Lecture 5: Linear Sequential Models

Vlad Niculae & André Martins



Deep Structured Learning Course, Fall 2019

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

Outline

1 Structured Prediction

Ø Generative Sequence Models

Markov Models Hidden Markov Models Unsupervised learning

3 Discriminative Sequence Models

Structured Perceptron Conditional Random Fields Structured SVM

So far, we looked at **classification**: simply predicting one-of-K classes.

How about more complicated output spaces?

McGrath	left	out	of	Ireland	World	Cup	squad
---------	------	-----	----	---------	-------	-----	-------

I	0	0	0	I	I	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

• Classify each word independently: $x_1 = \text{McGrath}; \quad y_1 = \text{I},$ $x_2 = \text{left}; \quad y_2 = 0, ...$

I	0	0	0	I	I	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

• Classify each word independently: $x_1 = \text{McGrath}; \quad y_1 = \text{I},$ $x_2 = \text{left}; \quad y_2 = 0, ...$

В	0	0	0	В	В	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

- Classify each word independently: $x_1 = McGrath; \quad y_1 = B,$ $x_2 = left; \quad y_2 = 0, ...$
- This labeling scheme is called **BIO** (beginning/inside/outside).

В	0	0	0	В	В	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

- Classify each word independently: $x_1 = McGrath; \quad y_1 = B,$ $x_2 = left; \quad y_2 = 0, ...$
- This labeling scheme is called **BIO** (beginning/inside/outside).
- Could add more context into *x*, e.g. a window:
 - $egin{aligned} & x_1 = (\$, \mathsf{McGrath}, \mathsf{left}); & y_1 = \mathsf{B}, \ & x_2 = (\mathsf{McGrath}, \mathsf{left}, \mathsf{out}); & y_2 = \mathsf{O}, \dots \end{aligned}$

В	0	0	0	В	В	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

- Classify each word independently: $x_1 = McGrath; \quad y_1 = B,$ $x_2 = left; \quad y_2 = 0, ...$
- This labeling scheme is called **BIO** (beginning/inside/outside).
- Could add more context into \boldsymbol{x} , e.g. a window: $\boldsymbol{x}_1 = (\$, \mathsf{McGrath}, \mathsf{left}); \qquad \boldsymbol{y}_1 = \mathsf{B},$
 - $x_2 = (McGrath, left, out); y_2 = 0, ...$
- But the predictions can't depend on one another!
 OI is not allowed; BB is allowed but relatively rare!
 ... what do we *really* want?

В	0	0	0	В	В	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

- Classify each word independently: this makes a big assumption! $x_1 = McGrath; \quad y_1 = B,$
- Really, we should think at sentence-level:

 $x_1 =$ McGrath left out of...; $y_1 =$ BOOOBBIO

- Can capture more statistics (BB is rare, BBBB even rarer, 00 is common...)
- Can we attempt this with our multi-class toolbox (e.g. naïve Bayes?)

В	0	0	0	В	В	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

- Classify each word independently: this makes a big assumption! $x_1 = McGrath; \quad y_1 = B,$
- Really, we should think at sentence-level:

 $x_1 =$ McGrath left out of...; $y_1 =$ BOOOBBIO

- Can capture more statistics (BB is rare, BBBB even rarer, 00 is common...)
- Can we attempt this with our multi-class toolbox (e.g. naïve Bayes?) How many sentences have the label B 0 0 0 B B I 0?

В	0	0	0	В	В	I	0
McGrath	left	out	of	Ireland	World	Cup	squad

- Classify each word independently: this makes a big assumption! $x_1 = McGrath; \quad y_1 = B,$
- Really, we should think at sentence-level:

 $x_1 =$ McGrath left out of...; $y_1 =$ BOOOBBIO

- Can capture more statistics (BB is rare, BBBB even rarer, 00 is common...)
- Can we attempt this with our multi-class toolbox (e.g. naïve Bayes?) How many sentences have the label B 0 0 0 B B I 0?

rank	sequence	count
1	0	1072
2	B O	663
3	B O O O O O O O	446
4	BOBO	378
5	0 B I 0 B 0 0	272
4856	B O O O B B I O	1

We are essentially treating each label sequence *y* as a distinct object, *ignoring its internal structure!*

We are essentially treating each label sequence y as a distinct object, *ignoring its internal structure!*

We've never seen sentences labelled 0 0 0 0 0 B, but we *have* seen sentences labelled 0 0 0 0 B. Surely that could help!

Structured Prediction

A framework for handling structured, constrained, inter-dependent outputs.

NLP

- Named Entity Recognition
- Machine Translation
- Syntactic Parsing



Speech Processing

- Speaker ID
- Speech Recognition



Computer Vision

- Object detection
- Segmentation



computational biology, robotics / planning, time series forecasting, etc.

Structured Prediction



Structured Prediction



Today, we talk about sequences.

Roadmap: Models for Structured Prediction

Binary/Multi-class	Structured Prediction
Naive Bayes	?
Logistic Regression	?
Perceptron	?
SVMs	?

Outline

O Structured Prediction

O Generative Sequence Models

Markov Models Hidden Markov Models Unsupervised learning

B Discriminative Sequence Models

Structured Perceptron

Conditional Random Fields

Structured SVM

Outline

Structured Prediction

O Generative Sequence Models

Markov Models

Hidden Markov Models

Unsupervised learning

B Discriminative Sequence Models

Structured Perceptron Conditional Random Fields Structured SVM

To begin, let's forget for a moment about x and focus on y. $P(y = \begin{bmatrix} 0 & 0 & 0 & B \end{bmatrix}) =?$

To begin, let's forget for a moment about x and focus on y. $P(y = \begin{bmatrix} 0 & 0 & 0 & B \end{bmatrix}) =?$ What are all possible labellings?

Vlad Niculae & André Martins (IST)

To begin, let's forget for a moment about x and focus on y. $P(y = \begin{bmatrix} 0 & 0 & 0 & B \end{bmatrix}) =?$ What are all possible labellings?

Let $\Sigma = \{B, I, 0\}.$ Then, $\mathfrak{Y} = \Sigma^{\star} = \Sigma \cup \Sigma^2 \cup \Sigma^3 \cup ...$

To begin, let's forget for a moment about x and focus on y. $P(y = \begin{bmatrix} 0 & 0 & 0 & B \end{bmatrix}) = ?$ What are all possible labellings? Let $\Sigma = \{B, I, 0\}$. Then, $\Im = \Sigma^* = \Sigma \cup \Sigma^2 \cup \Sigma^3 \cup ...$

Could set $\mathsf{P}(y) \propto \# y$

To begin, let's forget for a moment about x and focus on y. $P(y = [0 \ 0 \ 0 \ B]) = ?$ What are all possible labellings? Let $\Sigma = \{B, I, 0\}$. Then, $\mathcal{Y} = \Sigma^* = \Sigma \cup \Sigma^2 \cup \Sigma^3 \cup ...$ rank sequence count 1072 0 2 B O 663 3 B O O O O O O 446 Could set $\mathsf{P}(y) \propto \# y$ 4 BOBO 378 5 OBTOBOO 272 BOOOBBIO 4856 1

To begin, let's forget for a moment about x and focus on y. $P(y = [0 \ 0 \ 0 \ B]) = ?$ What are all possible labellings? Let $\Sigma = \{B, I, 0\}$. Then, $\mathcal{Y} = \Sigma^* = \Sigma \cup \Sigma^2 \cup \Sigma^3 \cup ...$ rank sequence count 1072 0 2 B O 663 3 B O O O O O O 446 Could set $\mathsf{P}(y) \propto \# y$ 4 B O B O 378 5 OBTOBOO 272 4856 B O O B B I O 1

Issues:

P(y) = 0 for **most** y! No sharing between sequences that are similar! Ignores the fact that $y = [y_1, \dots, y_L]$.

Lower extreme: Bag-of-words model

Let
$$\boldsymbol{y} = [y_1, \dots, y_L]$$
. $\mathsf{P}(\boldsymbol{y}) = \prod_{i=1}^L \mathsf{P}(y_i)$

- Also called "unigram" model
- Assumes every word is generated independently of other words therefore, abandons the structure of *y* entirely.

Lower extreme: Bag-of-words model

Let
$$\boldsymbol{y} = [y_1, \dots, y_L]$$
. $\mathsf{P}(\boldsymbol{y}) = \prod_{i=1}^L \mathsf{P}(y_i)$

- Also called "unigram" model
- Assumes every word is generated independently of other words therefore, abandons the structure of *y* entirely.
- Probability of a string is insensitive to word order:

$$P([0, B, 0]) = P([B, 0, 0]) = P(B) P(0)^{2}$$

Lower extreme: Bag-of-words model

Let
$$\boldsymbol{y} = [y_1, \dots, y_L]$$
. $\mathsf{P}(\boldsymbol{y}) = \prod_{i=1}^L \mathsf{P}(y_i)$

- Also called "unigram" model
- Assumes every word is generated independently of other words therefore, abandons the structure of *y* entirely.
- Probability of a string is insensitive to word order:

$$P([0, B, 0]) = P([B, 0, 0]) = P(B) P(0)^{2}$$

• How many parameters do we need to estimate and how?

Let
$$\boldsymbol{y} = [\texttt{start}, y_1, y_2, \dots, y_L, \texttt{end}].$$

 $\mathsf{P}(\boldsymbol{y}) =$

Let
$$\boldsymbol{y} = [\texttt{start}, y_1, y_2, \dots, y_L, \texttt{end}].$$

 $\mathsf{P}(\boldsymbol{y}) = \mathsf{P}(y_1 | \texttt{start})$
 $\cdot \mathsf{P}(y_2 | \texttt{start}, y_1)$
 $\cdot \dots$
 $\cdot \mathsf{P}(y_L | \texttt{start}, y_1, y_2, \dots, y_{L-1})$

 $\cdot \mathsf{P}(\mathsf{end}|\mathsf{start}, y_1, y_2, \dots, y_{L-1}, y_L)$

Let
$$\boldsymbol{y} = [\texttt{start}, y_1, y_2, \dots, y_L, \texttt{end}].$$

 $\mathsf{P}(\boldsymbol{y}) = \mathsf{P}(y_1 | \texttt{start})$
 $\cdot \mathsf{P}(y_2 | \texttt{start}, y_1)$
 $\cdot \dots$
 $\cdot \mathsf{P}(y_L | \texttt{start}, y_1, y_2, \dots, y_{L-1})$
 $\cdot \mathsf{P}(\texttt{end} | \texttt{start}, y_1, y_2, \dots, y_{L-1}, y_L)$
 $= \prod_{i=1}^{L+1} \mathsf{P}(y_i | \underbrace{y_1, \dots, y_{i-1}}_{\boldsymbol{y}'})$

Vlad Niculae & André Martins (IST)

Ρ

Let
$$\boldsymbol{y} = [\texttt{start}, y_1, y_2, \dots, y_L, \texttt{end}].$$

$$P(\boldsymbol{y}) = P(y_1 | \texttt{start})$$

$$\cdot P(y_2 | \texttt{start}, y_1)$$

$$\cdot \dots$$

$$\cdot P(y_L | \texttt{start}, y_1, y_2, \dots, y_{L-1})$$

$$\cdot P(\texttt{end} | \texttt{start}, y_1, y_2, \dots, y_{L-1}, y_L)$$

$$= \prod_{i=1}^{L+1} P(y_i | \underbrace{y_1, \dots, y_{i-1}}_{\boldsymbol{y'}})$$

• Each symbol y_i generated based on entire history y'.

Let
$$\boldsymbol{y} = [\texttt{start}, y_1, y_2, \dots, y_L, \texttt{end}].$$

 $\mathsf{P}(\boldsymbol{y}) = \mathsf{P}(y_1 | \texttt{start})$
 $\cdot \mathsf{P}(y_2 | \texttt{start}, y_1)$
 $\cdot \dots$
 $\cdot \mathsf{P}(y_L | \texttt{start}, y_1, y_2, \dots, y_{L-1})$
 $\cdot \mathsf{P}(\texttt{end} | \texttt{start}, y_1, y_2, \dots, y_{L-1}, y_L)$
 $= \prod_{i=1}^{L+1} \mathsf{P}(y_i | \underbrace{y_1, \dots, y_{i-1}}_{\boldsymbol{y}'})$

Each symbol y_i generated based on entire history y'.

Ρ

- Must estimate P(y|y') for every possible history y'!
 - ... and we're back where we started: same as counting sequences.

Let
$$\boldsymbol{y} = [\mathtt{start}, y_1, y_2, \dots, y_L, \mathtt{end}].$$

 $\mathsf{P}(\boldsymbol{y}) = \mathsf{P}(y_1 | \mathtt{start})$
 $\cdot \mathsf{P}(y_2 | \mathtt{start}, y_1)$
 $\cdot \dots$
 $\cdot \mathsf{P}(y_L | \mathtt{start}, y_1, y_2, \dots, y_{L-1})$
 $\cdot \mathsf{P}(\mathtt{end} | \mathtt{start}, y_1, y_2, \dots, y_{L-1}, y_L)$
 $= \prod_{i=1}^{L+1} \mathsf{P}(y_i | \underbrace{y_1, \dots, y_{i-1}}_{\boldsymbol{y}'})$

- Each symbol y_i generated based on entire history y'.
- Must estimate P(y|y') for every possible history y'!
 - ... and we're back where we started: same as counting sequences.
- Idea: condition only on the last *few* symbols.

Vlad Niculae & André Martins (IST)

Lecture 5: Linear Sequential Models
In-between: Markov Models

Let
$$\boldsymbol{y} = [\mathtt{start}, y_1, y_2, \dots, y_L, \mathtt{end}].$$

$$P(\boldsymbol{y}) = P(y_1 | \mathtt{start}) \cdot P(y_2 | y_1) \cdot \dots \cdot P(y_L | y_{L-1}) \cdot P(\mathtt{end} | y_L)$$

$$= \prod_{i=1}^{L+1} P(y_i | y_{i-1})$$

- Each symbol only depends on the previous word.
- We estimate transition probabilities P(y_i|y_{i-1});
 Including initial and final probabilities P(y₁|start) and P(end|y_L).
- Total number of parameters:

In-between: Markov Models

Let
$$\boldsymbol{y} = [\mathtt{start}, y_1, y_2, \dots, y_L, \mathtt{end}].$$

$$\mathsf{P}(\boldsymbol{y}) = \mathsf{P}(y_1 | \mathtt{start}) \cdot \mathsf{P}(y_2 | y_1) \cdot \dots \cdot \mathsf{P}(y_L | y_{L-1}) \cdot \mathsf{P}(\mathtt{end} | y_L)$$

$$= \prod_{i=1}^{L+1} \mathsf{P}(y_i | y_{i-1})$$

- Each symbol only depends on the previous word.
- We estimate transition probabilities P(y_i|y_{i-1});
 Including initial and final probabilities P(y₁|start) and P(end|y_L).
- Total number of parameters: $O(|\Sigma|^2)$.

$$\begin{array}{lll} P(B|\text{start}) = .393 & P(B|B) = .009 & P(B|I) = .003 & P(B|0) = .102 \\ P(I|\text{start}) = .0 & P(I|B) = .369 & P(I|I) = .178 & P(I|0) = .0 \\ P(0|\text{start}) = .607 & P(0|B) = .610 & P(0|I) = .779 & P(0|0) = .815 \\ P(\text{end}|\text{start}) = .0 & P(\text{end}|B) = .013 & P(\text{end}|I) = .040 & P(\text{end}|0) = .084 \end{array}$$

Aside: kth order Markov Models

Let
$$y = [\text{start}, y_1, y_2, ..., y_L, \text{end}].$$
 $P(y) = \prod_{i=1}^{L+1} P(y_i | y_{i-1}, ..., y_{i-k})$

- Each symbol depends on k previous symbols.
- Transition probabilities $P(y_i|y_{i-1}, \ldots, y_{i-k})$
- Total number of parameters: $O(|\Sigma|^{k+1})$
- Widely used in language modeling
 - Here, Σ = the vocabulary of English words.
 - Goal: next word prediction; P(w|"can we rely") =?

$$\mathsf{P}(y_i = \mathbf{b} | y_{i-1} = \mathbf{a}) = \frac{\#[\mathbf{a}, \mathbf{b}]}{\sum_{\mathbf{b}' \in \Sigma} \#[\mathbf{a}, \mathbf{b}']} = \frac{\#[\mathbf{a}, \mathbf{b}]}{\#\mathbf{a}}$$

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	P(I I) = .178	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

$$\mathsf{P}(y_i = \mathbf{b}|y_{i-1} = \mathbf{a}) = \frac{\#[\mathbf{a}, \mathbf{b}]}{\sum_{\mathbf{b}' \in \Sigma} \#[\mathbf{a}, \mathbf{b}']} = \frac{\#[\mathbf{a}, \mathbf{b}]}{\#\mathbf{a}}$$

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	P(I I) = .178	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

We never saw $\begin{bmatrix} 0 & I \end{bmatrix}$, so P(I|0) = 0. The model rules out **any sequence containing** $\begin{bmatrix} 0 & I \end{bmatrix}$!

$$\mathsf{P}([\texttt{start}, 0, 0, \texttt{I}, \texttt{B}, 0, \texttt{end}]) = \mathsf{P}(0|\texttt{start}) \cdot \mathsf{P}(0|0) \cdot \underbrace{\mathsf{P}(\texttt{I}|0)}_{=0} \cdot \ldots = 0$$

$$\mathsf{P}(y_i = \mathbf{b} | y_{i-1} = \mathbf{a}) = \frac{\#[\mathbf{a}, \mathbf{b}]}{\sum_{\mathbf{b}' \in \Sigma} \#[\mathbf{a}, \mathbf{b}']} = \frac{\#[\mathbf{a}, \mathbf{b}]}{\#\mathbf{a}}$$

 $\begin{array}{lll} P(B|\texttt{start}) = .393 & P(B|B) = .009 & P(B|I) = .003 & P(B|0) = .102 \\ P(I|\texttt{start}) = .0 & P(I|B) = .369 & P(I|I) = .178 & P(I|0) = .0 \\ P(0|\texttt{start}) = .607 & P(0|B) = .610 & P(0|I) = .779 & P(0|0) = .815 \\ P(\texttt{end}|\texttt{start}) = .0 & P(\texttt{end}|B) = .013 & P(\texttt{end}|I) = .040 & P(\texttt{end}|0) = .084 \end{array}$

We never saw $\begin{bmatrix} 0 & I \end{bmatrix}$, so P(I|0) = 0. The model rules out **any sequence containing** $\begin{bmatrix} 0 & I \end{bmatrix}$!

$$\mathsf{P}([\texttt{start}, 0, 0, \texttt{I}, \texttt{B}, 0, \texttt{end}]) = \mathsf{P}(0|\texttt{start}) \cdot \mathsf{P}(0|0) \cdot \underbrace{\mathsf{P}(\texttt{I}|0)}_{=0} \cdot \ldots = 0$$

Here, this is correct; other times it may be just due to insufficient data.

$$\mathsf{P}(y_i = \mathbf{b} | y_{i-1} = \mathbf{a}) = \frac{\#[\mathbf{a}, \mathbf{b}]}{\sum_{\mathbf{b}' \in \Sigma} \#[\mathbf{a}, \mathbf{b}']} = \frac{\#[\mathbf{a}, \mathbf{b}]}{\#\mathbf{a}}$$

 $\begin{array}{lll} P(B|\texttt{start}) = .393 & P(B|B) = .009 & P(B|I) = .003 & P(B|0) = .102 \\ P(I|\texttt{start}) = .0 & P(I|B) = .369 & P(I|I) = .178 & P(I|0) = .0 \\ P(0|\texttt{start}) = .607 & P(0|B) = .610 & P(0|I) = .779 & P(0|0) = .815 \\ P(\texttt{end}|\texttt{start}) = .0 & P(\texttt{end}|B) = .013 & P(\texttt{end}|I) = .040 & P(\texttt{end}|0) = .084 \end{array}$

We never saw $\begin{bmatrix} 0 & I \end{bmatrix}$, so P(I|0) = 0. The model rules out **any sequence containing** $\begin{bmatrix} 0 & I \end{bmatrix}$!

$$\mathsf{P}([\texttt{start}, 0, 0, \texttt{I}, \texttt{B}, 0, \texttt{end}]) = \mathsf{P}(0|\texttt{start}) \cdot \mathsf{P}(0|0) \cdot \underbrace{\mathsf{P}(\texttt{I}|0)}_{=0} \cdot \ldots = 0$$

Here, this is correct; other times it may be just due to insufficient data. Smoothing (one way): pretend we saw each possible transition once more

$$\mathsf{P}(y_i = \mathsf{b}|y_{i-1} = \mathsf{a}) = \frac{1 + \#[\mathsf{a},\mathsf{b}]}{\sum_{\mathsf{b}}, 1 + \#[\mathsf{a},\mathsf{b}']} = \frac{1 + \#[\mathsf{a},\mathsf{b}]}{|\Sigma| + \#\mathsf{a}}$$

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	$P(\mathbf{I} \mathbf{I}) = .178$	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	$P(\mathbf{I} \mathbf{B}) = .369$	$P(\mathbf{I} \mathbf{I}) = .178$	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

• Given a sequence, assess its likelihood under the model.

y = [start, B, 0, 0, end]P(y) = P(B|start) P(0|B) P(0|0) P(end|0) = .393 · .610 · .815 · .084 = .0164

- Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	$P(\mathbf{I} \mathbf{B}) = .369$	$P(\mathbf{I} \mathbf{I}) = .178$	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

• Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	$P(\mathbf{I} \mathbf{I}) = .178$	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

• Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	$P(\mathbf{I} \mathbf{I}) = .178$	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

• Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



- Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



$$\begin{array}{c|c} P(B|\texttt{start}) = .393 \\ P(I|\texttt{start}) = .0 \\ P(0|\texttt{start}) = .607 \\ P(\texttt{end}|\texttt{start}) = .0 \end{array} \begin{array}{c|c} P(B|B) = .009 \\ P(I|B) = .369 \\ P(I|I) = .178 \\ P(I|I) = .178 \\ P(0|I) = .178 \\ P(0|I) = .779 \\ P(0|0) = .815 \\ P(\texttt{end}|\texttt{start}) = .0 \end{array} \begin{array}{c|c} P(0|B) = .610 \\ P(0|I) = .779 \\ P(0|0) = .815 \\ P(\texttt{end}|I) = .040 \\ P(\texttt{end}|0) = .084 \end{array}$$

• Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



- Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



$$\begin{array}{c|c} P(B|\texttt{start}) = .393 \\ P(I|\texttt{start}) = .0 \\ P(0|\texttt{start}) = .607 \\ P(\texttt{end}|\texttt{start}) = .0 \end{array} \begin{array}{c|c} P(B|B) = .009 \\ P(I|B) = .369 \\ P(I|I) = .178 \\ P(I|I) = .178 \\ P(0|I) = .178 \\ P(0|I) = .779 \\ P(0|0) = .815 \\ P(\texttt{end}|\texttt{start}) = .0 \end{array} \begin{array}{c|c} P(0|B) = .610 \\ P(0|I) = .779 \\ P(0|0) = .815 \\ P(\texttt{end}|I) = .040 \\ P(\texttt{end}|0) = .084 \end{array}$$

• Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



- Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164



- Given a sequence, assess its likelihood under the model.

$$y = [start, B, 0, 0, end]$$

P(y) = P(B|start) P(0|B) P(0|0) P(end|0)
= .393 · .610 · .815 · .084 = .0164

• Sample sequences, going from left to right!



• Predict the most likely next symbol (like your phone's autocomplete)! (As above, but take max instead of a random sample.

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	P(I I) = .178	$P(\mathbf{I} 0) = .0$
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

• What is the probability that the second symbol in a sequence is B?

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	$P(\mathbf{I} \mathbf{I}) = .178$	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

What is the probability that the second symbol in a sequence is B?
 P([start,?,B]) = P(y₀ = start, y₂ = B) =?

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	P(I I) = .178	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

- What is the probability that the second symbol in a sequence is B?
 P([start,?,B]) = P(y₀ = start, y₂ = B) =?
- We must consider all possible choices for y_1 . Recall $P(b) = \sum_a P(a, b)!$

$$\begin{split} \mathsf{P}([\texttt{start},?,\texttt{B}]) &= \sum_{y \in \Sigma} \mathsf{P}([\texttt{start},y,\texttt{B}]) \\ &= \mathsf{P}([\texttt{start},\texttt{B},\texttt{B}]) + \mathsf{P}([\texttt{start},\texttt{I},\texttt{B}]) + \mathsf{P}([\texttt{start},0,\texttt{B}]) \\ &= .393 \cdot .009 + .0 \cdot .003 + .607 \cdot .102 = .065 \end{split}$$

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	P(I I) = .178	$P(\mathbf{I} 0) = .0$
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

- What is the probability that the second symbol in a sequence is B?
 P([start,?,B]) = P(y₀ = start, y₂ = B) =?
- We must consider all possible choices for y_1 . Recall $P(b) = \sum_a P(a, b)!$

$$\begin{split} \mathsf{P}([\texttt{start},?,B]) &= \sum_{y \in \Sigma} \mathsf{P}([\texttt{start},y,B]) \\ &= \mathsf{P}([\texttt{start},B,B]) + \mathsf{P}([\texttt{start},I,B]) + \mathsf{P}([\texttt{start},0,B]) \\ &= .393 \cdot .009 + .0 \cdot .003 + .607 \cdot .102 = .065 \end{split}$$

• How about the proba of **B** as *third* symbol?

P(B start) = .393	P(B B) = .009	P(B I) = .003	P(B 0) = .102
P(I start) = .0	P(I B) = .369	P(I I) = .178	P(I 0) = .0
P(0 start) = .607	P(0 B) = .610	P(0 I) = .779	P(0 0) = .815
P(end start) = .0	P(end B) = .013	P(end I) = .040	P(end 0) = .084

- What is the probability that the second symbol in a sequence is B?
 P([start,?,B]) = P(y₀ = start, y₂ = B) =?
- We must consider all possible choices for y_1 . Recall $P(b) = \sum_a P(a, b)!$

$$\begin{split} \mathsf{P}([\texttt{start},?,B]) &= \sum_{y \in \Sigma} \mathsf{P}([\texttt{start},y,B]) \\ &= \mathsf{P}([\texttt{start},B,B]) + \mathsf{P}([\texttt{start},I,B]) + \mathsf{P}([\texttt{start},0,B]) \\ &= .393 \cdot .009 + .0 \cdot .003 + .607 \cdot .102 = .065 \end{split}$$

• How about the proba of B as third symbol?

$$\begin{split} \mathsf{P}([\texttt{start},?,?B]) &= \sum_{y_1 \in \Sigma} \sum_{y_2 \in \Sigma} \mathsf{P}([\texttt{start},y_1,y_2,B]) = \mathsf{P}([\texttt{start},B,B,B]) + \mathsf{P}([\texttt{start},B,I,B]) + ... \\ (\texttt{regrouping}) &= \mathsf{P}([\texttt{start},?,B,B]) + \mathsf{P}([\texttt{start},?,I,B]) + \mathsf{P}([\texttt{start},?,0,B]) \\ &= \mathsf{P}([\texttt{start},?B]) \, \mathsf{P}(B|B) + \mathsf{P}([\texttt{start},?,I]) \, \mathsf{P}(B|I) + \mathsf{P}([\texttt{start},?,0]) \, \mathsf{P}(B|0) \\ &= \sum_{y_2 \in \Sigma} \mathsf{P}([\texttt{start},?,y_2]) \, \mathsf{P}(y_2|B) \quad (\texttt{A pattern we'll see again soon!}) \end{split}$$

What can we do with a MM?



• Can we do NER?



What can we do with a MM?



• Can we do NER? start B 0 0 end Halloween is coming

Not well! — a model of P(y) does not take x into account!

What can we do with a MM?





Not well! — a model of P(y) does not take x into account!

• We want a model that can take *x* into account, too.

Outline

Structured Prediction

O Generative Sequence Models

Markov Models

Hidden Markov Models

Unsupervised learning

B Discriminative Sequence Models

Structured Perceptron

Conditional Random Fields

Structured SVM

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models:



- Observe $m{x}$ and use Bayes' rule: $\mathsf{P}(m{y}|m{x}) \propto \mathsf{P}(m{y}) \mathsf{P}(m{x}|m{y})$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$
 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models: y P(y)x P(x|y)
- Observe $m{x}$ and use Bayes' rule: $\mathsf{P}(m{y}|m{x}) \propto \mathsf{P}(m{y}) \mathsf{P}(m{x}|m{y})$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$
 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

 $\begin{array}{ccccccc} & B & I & 0 \\ P(y) = .115 & .052 & .833 \\ P(\text{South}|y) = & .0051 & .0002 & .0 \\ P(\text{Africa}|y) = & .0003 & .0056 & .0 \\ P(\text{under}|y) = & .0 & .0 & .0009 \\ P(\text{Mountains}|y) = & .0 & .0002 & .0 \\ P(\text{of}|y) = & .0 & .0127 & .0212 \\ \end{array}$

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models: y P(y)x P(x|y)
- Observe $m{x}$ and use Bayes' rule: $\mathsf{P}(m{y}|m{x}) \propto \mathsf{P}(m{y}) \mathsf{P}(m{x}|m{y})$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$

 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

В Ι P(y) = .115.052 .833 P(South|y) = .0051.0002 .0 P(Africa|y) = .0003 .0056.0 P(under|y) = .0.0 .0009 P(Mountains|y) = .0.0002 .0 P(of|y) = .0.0127 .0212 *Y*1 Y2 Уз X2 *X*1 X3

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models: y P(y)x P(x|y)
- Observe $m{x}$ and use Bayes' rule: $\mathsf{P}(m{y}|m{x}) \propto \mathsf{P}(m{y}) \mathsf{P}(m{x}|m{y})$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$

 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

В Ι P(y) = .115 .052 .833 P(South|y) = .0051.0002 .0 P(Africa|y) = .0003 .0056.0 P(under|y) = .0.0 .0009 P(Mountains|y) = .0 .0002 .0 P(of|y) = .0.0127 .0212 *Y*1 Y2 Уз of Mountains Africa

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models: y P(y)x P(x|y)
- Observe $m{x}$ and use Bayes' rule: $\mathsf{P}(m{y}|m{x}) \propto \mathsf{P}(m{y}) \mathsf{P}(m{x}|m{y})$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$
 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

В Ι P(y) = .115 .052 .833 P(South|y) = .0051.0002 .0 P(Africa|y) = .0003 .0056.0 P(under|y) = .0.0 .0009 P(Mountains|y) = .0 .0002 .0 P(of|y) = .0.0127 .0212 Ι Τ of Mountains Africa

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models: y P(y)x P(x|y)
- Observe x and use Bayes' rule: $\mathsf{P}(y|x) \propto \mathsf{P}(y) \,\mathsf{P}(x|y)$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$
 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

$$P(y) = .115 .052 .833$$

$$P(South|y) = .0051 .0002 .0$$

$$P(Africa|y) = .0003 .0056 .0$$

$$P(under|y) = .0 .0 .0009$$

$$P(Mountains|y) = .0 .0002 .0$$

$$P(of|y) = .0 .0127 .0212$$
...
I
$$O$$
I
$$Africa$$

Why so bad?

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models: y P(y)x P(x|y)
- Observe $m{x}$ and use Bayes' rule: $\mathsf{P}(m{y}|m{x}) \propto \mathsf{P}(m{y}) \mathsf{P}(m{x}|m{y})$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$
 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

Why so bad? Smoothing doesn't fix it.

- We want $\mathsf{P}(m{y}|m{x})$ for prediction.
- We build two simpler models: y P(y)x P(x|y)
- Observe x and use Bayes' rule: $\mathsf{P}(y|x) \propto \mathsf{P}(y) \,\mathsf{P}(x|y)$
- For NER, a form of NB:

$$\mathsf{P}(\boldsymbol{y}) = \prod_{i} \mathsf{P}(y_{i})$$
 $\mathsf{P}(\boldsymbol{x}|\boldsymbol{y}) = \prod_{i} \mathsf{P}(x_{i}|y_{i})$

Why so bad? Smoothing doesn't fix it. $P([I \ 0 \ I]) = .225 \gg 0!$

Vlad Niculae & André Martins (IST)

Lecture 5: Linear Sequential Models
Hidden Markov Models

- Jointly model a sequence of observations x_i and hidden states y_i.
- States modeled by a first-order Markov model.
- Each observation is conditioned only on the corresponding state.



$$oldsymbol{x} = [x_1, \dots, x_L]; \quad oldsymbol{y} = [extstyle{start}, y_1, \dots, y_L, extstyle{end}].$$
 $\mathsf{P}(oldsymbol{x}, oldsymbol{y}) = \underbrace{\prod_{i=1}^{L+1} \mathsf{P}(y_i | y_{i-1})}_{\mathsf{P}(oldsymbol{y})} \cdot \underbrace{\prod_{i=1}^{L} \mathsf{P}(x_i | y_i)}_{\mathsf{P}(oldsymbol{x}|oldsymbol{y})}$

Estimating HMMs: Maximum Likelihood

maximize
$$\left(\prod_{(\boldsymbol{x},\boldsymbol{y})\in\mathcal{D}}\mathsf{P}(\boldsymbol{x},\boldsymbol{y})=\prod_{(\boldsymbol{x},\boldsymbol{y})\in\mathcal{D}}\prod_{i=1}^{L(\boldsymbol{x})+1}\mathsf{P}(y_i|y_{i-1})\cdot\prod_{i=1}^{L(\boldsymbol{x})}\mathsf{P}(x_i|y_i)\right)$$

A HMM is "just" a Markov model and an emission model :P

Estimating HMMs: Maximum Likelihood

maximize
$$\left(\prod_{(\boldsymbol{x},\boldsymbol{y})\in\mathcal{D}}\mathsf{P}(\boldsymbol{x},\boldsymbol{y})=\prod_{(\boldsymbol{x},\boldsymbol{y})\in\mathcal{D}}\prod_{i=1}^{L(\boldsymbol{x})+1}\mathsf{P}(y_i|y_{i-1})\cdot\prod_{i=1}^{L(\boldsymbol{x})}\mathsf{P}(x_i|y_i)\right)$$

A HMM is "just" a Markov model and an emission model :P

Transition probabilities: $P(y_i = b | y_{i-1} = a) = \frac{\#[a,b]}{\#a}$ just like the MM!P(B|start) = .393P(B|B) = .009P(B|I) = .003P(B|0) = .102P(I|start) = .0P(I|B) = .369P(I|I) = .178P(I|0) = .0P(0|start) = .607P(0|B) = .610P(0|I) = .779P(0|0) = .815P(end|start) = .0P(end|B) = .013P(end|I) = .040P(end|0) = .084

Estimating HMMs: Maximum Likelihood

maximize
$$\left(\prod_{(\boldsymbol{x},\boldsymbol{y})\in\mathcal{D}}\mathsf{P}(\boldsymbol{x},\boldsymbol{y}) = \prod_{(\boldsymbol{x},\boldsymbol{y})\in\mathcal{D}}\prod_{i=1}^{L(\boldsymbol{x})+1}\mathsf{P}(y_i|y_{i-1})\cdot\prod_{i=1}^{L(\boldsymbol{x})}\mathsf{P}(x_i|y_i)\right)$$

A HMM is "just" a Markov model and an emission model :P

Transition probabilities: $P(y_i = b | y_{i-1} = a) = \frac{\#[a,b]}{\#_a}$ just like the MM! $\begin{array}{ll} \mathsf{P}(\mathsf{B}|\texttt{start}) = .393 & \mathsf{P}(\mathsf{B}|\mathsf{B}) = .009 & \mathsf{P}(\mathsf{B}|\mathsf{I}) = .003 & \mathsf{P}(\mathsf{B}|\mathsf{0}) = .102 \\ \mathsf{P}(\mathsf{I}|\texttt{start}) = .0 & \mathsf{P}(\mathsf{I}|\mathsf{B}) = .369 & \mathsf{P}(\mathsf{I}|\mathsf{I}) = .178 & \mathsf{P}(\mathsf{I}|\mathsf{0}) = .0 \end{array}$ P(0|start) = .607 P(0|B) = .610 P(0|I) = .779 P(0|0) = .815P(end|start) = .0 P(end|B) = .013 P(end|I) = .040 P(end|0) = .084**Emission probabilities:** $P(x_i = w | y_i = a) = \frac{\# tag(w) = a}{\# a}$ just like NB! BIN $P(\text{South}|y) = .0051 .0002 \epsilon$ $P(Africa|y) = .0003 .0056 \epsilon$ $P(Mountains|y) = \epsilon$.0002 ϵ $P(of|y) = \epsilon$.0127 .0212

	В	I	0
P(South y) =	.0051	.0002	ϵ
P(Africa y) =	.0003	.0056	ϵ
P(Mountains y) =	ϵ	.0002	ϵ
P(of y) =	ϵ	.0127	.0212

	start	В	I	0
P(B y) =	.393	.009	.003	.102
P(I y) =	.0	.369	.178	.0
P(0 y) =	.607	.610	.779	.815
P(end y) =	.0	.013	.040	.084

	atoxt	D	т	0		В	T	U
- ())	Start	D	1	U	P(South v) =	.0051	.0002	ϵ
P(B y) =	.393	.009	.003	.102	P(Africaly) =	0002	0056	<i>c</i>
P(I v) =	.0	.369	.178	.0	F(AIIICa y) =	.0005	.0056	e
$P(\mathbf{n} \mathbf{v}) =$	607	610	770	015	P(Mountains y) =	ϵ	.0002	ϵ
$P(\mathbf{u} \mathbf{y}) =$.607	.610	.119	.615	P(of v) =	F	0127	0212
P(end y) =	.0	.013	.040	.084	((()))	0		

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B}, 0, \texttt{B}, \texttt{end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{O}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{O}) \,\mathsf{P}(\texttt{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{O}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$

	atomt	D	т	0	B I (U
- ())	start	Б	1	U	P(South v) = .0051 .0002	ϵ
P(B y) =	.393	.009	.003	.102	P(Africaly) = 0.002 0.056	c
$P(\mathbf{I} y) =$.0	.369	.178	.0	P(AIICa y) = .0003 .0030 0	E
P(n y) =	607	610	779	815	$P(Mountains y) = \epsilon$.0002 e	ε
$(\mathbf{u} \mathbf{y}) =$.001	.010	.115	.015	$P(of y) = \epsilon$.0127	.0212
P(end y) =	.0	.013	.040	.084		
$P(0 y) = P(\mathbf{end} y) =$.607 .0	.610 .013	.779 .040	.815 .084	$P(of y) = \epsilon$.0127 .	.0212

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



		D	т	0	B I U	J
	start	В	1	0	P(South y) = .0051 .0002 e	ε
P(B y) =	.393	.009	.003	.102	P(Africa v) = .0003 .0056	e
P(I y) =	.0	.369	.178	.0	P(Mountains y) = c = 0002	-
P(0 y) =	.607	.610	.779	.815	$P(MOUIIIaIIIS y) = \epsilon .0002 \epsilon$	5 0010
P(end v) =	0	013	040	084	$P(ot y) = \epsilon$.0127.	.0212
		.015				

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



	atoxt	D	т	0		В	T	U
- ()	start	Б	1	U	P(South v) =	.0051	.0002	ϵ
P(B y) =	.393	.009	.003	.102	P(Africaly) =	0002	0056	~
$P(\mathbf{I} v) =$.0	.369	.178	.0	F(Airica y) =	.0005	.0050	e
P(n y) =	607	610	770	015	P(Mountains y) =	ϵ	.0002	ϵ
$\Gamma(\mathbf{u} \mathbf{y}) =$.007	.010	.115	.015	P(of v) =	ϵ	.0127	.0212
P(end y) =	.0	.013	.040	.084				

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



-		ъ	т	0		В	1	U
D(T)	start.	Б	1	U	P(South y) =	.0051	.0002	ϵ
P(B y) =	393	.009	.003	.102	P(Africa y) =	.0003	.0056	ϵ
$P(1 \mathbf{y}) = .0$	0	.369	.178	.0	P(Mountains y) =	ϵ	.0002	ϵ
$P(\mathbf{U} \mathbf{y}) = .0$	60 <i>1</i>	.610	.119	.815	P(of y) =	ε	.0127	.0212
P(ena y) = .0	0	.013	.040	.084				

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



	atoxt	D	т	0		В	T	U
- ()	start	Б	1	U	P(South v) =	.0051	.0002	ϵ
P(B y) =	.393	.009	.003	.102	P(Africaly) =	0002	0056	~
$P(\mathbf{I} v) =$.0	.369	.178	.0	F(Airica y) =	.0005	.0050	e
P(n y) =	607	610	770	015	P(Mountains y) =	ϵ	.0002	ϵ
$\Gamma(\mathbf{u} \mathbf{y}) =$.007	.010	.115	.015	P(of v) =	ϵ	.0127	.0212
P(end y) =	.0	.013	.040	.084				

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



		_		-		В	1	0
5	start	В	I	0	P(South y) =	0051	0002	~
$P(\mathbf{R} y) =$	303	nna	003	102	r (Southy) =	.0051	.0002	e
$(\mathbf{u}_{ \mathbf{y} }) = .$.555	.005	.00.	.102	P(Africa v) -	0003	0056	F
P(T v) =	0	369	178	0	r (/ (/ (/ (/) /)) =	.0005	.0050	C
(-)		.000		.0	P(Mountains v) =	F	.0002	F
P(0 v) = .	.607	.610	.779	.815		c		
	_				P(ot v) =	ϵ	.0127	.0212
P(end y) = .	.0	.013	.04(.084				

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



	atomt	D	т	0	B I		U
- ())	Start	D	1	U	P(South v) = .0051 .00)02	ϵ
P(B y) =	.393	.009	.003	.102	P(Africaly) = 0.003 00	156	c .
$P(\mathbf{I} y) =$.0	.369	.178	.0	F(AIICa y) = .0003 .00	150	e
P(n y) =	607	610	779	Q15	$P(Mountains y) = \epsilon$.00)02	ϵ
	.001	.010	.115	.015	$P(of v) = \epsilon$.01	27	.0212
P(end y) =	.0	.013	.040	.084			

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



-

	-+ - m+	D	т	0			В	T	U
	Start	D	1	U	P(Sout	h v) =	.0051	.0002	ϵ
P(B y) = .	.393	.009	.003	.102	P(Afric	a(v) =	.0003	.0056	F
P(I y) = .	.0	.369	.178	.0	P(Mountain	(v) =	6	0002	c
P(0 y) = .	607	.610	.779	.815		s y =	c	0107	с 0010
P(end y) = .	0	.013	.040	.084	P(C	(עויי) =	e	.0127	.0212

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



stowt P	т	0	B I U	
Start D	1	U	$P(South v) = .0051 .0002 \epsilon$	
P(B y) = .393 .009	.003	.102	P(Africaly) = 0.003 0.056 c	
P(I v) = .0 .369	.178	.0	$\Gamma(AIIICa[y]) = .0003 .0030 e$	
P(0 y) = 607 = 610	779	815	$P(Mountains y) = \epsilon$.0002 ϵ	
P(0 y) = 1001	.115	.015	$P(of y) = \epsilon$.0127 .02	212
$P(\text{end} y) = .0 \qquad .013$.040	.084		

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



	D T	0	В	T	U
stari	вт	U	P(South y) = 0.051	0002	F
P(B y) = .393	.009 .003	.102	D(Africaly) = 0000	0050	c
P(T v) = 0	369 178	0	P(Africa y) = .0003	.0056	ϵ
$\Gamma(\underline{1} y) = .0$.505 .110	.0	$P(Mountains y) = \epsilon$.0002	ϵ
P(0 y) = .607	.610 .779	.815	P(of v) = c	0127	0212
P(end v) = .0	.013 .040	.084	$\Gamma(0 y) = \epsilon$.0121	.0212

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



stowt P	т	0	B I U	
Start D	1	U	$P(South v) = .0051 .0002 \epsilon$	
P(B y) = .393 .009	.003	.102	P(Africaly) = 0.003 0.056 c	
P(I v) = .0 .369	.178	.0	$\Gamma(AIIICa[y]) = .0003 .0030 e$	
P(0 y) = 607 = 610	779	815	$P(Mountains y) = \epsilon$.0002 ϵ	
P(0 y) = 1001	.115	.015	$P(of y) = \epsilon$.0127 .02	212
$P(\text{end} y) = .0 \qquad .013$.040	.084		

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \,\mathsf{P}(\mathsf{0}|\mathsf{B}) \,\mathsf{P}(\mathsf{B}|\mathsf{0}) \,\mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \,\mathsf{P}(\text{Mountains}|\mathsf{B}) \,\mathsf{P}(\mathsf{of}|\mathsf{0}) \,\mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



-

_	-++	D	т	0			В	T	U
	start	D	1	U	P(Sou	th v) =	.0051	.0002	ϵ
P(B y) = .	393	.009	.003	.102	P(Afri	ca(v) =	.0003	.0056	F
$P(\mathbf{I} y) = .$	0	.369	.178	.0	P(Mountai	ns(v) =	6	0002	c
P(0 y) = .	607	.610	.779	.815		(y) = (y)	C .	0107	с 0010
P(end y) = .	0	.013	.040	.084	P(o(y) =	e	.0127	.0212

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \, \mathsf{P}(\mathsf{0}|\mathsf{B}) \, \mathsf{P}(\mathsf{B}|\mathsf{0}) \, \mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \, \mathsf{P}(\text{Mountains}|\mathsf{B}) \, \mathsf{P}(\mathsf{of}|\mathsf{0}) \, \mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



stowt P	т	0	B I U	
Start D	1	U	$P(South v) = .0051 .0002 \epsilon$	
P(B y) = .393 .009	.003	.102	P(Africaly) = 0.003 0.056 c	
P(I v) = .0 .369	.178	.0	$\Gamma(AIIICa[y]) = .0003 .0030 e$	
P(0 y) = 607 = 610	779	815	$P(Mountains y) = \epsilon$.0002 ϵ	
P(0 y) = 1001	.115	.015	$P(of y) = \epsilon$.0127 .02	212
$P(\text{end} y) = .0 \qquad .013$.040	.084		

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \, \mathsf{P}(\mathsf{0}|\mathsf{B}) \, \mathsf{P}(\mathsf{B}|\mathsf{0}) \, \mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \, \mathsf{P}(\text{Mountains}|\mathsf{B}) \, \mathsf{P}(\mathsf{of}|\mathsf{0}) \, \mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



	D T	0	В	T	U
stari	вт	U	P(South y) = 0.051	0002	F
P(B y) = .393	.009 .003	.102	D(Africaly) = 0000	0050	c
P(T v) = 0	369 178	0	P(Africa y) = .0003	.0056	ϵ
$\Gamma(\underline{1} y) = .0$.505 .110	.0	$P(Mountains y) = \epsilon$.0002	ϵ
P(0 y) = .607	.610 .779	.815	P(of v) = c	0127	0212
P(end v) = .0	.013 .040	.084	$\Gamma(0 y) = \epsilon$.0121	.0212

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \, \mathsf{P}(\mathsf{0}|\mathsf{B}) \, \mathsf{P}(\mathsf{B}|\mathsf{0}) \, \mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \, \mathsf{P}(\text{Mountains}|\mathsf{B}) \, \mathsf{P}(\mathsf{of}|\mathsf{0}) \, \mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



stowt P	т	0	B I U	
Start D	1	U	$P(South v) = .0051 .0002 \epsilon$	
P(B y) = .393 .009	.003	.102	P(Africaly) = 0.003 0.056 c	
P(I v) = .0 .369	.178	.0	$\Gamma(AIIICa[y]) = .0003 .0030 e$	
P(0 y) = 607 = 610	779	815	$P(Mountains y) = \epsilon$.0002 ϵ	
P(0 y) = 1001 .000	.115	.015	$P(of y) = \epsilon$.0127 .02	212
$P(\text{end} y) = .0 \qquad .013$.040	.084		

Given a labeled sequence (x, y), assess its likelihood under the model.

 $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \, \mathsf{P}(\mathsf{0}|\mathsf{B}) \, \mathsf{P}(\mathsf{B}|\mathsf{0}) \, \mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \, \mathsf{P}(\text{Mountains}|\mathsf{B}) \, \mathsf{P}(\mathsf{of}|\mathsf{0}) \, \mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$



	atomt	D	т	0		В	1	U
- ())	Start	D	1	U	P(South v) =	.0051	.0002	ϵ
P(B y) =	.393	.009	.003	.102	P(Africaly) =	0002	0056	
P(I v) =	.0	.369	.178	.0	F(AIIICa y) =	.0005	.0056	e
D(0 y)	607	610	770	015	P(Mountains y) =	ϵ	.0002	ϵ
$P(\mathbf{u} \mathbf{y}) =$.607	.610	.119	.015	P(of v) =	F	0127	0212
P(end y) =	.0	.013	.040	.084		0		
					• • •			

Given a labeled sequence (x, y), assess its likelihood under the model.

- $\begin{aligned} \boldsymbol{x} &= [\text{Mountains, of, Africa}]; \quad \boldsymbol{y} &= [\texttt{start, B, 0, B, end}] \\ \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \mathsf{P}(\mathsf{B}|\texttt{start}) \, \mathsf{P}(\mathsf{0}|\mathsf{B}) \, \mathsf{P}(\mathsf{B}|\mathsf{0}) \, \mathsf{P}(\mathsf{end}|\mathsf{B}) \\ &\cdot \, \mathsf{P}(\text{Mountains}|\mathsf{B}) \, \mathsf{P}(\mathsf{of}|\mathsf{0}) \, \mathsf{P}(\mathsf{Africa}|\mathsf{B}) = 8 \cdot 10^{-17} \end{aligned}$
- Sample labelled sequences.



• Auto-complete equivalent: predict the most likely x_{k+1} given x_k, y_k .

Big HMM questions



- 1 What is the most likely label sequence *y*, given *x*?
- 2 What is the probability of x?
- 3 What is the probability of each assignment *y_i*, given *x*?
- What is the probability of each transition $y \rightarrow y'$, given x?
- 5 What sequence *y* minimizes the *Hamming cost?*

ł

- $\begin{array}{c|cccc} & B & I & 0 \\ P(South|y) = & .0051 & .0002 & \epsilon \\ P(Africa|y) = & .0003 & .0056 & \epsilon \\ P(Mountains|y) = & \epsilon & .0002 & \epsilon \\ P(of|y) = & \epsilon & .0127 & .0212 \\ & & & & \\ & & & \\$
- Can we predict the most likely label sequence y for a given x? arg max $_y$ P(y|x)

	start	B	т	Ο	D	T	U
	202	000	-002	102	P(South y) = .0051	.0002	ϵ
P(B y) =	.393	.009	.003	.102	P(Africav) = 0.003	0056	F
P(I y) =	.0	.369	.178	.0	D(Mountainslu) = 1	0000	
P(0 v) =	.607	.610	.779	.815	$P(Mountains y) = \epsilon$.0002	e
P(and y) =	0	012	040	0.04	$P(ot y) = \epsilon$.0127	.0212
r(enaly) =	.0	.015	.040	.004			

Can we predict the most likely label sequence y for a given x? $rg\max_y \mathsf{P}(y|x) = rg\max_y rac{\mathsf{P}(x,y)}{\mathsf{P}(x)} = rg\max_y \mathsf{P}(x,y)$

B I O

start	В	I	0		1	U
P(B v) = .393	.009	.003	.102	P(South y) = .0051	.0002	ϵ
P(T y) = 0	369	178	0	P(Africa y) = .0003	.0056	ϵ
$P(\mathbf{n} y) = .0$.505	.110	.0	$P(Mountains y) = \epsilon$.0002	ϵ
$P(\mathbf{u} \mathbf{y}) = .007$.010	.119	.015	$P(of y) = \epsilon$.0127	.0212
P(end y) = .0	.013	.040	.084			

• Can we predict the most likely label sequence y for a given x? arg max_y $P(y|x) = \arg \max_y \frac{P(x,y)}{P(x)} = \arg \max_y P(x,y)$ One algo: enumeration (correct, but prohibitive):

for $oldsymbol{y} \in \Sigma^L$ compute $P(oldsymbol{x},oldsymbol{y})$ Return the highest found.

0

- 4	D.	-	0		В	I	0
ST ST	Cart B	1	U	P(South y) =	.0051	.0002	ϵ
P(B y) = .35	93 .009	.003	.102	P(Africa y) =	0003	0056	F
P(I y) = .0	.369	.178	.0	P(Mountains y) =		0000	c c
P(0 y) = .6	07 .610	.779	.815	F(MOUIIIIIIIS y) =	e	.0002	e
P(end v) = .0	.013	.040	.084	P(of y) =	ϵ	.0127	.0212

• Can we predict the most likely label sequence y for a given x? arg max_y $P(y|x) = \arg \max_y \frac{P(x,y)}{P(x)} = \arg \max_y P(x,y)$ One algo: enumeration (correct, but prohibitive):

for
$$y_1 \in \Sigma$$
:
for $y_2 \in \Sigma$:
...
for $y_L \in \Sigma$:
compute $P(x, y)$
Return the highest found.

Fast idea: Greedy prediction!



Fast idea: Greedy prediction!



• $y_1 = \operatorname{arg} \max_{y \in \Sigma} \mathsf{P}(y| \texttt{start}) \mathsf{P}(x_1|y)$

Fast idea: Greedy prediction!



- $y_1 = \operatorname{arg} \max_{y \in \Sigma} \mathsf{P}(y|\operatorname{\mathtt{start}}) \mathsf{P}(x_1|y)$
- $y_2 = \operatorname{arg max}_{y \in \Sigma} \mathsf{P}(y|y_1) \mathsf{P}(x_2|y)$

Fast idea: Greedy prediction!



- $y_1 = \operatorname{arg} \max_{y \in \Sigma} \mathsf{P}(y| \texttt{start}) \mathsf{P}(x_1|y)$
- $y_2 = \operatorname{arg max}_{y \in \Sigma} \mathsf{P}(y|y_1) \mathsf{P}(x_2|y)$
- $y_3 = \arg \max_{y \in \Sigma} \mathsf{P}(y|y_2) \mathsf{P}(x_3|y) \mathsf{P}(\mathsf{end}|y)$

Fast idea: Greedy prediction!



- $y_1 = \arg \max_{y \in \Sigma} \mathsf{P}(y|\texttt{start}) \mathsf{P}(x_1|y)$
- $y_2 = \operatorname{arg max}_{y \in \Sigma} \mathsf{P}(y|y_1) \mathsf{P}(x_2|y)$
- $y_3 = \operatorname{arg\,max}_{y \in \Sigma} \mathsf{P}(y|y_2) \mathsf{P}(x_3|y) \mathsf{P}(\operatorname{end} |y)$

Is this algorithm correct? (Does it return $rg\max_y \mathsf{P}(x,y)$?)








Greedy prediction!



- The true arg max $_{m{y}}$ P $(m{x},m{y})=[m{0},m{B},m{0}]$
- We just got lucky the first time! What went wrong?

Greedy prediction!



- The true arg max $_{oldsymbol{y}}$ P $(oldsymbol{x},oldsymbol{y})=[extsf{0}, extsf{B}, extsf{0}]$
- We just got lucky the first time! What went wrong?
- We commit to a wrong label in the beginning, and can't go back.

Greedy prediction!



- The true arg max $_{oldsymbol{y}}$ P $(oldsymbol{x},oldsymbol{y})=[0, extsf{B},0]$
- We just got lucky the first time! What went wrong?
- We commit to a wrong label in the beginning, and can't go back.
- We should keep some memory of the past, to be able to undo.



$$P(x_1,...,x_i,y_1,...,y_i) = \prod_{k=1}^i P(x_k|y_k) P(y_k|y_{k-1})$$

At step *i*, let's assign a score to each possible choice for *y_i*. It will depend on the optimal labels *y*₁,..., *y_{i-1}*.

score_i(y) =
$$\max_{y_1, \dots, y_{i-1}} P(x_1, y_1, \dots, x_{i-1}, y_{i-1}, x_i, y_i = y)$$

But wait!



$$P(x_1,...,x_i,y_1,...,y_i) = \prod_{k=1}^i P(x_k|y_k) P(y_k|y_{k-1})$$

At step *i*, let's assign a score to each possible choice for *y_i*. It will depend on the optimal labels *y*₁,..., *y_{i-1}*.

score_i(y) =
$$\max_{y_1,...,y_{i-1}} P(x_1, y_1, ..., x_{i-1}, y_{i-1}, x_i, y_i = y)$$

But wait!

$$= \max_{y_1,...,y_{i-1}} \left(\mathsf{P}(x_i|y) \mathsf{P}(y|y_{i-1}) \mathsf{P}(x_1,y_1,\ldots,x_{i-1},y_{i-1}) \right)$$



$$P(x_1,...,x_i,y_1,...,y_i) = \prod_{k=1}^i P(x_k|y_k) P(y_k|y_{k-1})$$

At step *i*, let's assign a score to each possible choice for *y_i*. It will depend on the optimal labels *y*₁,..., *y_{i-1}*.

$$score_i(y) = \max_{y_1, \dots, y_{i-1}} P(x_1, y_1, \dots, x_{i-1}, y_{i-1}, x_i, y_i = y)$$

But wait!

$$= \max_{y_1, \dots, y_{i-1}} \left(\mathsf{P}(x_i | y) \mathsf{P}(y | y_{i-1}) \mathsf{P}(x_1, y_1, \dots, x_{i-1}, y_{i-1}) \right)$$

= $\mathsf{P}(x_i | y) \max_{y'} \left(\mathsf{P}(y | y') \max_{y_1, \dots, y_{i-2}} \mathsf{P}(x_1, y_1, \dots, x_{i-1}, y_{i-1} = y') \right)$



$$P(x_1,...,x_i,y_1,...,y_i) = \prod_{k=1}^i P(x_k|y_k) P(y_k|y_{k-1})$$

At step *i*, let's assign a score to each possible choice for *y_i*. It will depend on the optimal labels *y*₁,..., *y_{i-1}*.

$$score_i(y) = \max_{y_1, \dots, y_{i-1}} P(x_1, y_1, \dots, x_{i-1}, y_{i-1}, x_i, y_i = y)$$

But wait!

$$= \max_{y_1,...,y_{i-1}} \left(P(x_i|y) P(y|y_{i-1}) P(x_1, y_1, \dots, x_{i-1}, y_{i-1}) \right)$$

= $P(x_i|y) \max_{y'} \left(P(y|y') \max_{y_1,...,y_{i-2}} P(x_1, y_1, \dots, x_{i-1}, y_{i-1} = y') \right)$
= $P(x_i|y) \max_{y'} \left(P(y|y') \text{ score}_{i-1}(y') \right)$



$$P(x_1,...,x_i,y_1,...,y_i) = \prod_{k=1}^{i} P(x_k|y_k) P(y_k|y_{k-1})$$

At step *i*, let's assign a score to each possible choice for *y_i*. It will depend on the optimal labels *y*₁,..., *y_{i-1}*.

$$score_i(y) = \max_{y_1, \dots, y_{i-1}} P(x_1, y_1, \dots, x_{i-1}, y_{i-1}, x_i, y_i = y)$$

But wait!

$$= \max_{y_1, \dots, y_{i-1}} \left(P(x_i | y) P(y | y_{i-1}) P(x_1, y_1, \dots, x_{i-1}, y_{i-1}) \right)$$

= $P(x_i | y) \max_{y'} \left(P(y | y') \max_{y_1, \dots, y_{i-2}} P(x_1, y_1, \dots, x_{i-1}, y_{i-1} = y') \right)$
= $P(x_i | y) \max_{y'} \left(P(y | y') \text{ score}_{i-1}(y') \right)$

• We can compute the scores recursively!

Vlad Niculae & André Martins (IST)



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

At i = 1, the previous label can only be start:
 score1(B) = P(New|B) P(B|start)
 score1(I) = P(New|I) P(I|start)
 score1(0) = P(New|0) P(0|start)



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

• At i = 1, the previous label can only be start: $score_1(B) = P(New|B) P(B|start) = .006 \cdot .393$ $score_1(I) = P(New|I) P(I|start) = .001 \cdot .0$ $score_1(0) = P(New|0) P(0|start) = .0006 \cdot .607$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

- At *i* = 1, the previous label can only be start:
 score₁(B) = P(New|B) P(B|start) = .006 · .393
 score₁(I) = P(New|I) P(I|start) = .001 · .0
 score₁(0) = P(New|0) P(0|start) = .0006 · .607
- At *i* = 2:

$$score_{2}(B) = P(U.S.|B) \max_{y'} P(B|y') score_{1}(y')$$

= P(U.S.|B) max
$$\begin{cases} P(B|B) score_{1}(B) = .009 \cdot .002 = .00002 \\ P(B|I) score_{1}(I) = .003 \cdot .0 = .0 = .016 \cdot .00003 \\ P(B|O) score_{1}(O) = .102 \cdot .0003 = .00003 \end{cases}$$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

- At *i* = 1, the previous label can only be start:
 score₁(B) = P(New|B) P(B|start) = .006 · .393
 score₁(I) = P(New|I) P(I|start) = .001 · .0
 score₁(0) = P(New|0) P(0|start) = .0006 · .607
- At *i* = 2:

$$score_{2}(B) = P(U.S.|B) \max_{y'} P(B|y') score_{1}(y')$$

= P(U.S.|B) max
$$\begin{cases} P(B|B) score_{1}(B) = .009 \cdot .002 = .00002 \\ P(B|I) score_{1}(I) = .003 \cdot .0 = .0 = .016 \cdot .00003 \\ P(B|O) score_{1}(O) = .102 \cdot .0003 = .00003 \end{cases}$$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

- At *i* = 1, the previous label can only be start: score₁(B) = P(New|B) P(B|start) = .006 · .393 score₁(I) = P(New|I) P(I|start) = .001 · .0 score₁(0) = P(New|0) P(0|start) = .0006 · .607
- At *i* = 2:

$$score_{2}(B) = P(U.S.|B) \max_{y'} P(B|y') score_{1}(y')$$

= P(U.S.|B) max
$$\begin{cases} P(B|B) score_{1}(B) = .009 \cdot .002 = .00002 \\ P(B|I) score_{1}(I) = .003 \cdot .0 = .0 = .016 \cdot .00003 \\ P(B|O) score_{1}(O) = .102 \cdot .0003 = .00003 \end{cases}$$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

- At *i* = 1, the previous label can only be start:
 score₁(B) = P(New|B) P(B|start) = .006 · .393
 score₁(I) = P(New|I) P(I|start) = .001 · .0
 score₁(0) = P(New|0) P(0|start) = .0006 · .607
- At *i* = 2:

$$score_{2}(B) = P(U.S.|B) \max_{y'} P(B|y') score_{1}(y')$$

= P(U.S.|B) max
$$\begin{cases} P(B|B) score_{1}(B) = .009 \cdot .002 = .00002 \\ P(B|I) score_{1}(I) = .003 \cdot .0 = .0 = .016 \cdot .00003 \\ P(B|0) score_{1}(0) = .102 \cdot .0003 = .00003 \end{cases}$$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

• At the end we observe:

$$\max_{y} P(\text{end}|y) \operatorname{score}_{L}(y) = \max_{y} P(\text{end}|y) \max_{y_{1}, \dots, y_{L-1}} P(x_{1}, y_{1}, \dots, x_{L}, y_{L} = y)$$
$$= \max_{y} P(x, y)$$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

• At the end we observe:

$$\max_{y} P(\text{end}|y) \operatorname{score}_{L}(y) = \max_{y} P(\text{end}|y) \max_{y_{1}, \dots, y_{L-1}} P(x_{1}, y_{1}, \dots, x_{L}, y_{L} = y)$$
$$= \max_{y} P(x, y)$$

In this case,

$$\max \begin{cases} \mathsf{P}(\mathsf{end}|\mathsf{B}) \ \mathsf{score}_3(\mathsf{B}) = .013 \cdot 2e{\text{-}16} \\ \mathsf{P}(\mathsf{end}|\mathsf{I}) \ \mathsf{score}_3(\mathsf{I}) = .04 \cdot 2e{\text{-}14} \\ \mathsf{P}(\mathsf{end}|\mathsf{O}) \ \mathsf{score}_3(\mathsf{O}) = .084 \cdot 5e{\text{-}11} \end{cases} = 4e{\text{-}12}$$



$$score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$$

• At the end we observe:

$$\max_{y} P(\text{end}|y) \operatorname{score}_{L}(y) = \max_{y} P(\text{end}|y) \max_{y_{1}, \dots, y_{L-1}} P(x_{1}, y_{1}, \dots, x_{L}, y_{L} = y)$$
$$= \max_{y} P(x, y)$$

In this case,

$$\max \begin{cases} \mathsf{P}(\mathsf{end}|\mathsf{B}) \, \mathsf{score}_3(\mathsf{B}) = .013 \cdot 2e\text{-16} \\ \mathsf{P}(\mathsf{end}|\mathsf{I}) \, \mathsf{score}_3(\mathsf{I}) = .04 \cdot 2e\text{-14} \\ \mathsf{P}(\mathsf{end}|\mathsf{0}) \, \mathsf{score}_3(\mathsf{0}) = .084 \cdot 5e\text{-11} \end{cases} = 4e\text{-12}$$

• What is the *y* that gives this value?

Vlad Niculae & André Martins (IST)



 $score_i(y) = P(x_i|y) \max_{y'} P(y|y') score_i(y')$

• At the end we observe:

$$\max_{y} P(\text{end}|y) \operatorname{score}_{L}(y) = \max_{y} P(\text{end}|y) \max_{y_{1}, \dots, y_{L-1}} P(x_{1}, y_{1}, \dots, x_{L}, y_{L} = y)$$
$$= \max_{y} P(x, y)$$

In this case,

$$\max \begin{cases} \mathsf{P}(\mathsf{end}|\mathsf{B}) \ \mathsf{score}_3(\mathsf{B}) = .013 \cdot 2\mathsf{e}\text{-}16 \\ \mathsf{P}(\mathsf{end}|\mathsf{I}) \ \mathsf{score}_3(\mathsf{I}) = .04 \cdot 2\mathsf{e}\text{-}14 \\ \mathsf{P}(\mathsf{end}|\mathsf{O}) \ \mathsf{score}_3(\mathsf{O}) = .084 \cdot 5\mathsf{e}\text{-}11 \end{cases} = 4\mathsf{e}\text{-}12$$

• What is the y that gives this value? Backtrace remembering each max!

Vlad Niculae & André Martins (IST)

The Viterbi algorithm

input: sequence x_1, \ldots, x_L , emission and transition probabilities

Forward: compute scores recursively score₁(y) = P(y|start) · P(x₁|y) $\forall y \in \Sigma$ for i = 2 to L do for $y \in \Sigma$ do score_i(y) = P(x_i|y) · max_{y'} (P(y|y') · score_{i-1}(y')) backptr_i(y) = arg max_{y'} (P(y|y') · score_{i-1}(y'))

Backward: follow backpointers

$$p = \max_{y'} \left(P(end|y') \cdot score_L(y') \right)$$

$$\hat{y}_L = \arg \max_{y'} \left(P(end|y') \cdot score_L(y') \right)$$
for $i = L - 1$ down to 1 do
 $\hat{y}_i = backptr_{i+1}(\hat{y}_{i+1})$

output: the most likely sequence $\hat{y} = [\hat{y}_1, \dots, \hat{y}_L]$ and the joint likelihood $p = P(x, \hat{y})$

Viterbi with log-probabilities

- Notice how probabilities get tiny (1e-14) even for short sequences.
- Multiplying small numbers is not numerically robust. Fortunately,

 $u < v \iff \log u < \log v$ and $\log ab = \log a + \log b$

• We can equivalently find \hat{y} as:

$$\begin{aligned} \arg\max_{\boldsymbol{y}} \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) &= \arg\max_{\boldsymbol{y}} \mathsf{log} \, \mathsf{P}(\boldsymbol{x}, \boldsymbol{y}) \\ &= \arg\max_{\boldsymbol{y}} \mathsf{log} \left(\prod_{i=1}^{L+1} \mathsf{P}(y_i | y_{i-1}) \cdot \prod_{i=1}^{L} \mathsf{P}(x_i | y_i) \right) \\ &= \arg\max_{\boldsymbol{y}} \left(\sum_{i=1}^{L+1} \mathsf{log} \, \mathsf{P}(y_i | y_{i-1}) + \sum_{i=1}^{L} \mathsf{log} \, \mathsf{P}(x_i | y_i) \right) \end{aligned}$$

• Mathematically equivalent, but what good is mathematics when computers can't represent your numbers?

Vlad Niculae & André Martins (IST)

Lecture 5: Linear Sequential Models

Viterbi with log-probabilities (USE THIS ONE)

input: sequence *x*₁,...,*x*_{*L*}, emission and transition **log**-probabilities

Forward: compute scores recursively $\widetilde{score_1}(y) = \log P(y|\mathtt{start}) + \log P(x_1|y) \quad \forall y \in \Sigma$ for i = 2 to L do for $y \in \Sigma$ do $\widetilde{score_i}(y) = \log P(x_i|y) + \max_{y'} \left(\log P(y|y') + \widetilde{score_{i-1}}(y') \right)$ $\operatorname{backptr}_i(y) = \operatorname{arg max}_{y'} \left(\log P(y|y') + \widetilde{score_{i-1}}(y') \right)$

Backward: follow backpointers

$$l = \max_{y'} \left(\log P(end|y') + \widetilde{score}_{L}(y') \right)$$

$$\hat{y}_{L} = \arg \max_{y'} \left(\log P(end|y') + \widetilde{score}_{L}(y') \right)$$
for $i = L - 1$ down to 1 do

$$\hat{y}_{i} = \text{backptr}_{i+1}(\hat{y}_{i+1})$$

output: the most likely sequence $\hat{y} = [\hat{y}_1, \dots, \hat{y}_L]$ and the joint log-likelihood $l = \log P(x, \hat{y})$

note: $\widetilde{\text{score}}_i(y) = \log \operatorname{score}_i(y)$

IST. Fall 2019

35/89

Summing Up: Viterbi

- Computes the most likely sequence of tags: ${
 m arg\,max}_y \,{
 m P}(y|x)$
- This is called MAP (maximum a posteriori) decoding
- An instance of a dynamic programming algorithm: makes use of recurrence to reuse *partial solutions*.

Big HMM questions



- 1 What is the most likely label sequence *y*, given *x*?
- 2 What is the probability of x?
- 3 What is the probability of each assignment *y_i*, given *x*?
- What is the probability of each transition $y \rightarrow y'$, given x?
- 5 What sequence *y* minimizes the *Hamming cost?*



Treat y as unknown (missing). Marginal probability of x:

$$\mathsf{P}(x) = \sum_{oldsymbol{y}} \mathsf{P}(x, oldsymbol{y})$$



Treat y as unknown (missing). Marginal probability of x:

$$\mathsf{P}(x) = \sum_{oldsymbol{y}} \mathsf{P}(x, oldsymbol{y})$$

- Why?
 - Compare likelihood of observations $x^{(1)}, x^{(2)}$, e.g. language model, outlier detection...
 - Maximize this to learn HMM without supervision.

• Assess
$$\mathsf{P}(y|x) = rac{\mathsf{P}(x,y)}{\mathsf{P}(x)}$$
.



Treat y as unknown (missing). Marginal probability of x:

$$\mathsf{P}({m{x}}) = \sum_{{m{y}}} \mathsf{P}({m{x}},{m{y}})$$

- Why?
 - Compare likelihood of observations $x^{(1)}, x^{(2)},$ e.g. language model, outlier detection...
 - Maximize this to learn HMM without supervision.
 - Assess $\mathsf{P}(y|x) = \frac{\mathsf{P}(x,y)}{\mathsf{P}(x)}$.
- One algo: enumeration (correct, but prohibitive):

$$egin{array}{l} eta \leftarrow { extsf{0}} \ { extsf{for}} \ m{y} \in { extsf{\Sigma}}^{\scriptscriptstyle L} \end{array}$$
 for $m{y} \in { extsf{\Sigma}}^{\scriptscriptstyle L}$

$$p \leftarrow p + P(x, y)$$

Return p



- Treat $m{y}$ as unknown. Marginal probability of $m{x}$: $\mathsf{P}(m{x}) = \sum_{m{y}} \mathsf{P}(m{x},m{y})$
- Remember "what's the probability that the 3rd label is B?"



- Treat $m{y}$ as unknown. Marginal probability of $m{x}$: $\mathsf{P}(m{x}) = \sum_{m{y}} \mathsf{P}(m{x},m{y})$
- Remember "what's the probability that the 3rd label is B?"

$$\mathsf{P}(\boldsymbol{x}) = \sum_{\boldsymbol{y}} \mathsf{P}(\mathsf{end}|\boldsymbol{y}) \underbrace{\sum_{y_1, \dots, y_{L-1}} \mathsf{P}(x_1, \dots, x_L, y_1, \dots, y_{L-1}, y_L = \boldsymbol{y})}_{\alpha_L(\boldsymbol{y})}$$



- Treat y as unknown. Marginal probability of x: $\mathsf{P}(x) = \sum_{y} \mathsf{P}(x,y)$
- Remember "what's the probability that the 3rd label is B?"

$$P(\boldsymbol{x}) = \sum_{\boldsymbol{y}} P(\text{end}|\boldsymbol{y}) \underbrace{\sum_{\boldsymbol{y}_{1},\dots,\boldsymbol{y}_{L-1}} P(\boldsymbol{x}_{1},\dots,\boldsymbol{x}_{L},\boldsymbol{y}_{1},\dots,\boldsymbol{y}_{L-1},\boldsymbol{y}_{L} = \boldsymbol{y})}_{\boldsymbol{y}_{L}}$$
$$= \sum_{\boldsymbol{y}} P(\text{end}|\boldsymbol{y}) P(\boldsymbol{x}_{L}|\boldsymbol{y}) \underbrace{\sum_{\boldsymbol{y}'} P(\boldsymbol{y}|\boldsymbol{y}')}_{\boldsymbol{y}'} \underbrace{\sum_{\boldsymbol{y}_{1},\dots,\boldsymbol{y}_{L-2}} P(\boldsymbol{x}_{1},\dots,\boldsymbol{x}_{L-1},\boldsymbol{y}_{1},\dots,\boldsymbol{y}_{L-1} = \boldsymbol{y}')}_{\alpha_{L-1}(\boldsymbol{y}')}$$



- Treat y as unknown. Marginal probability of x: $\mathsf{P}(x) = \sum_{y} \mathsf{P}(x,y)$
- Remember "what's the probability that the 3rd label is B?"

$$P(x) = \sum_{y} P(end|y) \sum_{\substack{y_1, \dots, y_{L-1} \\ y_1 \in \mathbb{N}, y_1 \in \mathbb$$

• Recurrence: $\alpha_i(y) = \mathsf{P}(x_i|y) \sum_{y'} \mathsf{P}(y|y') \alpha_{i-1}(y')$

Notice a similarity?

Viterbi recurrence:

$$\operatorname{score}_i(y) = \operatorname{P}(x_i|y) \max_{y'} \operatorname{P}(y|y') \operatorname{score}_{i-1}(y')$$

Marginalization recurrence:

$$\alpha_i(y) = \mathsf{P}(x_i|y) \sum_{y'} \mathsf{P}(y|y') \, \alpha_{i-1}(y')$$





• At i = 1, the previous label can only be start:

$$\begin{aligned} &\alpha_1(B) = \mathsf{P}(\mathsf{New}|B) \, \mathsf{P}(B|\mathsf{start}) \\ &\alpha_1(I) = \mathsf{P}(\mathsf{New}|I) \, \mathsf{P}(I|\mathsf{start}) \\ &\alpha_1(0) = \mathsf{P}(\mathsf{New}|0) \, \mathsf{P}(0|\mathsf{start}) \end{aligned}$$



$$\alpha_i(y) = \mathsf{P}(x_i|y) \sum_{y_{i'}} \mathsf{P}(y|y_{i'}) \, \alpha_i(y_{i'})$$

• At i = 1, the previous label can only be start: $\alpha_1(B) = P(New|B) P(B|start) = .006 \cdot .393$ $\alpha_1(I) = P(New|I) P(I|start) = .001 \cdot .0$ $\alpha_1(0) = P(New|0) P(0|start) = .0006 \cdot .607$



Lecture 5: Linear Sequential Models




$$\alpha_{2}(B) = P(U.S.|B) \sum_{y'} P(B|y') \alpha_{1}(y')$$

= P(U.S.|B)sum
$$\begin{cases} P(B|B) \alpha_{1}(B) = .009 \cdot .002 = .00002 \\ P(B|I) \alpha_{1}(I) = .003 \cdot .0 = .0 = .016 \cdot .00005 \\ P(B|0) \alpha_{1}(0) = .102 \cdot .0003 = .00003 \end{cases}$$



$$\alpha_i(y) = \mathsf{P}(x_i|y) \sum_{y_{i'}} \mathsf{P}(y|y_{i'}) \, \alpha_i(y_{i'})$$

• At i = 1, the previous label can only be start: $\alpha_1(B) = P(New|B) P(B|start) = .006 \cdot .393$ $\alpha_1(I) = P(New|I) P(I|start) = .001 \cdot .0$ $\alpha_1(0) = P(New|0) P(0|start) = .0006 \cdot .607$ • At i = 2:

$$\begin{aligned} \alpha_{2}(B) &= \mathsf{P}(\mathsf{U.S.}|B) \sum_{y'} \mathsf{P}(B|y') \, \alpha_{1}(y') \\ &= \mathsf{P}(\mathsf{U.S.}|B) \mathsf{sum} \begin{cases} \mathsf{P}(B|B) \, \alpha_{1}(B) = .009 \cdot .002 &= .00002 \\ \mathsf{P}(B|I) \, \alpha_{1}(I) = .003 \cdot .0 &= .0 \\ \mathsf{P}(B|0) \, \alpha_{1}(0) = .102 \cdot .0003 &= .00003 \end{cases} \end{aligned}$$



$$\alpha_i(y) = \mathsf{P}(x_i|y) \sum_{y_{i'}} \mathsf{P}(y|y_{i'}) \, \alpha_i(y_{i'})$$

• At the end we can consider:

$$P(x) = \sum_{y} P(\text{end}|y) \alpha_L(y)$$



$$\alpha_i(y) = \mathsf{P}(x_i|y) \sum_{y_{i'}} \mathsf{P}(y|y_{i'}) \, \alpha_i(y_{i'})$$

• At the end we can consider:

$$P(\boldsymbol{x}) = \sum_{y} P(\texttt{end}|y) \, lpha_L(y)$$

In this case,

$$\sup \begin{cases} \mathsf{P}(\mathsf{end}|\mathsf{B}) \, \alpha_3(\mathsf{B}) = .013 \cdot 2\mathsf{e}\text{-}16 \\ \mathsf{P}(\mathsf{end}|\mathsf{I}) \, \alpha_3(\mathsf{I}) = .04 \cdot 2\mathsf{e}\text{-}14 \\ \mathsf{P}(\mathsf{end}|\mathsf{O}) \, \alpha_3(\mathsf{O}) = .084 \cdot 5\mathsf{e}\text{-}11 \end{cases} = 4\mathsf{e}\text{-}12$$

The Forward algorithm

input: sequence x_1, \ldots, x_l , emission and transition probabilities

Forward: compute
$$\alpha$$
 recursively
 $\alpha_1(y) = P(y|\text{start}) \cdot P(x_1|y) \quad \forall y \in \Sigma$
for $i = 2$ to L do
for $y \in \Sigma$ do
 $\alpha_i(y) = P(x_i|y) \cdot \sum_{y'} \left(P(y|y') \cdot \alpha_{i-1}(y') \right)$
 $p = \sum_{y'} \left(P(\text{end}|y') \cdot \alpha_L(y') \right)$

output: the marginal likelihood p = P(x)



$$\mathsf{P}(y_i = y | \boldsymbol{x}) = \frac{\mathsf{P}(y_i = y, \boldsymbol{x})}{\sum_{y'} \mathsf{P}(y_i = y', \boldsymbol{x})} = ?$$

- y_1, \ldots, y_{i-1} (in turn depending only on x_1, \ldots, x_{i-1})
- y_{i+1}, \ldots, y_L (in turn depending only on x_{i+1}, \ldots, x_L)



$$P(y_i = y | x) = \frac{P(y_i = y, x)}{\sum_{y'} P(y_i = y', x)} =?$$

- y_1, \ldots, y_{i-1} (in turn depending only on x_1, \ldots, x_{i-1})
- y_{i+1}, \ldots, y_L (in turn depending only on x_{i+1}, \ldots, x_L)

$$\mathsf{P}(y_{i} = y, x) = \sum_{y_{1}, \dots, y_{i-1}, y_{i+1}, \dots, y_{L}} \mathsf{P}(y_{1}, \dots, y_{i-1}, y_{i} = y, y_{i+1}, \dots, y_{L}, x)$$



$$\mathsf{P}(y_i = y | \boldsymbol{x}) = \frac{\mathsf{P}(y_i = y, \boldsymbol{x})}{\sum_{y'} \mathsf{P}(y_i = y', \boldsymbol{x})} = ?$$

- y_1, \ldots, y_{i-1} (in turn depending only on x_1, \ldots, x_{i-1})
- y_{i+1}, \ldots, y_L (in turn depending only on x_{i+1}, \ldots, x_L)

$$P(y_i = y, x) = \sum_{y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_L} P(y_1, \dots, y_{i-1}, y_i = y, y_{i+1}, \dots, y_L, x)$$

=
$$\sum_{y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_L} P(y_1, \dots, y_{i-1}, y_i = y, x_1, \dots, x_i) \cdot P(y_{i+1}, \dots, y_L, x_{i+1}, \dots, x_L | y_i = y)$$



$$\mathsf{P}(y_i = y | \boldsymbol{x}) = \frac{\mathsf{P}(y_i = y, \boldsymbol{x})}{\sum_{y'} \mathsf{P}(y_i = y', \boldsymbol{x})} = ?$$

- y_1, \ldots, y_{i-1} (in turn depending only on x_1, \ldots, x_{i-1})
- y_{i+1}, \ldots, y_L (in turn depending only on x_{i+1}, \ldots, x_L)

$$P(y_{i} = y, \boldsymbol{x}) = \sum_{y_{1}, \dots, y_{i-1}, y_{i+1}, \dots, y_{L}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, y_{i+1}, \dots, y_{L}, \boldsymbol{x})$$

$$= \sum_{y_{1}, \dots, y_{i-1}, y_{i+1}, \dots, y_{L}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i}) \cdot P(y_{i+1}, \dots, y_{L}, x_{i+1}, \dots, x_{L} | y_{i} = y)$$

$$= \underbrace{\sum_{y_{1}, \dots, y_{i-1}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i})}_{\alpha_{i}(y)} \cdot \sum_{y_{i+1}, \dots, y_{L}} P(y_{i+1}, \dots, y_{L}, x_{i+1}, \dots, x_{L} | y_{i} = y)}$$



$$\mathsf{P}(y_i = y | \boldsymbol{x}) = \frac{\mathsf{P}(y_i = y, \boldsymbol{x})}{\sum_{y'} \mathsf{P}(y_i = y', \boldsymbol{x})} = ?$$

- y_1, \ldots, y_{i-1} (in turn depending only on x_1, \ldots, x_{i-1})
- y_{i+1}, \ldots, y_L (in turn depending only on x_{i+1}, \ldots, x_L)

$$P(y_{i} = y, x) = \sum_{y_{1}, \dots, y_{i-1}, y_{i+1}, \dots, y_{L}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, y_{i+1}, \dots, y_{L}, x)$$

$$= \sum_{y_{1}, \dots, y_{i-1}, y_{i+1}, \dots, y_{L}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i}) \cdot P(y_{i+1}, \dots, y_{L}, x_{i+1}, \dots, x_{L} | y_{i} = y)$$

$$= \underbrace{\sum_{y_{1}, \dots, y_{i-1}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i})}_{\alpha_{i}(y)} \cdot \sum_{y_{i+1}, \dots, y_{L}} P(y_{i+1}, \dots, y_{L}, x_{i+1}, \dots, x_{L} | y_{i} = y)}_{\beta_{i}(y)}$$

Forward in reverse?!



• β looks a lot like running forward in the other direction...

$$P(x) = \sum_{y} P(y|\text{start}) P(x_1|y) \sum_{\substack{y_2, \dots, y_L \\ \beta_1(y)}} P(x_2, \dots, x_L, y_2, \dots, y_L|y_1 = y)$$

$$= \sum_{y} P(y|\text{start}) P(x_1|y) \sum_{y'} P(y'|y) P(x_2|y') \sum_{\substack{y_3, \dots, y_L \\ \beta_2(y')}} P(x_3, \dots, x_L, y_3, \dots, y_L|y_2 = y')$$

• Recurrence:
$$\beta_i(y) = \sum_{y'} \mathsf{P}(y'|y) \mathsf{P}(x_{i+1}|y') \beta_{i+1}(y')$$



$$(y_{i-1}) \xrightarrow{(y_{i+1})} \xrightarrow{(y_{i+1})} \xrightarrow{(y_{i+2})} \cdots$$

$$(x_{i-1}) \xrightarrow{(x_i)} \xrightarrow{(x_{i+1})} \xrightarrow{(x_{i+2})} P(y_i = y, y_{i+1} = y' | x) = \frac{P(y_i = y, y_{i+1} = y' x)}{P(x)}$$

$$\mathsf{P}(y_i = y, \boldsymbol{x}) = \sum_{\underbrace{y_1, \dots, y_{i-1}}} \mathsf{P}(y_1, \dots, y_{i-1}, y_i = y, x_1, \dots, x_i) \cdot \sum_{\underbrace{y_{i+1}, \dots, y_L}} \mathsf{P}(y_{i+1}, \dots, y_L, x_{i+1}, \dots, x_L | y_i = y)}_{\alpha_i(y)}$$

$$(y_{i-1}) \xrightarrow{(y_{i+1})} (y_{i+1}) \xrightarrow{(y_{i+2})} \cdots$$

$$(x_{i-1}) \xrightarrow{(x_i)} (x_{i+1}) \xrightarrow{(x_{i+2})} P(y_i = y, y_{i+1} = y' | x) = \frac{P(y_i = y, y_{i+1} = y' x)}{P(x)}$$

$$P(y_{i} = y, x) = \sum_{\substack{y_{1}, \dots, y_{i-1} \\ y_{i+1} = y'x}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i}) \cdot \sum_{\substack{y_{i+1}, \dots, y_{L} \\ y_{i+1}, \dots, y_{L}}} P(y_{i+1}, \dots, y_{L}, x_{i+1}, \dots, x_{L} | y_{i} = y)}$$

$$P(y_{i} = y, y_{i+1} = y'x) = \sum_{\substack{y_{1}, \dots, y_{i-1} \\ y_{i+2}, \dots, y_{L} \\ y_{i+2}, \dots, y_{L}}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i}) \cdot P(y' | y) \cdot P(x_{i+1} | y')$$

$$\sum_{\substack{y_{i+2}, \dots, y_{L} \\ \beta_{i+1}(y')}} P(y_{i+2}, \dots, y_{L}, x_{i+2}, \dots, x_{L} | y_{i+1} = y')}$$

$$(y_{i-1}) \xrightarrow{(y_{i+1})} (y_{i+1}) \xrightarrow{(y_{i+2})} \cdots$$

$$(x_{i-1}) \xrightarrow{(x_i)} (x_{i+1}) \xrightarrow{(x_{i+2})} P(y_i = y, y_{i+1} = y' | x) = \frac{P(y_i = y, y_{i+1} = y' x)}{P(x)}$$

$$P(y_{i} = y, x) = \sum_{\substack{y_{1}, \dots, y_{i-1} \\ (y_{i} = y, x_{i-1}) \\ (y_{i} = y, y_{i+1} = y'x)} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, y_{L}) P(y_{i+1}, \dots, y_{L}, x_{i+1}, \dots, x_{L}|y_{i} = y)$$

$$P(y_{i} = y, y_{i+1} = y'x) = \sum_{\substack{y_{1}, \dots, y_{i-1} \\ (y_{i} = y, y_{i-1}) \\ (y_{i+1}, \dots, y_{L}) \\ (y_{i+2}, \dots, y_{L}) P(y_{i+2}, \dots, y_{L}, x_{i+2}, \dots, x_{L}|y_{i+1} = y')}$$

$$P(y_{i} = y, y_{i+1} = y'x) = \sum_{\substack{y_{1}, \dots, y_{i-1} \\ (y_{i} = y, y_{i+1}) \\ (y_{i+2}, \dots, y_{L}) P(y_{i+2}, \dots, y_{L}, x_{i+2}, \dots, x_{L}|y_{i+1} = y')}$$

Independent of position (i, i + 1):

$$\mathsf{P}([y,y'],\boldsymbol{x}) =$$

Vlad Niculae & André Martins (IST)

$$(y_{i-1}) \xrightarrow{(y_{i+1})} (y_{i+1}) \xrightarrow{(y_{i+2})} \cdots$$

$$(x_{i-1}) \xrightarrow{(x_i)} (x_{i+1}) \xrightarrow{(x_{i+2})} P(y_i = y, y_{i+1} = y' | x) = \frac{P(y_i = y, y_{i+1} = y' x)}{P(x)}$$

$$P(y_{i} = y, x) = \sum_{\substack{y_{1}, \dots, y_{i-1} \\ \alpha_{i}(y)}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i}) \cdot \sum_{\substack{y_{i+1}, \dots, y_{L} \\ \beta_{i}(y)}} P(y_{i+1}, \dots, y_{L}, x_{i+1}, \dots, x_{L} | y_{i} = y)}{\beta_{i}(y)}$$

$$P(y_{i} = y, y_{i+1} = y'x) = \sum_{\substack{y_{1}, \dots, y_{i-1} \\ \alpha_{i}(y) \\ \vdots \\ y_{i+2}, \dots, y_{L}}} P(y_{1}, \dots, y_{i-1}, y_{i} = y, x_{1}, \dots, x_{i}) \cdot P(y' | y) \cdot P(x_{i+1} | y')$$

$$\sum_{\substack{y_{i+2}, \dots, y_{L} \\ \beta_{i+1}(y')}} P(y_{i+2}, \dots, y_{L}, x_{i+2}, \dots, x_{L} | y_{i+1} = y')}{\beta_{i+1}(y')}$$

Independent of position (i, i + 1):

$$\mathsf{P}([y,y'], \boldsymbol{x}) = \sum_{i=1}^{L-1} \mathsf{P}(y_i = y, y_{i+1} = y', \boldsymbol{x}) = \mathsf{P}(y'|y) \sum_{i=1}^{L-1} \alpha_i(y) \,\mathsf{P}(x_{i+1}|y') \,\beta_{i+1}(y')$$

Vlad Niculae & André Martins (IST)

The Forward-Backward algorithm

input: sequence x_1, \ldots, x_L , emission and transition probabilities

Forward: compute
$$\alpha$$
 recursively
 $\alpha_1(y) = P(y|\text{start}) \cdot P(x_1|y) \quad \forall y \in \Sigma$
for $i = 2$ to L do
for $y \in \Sigma$ do
 $\alpha_i(y) = P(x_i|y) \sum_{y'} \left(P(y|y') \cdot \alpha_{i-1}(y') \right)$

Forward: compute
$$\beta$$
 recursively
 $\beta_L(y) = P(end|y) \quad \forall y \in \Sigma$
for $i = L - 1$ down to 1 do
for $y \in \Sigma$ do
 $\beta_i(y) = \sum_{y'} \left(P(y'|y) \cdot P(x_{i+1}|y') \cdot \beta_{i+1}(y') \right)$

output: The marginal likelihood $P(x) = \sum_{y'} \alpha_i(y') \beta_i(y')$ for any *i*; Posterior unigram marginal probas $P(y_i = y | x) = \frac{\alpha_i(y) \beta_i(y)}{P(x)}$; Posterior transition marginal probas $P([y, y'] | x) = \frac{P(y'|y)}{P(x)} \sum_{i=1}^{L-1} \alpha_i(y) P(x_{i+1}|y') \beta_{i+1}(y')$. Viterbi recurrence:

$$\operatorname{score}_i(y) = \operatorname{P}(x_i|y) \max_{y'} \operatorname{P}(y|y') \operatorname{score}_{i-1}(y')$$

Forward recurrence:

$$\alpha_i(y) = \mathsf{P}(x_i|y) \sum_{y'} \mathsf{P}(y|y') \, \alpha_{i-1}(y')$$

Backward recurrence:

$$\beta_i(y) = \sum_{y'} \mathsf{P}(y'|y) \, \mathsf{P}(x_{i+1}|y') \, \beta_{i+1}(y')$$

In log-domain:

Viterbi recurrence:

$$\widetilde{\text{score}}_i(y) = \log P(x_i|y) + \max_{y'} \left(\log P(y|y') + \widetilde{\text{score}}_{i-1}(y') \right)$$

Forward recurrence:

$$\widetilde{\alpha}_i(y) = \log \mathsf{P}(x_i|y) + \log \sum_{y'} \exp\left(\log \mathsf{P}(y|y') + \widetilde{\alpha}_{i-1}(y')\right)$$

Backward recurrence:

$$\widetilde{\beta}_i(y) = \log \sum_{y'} \exp\left(\log \mathsf{P}(y'|y) + \log \mathsf{P}(x_{i+1}|y') + \widetilde{\beta}_{i+1}(y')\right)$$

note: $\widetilde{\alpha}_i(y) = \log \alpha_i(y)$, etc.

Log-sum-exp



The Forward-Backward algorithm in log-domain

input: sequence x_1, \ldots, x_L , emission and transition **log**-probabilities

Forward: compute
$$\tilde{\alpha}$$
 recursively
 $\tilde{\alpha}_{1}(y) = \log P(y|\text{start}) + \log P(x_{1}|y) \quad \forall y \in \Sigma$
for $i = 2$ to L do
for $y \in \Sigma$ do
 $\tilde{\alpha}_{i}(y) = \log P(x_{i}|y) + \log \sum_{y'} \exp\left(\log P(y|y') + \tilde{\alpha}_{i-1}(y')\right)$
Backward: compute $\tilde{\beta}$ recursively
 $\tilde{\beta}_{L}(y) = \log P(\text{end}|y) \quad \forall y \in \Sigma$
for $i = L - 1$ down to 1 do
for $y \in \Sigma$ do
 $\tilde{\beta}_{i}(y) = \log \sum_{y'} \exp\left(\log P(y'|y) + \log P(x_{i+1}|y') + \tilde{\beta}_{i+1}(y')\right)$
output: The log-marginal log $P(x) = \log \sum_{y'} \exp\left(\tilde{\alpha}_{i}(y') + \tilde{\beta}_{i}(y')\right)$ for any i ;
Posterior marginal log-probas: unigram
 $\log P(y_{i} = y|x) = \tilde{\alpha}_{i}(y) + \tilde{\beta}_{i}(y) - \log P(x)$;
and transition:
 $\log P([y, y']|x) = \log P(y'|y) + \log \sum_{i=1}^{L-1} \exp\left(\tilde{\alpha}_{i}(y) + \log P(x_{i+1}|y') + \tilde{\beta}_{i+1}(y')\right) - \log P(x)$.
 $\operatorname{note:} \tilde{\alpha}_{i}(y) = \log \alpha_{i}(y), \tilde{\beta}_{i}(y) = \log \beta_{i}(y)$

 $\log \beta_i(y)$

Big HMM questions



- 1 What is the most likely label sequence *y*, given *x*?
- 2 What is the probability of x?
- 3 What is the probability of each assignment *y_i*, given *x*?
- What is the probability of each transition $y \rightarrow y'$, given x?
- 5 What sequence *y* minimizes the *Hamming cost?*

Minimizing costs, a.k.a., risks

- Our HMM defines a distribution over labelings $\mathsf{P}(y|x)$.
- The HMM will be given a new sequence x to label, producing \hat{y} .
- We then observe the true y^{\star} , and incur a cost

 $\mathsf{cost}(\hat{m{y}}, m{y}^{\star}).$

How do we predict so as to minimize the expected cost

$$\hat{oldsymbol{y}} = \mathop{{\mathsf{arg\,min}}}\limits_{oldsymbol{y}} \mathbb{E}_{\mathsf{P}(oldsymbol{y}|oldsymbol{x})} \Big[\operatorname{\mathsf{cost}}(oldsymbol{y},oldsymbol{y}^{\star}) \Big]$$

Minimizing costs, a.k.a., risks

• Consider the sequence zero-one cost:

$$\mathsf{cost}_{01}(oldsymbol{y}, \hat{oldsymbol{y}}) = egin{cases} 1, & oldsymbol{y}
eq \hat{oldsymbol{y}} \ 0, & oldsymbol{y} = \hat{oldsymbol{y}} \end{cases}$$

The cost we may expect to pay is

$$\mathbb{E}_{\mathsf{P}(\boldsymbol{y}|\boldsymbol{x})}\Big[\operatorname{cost}_{01}(\boldsymbol{y}, \boldsymbol{y}^{\star})\Big] = \sum_{\boldsymbol{y} \neq \boldsymbol{y}^{\star}} \mathsf{P}(\boldsymbol{y}|\boldsymbol{x})$$

= 1 - P($\boldsymbol{y}^{\star}|\boldsymbol{x}$)

• We should return the sequence it assigns most probability to:

$$\hat{oldsymbol{y}} = rgmin_{oldsymbol{y}} \left(1 - \mathsf{P}(oldsymbol{y}|oldsymbol{x})
ight) = rgmax_{oldsymbol{y}} \mathsf{P}(oldsymbol{y}|oldsymbol{x})$$

Viterbi computes this sequence!

Minimizing costs, a.k.a., risks

• Now consider the Hamming (word-wise) cost:

$$ext{cost}_{H}(oldsymbol{y}, \hat{oldsymbol{y}}) = \sum_{i} ext{cost}_{01}(y_i, \hat{y}_i)$$

The expected cost we pay is

$$\mathbb{E}_{\mathsf{P}(\boldsymbol{y}|\boldsymbol{x})}\Big[\operatorname{cost}_{H}(\boldsymbol{y}, \boldsymbol{y}^{\star})\Big] = \sum_{i} \sum_{y} \operatorname{cost}_{01}(y, \hat{y}_{i}) \operatorname{\mathsf{P}}(y_{i} = y|\boldsymbol{x})$$
$$= \sum_{i} \left(1 - \operatorname{\mathsf{P}}(y_{i} = y_{i}^{\star}|\boldsymbol{x})\right)$$

• Posterior decoding: Get α, β (Forward-Backward); pick \hat{y} such that

$$\hat{y}_i = rg\max_{y} \mathsf{P}(y_i = y | \boldsymbol{x}) = rg\max_{y} \left(\widetilde{lpha}_i(y) + \widetilde{eta}_i(y) \right)$$

• Exercise: This can be extended for any cost of the form $cost(\boldsymbol{y}, \hat{\boldsymbol{y}}) = \sum_{i=1}^{L} c_u(y_i, \hat{y}_i) + \sum_{i=1}^{L-1} c_t(y_i, y_{i+1}, \hat{y}_i, \hat{y}_{i+1})$

Vlad Niculae & André Martins (IST)

Big HMM questions



- 1 What is the most likely label sequence *y*, given *x*?
- 2 What is the probability of x?
- 3 What is the probability of each assignment *y_i*, given *x*?
- What is the probability of each transition $y \rightarrow y'$, given x?
- 5 What sequence *y* minimizes the *Hamming cost?*

Two important algorithms for sequential models

All of these big questions are solved by these two similar algorithms.



Vlad Niculae & André Martins (IST)

- So far y was a sequence of labels, and x a sequence of words, so emission probabilities P(x_i|y_i) are just tables.
- Everything works with other choices for $P(x_i|y_i)$. Examples:

- So far y was a sequence of labels, and x a sequence of words, so emission probabilities P(x_i|y_i) are just tables.
- Everything works with other choices for $P(x_i|y_i)$. Examples:
 - x_i are sentences. Example: sentence-level review sentiment. x

 x_1 : I bought this knife set last year. y_1 : neutral x_2 : I was pleasantly surprised with it. y_2 : positive x_3 : They're still sharp. y_3 : positive

Naïve Bayes assumption for emissions: $P(x_i|y_i) = \prod_{w \in x_i} P(w|y_i)$

- So far y was a sequence of labels, and x a sequence of words, so emission probabilities $P(x_i|y_i)$ are just tables.
- Everything works with other choices for $P(x_i|y_i)$. Examples:
 - x_i are sentences. Example: sentence-level review sentiment. \boldsymbol{x}
 - x_1 : I bought this knife set last year. y_1 : neutral x_2 : I was pleasantly surprised with it. y_2 : positive x₃: They're still sharp. y₃: positive

Naïve Bayes assumption for emissions: $P(x_i|y_i) = \prod_{w \in x_i} P(w|y_i)$



Gaussian emissions: $\mathsf{P}(\mathbf{x}_i|\mathbf{y}_i) = \mathcal{N}(\mu_{\mathbf{y}_i}, \mathbf{S}_{\mathbf{y}_i})$

- So far y was a sequence of labels, and x a sequence of words, so emission probabilities $P(x_i|y_i)$ are just tables.
- Everything works with other choices for $P(x_i|y_i)$. Examples:
 - x_i are sentences. Example: sentence-level review sentiment. \boldsymbol{x}
 - x_1 : I bought this knife set last year. y_1 : neutral x_2 : I was pleasantly surprised with it. y_2 : positive x₃: They're still sharp. y₃: positive

Naïve Bayes assumption for emissions: $P(x_i|y_i) = \prod_{w \in x_i} P(w|y_i)$



Gaussian emissions: $\mathsf{P}(\mathbf{x}_i|\mathbf{y}_i) = \mathcal{N}(\mu_{\mathbf{y}_i}, \mathbf{S}_{\mathbf{y}_i})$

Or, turn to feature-based discriminative models (later.)

4000 3000 requency (Hz)

2000

1000

Outline

O Structured Prediction

2 Generative Sequence Models

Markov Models

Hidden Markov Models

Unsupervised learning

B Discriminative Sequence Models

Structured Perceptron

Conditional Random Fields

Structured SVM

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

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

Assumes all we have is *unlabeled* training data $\{x^{(1)}, \dots, x^{(N)}\}$ Maximum Likelihood Estimation with incomplete data!

maximize
$$\prod_{n=1}^{N} \mathsf{P}(\boldsymbol{x}^{(n)}) = \prod_{n=1}^{N} \sum_{\boldsymbol{y}} \mathsf{P}(\boldsymbol{x}^{(n)}, \boldsymbol{y}^{(n)})$$

Algorithm: Expectation-Maximization (EM).

Expectation-Maximization (Baum-Welch)

- If we knew y, could fit HMM by counting and normalizing
- If we knew the model parameters, we could estimate the posterior marginal probabilities (soft counts) $P(y_i | x)$ and $P(y_{i-1}, y_i | x)$
- This is a chicken-and-egg problem!
Initialize HMM at random.

Alternate:

- E-step: Get a soft-labelling of x from current model, keep track of soft counts P(y_i | x) and P(y_{i-1}, y_i | x)
 - Forward-Backward for each data point.
- M-step: Do a "supervised" update of the HMM
 - Count & normalize the soft-labels!

Guarantees improvement, but the problem has multiple optima.

Summary of HMMs

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

Binary/Multi-class	Structured Prediction
Naive Bayes	HMMs
Perceptron	?
Logistic Regression	?
SVMs	?

Outline

Structured Prediction

Ø Generative Sequence Models

Markov Models Hidden Markov Mode

Unsupervised learning

3 Discriminative Sequence Models

Structured Perceptron Conditional Random Fields Structured SVM

Generative models

- Model the joint $\mathsf{P}(x,y)$ by choosing $\mathsf{P}(y),\mathsf{P}(x|y)$
 - E.g., Naïve Bayes
 - Easy to fit: just count, one pass over the data.
- Make predictions: ${
 m arg\,max}_y \,{
 m P}(y|x) = {
 m arg\,max}_y \,{
 m P}(x,y)$
- Use Bayes' rule to get $\mathsf{P}(y|x)$

Generative models

- Model the joint $\mathsf{P}(x,y)$ by choosing $\mathsf{P}(y),\mathsf{P}(x|y)$
 - E.g., Naïve Bayes
 - Easy to fit: just count, one pass over the data.
- Make predictions: ${
 m arg\,max}_{m y}\,{
 m P}(m y|m x)={
 m arg\,max}_{m y}\,{
 m P}(x,m y)$
- Use Bayes' rule to get $\mathsf{P}(y|x)$

Discriminative models

- Directly try to assign a **score** to each class
 - linear, feature-driven: $s_{m{y}} = m{w} \cdot \phi(m{x},m{y})$,
 - or a neural network: $s_y = f(y; x)$

Generative models

- Model the joint $\mathsf{P}(x,y)$ by choosing $\mathsf{P}(y),\mathsf{P}(x|y)$
 - E.g., Naïve Bayes
 - Easy to fit: just count, one pass over the data.
- Make predictions: ${
 m arg\,max}_{m y}\,{
 m P}(m y|m x)={
 m arg\,max}_{m y}\,{
 m P}(x,m y)$
- Use Bayes' rule to get $\mathsf{P}(y|x)$

Discriminative models

- Directly try to assign a score to each class
 - linear, feature-driven: $s_{m{y}} = m{w} \cdot \phi(m{x},m{y}),$
 - or a neural network: $s_y = f(y; x)$
- Make predictions: arg max_y s_y

Generative models

- Model the joint $\mathsf{P}(x,y)$ by choosing $\mathsf{P}(y),\mathsf{P}(x|y)$
 - E.g., Naïve Bayes
 - Easy to fit: just count, one pass over the data.
- Make predictions: ${
 m arg\,max}_y \,{
 m P}(y|x) = {
 m arg\,max}_y \,{
 m P}(x,y)$
- Use Bayes' rule to get $\mathsf{P}(y|x)$

Discriminative models

- Directly try to assign a score to each class
 - linear, feature-driven: $s_{m{y}} = m{w} \cdot \phi(m{x},m{y})$,
 - or a neural network: $s_y = f(y; x)$
- Make predictions: arg max_y s_y
- Harder to train, needs iterative optimization
 - Perceptron: Try to make $s_{y^{true}} > s_{y}$
 - Logistic regression: Define $P(y|x) \propto \exp(s_y)$, maximize $P(y^{true}|x)$.
 - SVM: Try to make $s_{y^{true}} > 1 + s_{y}$

Generative models

- Model the joint P(x, y) by choosing P(y), P(x|y)
 - E.g., Naïve Bayes
 - Easy to fit: just count, one pass over the data.
- Make predictions: arg max_y P(y|x) = arg max_y P(x, y)
 Use Bayes' rule to get P(y|x)

Discriminative models

- Directly try to assign a score to each class
 - linear, feature-driven: $s_u = w \cdot \phi(x, y)$,
 - or a neural network: $s_{y} = f(y; x)$
- Make predictions: $\arg \max_{u} s_{u}$
- Harder to train, needs iterative optimization
 - Perceptron: Try to make $s_{u^{true}} > s_{u}$
 - Logistic regression: Define $P(y|x) \propto \exp(s_y)$, maximize $P(y^{true}|x)$.
 - SVM: Try to make $s_{u^{true}} > 1 + s_{u}$

In HMM, with Viterbi resp. Forward-Backward

Does the same trick work?

Outline

Structured Prediction

Ø Generative Sequence Models

Markov Models

Hidden Markov Models

Unsupervised learning

Oiscriminative Sequence Models

Structured Perceptron

Conditional Random Fields

Structured SVM

Recall the simple & powerful Perceptron algorithm.

- Process one pair (x, y^{true}) at each round
 - Take x; predict a sequence \widehat{y} .
 - If prediction is correct, proceed. If not, adjust.

input: labeled data ${\mathcal D}$

initialize $m{w}$

repeat

get new training example (x, y^{true}) predict $\widehat{y} = \arg \max_{y \in \mathcal{Y}} f(y; x)$ if $\widehat{y} \neq y^{true}$ then update $w \leftarrow w + \nabla_w f(y^{true}; x) - \nabla_w f(\widehat{y}; x)$ until max. epochs output: model weights w

Recall the simple & powerful Perceptron algorithm.

- Process one pair (x, y^{true}) at each round
 - Take x; predict a sequence \widehat{y} .
 - If prediction is correct, proceed. If not, adjust.

input: labeled data ${\mathcal D}$

initialize $m{w}$

repeat

get new training example (x, y^{true}) predict $\widehat{y} = \arg \max_{y \in \mathcal{Y}} f(y; x)$ if $\widehat{y} \neq y^{true}$ then update $w \leftarrow w + \nabla_w f(y^{true}; x) - \nabla_w f(\widehat{y}; x)$ until max. epochs output: model weights w

Recall the simple & powerful Perceptron algorithm.

- Process one pair (x, y^{true}) at each round
 - Take x; predict a sequence \widehat{y} .
 - If prediction is correct, proceed. If not, adjust.

input: labeled data ${\mathcal D}$

initialize $m{w}$

repeat

get new training example (x, y^{true}) predict $\widehat{y} = \arg \max_{y \in \mathcal{Y}} f(y; x)$ if $\widehat{y} \neq y^{true}$ then update $w \leftarrow w + \nabla_w f(y^{true}; x) - \nabla_w f(\widehat{y}; x)$ until max. epochs output: model weights w

Why couldn't y be an entire sequence here?

Recall the simple & powerful Perceptron algorithm.

- Process one pair (x, y^{true}) at each round
 - Take x; predict a sequence \widehat{y} .
 - If prediction is correct, proceed. If not, adjust.

input: labeled data ${\mathcal D}$

initialize $m{w}$

repeat

get new training example (x, y^{true}) predict $\widehat{y} = \arg \max_{y \in \mathcal{Y}} f(y; x)$ if $\widehat{y} \neq y^{true}$ then update $w \leftarrow w + \nabla_w f(y^{true}; x) - \nabla_w f(\widehat{y}; x)$ until max. epochs output: model weights w

Why couldn't y be an entire sequence here? Mathematically, all is cool. Algorithmically...

Vlad Niculae & André Martins (IST)

Lecture 5: Linear Sequential Models

Recall the simple & powerful Perceptron algorithm.

- Process one pair $(x,y^{ ext{true}})$ at each round
 - Take x; predict a sequence \widehat{y} .
 - If prediction is correct, proceed. If not, adjust.

input: labeled data ${\mathcal D}$

initialize $m{w}$

repeat

get new training example (x, y^{true}) predict $\widehat{y} = \arg \max_{y \in \mathcal{Y}} f(y; x)$ if $\widehat{y} \neq y^{true}$ then update $w \leftarrow w + \nabla_w f(y^{true}; x) - \nabla_w f(\widehat{y}; x)$ until max. epochs output: model weights w

Why couldn't y be an entire sequence here? Mathematically, all is cool. Algorithmically...

Vlad Niculae & André Martins (IST)

Lecture 5: Linear Sequential Models

Everything would be fine if we had a way to calculate

 $rg\max_{oldsymbol{y}\in \mathbb{Y}} f(oldsymbol{y};oldsymbol{x})$

This would work, but prohibitive:.

 $\begin{array}{l} \textbf{for } \boldsymbol{y} = [\texttt{start}, y_1, \dots, y_L, \texttt{end}] \in \boldsymbol{\mathcal{Y}} \ \textbf{do} \\ \text{compute } f(\boldsymbol{y}; \boldsymbol{x}) & \# \text{ e.g. } \boldsymbol{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}), \text{ or nnet forward pass} \\ \text{return the highest scoring } \boldsymbol{y} \end{array}$

Everything would be fine if we had a way to calculate

 $rg\max_{oldsymbol{y}\in \mathbb{Y}} f(oldsymbol{y};oldsymbol{x})$

This would work, but prohibitive:.

 $\begin{array}{l} \textbf{for } \boldsymbol{y} = [\texttt{start}, y_1, \dots, y_L, \texttt{end}] \in \boldsymbol{\mathcal{Y}} \ \textbf{do} \\ \text{compute } f(\boldsymbol{y}; \boldsymbol{x}) & \# \texttt{e.g.} \ \boldsymbol{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}), \texttt{or nnet forward pass} \\ \text{return the highest scoring } \boldsymbol{y} \end{array}$

This looks similar to the problem Viterbi solved in HMMs!

 $\argmax_{\boldsymbol{y} \in \boldsymbol{\mathbb{Y}}} \mathsf{P}(\boldsymbol{y}, \boldsymbol{x})$

Everything would be fine if we had a way to calculate

 $rg\max_{oldsymbol{y}\in \mathbb{Y}} f(oldsymbol{y};oldsymbol{x})$

This would work, but prohibitive:.

 $\begin{array}{l} \textbf{for } \boldsymbol{y} = [\texttt{start}, y_1, \dots, y_L, \texttt{end}] \in \boldsymbol{\mathcal{Y}} \, \textbf{do} \\ \text{compute } f(\boldsymbol{y}; \boldsymbol{x}) & \# \texttt{e.g.} \, \boldsymbol{w} \cdot \boldsymbol{\phi}(\boldsymbol{x}, \boldsymbol{y}), \texttt{or nnet forward pass} \\ \text{return the highest scoring } \boldsymbol{y} \end{array}$

This looks similar to the problem Viterbi solved in HMMs!

 $\argmax_{\boldsymbol{y} \in \boldsymbol{\mathbb{Y}}} \mathsf{P}(\boldsymbol{y}, \boldsymbol{x})$

What magic made Viterbi work there? - can we replicate it?

The HMM magic formula: decomposition into parts

Viterbi was able to efficiently compute

 $rg\max_{oldsymbol{y}\in rak{Y}} \log \mathsf{P}(oldsymbol{y},oldsymbol{x})$

because of the decomposition into parts:

$$\log P(\boldsymbol{y}, \boldsymbol{x}) = \sum_{i=1}^{L} \underbrace{\log P(x_i | y_i)}_{\text{emission log-proba}} + \sum_{i=1}^{L+1} \underbrace{\log P(y_i | y_{i-1})}_{\text{transition log-proba}}$$

We should try to design our scorer f such that



The HMM magic formula: decomposition into parts

Viterbi was able to efficiently compute

 $rg\max_{oldsymbol{y}\in rak{Y}} \log \mathsf{P}(oldsymbol{y},oldsymbol{x})$

because of the decomposition into parts:

$$\log \mathsf{P}(\boldsymbol{y}, \boldsymbol{x}) = \sum_{i=1}^{L} \underbrace{\log \mathsf{P}(x_i | y_i)}_{\text{emission log-proba}} + \sum_{i=1}^{L+1} \underbrace{\log \mathsf{P}(y_i | y_{i-1})}_{\text{transition log-proba}}$$

We should try to design our scorer f such that

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} \underbrace{f_i^{(u)}(y_i; \boldsymbol{x})}_{\text{unary score}} + \sum_{i=2}^{L} \underbrace{f_i^{(t)}(y_i, y_{i-1}; \boldsymbol{x})}_{\text{transition score}}$$

Given *x*, unary scores form a $|\Sigma| \times L$ array $s_{y,i}^{(u)}$ and transition scores form a $|\Sigma| \times |\Sigma| \times L$ array $s_{y',y,i}^{(t)}$

Transitions from start and to end?

The HMM magic formula: decomposition into parts

Viterbi was able to efficiently compute

 $rg\max_{oldsymbol{y}\in rak{Y}} \log \mathsf{P}(oldsymbol{y},oldsymbol{x})$

because of the decomposition into parts:

$$\log \mathsf{P}(\boldsymbol{y}, \boldsymbol{x}) = \sum_{i=1}^{L} \underbrace{\log \mathsf{P}(x_i | y_i)}_{\text{emission log-proba}} + \sum_{i=1}^{L+1} \underbrace{\log \mathsf{P}(y_i | y_{i-1})}_{\text{transition log-proba}}$$

We should try to design our scorer f such that

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} \underbrace{f_{i}^{(u)}(y_{i}; \boldsymbol{x})}_{\text{unary score}} + \sum_{i=2}^{L} \underbrace{f_{i}^{(t)}(y_{i}, y_{i-1}; \boldsymbol{x})}_{\text{transition score}}$$

Given x, unary scores form a $|\Sigma| \times L$ array $s_{y,i}^{(u)}$ and transition scores form a $|\Sigma| \times |\Sigma| \times L$ array $s_{y',y,i}^{(t)}$

Transitions from start and to end? Can be added to $s_{v,1}^{(u)}$, $s_{v,L}^{(u)}$

Score-based Viterbi

input: Unary scores $s^{(u)}$ ($|\Sigma| \times L$ array) Transition scores $s^{(t)}$ ($|\Sigma| \times |\Sigma| \times (L - 1)$ array)

Forward: compute scores recursively $\widetilde{score}_1(y) = s_{y,1}^{(u)} \quad \forall y \in \Sigma$ for i = 2 to L do for $y \in \Sigma$ do $\widetilde{score}_i(y) = s_{y,i}^{(u)} + \max_{y'} \left(s_{y',y,i}^{(t)} + \widetilde{score}_{i-1}(y') \right)$ backptr_i(y) = arg max_{y'} $\left(s_{y',y,i}^{(t)} + \widetilde{score}_{i-1}(y') \right)$

Backward: follow backpointers $f^* = \max_{y'} \widetilde{\text{score}}_L(y')$ $\widehat{y}_L = \arg \max_{y'} \widetilde{\text{score}}_L(y')$

for
$$i = L - 1$$
 down to 1 d
 $\widehat{y}_i = backptr_{i+1}(\widehat{y}_{i+1})$

output: The highest-scoring $\widehat{y} = [\widehat{y}_1, \dots, \widehat{y}_L]$ and its total score f^* .

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}$$

- More expressive than an HMM: can look at entire *x*. Useful cases:
- We can simulate an HMM: $W^{(t)} \in \mathbb{R}^{|\Sigma| \times |\Sigma|}, \quad W^{(u)} \in \mathbb{R}^{|\Sigma| \times |V|};$

•
$$s_{y',y,i}^{(t)} = w_{y',y}^{(t)}$$
 (ignore x and i)
• $s_{y,i}^{(u)} = w_{y,x_i}^{(u)}$

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}$$

- More expressive than an HMM: can look at entire *x*. Useful cases:
- We can simulate an HMM: $W^{(t)} \in \mathbb{R}^{|\Sigma| \times |\Sigma|}, \quad W^{(u)} \in \mathbb{R}^{|\Sigma| \times |V|};$

•
$$s_{y',y,i}^{(t)} = w_{y',y}^{(t)}$$
 (ignore x and i)
• $s_{y,i}^{(u)} = w_{y,x_i}^{(u)}$

- Transitions as above; linear / neural model for unary scores
 - A word-level classifier augmented with transition scores.

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}$$

- More expressive than an HMM: can look at entire x. Useful cases:
- We can simulate an HMM: $W^{(t)} \in \mathbb{R}^{|\Sigma| imes |\Sigma|}, \quad W^{(u)} \in \mathbb{R}^{|\Sigma| imes |V|};$

•
$$s_{y',y,i}^{(t)} = w_{y',y}^{(t)}$$
 (ignore x and i)
• $s_{y,i}^{(u)} = w_{y,x_i}^{(u)}$

- Transitions as above; linear / neural model for unary scores
 - A word-level classifier augmented with transition scores.
 - Linear scores: $s_{y,i}^{(u)} = w^{(u)} \cdot \phi^{(u)}((x,i),y)$ Unary features:
 - y = B, x_i capitalized
 - $y = B, x_{i-1}$ indefinite article (a/an)

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}$$

- More expressive than an HMM: can look at entire *x*. Useful cases:
- We can simulate an HMM: $W^{(t)} \in \mathbb{R}^{|\Sigma| \times |\Sigma|}, \quad W^{(u)} \in \mathbb{R}^{|\Sigma| \times |V|};$

•
$$s_{y',y,i}^{(t)} = w_{y',y}^{(t)}$$
 (ignore x and i)
• $s_{y,i}^{(u)} = w_{y,x_i}^{(u)}$

- Transitions as above; linear / neural model for unary scores
 - A word-level classifier augmented with transition scores.
 - Linear scores: $\mathbf{s}_{y,i}^{(u)} = \mathbf{w}^{(u)} \cdot \phi^{(u)}((\mathbf{x},i),y)$ Unary features:
 - y = B, x_i capitalized
 - $y = B, x_{i-1}$ indefinite article (a/an)
 - Neural unary scores (example)
 - 1 Encode each word into a vector $h_i = g(x, i)$.

2 apply "output layer":
$$s_{y,i}^{(u)} = (\boldsymbol{W}^{(out)}\boldsymbol{h}_i + \boldsymbol{b}^{(out)})_y$$

Decomposable feature-based scorers

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} s_{y_{i},i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_{i},i}^{(t)}$$

Lots of NLP literature uses feature representations for everything.

$$\begin{split} s_{y,i}^{(u)} &= \boldsymbol{w}^{(u)} \cdot \phi^{(u)}((x,i),y) \\ s_{y',y,i}^{(t)} &= \boldsymbol{w}^{(t)} \cdot \phi^{(t)}((x,i),(y',y)) \\ \text{or, compactly,} \quad f(\boldsymbol{y},\boldsymbol{x}) &= \boldsymbol{w} \cdot \phi(\boldsymbol{x},\boldsymbol{y}) \\ \text{where} \quad \phi(\boldsymbol{y},\boldsymbol{x}) &= \mathsf{cat} \left[\sum_{i} \phi^{(u)}((x,i),y_{i}), \ \sum_{i} \phi^{(t)}((x,i),y_{i+1}) \right] \\ \boldsymbol{w} &= \mathsf{cat} \left[\boldsymbol{w}^{(u)}, \ \boldsymbol{w}^{(t)} \right] \end{split}$$

• Unary features: just like in multi-class. Transition features:

$$y = B, y' = 0, i = 3$$

•
$$y = B, y' = 0, x_i$$
 capitalized, $x_{i-1} =$ "from"

•
$$y = B, y' = 0, x_{last} = "?"$$

Structured Perceptron

input: labeled data ${\mathcal D}$

initialize $m{w}$

repeat

get new training example (x, y^{true}) compute unary and transition scores, $s^{(u)}, s^{(t)}$. predict $\widehat{y} = \arg \max_{y \in \mathbb{Y}} f(y; x)$ using Viterbi. if $\widehat{y} \neq y^{true}$ then update $w \leftarrow w + \nabla_w f(y^{true}; x) - \nabla_w f(\widehat{y}; x)$ until max. epochs output: model weights w

• If linear
$$f(y;x) = w \cdot \phi(x,y)$$
, then $abla_w f(y;x) = \phi(x,y)$.

• If neural, $\nabla_{\boldsymbol{w}} f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} \nabla_{\boldsymbol{w}} s_{y_{i,i}}^{(u)} + \sum_{i=2}^{L} \nabla_{\boldsymbol{w}} s_{y_{i-1},y_{i,i}}^{(t)}$ from autodiff.

Outline

Structured Prediction

Ø Generative Sequence Models

Markov Models

Hidden Markov Models

Unsupervised learning

3 Discriminative Sequence Models

Structured Perceptron

Conditional Random Fields

Structured SVM

What if we want a **discriminative probabilistic model**? i.e., one that gives P(y|x) and not just a prediction \hat{y} ? For multi-class, we had **logistic regression**.

$$\mathsf{P}(oldsymbol{y}|oldsymbol{x}) = rac{\mathsf{exp}\, \mathsf{s}_{oldsymbol{y}}}{\sum_{oldsymbol{y}}' \mathsf{exp}\, \mathsf{s}_{oldsymbol{y}}'}$$

$$\mathsf{P}(oldsymbol{y}|oldsymbol{x}) = rac{\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}}}{\sum_{oldsymbol{y}}'\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}'}}$$

To learn, we maximize

$$\log \mathsf{P}(y^{\mathsf{true}}|x) = \mathsf{s}_{y^{\mathsf{true}}} - \log \sum_{y} \exp \mathsf{s}_{y}$$

with gradient descent, noting that

$$abla \log \mathsf{P}(oldsymbol{y}^{\mathsf{true}} | oldsymbol{x}) =
abla \mathsf{s}_{oldsymbol{y}^{\mathsf{true}}} - \mathbb{E}_{oldsymbol{y}}
abla \mathsf{s}_{oldsymbol{y}}.$$

$$\mathsf{P}(oldsymbol{y}|oldsymbol{x}) = rac{\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}}}{\sum_{oldsymbol{y}}'\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}'}}$$

To learn, we maximize

$$\log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}} | \boldsymbol{x}) = \boldsymbol{s}_{\boldsymbol{y}^{\mathsf{true}}} - \log \sum_{\boldsymbol{y}} \exp \boldsymbol{s}_{\boldsymbol{y}}$$

with gradient descent, noting that

$$abla \log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) =
abla \mathsf{s}_{oldsymbol{y}^{\mathsf{true}}} - \mathbb{E}_{oldsymbol{y}}
abla \mathsf{s}_{oldsymbol{y}}.$$

In particular, for linear models,

$$abla_{m{w}} \log \mathsf{P}(m{y}^{\mathsf{true}}|m{x}) = \phi(m{x},m{y}^{\mathsf{true}}) - \mathbb{E}\phi(m{x},m{y})$$

$$\mathsf{P}(oldsymbol{y}|oldsymbol{x}) = rac{\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}}}{\sum_{oldsymbol{y}}'\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}'}}$$

To learn, we maximize

$$\log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}} | \boldsymbol{x}) = \mathsf{s}_{\boldsymbol{y}^{\mathsf{true}}} - \log \sum_{\boldsymbol{y}} \exp \mathsf{s}_{\boldsymbol{y}}$$

with gradient descent, noting that

$$abla \log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) =
abla {\mathsf{s}}_{oldsymbol{y}^{\mathsf{true}}} - \mathbb{E}_{oldsymbol{y}}
abla {\mathsf{s}}_{oldsymbol{y}}.$$

In particular, for linear models,

$$abla_w \log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) = \phi(oldsymbol{x},oldsymbol{y}^{\mathsf{true}}) - \mathbb{E}\phi(oldsymbol{x},oldsymbol{y})$$

For neural nets, grad wrt. score vector

$$egin{aligned}
abla_s \log \mathsf{P}(m{y}^{\mathsf{true}} | m{x}) &= m{e}_{m{y}^{\mathsf{true}}} - \mathbb{E}m{e}_{m{y}} \ &= m{e}_{m{y}^{\mathsf{true}}} - \mathsf{softmax}(s) \end{aligned}$$

$$\mathsf{P}(oldsymbol{y}|oldsymbol{x}) = rac{\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}}}{\sum_{oldsymbol{y}}'\mathsf{exp}\,\mathsf{s}_{oldsymbol{y}'}}$$

To learn, we maximize

$$\log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}} | \boldsymbol{x}) = \mathsf{s}_{\boldsymbol{y}^{\mathsf{true}}} - \log \sum_{\boldsymbol{y}} \exp \mathsf{s}_{\boldsymbol{y}}$$

with gradient descent, noting that

$$abla \log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) =
abla {\mathsf{s}}_{oldsymbol{y}^{\mathsf{true}}} - \mathbb{E}_{oldsymbol{y}}
abla {\mathsf{s}}_{oldsymbol{y}}.$$

In particular, for linear models,

$$abla_{m{w}} \log \mathsf{P}(m{y}^{\mathsf{true}}|m{x}) = \phi(m{x},m{y}^{\mathsf{true}}) - \mathbb{E}\phi(m{x},m{y})$$

For neural nets, grad wrt. score vector

$$egin{aligned}
abla_s \log \mathsf{P}(m{y}^{\mathsf{true}} | m{x}) &= m{e}_{m{y}^{\mathsf{true}}} - \mathbb{E}m{e}_{m{y}} \ &= m{e}_{m{y}^{\mathsf{true}}} - \mathsf{softmax}(s) \end{aligned}$$

For sequence models, Forward-Backward computes what we need!

Vlad Niculae & André Martins (IST)

Conditional Random Fields

Discriminative (non-generative) structured models.

$$\log \mathsf{P}(\boldsymbol{y}|\boldsymbol{x}) = f(\boldsymbol{y};\boldsymbol{x}) - \log \sum_{\boldsymbol{y}'} \exp f(\boldsymbol{y}';\boldsymbol{x})$$

Given a decomposable sequence scorer,

$$f(\boldsymbol{y}; \boldsymbol{x}) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}$$

Forward-Backward computes

- Normalizer / log-partition function $\log \sum_{y'} \exp f(y'; x) = \log Z$.
- Unary and transition posterior marginals: log P($y_i = y | x$), log P($y_i = y, y_{i+1} = y' | x$).

Score-based Forward-Backward

input: Unary scores $s^{(u)}$ ($|\Sigma| \times L$ array) Transition scores $s^{(t)}$ ($|\Sigma| \times |\Sigma| \times (L-1)$ array)

Forward: compute $\widetilde{\alpha}$ recursively $\widetilde{\alpha}_1(y) = s_{y,1}^{(u)} \quad \forall y \in \Sigma$ for i = 2 to L do for $y \in \Sigma$ do $\widetilde{\alpha}_i(y) = s_{y,i}^{(u)} + \log \sum_{y'} \exp\left(s_{y',y,i}^{(t)} + \widetilde{\alpha}_{i-1}(y')\right)$

Backward: compute $\widetilde{\beta}$ recursively $\widetilde{\beta}_{L}(y) = 0$ for i = L - 1 down to 1 do for $y \in \Sigma$ do $\widetilde{\beta}_{i}(y) = \log \sum_{y'} \exp \left(s_{y,y',i}^{(t)} + s_{y',i+1}^{(u)} + \widetilde{\beta}_{i+1}(y') \right)$

output: The log-partition function $\log Z = \log \sum_{y'} \exp \left(\widetilde{\alpha}_i(y') + \widetilde{\beta}_i(y') \right)$ for any *i*; Posterior unigram and transition marginal log-probas: $\log P(y_i = y | \boldsymbol{x}) = \widetilde{\alpha}_i(y) + \widetilde{\beta}_i(y) - \log Z;$ $\log P(y_i = y, y_{i+1} = y' | \boldsymbol{x}) = \widetilde{\alpha}_i(y) + \widetilde{\beta}_{i+1}(y') + s_{y',i+1}^{(u)} + s_{y,y',i+1}^{(t)} - \log Z.$
Training a linear model CRF

We want to maximize

$$\log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) = oldsymbol{w} \cdot \phi(oldsymbol{x},oldsymbol{y}^{\mathsf{true}}) - \log \sum_{oldsymbol{y}} \exp oldsymbol{w} \cdot \phi(oldsymbol{x},oldsymbol{y})$$

where

$$egin{aligned} \phi(oldsymbol{y},oldsymbol{x}) &= \mathsf{cat}\left[\sum_i \phi^{(u)}((oldsymbol{x},i),y_i), \ \sum_i \phi^{(t)}((oldsymbol{x},i),y_i,y_{i+1})
ight] \ oldsymbol{w} &= \mathsf{cat}\left[oldsymbol{w}^{(u)}, \ oldsymbol{w}^{(t)}
ight]. \end{aligned}$$

Training a linear model CRF

We want to maximize

$$\log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) = oldsymbol{w} \cdot \phi(oldsymbol{x},oldsymbol{y}^{\mathsf{true}}) - \log \sum_{oldsymbol{y}} \exp oldsymbol{w} \cdot \phi(oldsymbol{x},oldsymbol{y})$$

where

$$egin{aligned} \phi(oldsymbol{y},oldsymbol{x}) &= \mathsf{cat}\left[\sum_i \phi^{(u)}((oldsymbol{x},i),y_i), \ \sum_i \phi^{(t)}((oldsymbol{x},i),y_i,y_{i+1})
ight] \ oldsymbol{w} &= \mathsf{cat}\left[oldsymbol{w}^{(u)}, \ oldsymbol{w}^{(t)}
ight]. \end{aligned}$$

Similar to logistic regression, we get

$$\begin{aligned} \nabla_{\boldsymbol{w}^{(u)}} \log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}} | \boldsymbol{x}) &= \sum_{i} \Big(\phi^{(u)}((\boldsymbol{x}, i), y_{i}^{\mathsf{true}}) - \mathbb{E}_{\boldsymbol{y}} \phi^{(u)}((\boldsymbol{x}, i), y_{i}) \Big) \\ &= \sum_{i} \Big(\phi^{(u)}((\boldsymbol{x}, i), y_{i}^{\mathsf{true}}) - \sum_{\boldsymbol{y}} \phi^{(u)}((\boldsymbol{x}, i), \boldsymbol{y}) \,\mathsf{P}(\boldsymbol{y}_{i} = \boldsymbol{y} | \boldsymbol{x}) \Big) \end{aligned}$$

Training a linear model CRF

We want to maximize

$$\log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) = oldsymbol{w} \cdot \phi(oldsymbol{x},oldsymbol{y}^{\mathsf{true}}) - \log \sum_{oldsymbol{y}} \exp oldsymbol{w} \cdot \phi(oldsymbol{x},oldsymbol{y})$$

where

$$egin{aligned} \phi(oldsymbol{y},oldsymbol{x}) &= \mathsf{cat}\left[\sum_i \phi^{(u)}((oldsymbol{x},i),y_i), \ \sum_i \phi^{(t)}((oldsymbol{x},i),y_i,y_{i+1})
ight] \ oldsymbol{w} &= \mathsf{cat}\left[oldsymbol{w}^{(u)}, \ oldsymbol{w}^{(t)}
ight]. \end{aligned}$$

Similar to logistic regression, we get

$$\begin{aligned} \nabla_{\boldsymbol{w}^{(u)}} \log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}} | \boldsymbol{x}) &= \sum_{i} \Big(\phi^{(u)}((\boldsymbol{x}, i), y_{i}^{\mathsf{true}}) - \mathbb{E}_{\boldsymbol{y}} \phi^{(u)}((\boldsymbol{x}, i), y_{i}) \Big) \\ &= \sum_{i} \Big(\phi^{(u)}((\boldsymbol{x}, i), y_{i}^{\mathsf{true}}) - \sum_{\boldsymbol{y}} \phi^{(u)}((\boldsymbol{x}, i), \boldsymbol{y}) \,\mathsf{P}(\boldsymbol{y}_{i} = \boldsymbol{y} | \boldsymbol{x}) \Big) \end{aligned}$$

and, for transitions,

$$\begin{aligned} \nabla_{\boldsymbol{w}^{(t)}} \log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}} | \boldsymbol{x}) &= \sum_{i} \Big(\phi^{(t)}((\boldsymbol{x}, i), y_{i-1}^{\mathsf{true}}, y_{i}^{\mathsf{true}}) \\ &- \sum_{\boldsymbol{y}, \boldsymbol{y}'} \phi^{(t)}((\boldsymbol{x}, i), \boldsymbol{y}, \boldsymbol{y}') \,\mathsf{P}(\boldsymbol{y}_{i-1} = \boldsymbol{y}, \boldsymbol{y}_{i} = \boldsymbol{y}' | \boldsymbol{x}) \Big) \end{aligned}$$

Training a neural CRF

We want to maximize

$$\log \mathsf{P}(m{y}^{\mathsf{true}}|m{x}) = f(m{x},m{y}) - \log \sum_{m{y}} \exp f(m{x},m{y})$$

where

$$f(x, y) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}$$

Training a neural CRF

We want to maximize

$$\log \mathsf{P}(oldsymbol{y}^{ ext{true}}|oldsymbol{x}) = f(oldsymbol{x},oldsymbol{y}) - \log \sum_{oldsymbol{y}} \exp f(oldsymbol{x},oldsymbol{y})$$

where

$$f(\boldsymbol{x}, \boldsymbol{y}) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}.$$

Note that

$$\frac{\partial f(\boldsymbol{x}, \boldsymbol{y})}{\partial s_{\boldsymbol{y}, i}^{(u)}} = \begin{cases} 1, & y_i = y \\ 0, & \text{otherwise} \end{cases} = [[y_i = y]], \qquad \frac{\partial f(\boldsymbol{x}, \boldsymbol{y})}{\partial s_{\boldsymbol{y}, \boldsymbol{y}', i}^{(t)}} = [[y_{i-1} = y, y_i = y']] \end{cases}$$

Training a neural CRF

We want to maximize

$$\log \mathsf{P}(oldsymbol{y}^{\mathsf{true}}|oldsymbol{x}) = f(oldsymbol{x},oldsymbol{y}) - \log \sum_{oldsymbol{y}} \exp f(oldsymbol{x},oldsymbol{y})$$

where

$$f(x, y) = \sum_{i=1}^{L} s_{y_i,i}^{(u)} + \sum_{i=2}^{L} s_{y_{i-1},y_i,i}^{(t)}$$

Note that

$$\frac{\partial f(\boldsymbol{x}, \boldsymbol{y})}{\partial s_{\boldsymbol{y}, i}^{(u)}} = \begin{cases} 1, & y_i = y \\ 0, & \text{otherwise} \end{cases} = [[y_i = y]], \qquad \frac{\partial f(\boldsymbol{x}, \boldsymbol{y})}{\partial s_{\boldsymbol{y}, \boldsymbol{y}', i}^{(t)}} = [[y_{i-1} = y, y_i = y']] \end{cases}$$

Therefore,

$$\frac{\partial \log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}}|\boldsymbol{x})}{\partial s_{y,i}^{(u)}} = [[y_i^{\mathsf{true}} = y]] - \mathsf{P}(y_i = y|\boldsymbol{x})$$
$$\frac{\partial \log \mathsf{P}(\boldsymbol{y}^{\mathsf{true}}|\boldsymbol{x})}{\partial s_{y,y',i}^{(t)}} = [[y_{i-1}^{\mathsf{true}} = y, y_i^{\mathsf{true}} = y']] - \mathsf{P}(y_{i-1} = y, y_i = y'|\boldsymbol{x})$$

Visually:



Outline

Structured Prediction

Ø Generative Sequence Models

Markov Models Hidden Markov Model

Unsupervised learning

3 Discriminative Sequence Models

Structured Perceptron Conditional Random Field

Structured SVM

SVMs are non-probabilistic max-margin classifiers.

The soft-margin perspective: minimizing a "hinge loss":

$$L(w; (x, y^{\mathsf{true}})) = \max_{oldsymbol{y}} w \cdot \phi(x, oldsymbol{y}) + [[oldsymbol{y}
otin v^{\mathsf{true}}]] - w \cdot \phi(x, oldsymbol{y}^{\mathsf{true}})$$

Intuition: the score of the correct class must be greater than the score of wrong classes by at least 1.

Structured SVM

$$L(oldsymbol{w};(oldsymbol{x},oldsymbol{y}^{\mathsf{true}})) = \max_{oldsymbol{y}} f(oldsymbol{x},oldsymbol{y}) + \operatorname{cost}(oldsymbol{y},oldsymbol{y}^{\mathsf{true}}) - f(oldsymbol{x},oldsymbol{y}^{\mathsf{true}})$$

where $cost(\boldsymbol{y}, \boldsymbol{y}^{true}) = \sum_{i=1}^{L} [[y_i \neq y_i^{true}]]$ is the Hamming cost.

Structured SVM

$$L(oldsymbol{w};(oldsymbol{x},oldsymbol{y}^{ extsf{true}})) = \max_{oldsymbol{y}} f(oldsymbol{x},oldsymbol{y}) + extsf{cost}(oldsymbol{y},oldsymbol{y}^{ extsf{true}}) - f(oldsymbol{x},oldsymbol{y}^{ extsf{true}})$$

where $cost(\boldsymbol{y}, \boldsymbol{y}^{true}) = \sum_{i=1}^{L} [[y_i \neq y_i^{true}]]$ is the Hamming cost.

Cost-augmented decoding:

Finding $\arg \max_y f(x, y) + \cot(y, y^{true})$ can be done by Viterbi with adjusted unary scores

$$\tilde{s}_{y,i}^{(u)} = \begin{cases} s_{y,i}^{(u)} & y_i = y_i^{\text{true}} \\ s_{y,i}^{(u)} + 1 & y_i \neq y_i^{\text{true}} \end{cases}$$

Intuition: We give a boost to the wrong labels, and we want the global prediction to still be correct.

Similar algorithm to Structured Perceptron.

Vlad Niculae & André Martins (IST)

Lecture 5: Linear Sequential Models

Structured SVM

input: labeled data \mathcal{D} , learning rate η initialize \boldsymbol{w}

repeat

get new training example (x, y^{true}) compute unary and transition scores, $s^{(u)}$, $s^{(t)}$ predict $\hat{y} = \arg \max_{y \in \mathcal{Y}} f(y; x) + \operatorname{cost}(y, y^{true})$ (cost-augmented Viterbi) if $\hat{y} \neq y$ then update $w \leftarrow w + \eta (\nabla_w f(y^{true}; x) - \nabla_w f(\hat{y}; x))$ until max. epochs output: model weights w

- If linear $f(m{y};m{x}) = m{w}\cdot \phi(m{x},m{y})$, then $abla_{m{w}}f(m{y};m{x}) = \phi(m{x},m{y})$.
- If neural, $\nabla_{w} f(y; x) = \sum_{i=1}^{L} \nabla_{w} s_{y_{i,i}}^{(u)} + \sum_{i=2}^{L} \nabla_{w} s_{y_{i-1},y_{i,i}}^{(t)}$ from autodiff.

Visually:

Unlike CRF, subgradient updates for Perceptron/SSVM look like:



Structured prediction models

Binary/Multi-class Structured Prediction

Naive Bayes Perceptron Logistic Regression SVMs HMMs Structured Perceptron CRF Structured SVM

Discriminative vs HMM

- HMM is less expressive, but fast to train.
- Discriminative models are powerful, but require iterative training.
- HMM and CRF have probabilistic interpretations. (Useful when posterior analyses are desirable)
- Structured SVM is a good classifier, can be extended to other cost functions.
- Perceptron and Structured SVM updates are sparse; (may be faster for very large Σ.)
- For structures more complicated than sequences, we may not have a Forward-Backward equivalent, but we may be able to approximate arg maxy f(y; x).