i)$
- 2 Alternate between sampling the next state and the current word from the current state

Three Problems in HMMs

Let $\mathbf{x} = \text{START}, \mathbf{x}_1, \dots, \mathbf{x}_L, \text{STOP}$ denote a word sequence.

- 1 Given \mathbf{x} , compute the *most likely sequence of states*
 $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x})$
- 2 Given \mathbf{x} , compute the *likelihood* $\mathbb{P}(\mathbf{x})$ and the *posterior state probabilities* $\mathbb{P}(\mathbf{y}_i | \mathbf{x})$
- 3 Given training sequences of word and states (or just words), *estimate the emission and transition probabilities.*

Three Problems in HMMs

Let $\mathbf{x} = \text{START}, \mathbf{x}_1, \dots, \mathbf{x}_L, \text{STOP}$ denote a word sequence.

- 1 Given \mathbf{x} , compute the *most likely sequence of states*
 $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x})$
- 2 Given \mathbf{x} , compute the *likelihood* $\mathbb{P}(\mathbf{x})$ and the *posterior state probabilities* $\mathbb{P}(\mathbf{y}_i | \mathbf{x})$
- 3 Given training sequences of word and states (or just words), *estimate the emission and transition probabilities.*

We'll see that:

- Problems 1 and 2 can be solved with **dynamic programming**
- Problem 3 can be solved by **counting and normalizing**

Three Problems in HMMs

Let $\mathbf{x} = \text{START}, \mathbf{x}_1, \dots, \mathbf{x}_L, \text{STOP}$ denote a word sequence.

- 1 Given \mathbf{x} , compute the *most likely sequence of states*
 $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x})$
- 2 Given \mathbf{x} , compute the *likelihood* $\mathbb{P}(\mathbf{x})$ and the *posterior state probabilities* $\mathbb{P}(\mathbf{y}_i | \mathbf{x})$
- 3 Given training sequences of word and states (or just words), *estimate the emission and transition probabilities.*

We'll see that:

- Problems 1 and 2 can be solved with **dynamic programming**
- Problem 3 can be solved by **counting and normalizing**

Dynamic Programming

A central concept applicable to many structured prediction problems

Key idea: use a table data structure to store partial quantities that will be reused many times

Optimal substructure: best solution to a problem relies on best solutions to its (similar-looking) subproblems

Overlapping subproblems: reuse a small number of quantities many times

Examples:

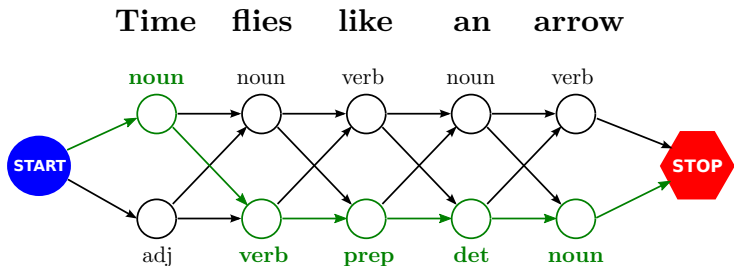
- Viterbi
- Levenshtein distance (a.k.a. edit distance)
- Dijkstra's shortest path algorithm
- Bellman equation in optimal control theory

Problem 1: Most Likely State Sequence

$$\begin{aligned}\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_L &= \arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x}) \\ &= \arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_L, \mathbf{y}_L, \text{STOP}) \\ &= \arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \prod_{i=1}^L \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i)\end{aligned}$$

- **Combinatorial problem: need to search over $|\Lambda|^L$ possibilities.**

Viterbi Trellis



- Each word's state depends on the word and on the nearby labels
- Key fact: given adjacent labels, the others do not matter

The Viterbi algorithm

- Key idea: recurrence
- Traverse the sequence left-to-right and compute the best way of reaching the i th word given the possible decisions at the $(i - 1)$ th word
- Fill a table along the way, then backtrack to recover the best sequence of states

$$V(i, \mathbf{y}_i) = \max_{\mathbf{y}_{i-1} \in \Lambda} \left(\mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \times V(i-1, \mathbf{y}_{i-1}) \right)$$

$$\psi(i, \mathbf{y}_i) = \arg \max_{\mathbf{y}_{i-1} \in \Lambda} \left(\mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \times V(i-1, \mathbf{y}_{i-1}) \right)$$

The Viterbi algorithm

input: sequence x_1, \dots, x_L , emission/transition probabilities

Forward pass: incrementally fill the table

$$V(1, \mathbf{y}_1) = \mathbb{P}(\mathbf{y}_1 | \text{START}) \times \mathbb{P}(x_1 | \mathbf{y}_1) \quad \forall \mathbf{y}_1 \in \Lambda$$

for $i = 2$ **to** L **do**

for $\mathbf{y}_i \in \Lambda$ **do**

$$V(i, \mathbf{y}_i) = \max_{\mathbf{y}_{i-1}} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(x_i | \mathbf{y}_i) \times V(i-1, \mathbf{y}_{i-1})$$

$$\psi(i, \mathbf{y}_i) = \arg \max_{\mathbf{y}_{i-1}} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(x_i | \mathbf{y}_i) \times V(i-1, \mathbf{y}_{i-1})$$

Backward pass: follow backpointers

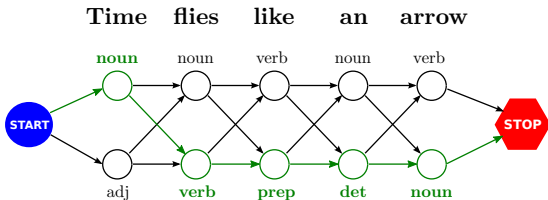
$$\hat{\mathbf{y}}_L = \arg \max_{\mathbf{y}_L} \mathbb{P}(\text{STOP} | \mathbf{y}_L) \times V(L, \mathbf{y}_L)$$

for $i = L - 1$ **to** 1 **do**

$$\hat{\mathbf{y}}_i = \psi(i + 1, \hat{\mathbf{y}}_{i+1})$$

output: the best state sequence $\hat{\mathbf{y}}_1, \dots, \hat{\mathbf{y}}_L$.

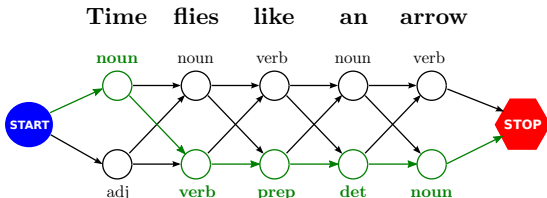
Example – $i = 1$



$$\begin{aligned}V(1, \text{noun}) &= \mathbb{P}(\text{noun}|\text{START}) \times \mathbb{P}(\text{"Time"}|\text{noun}) \\ &= 0.4 \times 0.1 = \mathbf{0.04}\end{aligned}$$

$$\begin{aligned}V(1, \text{adj}) &= \mathbb{P}(\text{adj}|\text{START}) \times \mathbb{P}(\text{"Time"}|\text{adj}) \\ &= 0.2 \times 0.05 = \mathbf{0.01}\end{aligned}$$

Example – $i = 2$

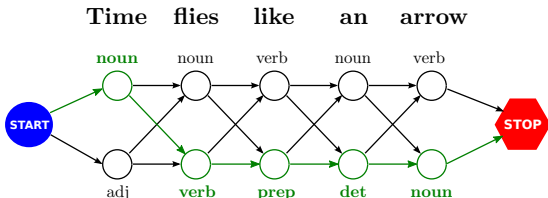


$$\begin{aligned} V(2, \text{noun}) &= \max\{\mathbb{P}(\text{noun}|\text{noun}) \times \mathbb{P}(\text{"flies"}|\text{noun}) \times V(1, \text{noun}), \\ &\quad \mathbb{P}(\text{noun}|\text{adj}) \times \mathbb{P}(\text{"flies"}|\text{noun}) \times V(1, \text{adj})\} \\ &= \max\{0.1 \times 0.1 \times 0.04, 0.8 \times 0.1 \times 0.01\} = 0.0008 \end{aligned}$$

$$\begin{aligned} \psi(2, \text{noun}) &= \text{adj} \\ V(2, \text{verb}) &= \max\{\mathbb{P}(\text{verb}|\text{noun}) \times \mathbb{P}(\text{"flies"}|\text{verb}) \times V(1, \text{noun}), \\ &\quad \mathbb{P}(\text{verb}|\text{adj}) \times \mathbb{P}(\text{"flies"}|\text{verb}) \times V(1, \text{adj})\} \\ &= \max\{0.4 \times 0.2 \times 0.04, 0.01 \times 0.2 \times 0.01\} = 0.0032 \end{aligned}$$

$$\psi(2, \text{verb}) = \text{noun}$$

Example – $i = 3$



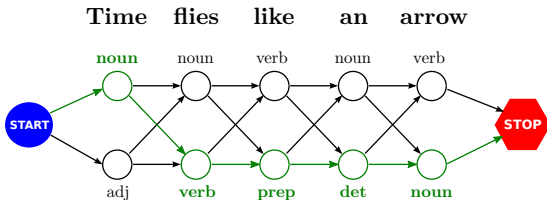
$$\begin{aligned} V(3, \text{verb}) &= \max\{\mathbb{P}(\text{verb}|\text{noun}) \times \mathbb{P}(\text{"like"}|\text{verb}) \times V(2, \text{noun}), \\ &\quad \mathbb{P}(\text{verb}|\text{verb}) \times \mathbb{P}(\text{"like"}|\text{verb}) \times V(2, \text{verb})\} \\ &= \max\{0.6 \times 0.1 \times 0.0008, 0.01 \times 0.1 \times 0.0032\} = 4.8 \times 10^{-5} \end{aligned}$$

$$\psi(3, \text{verb}) = \text{noun}$$

$$\begin{aligned} V(3, \text{prep}) &= \max\{\mathbb{P}(\text{prep}|\text{noun}) \times \mathbb{P}(\text{"like"}|\text{prep}) \times V(2, \text{noun}), \\ &\quad \mathbb{P}(\text{prep}|\text{verb}) \times \mathbb{P}(\text{"like"}|\text{prep}) \times V(2, \text{verb})\} \\ &= \max\{0.2 \times 0.5 \times 0.0008, 0.2 \times 0.5 \times 0.0032\} = 0.00032 \end{aligned}$$

$$\psi(3, \text{prep}) = \text{verb}$$

Example – $i = 4$



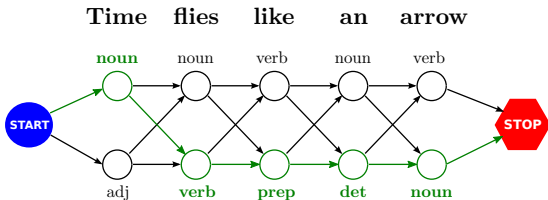
$$\begin{aligned}
 V(4, \text{noun}) &= \max\{\mathbb{P}(\text{noun}|\text{verb}) \times \mathbb{P}(\text{"an"}|\text{noun}) \times V(3, \text{verb}), \\
 &\quad \mathbb{P}(\text{noun}|\text{prep}) \times \mathbb{P}(\text{"an"}|\text{noun}) \times V(3, \text{prep})\} \\
 &= \max\{0.4 \times 0.001 \times 4.8 \times 10^{-5}, 0.2 \times 0.001 \times 0.00032\} = 6.4 \times 10^{-8}
 \end{aligned}$$

$$\psi(4, \text{noun}) = \text{prep}$$

$$\begin{aligned}
 V(4, \text{det}) &= \max\{\mathbb{P}(\text{det}|\text{verb}) \times \mathbb{P}(\text{"an"}|\text{det}) \times V(3, \text{verb}), \\
 &\quad \mathbb{P}(\text{det}|\text{prep}) \times \mathbb{P}(\text{"an"}|\text{det}) \times V(3, \text{prep})\} \\
 &= \max\{0.3 \times 0.5 \times 4.8 \times 10^{-5}, 0.5 \times 0.5 \times 0.00032\} = 8 \times 10^{-5}
 \end{aligned}$$

$$\psi(4, \text{det}) = \text{prep}$$

Example – $i = 5$



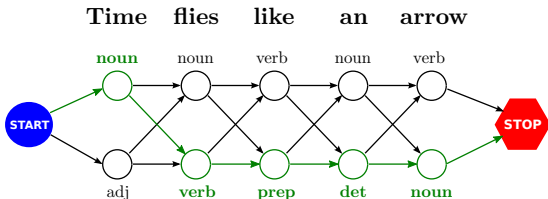
$$\begin{aligned}
 V(5, \text{verb}) &= \max\{\mathbb{P}(\text{verb}|\text{noun}) \times \mathbb{P}(\text{"arrow"}|\text{verb}) \times V(4, \text{noun}), \\
 &\quad \mathbb{P}(\text{verb}|\text{det}) \times \mathbb{P}(\text{"arrow"}|\text{verb}) \times V(4, \text{det})\} \\
 &= \max\{0.3 \times 0.01 \times 6.4 \times 10^{-8}, 0.01 \times 0.01 \times 8 \times 10^{-5}\} = 8 \times 10^{-9}
 \end{aligned}$$

$$\psi(5, \text{verb}) = \text{det}$$

$$\begin{aligned}
 V(5, \text{noun}) &= \max\{\mathbb{P}(\text{noun}|\text{noun}) \times \mathbb{P}(\text{"arrow"}|\text{noun}) \times V(4, \text{noun}), \\
 &\quad \mathbb{P}(\text{noun}|\text{det}) \times \mathbb{P}(\text{"arrow"}|\text{noun}) \times V(4, \text{det})\} \\
 &= \max\{0.1 \times 0.1 \times 6.4 \times 10^{-8}, 0.9 \times 0.1 \times 8 \times 10^{-5}\} = 7.2 \times 10^{-6}
 \end{aligned}$$

$$\psi(5, \text{noun}) = \text{det}$$

Example – Backward Pass



$$\begin{aligned} \text{Probability} &= \max\{\mathbb{P}(\text{STOP}|\text{verb}) \times V(5, \text{verb}), \\ &\quad \mathbb{P}(\text{STOP}|\text{noun}) \times V(5, \text{noun})\} \\ &= \max\{0.2 \times 8 \times 10^{-9}, 0.5 \times 7.2 \times 10^{-6}\} = 3.6 \times 10^{-6} \end{aligned}$$

$$\hat{y}_5 = \text{noun}$$

$$\hat{y}_4 = \psi(5, \text{noun}) = \text{det}$$

$$\hat{y}_3 = \psi(4, \text{det}) = \text{prep}$$

$$\hat{y}_2 = \psi(3, \text{prep}) = \text{verb}$$

$$\hat{y}_1 = \psi(2, \text{verb}) = \text{noun}$$

Preventing Underflowing: Computation in the Log-Domain

- You may have noticed that probabilities get smaller and smaller as we go through the several time steps
- This may cause underflows and numerical instability
- Easy solution: instead of multiplying probabilities, sum log-probabilities!

$$\log(\exp(a) \times \exp(b)) = a + b$$

- To add probabilities, implement a log-sum function:

$$\log(\exp(a) + \exp(b)) = a + \log(1 + \exp(b - a))$$

Summing Up: Viterbi

- An instance of a dynamic programming algorithm
- Computes the **most likely sequence of tags**
- In other words, the structured output that maximizes the total score/probability
- This is called **MAP** (*maximum a posteriori*) **decoding**
- Later we'll talk about **marginal decoding** as a way of computing the structured output that minimizes the total risk

Recap: Three Problems in HMMs

Let $\mathbf{x} = \text{START}, \mathbf{x}_1, \dots, \mathbf{x}_L, \text{STOP}$ denote a word sequence.

- 1 Given \mathbf{x} , compute the *most likely sequence of states*
 $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x})$
- 2 Given \mathbf{x} , compute the *likelihood* $\mathbb{P}(\mathbf{x})$ and the *posterior state probabilities* $\mathbb{P}(\mathbf{y}_i | \mathbf{x})$
- 3 Given training sequences of word and states (or just words), *estimate the emission and transition probabilities.*

Recap: Three Problems in HMMs

Let $\mathbf{x} = \text{START}, \mathbf{x}_1, \dots, \mathbf{x}_L, \text{STOP}$ denote a word sequence.

- 1 Given \mathbf{x} , compute the *most likely sequence of states*
 $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x})$
- 2 Given \mathbf{x} , compute the *likelihood* $\mathbb{P}(\mathbf{x})$ and the *posterior state probabilities* $\mathbb{P}(\mathbf{y}_i | \mathbf{x})$
- 3 Given training sequences of word and states (or just words), *estimate the emission and transition probabilities.*

Problem 2: Likelihood and State Posteriors

- Likelihood:

$$\begin{aligned}\mathbb{P}(\mathbf{x}) &= \sum_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_L, \mathbf{y}_L, \text{STOP}) \\ &= \sum_{\mathbf{y}_1, \dots, \mathbf{y}_L} \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \prod_{i=1}^L \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i)\end{aligned}$$

- Posterior state probabilities:

$$\begin{aligned}\mathbb{P}(\mathbf{y}_i | \mathbf{x}) &= \frac{\sum_{\mathbf{y}_1, \dots, \mathbf{y}_{i-1}, \mathbf{y}_{i+1}, \dots, \mathbf{y}_L} \mathbb{P}(\text{START}, \mathbf{x}_1, \mathbf{y}_1, \dots, \mathbf{x}_L, \mathbf{y}_L, \text{STOP})}{\mathbb{P}(\mathbf{x})} \\ &= \frac{\sum_{\mathbf{y}_1, \dots, \mathbf{y}_{i-1}, \mathbf{y}_{i+1}, \dots, \mathbf{y}_L} \prod_{j=1}^{L+1} \mathbb{P}(\mathbf{y}_j | \mathbf{y}_{j-1}) \times \prod_{j=1}^L \mathbb{P}(\mathbf{x}_j | \mathbf{y}_j)}{\mathbb{P}(\mathbf{x})}\end{aligned}$$

- Both involve marginalizing out the hidden states
- Huge summations ($O(|\Lambda|^L)$ terms).

The Forward-Backward Algorithm

$$\begin{aligned}\mathbb{P}(\mathbf{x}) &= \sum_{\mathbf{y}_1, \dots, \mathbf{y}_L} \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \prod_{i=1}^L \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \\ &= \sum_{\mathbf{y}_L} \mathbb{P}(\text{STOP} | \mathbf{y}_L) \sum_{\mathbf{y}_{L-1}} \mathbb{P}(\mathbf{y}_L | \mathbf{y}_{L-1}) \mathbb{P}(\mathbf{x}_L | \mathbf{y}_L) \dots \sum_{\mathbf{y}_1} \mathbb{P}(\mathbf{y}_2 | \mathbf{y}_1) \mathbb{P}(\mathbf{x}_2 | \mathbf{y}_2) (\mathbb{P}(\mathbf{y}_1 | \text{START}) \mathbb{P}(\mathbf{x}_1 | \mathbf{y}_1))\end{aligned}$$

- **Key idea: interleave summations with products!**

The Forward-Backward Algorithm

$$\begin{aligned}\mathbb{P}(\mathbf{x}) &= \sum_{\mathbf{y}_1, \dots, \mathbf{y}_L} \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \prod_{i=1}^L \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \\ &= \sum_{\mathbf{y}_L} \mathbb{P}(\text{STOP} | \mathbf{y}_L) \sum_{\mathbf{y}_{L-1}} \mathbb{P}(\mathbf{y}_L | \mathbf{y}_{L-1}) \mathbb{P}(\mathbf{x}_L | \mathbf{y}_L) \dots \sum_{\mathbf{y}_1} \mathbb{P}(\mathbf{y}_2 | \mathbf{y}_1) \mathbb{P}(\mathbf{x}_2 | \mathbf{y}_2) (\mathbb{P}(\mathbf{y}_1 | \text{START}) \mathbb{P}(\mathbf{x}_1 | \mathbf{y}_1))\end{aligned}$$

- **Key idea: interleave summations with products!**
- **Forward variables:**

$$\alpha(i, \mathbf{y}_i) = \sum_{\mathbf{y}_{i-1} \in \Lambda} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \times \alpha(i-1, \mathbf{y}_{i-1})$$

The Forward-Backward Algorithm

$$\begin{aligned}\mathbb{P}(\mathbf{x}) &= \sum_{\mathbf{y}_1, \dots, \mathbf{y}_L} \prod_{i=1}^{L+1} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \prod_{i=1}^L \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \\ &= \sum_{\mathbf{y}_L} \mathbb{P}(\text{STOP} | \mathbf{y}_L) \sum_{\mathbf{y}_{L-1}} \mathbb{P}(\mathbf{y}_L | \mathbf{y}_{L-1}) \mathbb{P}(\mathbf{x}_L | \mathbf{y}_L) \dots \sum_{\mathbf{y}_1} \mathbb{P}(\mathbf{y}_2 | \mathbf{y}_1) \mathbb{P}(\mathbf{x}_2 | \mathbf{y}_2) (\mathbb{P}(\mathbf{y}_1 | \text{START}) \mathbb{P}(\mathbf{x}_1 | \mathbf{y}_1))\end{aligned}$$

- **Key idea: interleave summations with products!**
- **Forward variables:**

$$\alpha(i, \mathbf{y}_i) = \sum_{\mathbf{y}_{i-1} \in \Lambda} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \times \alpha(i-1, \mathbf{y}_{i-1})$$

- **Forward algorithm:** same as Viterbi, but replacing max with \sum
- **Backward algorithm:** same rationale, but right-to-left

The Forward-Backward Algorithm

input: sequence $\mathbf{x}_1, \dots, \mathbf{x}_L$, emission/transition probabilities

Forward pass: compute the forward probabilities

$$\alpha(1, \mathbf{y}_1) = \mathbb{P}(\mathbf{y}_1 | \text{START}) \times \mathbb{P}(\mathbf{x}_1 | \mathbf{y}_1) \quad \forall \mathbf{y}_1 \in \Lambda$$

for $i = 2$ **to** L **do**

for $\mathbf{y}_i \in \Lambda$ **do**

$$\alpha(i, \mathbf{y}_i) = \sum_{\mathbf{y}_{i-1}} \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \times \alpha(i-1, \mathbf{y}_{i-1})$$

Backward pass: compute the backward probabilities

$$\beta(L, \mathbf{y}_L) = \mathbb{P}(\text{STOP} | \mathbf{y}_L) \quad \forall \mathbf{y}_L \in \Lambda$$

for $i = L-1$ **to** 1 **do**

for $\mathbf{y}_i \in \Lambda$ **do**

$$\beta(i, \mathbf{y}_i) = \sum_{\mathbf{y}_{i+1}} \mathbb{P}(\mathbf{y}_{i+1} | \mathbf{y}_i) \times \mathbb{P}(\mathbf{x}_{i+1} | \mathbf{y}_{i+1}) \times \beta(i+1, \mathbf{y}_{i+1})$$

output: Posterior unigram probabilities $\mathbb{P}(\mathbf{y}_i | \mathbf{x}) \propto \alpha(i, \mathbf{y}_i) \times \beta(i, \mathbf{y}_i)$

Posterior bigram probabilities

$$\mathbb{P}(\mathbf{y}_i, \mathbf{y}_{i-1} | \mathbf{x}) \propto \alpha(i-1, \mathbf{y}_{i-1}) \times \mathbb{P}(\mathbf{y}_i | \mathbf{y}_{i-1}) \times \mathbb{P}(\mathbf{x}_i | \mathbf{y}_i) \times \beta(i, \mathbf{y}_i)$$

$$\text{Likelihood } \mathbb{P}(\mathbf{x}) = \sum_{\mathbf{y}_i} \alpha(i, \mathbf{y}_i) \beta(i, \mathbf{y}_i).$$

Decoding: Dynamic Programming

- The Viterbi and forward-backward algorithms are two instances of dynamic programming algorithms
- Many other problems in NLP (e.g. parsing with context-free grammars) can also be solved with dynamic programming
- Examples: inside-outside, CKY, Eisner's, belief propagation, ...
- We'll see this later in another lecture.

Minimum Risk Decoding

- Recall that the forward-backward algorithm allow us to compute the **posterior marginal probabilities**:

$$\mathbb{P}(\mathbf{y}_i | \mathbf{x}) \propto \alpha(i, \mathbf{y}_i)\beta(i, \mathbf{y}_i)$$

- This is called **marginal decoding** (a.k.a. marginal inference)
- We can now compute, for each position i , the label $\hat{\mathbf{y}}_i$ that maximizes this posterior probability

$$\hat{\mathbf{y}}_i = \arg \max_{\mathbf{y}_i} \mathbb{P}(\mathbf{y}_i | \mathbf{x})$$

- This is called **minimum risk decoding**
- Does this give the same sequence as Viterbi decoding?

Minimum Risk Decoding

- Recall that the forward-backward algorithm allow us to compute the **posterior marginal probabilities**:

$$\mathbb{P}(\mathbf{y}_i | \mathbf{x}) \propto \alpha(i, \mathbf{y}_i)\beta(i, \mathbf{y}_i)$$

- This is called **marginal decoding** (a.k.a. marginal inference)
- We can now compute, for each position i , the label $\hat{\mathbf{y}}_i$ that maximizes this posterior probability

$$\hat{\mathbf{y}}_i = \arg \max_{\mathbf{y}_i} \mathbb{P}(\mathbf{y}_i | \mathbf{x})$$

- This is called **minimum risk decoding**
- Does this give the same sequence as Viterbi decoding? **No.**

Viterbi or Minimum Risk Decoding?

- They will not, in general, give the same label sequence.
- Sometimes one works better, sometimes the other.
- Posterior decoding can give a label sequence that itself gets zero probability!
- We can motivate each by the **cost** they minimize.

Imagine the following game:

- The HMM will be given a new sequence x , which we must label.
- Our predicted label sequence \hat{y} will be compared to the true one, y .
- Depending on how badly we do, we will pay a fine.
- We want to minimize the expected cost (a.k.a. the **risk**).
- Our best strategy will depend on how the cost is defined!

All-or-Nothing Cost

- Suppose we pay 1€ if we get the label sequence wrong, i.e., $\hat{y} \neq y$.
- Otherwise we pay nothing.
- What should we do?
- We should use the most probable whole sequence \hat{y} .
- i.e., we should do **Viterbi decoding**.

Hamming Cost

- Alternately, suppose we pay 0.1€ for every word that we label incorrectly.
- This is more forgiving, and suggests that we focus on reasoning about each word without worrying about the coherence of the whole sequence.
- What should we do?
- We should use the most probable label \hat{y}_i for each word.
- i.e., we should do **minimum risk decoding**.

Recap: Three Problems in HMMs

Let $\mathbf{x} = \text{START}, \mathbf{x}_1, \dots, \mathbf{x}_L, \text{STOP}$ denote a word sequence.

- 1 Given \mathbf{x} , compute the *most likely sequence of states*
 $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x})$
- 2 Given \mathbf{x} , compute the *likelihood* $\mathbb{P}(\mathbf{x})$ and the *posterior state probabilities* $\mathbb{P}(\mathbf{y}_i | \mathbf{x})$
- 3 Given training sequences of word and states (or just words), *estimate the emission and transition probabilities.*

Recap: Three Problems in HMMs

Let $\mathbf{x} = \text{START}, \mathbf{x}_1, \dots, \mathbf{x}_L, \text{STOP}$ denote a word sequence.

- 1 Given \mathbf{x} , compute the *most likely sequence of states*
 $\arg \max_{\mathbf{y}_1, \dots, \mathbf{y}_L} \mathbb{P}(\mathbf{y}_1, \dots, \mathbf{y}_L | \mathbf{x})$
- 2 Given \mathbf{x} , compute the *likelihood* $\mathbb{P}(\mathbf{x})$ and the *posterior state probabilities* $\mathbb{P}(\mathbf{y}_i | \mathbf{x})$
- 3 Given training sequences of word and states (or just words), *estimate the emission and transition probabilities.*

Problem 3: Estimating the Model Parameters

How to estimate the emission probabilities $\mathbb{P}(x_i|y_i)$ and the transition probabilities $\mathbb{P}(y_i|y_{i-1})$?

We will look at two cases:

- 1 **Supervised learning:** assumes we have *labeled* training data $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(N)}, \mathbf{y}^{(N)})\}$
- 2 **Unsupervised learning:** assumes all we have is *unlabeled* training data $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$

HMM Supervised Learning

Assumes we have *labeled* training data $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(N)}, \mathbf{y}^{(N)})\}$

Simply use Maximum Likelihood Estimation for the complete data!

- Choose parameters that maximize $\prod_{n=1}^N \mathbb{P}(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$
- Closed form solution: **just count events and normalize**
 - How many times two states \mathbf{y}_{i-1} and \mathbf{y}_i occur consecutively \implies
Estimate transition probabilities $P(\mathbf{y}_i | \mathbf{y}_{i-1})$
 - How many times a state \mathbf{y}_i and a symbol x_i occur together \implies
Estimate emission probabilities $P(x_i | \mathbf{y}_i)$
- This is the “structured” version of Naive Bayes!

HMM Unsupervised Learning

Assumes all we have is *unlabeled* training data $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$

This is a lot harder!

We never get to see any label \mathbf{y} , so the best we can aim is to learn a good HMM up to a permutation of the labels!

But which objective function should we try to optimize?

Common choice: Maximum Likelihood Estimation with incomplete data!

- Involves maximizing the marginal likelihood

$$\prod_{n=1}^N \mathbb{P}(\mathbf{x}^{(n)}) = \prod_{n=1}^N \sum_{\mathbf{y}^{(n)}} \mathbb{P}(\mathbf{x}^{(n)}, \mathbf{y}^{(n)})$$

- However, this is usually prohibitive because it requires marginalizing out $\mathbf{y}^{(n)}$
- “Standard” solution: **Expectation-Maximization (EM) algorithm**

Expectation-Maximization (Baum-Welch)

Expectation-Maximization (EM) is a general inference technique, applicable beyond HMMs

- It works for any graphical model mixing observed and latent variables
- Applicable to any generative model with latent variables (discrete or continuous)!
- You might have seen it for mixtures of Gaussians and other clustering algorithms (e.g. soft clustering)
- In the case of HMMs, this is also called the **Baum-Welch algorithm**.

Expectation-Maximization (Baum-Welch)

Recall that if we observed the hidden variables \mathbf{y} during training, the parameters could be estimated easily by counting and normalizing

- How many times two states \mathbf{y}_{i-1} and \mathbf{y}_i occur consecutively \implies
Estimate transition probabilities $P(\mathbf{y}_i | \mathbf{y}_{i-1})$
- How many times a state \mathbf{y}_i and a symbol x_i occur together \implies
Estimate emission probabilities $P(x_i | \mathbf{y}_i)$

Unfortunately we don't observe the \mathbf{y} 's...

... but **if we knew** the model parameters, we could estimate the posterior marginal probabilities $P(\mathbf{y}_i | \mathbf{x})$ and $P(\mathbf{y}_{i-1}, \mathbf{y}_i | \mathbf{x})$, which would give us **fractional counts**

This is a chicken-and-egg problem!

Expectation-Maximization (Baum-Welch)

EM initializes the model parameters randomly

Then it alternates between these two steps:

- E-step: apply the current model to each training sequence x to obtain $P(\mathbf{y}_i | x)$ and $P(\mathbf{y}_{i-1}, \mathbf{y}_i | x)$
 - We can do this with the **forward-backward algorithm!**
- M-step: update the model
 - We can do this by counting events (fractional counts) and normalizing!

This is guaranteed to **improve the observed data likelihood at every iteration** (try to prove this at home!)

But MLE with incomplete data is a non-convex problem, so we can get to a local minimum.

Alternatives to EM

Also tackling MLE:

- Gradient descent: tackles MLE directly, also guaranteed to converge to a local minimum, also reutilizes existing toolboxes
- SGD: similar, but does updates after observing each sequence

Method of moments:

- Instead of tackling MLE, try to match observed moments with expected moments under the model
- Collect high-order moments (e.g. triplets of consecutive observed symbols) and solve an equation to obtain the model parameters
- Can be done with **spectral learning**, which boils down to a SVD matrix factorization
- Also doable with **anchor learning**, which boils down to non-negative matrix factorization

Semi-Supervised Learning of HMMs

Can also be done with EM, by combining sequences where labels are observed with sequence where they are not observed.

Boils down to combining true counts with fractional (expected) counts.

Summary of HMMs

- Representation? **Directed sequence model.**
- Assumptions? **Markov assumption on states; words are conditionally independent given the state.**
- Decoding/Inference? **Viterbi/forward-backward algorithms.**
- Learning the parameters? **Maximum likelihood (count and normalize)** for the supervised case, EM for the unsupervised case.

Roadmap: Models for Structured Prediction

Binary/Multi-class	Structured Prediction
--------------------	-----------------------

Naive Bayes	HMMs
Logistic Regression	?
Perceptron	?
SVMs	?

Related: Class-Based Language Models

- Proposed by Brown et al. (1992)
- Similar to HMMs where the observed symbols are **words** and the states are **word classes**
- Detail: we assume that each word belongs to **one and only one class**
 - i.e. the emission probabilities $P(x_i | y_i)$ are hard-clustered around classes; $P(x_i | y_i) = 0$ if word x_i does not belong to class y_i
- We also assume classes have a hierarchy (from coarse-grained to fine-grained classes)
- Each class is represented with a bit-string, number of bits determines how fine-grained it is
- Classes are learned with agglomerative clustering
- This is known as **Brown clustering**.

Examples of Brown clusters (from Twitter data)

Path	Terms
001010110	never neva nvr gladly nevr #never nevrerr nver nevrerrr nevaa
001010111	ever eva evar evr everrr everr everrrr evah everrrrr everrrrr
01000010	does duz doess does sayeth doez doesss d0es deos

Path	Terms
0100	Monday
010100	Sunday
010101	Friday
0101100	Thursday
01011010	Saturday

(from http://www.cs.cmu.edu/~ark/TweetNLP/cluster_viewer.html)

Brown clustering

Brown clustering is a way of **representation learning**

It obtains **discrete word representations** (binary vectors) via unsupervised hierarchical clustering

It was extremely popular in NLP before the age of neural networks

We'll look at other ways of learning word representations (continuous ones) in the next class

① Structured Prediction

② Sequence Models

Markov Models

Hidden Markov Models

Conditional Random Fields

Structured Perceptron

Structured Support Vector Machines

Conditional Random Fields (Lafferty et al., 2001)

HMMs have two drawbacks in practical NLP problems:

- 1 Words being conditionally independent is too strong an assumption
- 2 Can't have coarser features than words (e.g capitalization, suffixes, etc.) which often help

Conditional Random Fields (Lafferty et al., 2001)

HMMs have two drawbacks in practical NLP problems:

- ① Words being conditionally independent is too strong an assumption
- ② Can't have coarser features than words (e.g. capitalization, suffixes, etc.) which often help

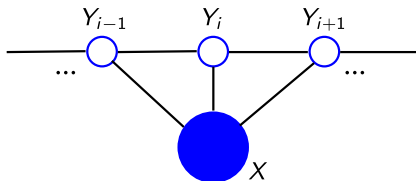
CRFs overcome these drawbacks—they're a structured variant of logistic regression!

Also used in vision (Kumar and Hebert, 2003; Quattoni et al., 2004)

Conditional Random Fields

Similar factorization as HMMs, but globally normalized!

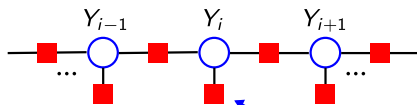
$$\mathbb{P}_{\mathbf{W}}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{W}, \mathbf{x})} \exp \left(\underbrace{\sum_i \mathbf{w}_{y_i} \cdot \phi_i(\mathbf{x})}_{\text{nodes}} + \underbrace{\sum_i \mathbf{w}_{y_i, y_{i-1}} \cdot \phi_{i, i-1}(\mathbf{x})}_{\text{edges}} \right)$$



Conditional Random Fields

Similar factorization as HMMs, but globally normalized!

$$\mathbb{P}_{\mathbf{W}}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{W}, \mathbf{x})} \exp \left(\underbrace{\sum_i \mathbf{w}_{y_i} \cdot \phi_i(\mathbf{x})}_{\text{nodes}} + \underbrace{\sum_i \mathbf{w}_{y_i, y_{i-1}} \cdot \phi_{i, i-1}(\mathbf{x})}_{\text{edges}} \right)$$



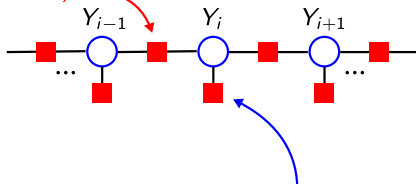
$$\psi_i(y_i) := \exp(w \cdot \phi_i(x, y_i)) \text{ (unary potentials)}$$

Conditional Random Fields

Similar factorization as HMMs, but globally normalized!

$$\mathbb{P}_{\mathbf{W}}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z(\mathbf{W}, \mathbf{x})} \exp \left(\underbrace{\sum_i \mathbf{w}_{y_i} \cdot \phi_i(\mathbf{x})}_{\text{nodes}} + \underbrace{\sum_i \mathbf{w}_{y_i, y_{i-1}} \cdot \phi_{i,i-1}(\mathbf{x})}_{\text{edges}} \right)$$

$\psi_{i,i-1}(y_i, y_{i-1}) := \exp(w \cdot \phi_{i,i-1}(x, y_i, y_{i-1}))$
(pairwise potentials)



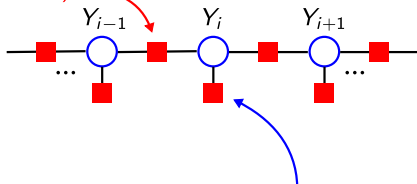
$\psi_i(y_i) := \exp(w \cdot \phi_i(x, y_i))$ (unary potentials)

Conditional Random Fields

Similar factorization as HMMs, but globally normalized!

$$\begin{aligned}\mathbb{P}_{\mathbf{W}}(\mathbf{y}|\mathbf{x}) &= \frac{1}{Z(\mathbf{W}, \mathbf{x})} \exp \left(\underbrace{\sum_i \mathbf{w}_{y_i} \cdot \phi_i(\mathbf{x})}_{\text{nodes}} + \underbrace{\sum_i \mathbf{w}_{y_i, y_{i-1}} \cdot \phi_{i,i-1}(\mathbf{x})}_{\text{edges}} \right) \\ &\propto \prod_i \psi_i(\mathbf{y}_i) \prod_i \psi_{i,i-1}(\mathbf{y}_i, \mathbf{y}_{i-1})\end{aligned}$$

$\psi_{i,i-1}(y_i, y_{i-1}) := \exp(w \cdot \phi_{i,i-1}(x, y_i, y_{i-1}))$
(pairwise potentials)



$\psi_i(y_i) := \exp(w \cdot \phi_i(x, y_i))$ (unary potentials)

Learning and Decoding in CRFs

- HMMs are a *generative* model: they model $\mathbb{P}(\mathbf{x}, \mathbf{y})$
- CRFs are a *discriminative* model: they model $\mathbb{P}(\mathbf{y}|\mathbf{x})$

Learning and Decoding in CRFs

- HMMs are a *generative* model: they model $\mathbb{P}(\mathbf{x}, \mathbf{y})$
- CRFs are a *discriminative* model: they model $\mathbb{P}(\mathbf{y}|\mathbf{x})$
- We can't compute $\mathbb{P}(\mathbf{x})$ with a CRF...
- ... But we can compute the most likely sequence of states and the posterior state probabilities
- **Just use the standard Viterbi and forward-backward algorithms (they work with unnormalized distributions too!)**

Learning and Decoding in CRFs

- Learning the parameters w (with supervision) is done by maximizing *conditional log-likelihood*:

$$\max_{\mathbf{w}} \frac{\lambda}{2} \|\mathbf{W}\|^2 + \sum_t \log \mathbb{P}_{\mathbf{w}}(\mathbf{y}^{(t)} | \mathbf{x}^{(t)})$$

- **A convex optimization problem, typically solved with L-BFGS or stochastic gradient descent**

Recap: Logistic Regression

- Define conditional probability

$$P_{\mathbf{W}}(\mathbf{y}|\mathbf{x}) = \frac{\exp(\mathbf{w}_{\mathbf{y}} \cdot \phi(\mathbf{x}))}{Z_{\mathbf{x}}}$$

- Set weights to maximize conditional log-likelihood of training data:

$$\mathbf{W} = \operatorname{argmax}_{\mathbf{W}} \sum_t \log P_{\mathbf{W}}(\mathbf{y}_t|\mathbf{x}_t) = \operatorname{argmin}_{\mathbf{W}} \sum_t L(\mathbf{W}; (\mathbf{x}_t, \mathbf{y}_t))$$

- Can find the gradient and run gradient descent (or any gradient-based optimization algorithm)

$$\nabla_{\mathbf{W}} L(\mathbf{W}; (\mathbf{x}, \mathbf{y})) = \sum_{\mathbf{y}'} P_{\mathbf{W}}(\mathbf{y}'|\mathbf{x}) \mathbf{e}_{\mathbf{y}'} \phi(\mathbf{x})^{\top} - \mathbf{e}_{\mathbf{y}} \phi(\mathbf{x})^{\top}$$

Structural Decomposition

For conditional random fields, the features and their expectation decompose according to the sequential structure:

$$W = \left[\overbrace{W_{\text{unig}}}^{[\dots w_{\mathbf{y}} \dots]}, \overbrace{W_{\text{big}}}^{[\dots w_{\mathbf{y}, \mathbf{y}_{\text{prev}}} \dots]} \right]$$

$$\phi_{\text{unig}}(\mathbf{x}) = \sum_i \phi_i(\mathbf{x})$$

$$\phi_{\text{big}}(\mathbf{x}) = \sum_i \phi_{i, i-1}(\mathbf{x})$$

$$\sum_{\mathbf{y}'} P_W(\mathbf{y}' | \mathbf{x}) e_{\mathbf{y}'} \phi_{\text{unig}}(\mathbf{x})^\top = \sum_i \sum_{\mathbf{y}'_i} P_W(\mathbf{y}'_i | \mathbf{x}) e_{\mathbf{y}'_i} \phi_i(\mathbf{x})^\top$$

$$\sum_{\mathbf{y}'} P_W(\mathbf{y}' | \mathbf{x}) e_{\mathbf{y}'} \phi_{\text{big}}(\mathbf{x})^\top = \sum_i \sum_{\mathbf{y}'_{i-1}, \mathbf{y}'_i} P_W(\mathbf{y}'_i, \mathbf{y}'_{i-1} | \mathbf{x}) e_{\mathbf{y}'_i \wedge \mathbf{y}'_{i-1}} \phi_{i, i-1}(\mathbf{x})^\top$$

All quantities can be obtained with the **forward-backward algorithm!**

Conditional Random Fields

- Representation? **Undirected sequence model.**
- Assumptions? **Markov assumption on states; no assumption of the words whatsoever.**
- Decoding/Inference? **Viterbi/forward-backward algorithms.**
- Learning the parameters? **Maximum conditional likelihood (convex optimization).**

Features for POS tagging (Ratnaparkhi, 1999)

Tag trigram	$y_{i-2} \wedge y_{i-1} \wedge y_i$
Tag bigram	$y_{i-1} \wedge y_i$
Word	$x_i \wedge y_i$
Previous word	$x_{i-1} \wedge y_i$
Two words before	$x_{i-2} \wedge y_i$
Next word	$x_{i+1} \wedge y_i$
Two words after	$x_{i+2} \wedge y_i$
Word has upper case	$\text{HASUPPER}(x_i) \wedge y_i$
Word has digit	$\text{HASDIGIT}(x_i) \wedge y_i$
Word has hyphen	$\text{HASHYPHEN}(x_i) \wedge y_i$
Word prefixes (up to 4 characters)	$\text{PREFIX}(x_i) \wedge y_i$
Word suffixes (up to 4 characters)	$\text{SUFFIX}(x_i) \wedge y_i$

- **Accuracies are $\sim 97\%$ in the Penn Treebank**

Features for NER (Ratinov and Roth, 2009)

Similar, plus some more features:

- Word shapes: “Leuven” → Aa+
- Flags for all-digits, all-upper, first-upper
- Same for previous/next word and two words before/after
- Bigram tag features also conjoined with words
- Gazetteer features (e.g. entities found in Wikipedia)

Standard evaluation metric is F_1 *measure* (harmonic mean of precision and recall)

- $F_1 \sim 80\%$ in the **English CoNLL 2003 dataset**

Roadmap: Models for Structured Prediction

Binary/Multi-class

Structured Prediction

Naive Bayes
Logistic Regression
Perceptron
SVMs

HMMs
CRFs
?
?

① Structured Prediction

② Sequence Models

Markov Models

Hidden Markov Models

Conditional Random Fields

Structured Perceptron

Structured Support Vector Machines

Structured Perceptron (Collins, 2002)

- The perceptron algorithm has their structured counterpart too
- The prediction $\hat{\mathbf{y}}$ according to the current model \mathbf{W} is now the most likely sequence of states (computed using the Viterbi algorithm)
- Everything else still applies (similar updates, similar mistake bound)

Recap: Perceptron

input: labeled data \mathcal{D}

initialize $\mathbf{W} = \mathbf{0}$

repeat

 get new training example $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$

 predict $\hat{\mathbf{y}} = \arg \max_{\mathbf{y}' \in \mathcal{Y}} \mathbf{w}_{\mathbf{y}'} \cdot \phi(\mathbf{x})$

if $\hat{\mathbf{y}} \neq \mathbf{y}$ **then**

 update $\mathbf{w}_{\mathbf{y}} \leftarrow \mathbf{w}_{\mathbf{y}} + \phi(\mathbf{x})$

 update $\mathbf{w}_{\hat{\mathbf{y}}} \leftarrow \mathbf{w}_{\hat{\mathbf{y}}} - \phi(\mathbf{x})$

until maximum number of epochs

output: model weights \mathbf{W}

Structured Perceptron

input: labeled data \mathcal{D}

initialize $\mathbf{W} = \mathbf{0}$

repeat

get new training example $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}$

predict $\hat{\mathbf{y}} = \arg \max_{\mathbf{y}' \in \mathcal{Y}} \sum_i \mathbf{w}_{\mathbf{y}'_i} \cdot \phi_i(\mathbf{x}) + \sum_i \mathbf{w}_{\mathbf{y}'_i, \mathbf{y}'_{i-1}} \cdot \phi_{i,i-1}(\mathbf{x})$

for $i = 1, \dots, L + 1$ **do**

if $i \leq L$ and $\hat{\mathbf{y}}_i \neq \mathbf{y}_i$ **then**

update $\mathbf{w}_{\mathbf{y}_i} \leftarrow \mathbf{w}_{\mathbf{y}_i} + \phi_i(\mathbf{x}_i)$

update $\mathbf{w}_{\hat{\mathbf{y}}_i} \leftarrow \mathbf{w}_{\hat{\mathbf{y}}_i} - \phi_i(\mathbf{x}_i)$

if $\hat{\mathbf{y}}_i \neq \mathbf{y}_i$ or $\hat{\mathbf{y}}_{i-1} \neq \mathbf{y}_{i-1}$ **then**

update $\mathbf{w}_{\mathbf{y}_i, \mathbf{y}_{i-1}} \leftarrow \mathbf{w}_{\mathbf{y}_i, \mathbf{y}_{i-1}} + \phi_{i,i-1}(\mathbf{x}_i)$

update $\mathbf{w}_{\hat{\mathbf{y}}_i, \hat{\mathbf{y}}_{i-1}} \leftarrow \mathbf{w}_{\hat{\mathbf{y}}_i, \hat{\mathbf{y}}_{i-1}} - \phi_{i,i-1}(\mathbf{x}_i)$

until maximum number of epochs

output: model weights \mathbf{W}

Structured Perceptron

- Representation? **Undirected sequence model.**
- Assumptions? **Markov assumption on states; no assumption of the words whatsoever.**
- Decoding/Inference? **Viterbi.**
- Learning the parameters? **Perceptron updates.**

① Structured Prediction

② Sequence Models

Markov Models

Hidden Markov Models

Conditional Random Fields

Structured Perceptron

Structured Support Vector Machines

Structured Support Vector Machines (Taskar et al., 2003; Altun et al., 2003; Tsochantaridis et al., 2004)

- SVMs have their structured counterpart too
- Similar idea to CRFs, but maximize margin
- No probabilistic interpretation anymore
- Can still be used to compute the most likely sequence of states (using the Viterbi algorithm)
- The loss function is **structured hinge loss**: a convex surrogate to **Hamming loss** (the fraction of incorrect states)

Recap: SVMs

Hinge loss:

$$\begin{aligned}L((\mathbf{x}, \mathbf{y}); \mathbf{W}) &= \max(0, 1 + \max_{\mathbf{y}' \neq \mathbf{y}} \mathbf{w}_{\mathbf{y}'} \cdot \phi(\mathbf{x}) - \mathbf{w}_{\mathbf{y}} \cdot \phi(\mathbf{x})) \\ &= \left(\max_{\mathbf{y}' \in \mathcal{Y}} \mathbf{w}_{\mathbf{y}'} \cdot \phi(\mathbf{x}) + [[\mathbf{y}' \neq \mathbf{y}]] \right) - \mathbf{w}_{\mathbf{y}} \cdot \phi(\mathbf{x})\end{aligned}$$

A **subgradient** of the hinge is

$$\partial_{\mathbf{W}} L((\mathbf{x}, \mathbf{y}); \mathbf{W}) \ni \mathbf{e}_{\hat{\mathbf{y}}} \phi(\mathbf{x})^\top - \mathbf{e}_{\mathbf{y}} \phi(\mathbf{x})^\top$$

where

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}' \in \mathcal{Y}} \mathbf{w}_{\mathbf{y}'} \cdot \phi(\mathbf{x}) + [[\mathbf{y}' \neq \mathbf{y}]]$$

Structured Support Vector Machines

Key idea: replace the 0/1 cost $[[\mathbf{y}' \neq \mathbf{y}]]$ by a **Hamming cost**:

$$\rho(\mathbf{y}', \mathbf{y}) = \sum_i [[\mathbf{y}'_i \neq \mathbf{y}_i]]$$

In other words: the “margin” counts how many labels are incorrect

Cost-augmented decoding:

$$\begin{aligned}\hat{\mathbf{y}} &= \arg \max_{\mathbf{y}' \in \mathcal{Y}} \mathbf{w}_{\mathbf{y}'} \cdot \phi(\mathbf{x}) + \rho(\mathbf{y}', \mathbf{y}) \\ &= \arg \max_{\mathbf{y}' \in \mathcal{Y}} \sum_i \left(\mathbf{w}_{\mathbf{y}'_i} \cdot \phi_i(\mathbf{x}) + [[\mathbf{y}'_i, \mathbf{y}_i]] \right) + \sum_i \mathbf{w}_{\mathbf{y}'_i, \mathbf{y}'_{i-1}} \cdot \phi_{i,i-1}(\mathbf{x})\end{aligned}$$

The Hamming cost is **decomposable**! Hence cost-augmented decoding can still be done with the Viterbi algorithm

A subgradient can then be computed as usual.

Structured Support Vector Machines

- Representation? **Undirected sequence model.**
- Assumptions? **Markov assumption on states; no assumption of the words whatsoever.**
- Decoding/Inference? **Viterbi algorithm.**
- Learning the parameters? **Minimize structured hinge loss (convex, non-smooth optimization).**

Roadmap: Models for Structured Prediction

Binary/Multi-class

Naive Bayes
Logistic Regression
Perceptron
SVMs

Structured Prediction

HMMs
CRFs
Structured Perceptron
Structured SVMs

Neural Sequence Models

So far, all models discussed are structured but **linear**

But we know that neural networks are more powerful than linear models...

Can we mix a neural network and a CRF?

Neural Sequence Models

So far, all models discussed are structured but **linear**

But we know that neural networks are more powerful than linear models...

Can we mix a neural network and a CRF? **Yes.**

Neural Sequence Models

Just append a CRF output layer to your favorite neural network!

For example:

- a shared feedforward neural network to each element of the sequence to compute unigram scores
- another shared feedforward to compute bigram scores (or just a label bigram indicator)
- an output layer with a sequential structure using those unigram and bigram scores

Use Forward-Backward to compute the gradient in the last layer; then gradient backpropagation will make all the work for us

Other configurations are possible instead of a shared feedforward: e.g. a CNN, a BILSTM, ... (will learn those later)

Works for a structured SVM output layer too.

References I

- Altun, Y., Tsochantaridis, I., and Hofmann, T. (2003). Hidden Markov support vector machines. In *Proc. of International Conference of Machine Learning*.
- Brants, T. (2000). Tnt: a statistical part-of-speech tagger. In *Proceedings of the sixth conference on Applied natural language processing*, pages 224–231. Association for Computational Linguistics.
- Brown, P. F., Desouza, P. V., Mercer, R. L., Pietra, V. J. D., and Lai, J. C. (1992). Class-based n-gram models of natural language. *Computational linguistics*, 18(4):467–479.
- Collins, M. (2002). Discriminative training methods for hidden Markov models: theory and experiments with perceptron algorithms. In *Proc. of Empirical Methods for Natural Language Processing*.
- Kumar, S. and Hebert, M. (2003). Discriminative random fields: A discriminative framework for contextual interaction in classification. In *Computer Vision, 2003. Proceedings. Ninth IEEE International Conference on*, pages 1150–1157. IEEE.
- Lafferty, J., McCallum, A., and Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proc. of International Conference of Machine Learning*.
- Quattoni, A., Collins, M., and Darrell, T. (2004). Conditional random fields for object recognition. In *Advances in neural information processing systems*, pages 1097–1104.
- Ratinov, L. and Roth, D. (2009). Design challenges and misconceptions in named entity recognition. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning*, pages 147–155. Association for Computational Linguistics.
- Ratnaparkhi, A. (1999). Learning to parse natural language with maximum entropy models. *Machine Learning*, 34(1):151–175.
- Taskar, B., Guestrin, C., and Koller, D. (2003). Max-margin Markov networks. In *Proc. of Neural Information Processing Systems*.
- Toutanova, K., Klein, D., Manning, C. D., and Singer, Y. (2003). Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proc. of the North American Chapter of the Association for Computational Linguistics*, pages 173–180.
- Tsochantaridis, I., Hofmann, T., Joachims, T., and Altun, Y. (2004). Support vector machine learning for interdependent and structured output spaces. In *Proc. of International Conference of Machine Learning*.
- Zhang, T. and Johnson, D. (2003). A robust risk minimization based named entity recognition system. In *Proceedings of the seventh conference on Natural language learning at HLT-NAACL 2003-Volume 4*, pages 204–207. Association for Computational Linguistics.