Predicting Sequences: Hidden Markov Models

CS 6355: Structured Prediction



Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences

Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences

Sequences

- Sequences of states
 - Text is a sequence of words or even letters
 - A video is a sequence of frames
- Even with a finite set of states, the set of unique state *sequences* is infinite
- Our goal (for now): Define probability distributions over sequences
- If x₁, x₂, ..., x_n is a sequence that has n tokens, we want to be able to define P(x₁, x₂, ..., x_n)
 ...for all values of n

A history-based model

$$P(x_1, x_2, \cdots, x_n) = \prod_{i=1}^n P(x_i \mid x_1, x_2, \cdots, x_{i-1})$$

Each token is dependent on every token that came before it

- Simple conditioning
- Each $P(x_i | x_1, x_2, \dots, x_{i-1})$ is a multinomial probability distribution over the tokens



It was a bright cold day in April.

P(It was a bright cold day in April) =

It was a bright cold day in April.

P(It was a bright cold day in April) =

 $P(\mathrm{It}) imes$ imes Probability of a word starting a sentence

It was a bright cold day in April.

P(It was a bright cold day in April) = $P(\text{It}) \times \longleftarrow \text{Probability of a word starting a sentence}$ $P(\text{was}|\text{It}) \times \longleftarrow \text{Probability of a word following "It"}$

It was a bright cold day in April.

It was a bright cold day in April.

P(It was a bright cold day in April) = $P(\text{It}) \times \longleftarrow$ Probability of a word starting a sentence $P(\text{was}|\text{It}) \times \longleftarrow$ Probability of a word following "It"

 $P(a|It was) \times$ — Probability of a word following "It was"

 $P(\text{bright}|\text{It was a}) imes extsf{formula}$ Probability of a word following "It was a"

It was a bright cold day in April.

P(It was a bright cold day in April) =

 $P(\mathrm{It}) \times \longleftarrow Probability of a word starting a sentence$ $P(\mathrm{was}|\mathrm{It}) \times \longleftarrow Probability of a word following "It"$ $P(\mathrm{a}|\mathrm{It} \mathrm{was}) \times \longleftarrow Probability of a word following "It was"$ $P(\mathrm{bright}|\mathrm{It} \mathrm{was} \mathrm{a}) \times \longleftarrow Probability of a word following "It was a"$ $P(\mathrm{cold}|\mathrm{It} \mathrm{was} \mathrm{a} \mathrm{bright}) \times$ $P(\mathrm{day}|\mathrm{It} \mathrm{was} \mathrm{a} \mathrm{bright} \mathrm{cold}) \times \cdots$

What's the problem with this strategy?

A history-based model

$$P(x_1, x_2, \cdots, x_n) = \prod_{i=1}^n P(x_i \mid x_1, x_2, \cdots, x_{i-1})$$

Each token is dependent on every token that came before it

- Simple conditioning
- Each $P(x_i | x_1, x_2, \dots, x_{i-1})$ is a multinomial probability distribution over the tokens

What's the problem with this strategy?

- How many parameters do we have?
 - Grows with the size of the sequence!

Solution: Lose the history

Make a modeling assumption: *The first-order Markov assumption*

The state of the system at any time is *independent* of the full sequence history *given the previous state*

$$P(x_i \mid x_1, x_2, \cdots, x_{i-1}) = P(x_i \mid x_{i-1})$$

Solution: Lose the history

Make a modeling assumption: *The first-order Markov assumption*

The state of the system at any time is *independent* of the full sequence history *given the previous state*

$$P(x_i \mid x_1, x_2, \cdots, x_{i-1}) = P(x_i \mid x_{i-1})$$

This allows us to simplify

$$P(x_1, x_2, x_3, \cdots, x_n) = \prod_i P(x_i \mid x_1, x_2 \cdots, x_{i-1})$$

These dependencies are ignored

Solution: Lose the history

Make a modeling assumption: *The first-order Markov assumption*

The state of the system at any time is *independent* of the full sequence history *given the previous state*

$$P(x_i \mid x_1, x_2, \cdots, x_{i-1}) = P(x_i \mid x_{i-1})$$

This allows us to simplify

$$P(x_1, x_2, x_3, \dots, x_n) = \prod_i P(x_i \mid x_{i-1})$$

First-order Markov models

Defined by two sets of probabilities

1. The initial state distribution: The probability that a sequence starts at a certain state $j: P(x_1 = \text{state}_i)$

2. The state transition distribution: The probability that the system will transition to a state k at some step if it was at a state j at the previous step: $P(x_{t+1} \text{state}_k \mid x_t = \text{state}_j)$

Example: Another language model

It was a bright cold day in April

 $P(\text{It was a bright cold day in April}) = P(\text{It}) \times \qquad \qquad Probability of a word starting a sentence} \\ P(\text{was}|\text{It}) \times \qquad \qquad Probability of a word following "It"} \\ P(a|\text{was}) \times \qquad \qquad Probability of a word following "was"} \\ P(bright|a) \times \qquad \qquad Probability of a word following "a"} \\ P(cold|bright) \times \\ P(day|cold) \times \cdots$

Example: Another language model

It was a bright cold day in April

 $P(\text{It was a bright cold day in April}) = P(\text{It}) \times \qquad \qquad Probability of a word starting a sentence} \\ P(\text{was}|\text{It}) \times \qquad \qquad Probability of a word following "It"} \\ P(a|\text{was}) \times \qquad \qquad Probability of a word following "was"} \\ P(bright|a) \times \qquad \qquad Probability of a word following "a"} \\ P(cold|bright) \times \\ P(day|cold) \times \cdots$

If there are K tokens/states, how many parameters do we need?

Example: Another language model

It was a bright cold day in April

 $P(\text{It was a bright cold day in April}) = P(\text{It}) \times \qquad \qquad Probability of a word starting a sentence} \\ P(\text{was}|\text{It}) \times \qquad \qquad Probability of a word following "It"} \\ P(a|\text{was}) \times \qquad \qquad Probability of a word following "was"} \\ P(bright|a) \times \qquad \qquad Probability of a word following "a"} \\ P(cold|bright) \times \\ P(day|cold) \times \cdots$

If there are K tokens/states, how many parameters do we need? $O(K^2)$

Example: The weather

Three states: rain, cloudy, sunny



mth order Markov Model

A generalization of the first order Markov Model

- Each state is only dependent on m previous states
- More parameters
- But still less than storing entire history

mth order Markov Model

A generalization of the first order Markov Model

- Each state is only dependent on m previous states
- More parameters
- But still less than storing entire history

Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences

Hidden Markov Model

- Discrete Markov Model:
 - States follow a Markov chain
 - Each state is an observation
- Hidden Markov Model:
 - States follow a Markov chain
 - States are not observed
 - Each state stochastically emits an observation

Given a sentence, find parts of speech of all the words

The Fed raises interest rates

Given a sentence, find parts of speech of all the words



Given a sentence, find parts of speech of all the words



Given a sentence, find parts of speech of all the words



If these were the only options allowed, we will have $1 \times 2 \times 2 \times 2 \times 2 = 16$ possible output sequences



Each edge here is associated with a **transition probability**



Each edge here is associated with a **transition probability**

P(The | Determiner) = 0.5 P(A | Determiner) = 0.3 P(An | Determiner) = 0.1 P(Fed | Determiner) = 0

Emissions

P(Fed| Noun) = 0.001 P(raises| Noun) = 0.04 P(interest| Noun) = 0.07 P(The| Noun) = 0

Emission probabilities: Given that the system is in a certain state, these are probabilities that it will emit a certain observation

...



Each edge here is associated with a **transition probability**

Emissions

...

P(Fed | Noun) = 0.001 P(raises | Noun) = 0.04 P(interest | Noun) = 0.07 P(The | Noun) = 0

Emission probabilities: Given that the system is in a certain state, these are probabilities that it will emit a certain observation

Initial

P(Determiner) = 0.9 P(Noun) = 0.08 P(Verb) = 0.02 *Initial probabilities*: What is the probability that the sequence starts in a certain state?

...



P(The | Determiner) = 0.5 P(A | Determiner) = 0.3 P(An | Determiner) = 0.1 P(Fed | Determiner) = 0

Emissions

P(Fed | Noun) = 0.001 P(raises | Noun) = 0.04 P(interest | Noun) = 0.07 P(The | Noun) = 0

...



P(The | Determiner) = 0.5 P(A | Determiner) = 0.3 P(An | Determiner) = 0.1 P(Fed | Determiner) = 0

Emissions

P(Fed | Noun) = 0.001 P(raises | Noun) = 0.04 P(interest | Noun) = 0.07 P(The | Noun) = 0

...



P(The | Determiner) = 0.5 P(A | Determiner) = 0.3 P(An | Determiner) = 0.1 P(Fed | Determiner) = 0

Emissions

P(Fed | Noun) = 0.001 P(raises | Noun) = 0.04 P(interest | Noun) = 0.07 P(The | Noun) = 0



Emissions

P(Fed| Noun) = 0.001 P(raises| Noun) = 0.04 P(interest| Noun) = 0.07 P(The| Noun) = 0



Emissions

P(Fed | Noun) = 0.001 P(raises | Noun) = 0.04 P(interest | Noun) = 0.07 P(The | Noun) = 0












Joint model over states and observations

- Notation
 - Number of states = K
 - Number of possible observations for any state = M
 - π : Initial probability over states (K 1 numbers)
 - A: Transition probabilities ($K \times K$ matrix)
 - B: Emission probabilities ($K \times M$ matrix)

Joint model over states and observations

• Notation

- Number of states = K
- Number of possible observations for any state = M
- π : Initial probability over states (K 1 numbers)
- A: Transition probabilities ($K \times K$ matrix)
- B: Emission probabilities ($K \times M$ matrix)
- Probability of states and observations
 - Denote states by y_1, y_2, \cdots and observations by x_1, x_2, \cdots

$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$
$$= \pi_{y_1} \prod_{i=1}^{n-1} A_{y_i, y_{i+1}} \prod_{i=1}^n B_{y_i, x_i}$$

44

Example: Named Entity Recognition

Goal: To identify persons, locations and organizations in text



Example: Named Entity Recognition

Goal: To identify persons, locations and organizations in text



Numerous other applications

- Speech recognition
 - Input: Speech signal
 - Output: Sequence of words
- NLP applications
 - Information extraction
 - Text chunking
- Computational biology
 - Aligning protein sequences
 - Labeling nucleotides in a sequence as exons, introns, etc.

Questions?

Three questions for HMMs

[Rabiner 1999]

- 1. Given an observation sequence x_1, x_2, \dots, x_n and a model (π, A, B) , how to efficiently calculate the probability of the observation?
- 2. Given an observation sequence x_1, x_2, \dots, x_n and a model (π, A, B) , how to efficiently calculate the most probable state sequence?
- 3. How to calculate (π, A, B) from observations?

Three questions for HMMs

[Rabiner 1999]

- 1. Given an observation sequence x_1, x_2, \dots, x_n and a model (π, A, B) , how to efficiently calculate the probability of the observation?
- 2. Given an observation sequence x_1, x_2, \dots, x_n and a model (π, A, B) , how to efficiently calculate the most probable state sequence?
- 3. How to calculate (π, A, B) from observations?

Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences

Most likely state sequence

- Input:
 - A hidden Markov model (π , A, B)
 - An observation sequence $\mathbf{x} = (x_1, x_2, \dots, x_n)$
- Output: A state sequence y = (y₁, y₂, …, y_n) that corresponds to argmax P(y | x, π, A, B)
 y
 Maximum *a posteriori* inference (MAP inference)
- Computationally: combinatorial optimization

MAP inference

- We want to find $\underset{y}{\operatorname{argmax}} P(\mathbf{y} \mid \mathbf{x}, \pi, A, B)$
- We have defined

$$P(x_1, x_2, \dots, x_n, y_1, y_2, \dots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1} \mid y_i) \prod_{i=1}^n P(x_i \mid y_i)$$

- But, $P(\mathbf{y} | \mathbf{x}, \pi, A, B) \propto P(\mathbf{x}, \mathbf{y} | \pi, A, B)$ - And we don't care about $P(\mathbf{x})$ we are maximizing over
 - And we don't care about $P(\mathbf{x})$ we are maximizing over \mathbf{y}
- That is

 $\underset{y}{\operatorname{argmax}} P(\mathbf{y} \mid \mathbf{x}, \pi, A, B) = \underset{y}{\operatorname{argmax}} P(\mathbf{x}, \mathbf{y} \mid \pi, A, B)$

How many possible sequences?

The	Fed	raises	interest	rates
	Suppose each word	allows only the fo	ollowing tags	
Determiner	Verb	Verb	Verb	Verb
	Noun	Noun	Noun	Noun
1	2	2	2	2

In this simple case, $1 \times 2 \times 2 \times 2 \times 2 = 16$ possible sequences exist

How many possible sequences?

Observations **X**₁ X_2 Xn ... Suppose each observation allows any of the following k states S_1 S_1 S_1 ... **S**₂ S_2 S_2 **S**₃ S_2 **S**₃ . • . S_K S_K S_K

Output: One state per observation $y_i = s_j$ Kⁿ possible sequences to consider for $\operatorname{argmax} P(\mathbf{y} | \mathbf{x}, \pi, A, B)$

Naïve approaches

- 1. Try out every sequence
 - Score the sequence \mathbf{y} as $P(\mathbf{y} | \mathbf{x}, \pi, A, B)$
 - Return the highest scoring one
 - Correct, but slow, O(Kⁿ)
- 2. Greedy search
 - Construct the output left to right
 - For each i, elect the best y_i using y_{i-1} and x_i
 - Incorrect but fast, O(n)

Solution: Use the independence assumptions

Take advantage of the first order Markov assumption

The state for any observation is only influenced by the previous state, the next state and the observation itself

Given the adjacent labels, the others do not matter

Suggests a recursive algorithm

$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

What we want: An assignment to all the y_i 's that maximizes this product



$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

 $\max_{y_1, y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$





Initial probability



Deriving the recursive algorithm $\frac{n-1}{n}$









$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

 $\max_{y_1, y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$ = $\max_{y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$



$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$



The only terms that depend on y_1



$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

 $\max_{y_1, y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$

- $= \max_{y_2, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$ $= \max_{y_2, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$



$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

 $P(y_n|y_{n-1})P(x_n|y_n)\cdots P(y_2|y_1)P(x_2|y_2)P(y_1)P(x_1|y_1)$ max y_1, y_2, \cdots, y_n

- $= \max_{y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$
- $= \max_{y_2, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$ $= \max_{y_3, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_2} P(y_3 | y_2) P(x_3 | y_3) \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$



$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

 $\max_{y_1, y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$

- $= \max_{y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$
- $= \max_{y_2, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$

$$= \max_{y_3, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_2} P(y_3 | y_2) P(x_3 | y_3) \max_{y_1} \frac{P(y_2 | y_1) P(x_2 | y_2)}{1} \operatorname{score}_1(y_1)$$

Only terms that depend on y₂



$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

 $\max_{y_1, y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$

- $= \max_{y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$
- $= \max_{y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$
- $= \max_{y_3, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_2} P(y_3 | y_2) P(x_3 | y_3) \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$

$$= \max_{y_3, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_2} P(y_3 | y_2) P(x_3 | y_3) \operatorname{score}_2(y_2)$$



Abstract away the score for all decisions till here into score

$$P(x_1, x_2, \cdots, x_n, y_1, y_2, \cdots, y_n) = P(y_1) \prod_{i=1}^{n-1} P(y_{i+1}|y_i) \prod_{i=1}^n P(x_i|y_i)$$

 $\max_{y_1, y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$

- $= \max_{y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) P(y_1) P(x_1 | y_1)$
- $= \max_{y_2, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$
- $= \max_{y_3, \cdots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_2} P(y_3 | y_2) P(x_3 | y_3) \max_{y_1} P(y_2 | y_1) P(x_2 | y_2) \operatorname{score}_1(y_1)$

$$= \max_{y_3, \dots, y_n} P(y_n | y_{n-1}) P(x_n | y_n) \cdots \max_{y_2} P(y_3 | y_2) P(x_3 | y_3) \operatorname{score}_2(y_2)$$



Abstract away the score for all decisions till here into score

$$P(x_{1}, x_{2}, \cdots, x_{n}, y_{1}, y_{2}, \cdots y_{n}) = P(y_{1}) \prod_{i=1}^{n-1} P(y_{i+1}|y_{i}) \prod_{i=1}^{n} P(x_{i}|y_{i})$$

$$\max_{y_{1}, y_{2}, \cdots, y_{n}} P(y_{n}|y_{n-1})P(x_{n}|y_{n}) \cdots \max_{y_{1}} P(y_{2}|y_{2})P(y_{1})P(x_{1}|y_{1})$$

$$= \max_{y_{2}, \cdots, y_{n}} P(y_{n}|y_{n-1})P(x_{n}|y_{n}) \cdots \max_{y_{1}} P(y_{2}|y_{1})P(x_{2}|y_{2})\operatorname{score}_{1}(y_{1})$$

$$= \max_{y_{3}, \cdots, y_{n}} P(y_{n}|y_{n-1})P(x_{n}|y_{n}) \cdots \max_{y_{2}} P(y_{3}|y_{2})P(x_{3}|y_{3}) \max_{y_{1}} P(y_{2}|y_{1})P(x_{2}|y_{2})\operatorname{score}_{1}(y_{1})$$

$$= \max_{y_{3}, \cdots, y_{n}} P(y_{n}|y_{n-1})P(x_{n}|y_{n}) \cdots \max_{y_{2}} P(y_{3}|y_{2})P(x_{3}|y_{3})\operatorname{score}_{2}(y_{2})$$

$$\vdots$$

$$= \max_{y_{3}, \cdots, y_{n}} P(y_{n}|y_{n-1})P(x_{n}|y_{n}) \cdots \max_{y_{2}} P(y_{3}|y_{2})P(x_{3}|y_{3})\operatorname{score}_{2}(y_{2})$$

$$\vdots$$

Abstract away the score for all decisions till here into score

 $\begin{pmatrix} x_1 \end{pmatrix}$



$$= \max_{y_n} \operatorname{score}_n(y_n)$$

Viterbi algorithm

Max-product algorithm for first order sequences

1. Initial: For each state s, calculate $score_1(s) = P(s)P(x_1 | s)$

2. Recurrence: For i = 2 to n, for every state s, calculate $score_i(s) = \max_{y_{i-1}} P(s | y_{i-1}) P(x_i | s) score_{i-1}(y_{i-1})$

3. At the final state: calculate $\max_{y_{i-1}} P(y, x \mid \pi, A, B) = \max_{s} \frac{score_n(s)}{s}$

Viterbi algorithm

Max-product algorithm for first order sequences

 π : Initial probabilities A: Transitions B: Emissions

1. Initial: For each state s, calculate $score_1(s) = P(s)P(x_1 | s) = \pi_s B_{x_1,s}$

- 2. Recurrence: For i = 2 to n, for every state s, calculate $score_{i}(s) = \max_{y_{i-1}} P(s \mid y_{i-1}) P(x_{i} \mid s) score_{i-1}(y_{i-1})$ $= \max_{y_{i-1}} A_{y_{i-1,s}} B_{s,x_{i}} score_{i-1}(y_{i-1})$
- 3. At the final state: calculate $\max_{y_{i-1}} P(y, x \mid \pi, A, B) = \max_{s} \frac{score_n(s)}{s}$
Viterbi algorithm

Max-product algorithm for first order sequences

 π : Initial probabilities A: Transitions B: Emissions

1. Initial: For each state s, calculate $score_1(s) = P(s)P(x_1 | s) = \pi_s B_{x_1,s}$

2. Recurrence: For i = 2 to n, for every state s, calculate $score_{i}(s) = \max_{y_{i-1}} P(s \mid y_{i-1}) P(x_{i} \mid s) score_{i-1}(y_{i-1})$ $= \max_{y_{i-1}} A_{y_{i-1,s}} B_{s,x_{i}} score_{i-1}(y_{i-1})$

3. At the final state: calculate $\max_{y_{i-1}} P(y, x \mid \pi, A, B) = \max_{s} score_{n}(s)$

This only calculates the max. To get final answer (*argmax*):

- keep track of which state corresponds to the max at each step
- build the answer using these back pointers

Viterbi algorithm

Max-product algorithm for first order sequences

 π : Initial probabilities A: Transitions B: Emissions

1. Initial: For each state s, calculate $score_1(s) = P(s)P(x_1 | s) = \pi_s B_{x_1,s}$

- 2. Recurrence: For i = 2 to n, for every state s, calculate $score_{i}(s) = \max_{y_{i-1}} P(s \mid y_{i-1}) P(x_{i} \mid s) score_{i-1}(y_{i-1})$ $= \max_{y_{i-1}} A_{y_{i-1,s}} B_{s,x_{i}} score_{i-1}(y_{i-1})$
- 3. At the final state: calculate $\max_{y_{i-1}} P(y, x \mid \pi, A, B) = \max_{s} score_{n}(s)$

This only calculates the max. To get final answer (argmax):

- keep track of which state corresponds to the max at each step
- build the answer using these back pointers

Questions?⁷⁴

General idea

- Dynamic programming
 - The best solution for the full problem relies on best solution to sub-problems
 - Memoize partial computation
- Examples
 - Viterbi algorithm
 - Dijkstra's shortest path algorithm
 - ...

Viterbi algorithm as best path

Goal: To find the highest scoring path in this trellis



Viterbi algorithm as best path

Goal: To find the highest scoring path in this trellis



Complexity of inference

- Complexity parameters
 - Input sequence length: n
 - Number of states: K
- Memory
 - Storing the table: nK (scores for all states at each position)
- Runtime
 - At each step, go over pairs of states
 - O(nK²)

Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences

Assume that we know the number of states in the HMM

Two possible scenarios

Assume that we know the number of states in the HMM

Two possible scenarios

 We are given a data set D = {<x_i, y_i>} of sequences labeled with states

And we have to learn the parameters of the HMM (π , A, B)

Assume that we know the number of states in the HMM

Two possible scenarios

 We are given a data set D = {<x_i, y_i>} of sequences labeled with states

And we have to learn the parameters of the HMM (π , A, B)

Supervised learning with complete data

Assume that we know the number of states in the HMM

Two possible scenarios

 We are given a data set D = {<x_i, y_i>} of sequences labeled with states

And we have to learn the parameters of the HMM (π, A, B) Supervised learning with complete data

2. We are given only a collection of sequences $D = \{x_i\}$ And we have to learn the parameters of the HMM (π, A, B)

Assume that we know the number of states in the HMM

Two possible scenarios

 We are given a data set D = {<x_i, y_i>} of sequences labeled with states

And we have to learn the parameters of the HMM (π, A, B) Supervised learning with complete data

2. We are given only a collection of sequences $D = \{x_i\}$ And we have to learn the parameters of the HMM (π, A, B) Unsupervised learning, with incomplete data

EM algorithm and its siblings: a subsequent lecture

Assume that we know the number of states in the HMM

Two possible scenarios

We are given a data set D = {<x_i, y_i>} of sequences labeled with states
And we have to learn the parameters of the HMM (π, A, B)

Supervised learning with complete data

2. We are given only a collection of sequences $D = \{x_i\}$ And we have to learn the parameters of the HMM (π, A, B) Unsupervised learning, with incomplete data

EM algorithm and its siblings: a subsequent lecture

We are given a dataset $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$

Each x_i is a sequence of observations and y_i is a sequence of states that correspond to x_i

Goal: Learn initial, transition, emission distributions (π, A, B)

We are given a dataset $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$

Each x_i is a sequence of observations and y_i is a sequence of states that correspond to x_i

Goal: Learn initial, transition, emission distributions (π, A, B)

• How do we learn the parameters of the HMM?

We are given a dataset $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$

Each x_i is a sequence of observations and y_i is a sequence of states that correspond to x_i

Goal: Learn initial, transition, emission distributions (π, A, B)

• How do we learn the parameters of the HMM?

- The maximum likelihood principle Where have we seen this before?

We are given a dataset $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$

- Each \mathbf{x}_i is a sequence of observations and \mathbf{y}_i is a sequence of states that correspond to \mathbf{x}_i

Goal: Learn initial, transition, emission distributions (π, A, B)

How do we learn the parameters of the HMM?
The maximum likelihood principle

$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B)$$

We are given a dataset $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$

- Each \mathbf{x}_i is a sequence of observations and \mathbf{y}_i is a sequence of states that correspond to \mathbf{x}_i

Goal: Learn initial, transition, emission distributions (π, A, B)

How do we learn the parameters of the HMM?
The maximum likelihood principle

$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_{i} P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

We are given a dataset $D = \{\langle \mathbf{x}_i, \mathbf{y}_i \rangle\}$

- Each \mathbf{x}_i is a sequence of observations and \mathbf{y}_i is a sequence of states that correspond to \mathbf{x}_i

Goal: Learn initial, transition, emission distributions (π, A, B)

How do we learn the parameters of the HMM?
The maximum likelihood principle

$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_{i} \left[P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B) \right]$$

And we know how to write this in terms of the parameters of the HMM

$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_{i} P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

 (π, A, B) can be estimated separately just by counting – Makes learning simple and fast

$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_{i} P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

 (π, A, B) can be estimated separately just by counting - Makes learning simple and fast

Initial
$$\pi_s = \frac{\operatorname{count}(\operatorname{start} \to s)}{n}$$
 probabilities



$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_{i} P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

 (π, A, B) can be estimated separately just by counting

- Makes learning simple and fast



$$(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_{i} P(\mathbf{x}_i, \mathbf{y}_i | \pi, A, B)$$

 (π, A, B) can be estimated separately just by counting

Makes learning simple and fast



Supervised learning details $(\hat{\pi}, \hat{A}, \hat{B}) = \max_{\pi, A, B} P(D|\pi, A, B) = \max_{\pi, A, B} \prod_{i} P(\mathbf{x}_{i}, \mathbf{y}_{i}|\pi, A, B)$

 (π, A, B) can be estimated separately just by counting

- Makes learning simple and fast

Initial
probabilities
$$\pi_s = \frac{\operatorname{count}(\operatorname{start} \to s)}{n}$$
Transition
probabilities $A_{s',s} = \frac{\operatorname{count}(s \to s')}{\operatorname{count}(s)}$ Emission
probabilities $B_{s,x} = \frac{\operatorname{count}\begin{pmatrix}s\\\downarrow\\x\end{pmatrix}}{\operatorname{count}(s)}$

Exercise: Derive these using derivatives of the log likelihood. Requires Lagrangian multipliers.

Priors and smoothing

- Maximum likelihood estimation works best with lots of annotated data
 - Never the case
- Priors inject information about the probability distributions
 - Dirichlet priors for multinomial distributions
- Effectively additive smoothing
 - Add small constants to the counts

Hidden Markov Models summary

- Predicting sequences
 - As many output states as observations
- Markov assumption helps decompose the score
- Several algorithmic questions
 - Most likely state
 - Learning parameters
 - Supervised, Unsupervised
 - Probability of an observation sequence
 - Sum over all assignments to states, replace max with sum in Viterbi
 - Probability of state for each observation
 - Sum over all assignments to all other states

Questions?

Outline

- Sequence models
- Hidden Markov models
 - Inference with HMM
 - Learning
- Conditional Models and Local Classifiers
- Global models
 - Conditional Random Fields
 - Structured Perceptron for sequences