#### Hidden Markov Model

#### course based on Jurafsky and Martin [2009, Chap.6]



#### UPPSALA UNIVERSITET MARIE DUBREMETZ marie.dubremetz@lingfil.uu.se

Uppsala, 2015

# Why Markov Models? Historical Introduction Formal Introduction



### **Table of Contents**





### **Table of Contents**





# So far you have learnt how to compute ngram probabilities.(MLE+smoothing)

However, applications, like PoS tagging, require to compute and combine even more probabilities over the ngram probabilities. For that you need to learn the concept of Hidden Markov Models. Here is an history of how we came up with this mathematical model : https://www.youtube.com/watch?v=o-jdJxXL\_W4

So far you have learnt how to compute ngram probabilities.(MLE+smoothing) However, applications, like PoS tagging, require to compute and combine even more probabilities over the ngram probabilities. For that you need to learn the concept of Hidden Markov Models Here is an history of how we came up with this mathematical model, https://www.woutube.com/watch?v=o=id.lvXI\_\_W4

So far you have learnt how to compute ngram probabilities.(MLE+smoothing) However, applications, like PoS tagging, require to compute and combine even more probabilities over the ngram probabilities. For that you need to learn the concept of Hidden Markov Models. Here is an history of how we came up with this mathematical model : https://www.youtube.com/watch?v=o-jdJxXL\_W4

To sum up and relate to our general problems : Markov has demonstrated that you can build a conceptual machine that represents the generation of a succession of events even if those events are themselves dependent of another succession of other events happening in the background.

Instead of using it on meteorological events, letters or pearls we will use it on PoS Tags and words.

### **Table of Contents**







### Hidden Markov Models

- Markov models are probabilistic sequence models used for problems such as:
  - 1. Speech recognition
  - 2. Spell checking
  - 3. Part-of-speech tagging
  - 4. Named entity recognition
- A Markov model runs through a sequence of states emitting observable signals
- If the state sequence cannot be determined from the observation sequence, the model is said to be hidden





#### **Markov Assumptions**

State transitions are assumed to be independent of everything except the current state:

$$P(q_1,\ldots,q_n)=\prod_{i=1}^n P(q_i\mid q_{i-1})$$

Signal emissions are assumed to be independent of everything except the current state:

$$P(q_1,\ldots,q_n,s_1,\ldots,s_n)=P(q_1,\ldots,q_n)\prod_{i=1}^n P(s_i\mid q_i)$$

NB: subscripts on states and signals refer to sequence positions

### More Formally

$$Q = q_1 q_2 \dots q_N$$
$$A = a_{01} a_{02} \dots a_{n1} \dots a_{nn}$$

$$O = o_1 o_2 \dots o_N$$

$$B = b_i(o_t)$$

 $q_0, q_{end}$ 

#### a set of **states**

a **transition probability matrix** *A*, each  $a_{ij}$  representing the probability of moving from state *i* to state *j*, s.t.  $\sum_{j=1}^{n} a_{ij} = 1 \quad \forall i$ 

a set of **observations**, each one drawn from a vocabulary  $V = v_1, v_2, ..., v_V$ .

A set of **observation likelihoods:**, also called **emission probabilities**, each expressing the probability of an observation  $o_t$  being generated from a state *i*.

a special **start and end state** which are not associated with observation.

#### Exercise

The next slide describes a HMM

According to the formalism, color/circle what represents :

- A in red
- B in green
- V in dotted green
- Q in blue

#### • $q_0$ and $q_{end}$ (if not missing) in dotted blue

$Q = q_1 q_2 \dots q_N$	a set of <b>states</b>
$A = a_{01}a_{02}\ldots a_{n1}\ldots a_{nn}$	a <b>transition probability matrix</b> <i>A</i> , each $a_{ij}$ representing the probability of moving from state <i>i</i> to state <i>j</i> , s.t. $\sum_{j=1}^{n} a_{ij} = 1  \forall i$
$O = o_1 o_2 \dots o_N$	a set of <b>observations</b> , each one drawn from a vo- cabulary $V = v_1, v_2,, v_V$ .
$B = b_i(o_t)$	A set of <b>observation likelihoods:</b> , also called <b>emission probabilities</b> , each expressing the probability of an observation $o_t$ being generated from a state <i>i</i> .
$q_0, q_{end}$	a special <b>start and end state</b> which are not asso- ciated with observation.



### A Simple First-Order HMM for Tagging



### Tasks on HMM

There is three kinds of problems when you deal with Hidden Markov Models :

**Problem 1 (Computing Likelihood):** Given an HMM  $\lambda = (A, B)$  and an observation sequence *O*, determine the likelihood  $P(O|\lambda)$ .

**Problem 2 (Decoding):** Given an observation sequence *O* and an HMM  $\lambda = (A, B)$ , discover the best hidden state sequence *Q*.

**Problem 3 (Learning):** Given an observation sequence *O* and the set of states in the HMM, learn the HMM parameters *A* and *B*.

Write this down for the quiz!

#### Question

Can you find which sentence answers to which HMM Problem?

- The probability of this HMM to generate the sequence of observation O=[he,eats,cakes] is 0.34859
- One transition probabilities are A=[0.6, 0.4...] and the observations likelyhood are B=[0.9,0.1...]
- The most probable sequence of states that has generated O, is the state sequence Q=[Start, Noun, Verb, Adj, End]
- a Decoding
- b Likelihood
- c Learning

### **Table of Contents**

# Why Markov Models? Historical Introduction Formal Introduction



#### Link to the slides on Viterbi :

The following slides (with a blue background) are extracted from this course :

http://courses.washington.edu/ling570/gina\_fall11/
slides/ling570\_class12\_viterbi.pdf

(Slides designed by Fei Xia and Gina-Anne Levow used with their kind authorisation.)

We will learn formally what is the Viterbi algorithm. But don't hesitate to look at this link if you want more details about Viterbi implementation. Or if you would like to continue beyond what we will do in the next exercise.

### Three Problems for HMMs

- Likelihood:
  - Find the probability of an observation sequence given a model
    - Forward algorithm
- Decoding:
  - Find the most likely path through a model given an observed sequence
    - Viterbi algorithm
- Learning:
  - Find the most likely model (parameters) given an observed sequence
    - Supervised (MLE) or unsupervised Baum-Welch

- Have complete model of ngram POS tagging
  - Need to compute

$$\hat{t}_1^n = \operatorname*{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

Possible approach:

Have complete model of ngram POS tagging

• Need to compute

$$\hat{t}_1^n = \operatorname*{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- Possible approach:
  - Enumerate all paths through HMM, pick highest score
  - Good idea?

- Have complete model of ngram POS tagging
  - Need to compute

$$\hat{t}_1^n = \operatorname*{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- Possible approach:
  - Enumerate all paths through HMM, pick highest score
  - Good idea? No. Why?

- Have complete model of ngram POS tagging
  - Need to compute

$$\hat{t}_1^n = \operatorname*{argmax}_{t_1^n} P(t_1^n | w_1^n)$$

- Possible approach:
  - Enumerate all paths through HMM, pick highest score
  - Good idea? No. Why?
    - Computationally intractable
  - Dynamic programming can help!

# Example

• time flies like an arrow

# HMM Model

\start_state	\emission
0	BOS <s> 1.0</s>
\transition 0 BOS 1.0	N time 0.1
BOS N 0.5 BOS DT 0.4	V time 0.1
BOS V 0.1 DT N 1.0	N flies 0.1
N N 0.2 N V 0.7	V flies 0.2
N P 0.1	V like 0.2
V N 0.4	P like 0.1
V V 0.1	DT an 0.3
P DT 0.6 P N 0.4	N arrow 0.1



















- Find best hidden state sequence given observations and model
- Each cell represents:
  - Probability of being in state j after first t observations, passing through most probable sequence in model λ

• Formally,  $v_t(j) = \max_{q_0,q_1,q_2,...,q_{t-1}} P(q_0,q_1,...q_{t-1},o_1,o_2,...o_t,q_t = j \mid \lambda)$ 

# Viterbi

• Initialization:  $v_1(j) = a_{0j}b_j(o_1)$  $bt_1(j) = 0$ 

### Viterbi

• Initialization:  $v_1(j) = a_{0j}b_j(o_1)$   $bt_1(j) = 0$ • Recursion:  $v_t(j) = \max_{i=1}^N v_{t-1}(i)a_{ij}b_j(o_t), 1 \le j \le N, 1 < t \le T$  $bt_t(j) = \operatorname*{argmax}_{i=1}^N v_{t-1}(i)a_{ij}b_j(o_t), 1 \le j \le N, 1 < t \le T$ 

### Viterbi

- Initialization:  $v_1(j) = a_{0j}b_j(o_1)$  $bt_1(j) = 0$
- Recursion:  $v_t(j) = \max_{i=1}^{N} v_{t-1}(i)a_{ij}b_j(o_t), 1 \le j \le N, 1 < t \le T$   $bt_t(j) = rgmax_{i=1}^{N} v_{t-1}(i)a_{ij}b_j(o_t), 1 \le j \le N, 1 < t \le T$ • Termination:  $P^* = v_t(q_F) = \max_{i=1}^{N} v_T(i)a_{iF}$

$$q_T^* = bt_T(q_F) = \operatorname*{argmax}_{i=1} v_T(i)a_{iF}$$

	1	2	3	4	5	
Ν	0					
V	0					
Ρ	0					
D	0					
BOS	P(BOS 0)* P( <s> BOS) =1.0</s>					

	1	2	3	4	5	
N	0	[BOS,1]* P(N BOS)* P(time N) =				
V	0					
Ρ	0					
D	0					
BOS	1.0 0					

	1	2	3	4	5	
N	0	[BOS,1]* P(N BOS)* P(time N) =				
V	0	[BOS,1]* P(V BOS)* P(time V) =				
Ρ	0					
D	0					
BOS	1.0 0					

	1	2	3	4	5	
N	0	[BOS,1]* P(N BOS)* P(time N) =0.05				
V	0	[BOS,1]* P(V BOS)* P(time V) =0.01				
Ρ	0	0				
D	0	0				
BOS	1.0 0	0				

	1	2	3	4	5	
N	0	0.05 BOS				
V	0	0.01 BOS				
Ρ	0	0				
D	0	0				
BOS	1.0 0	0				

	1	2	3	4	5	
Ν	0	0.05 BOS	max(([N,2]*P(N N),[V, 2]*P(N V))*P(flies N)=			
V	0	0.01 BOS	max(([V,2]*P(V V),[N, 2]*P(V N))*P(flies V)=			
Ρ	0	0				
D	0	0				
BOS	1.0 0	0				

#### Exercise

Make the calculations of the 3rd column.

#### After computing 3rd column what do we learn?

At column 3 we know that we should not consider other sequences than the one starting by O,BOS,N,N or O,BOS,N,V. For instance we don't need to compute O,BOS,V,N,N...

- Create an array
  - With columns corresponding to inputs
  - Rows corresponding to possible states

- Create an array
  - With columns corresponding to inputs
  - Rows corresponding to possible states
- Sweep through the array in one pass
  - Fill columns left-to-right based on transition/emission

- Create an array
  - With columns corresponding to inputs
  - Rows corresponding to possible states
- Sweep through the array in one pass
  - Fill columns left-to-right based on transition/emission
- Dynamic programming key

- Create an array
  - With columns corresponding to inputs
  - Rows corresponding to possible states
- Sweep through the array in one pass
  - Fill columns left-to-right based on transition/emission
- Dynamic programming key
  - Store maximum probability of path to each cell

- Create an array
  - With columns corresponding to inputs
  - Rows corresponding to possible states
- Sweep through the array in one pass
  - Fill columns left-to-right based on transition/emission
- Dynamic programming key
  - Store maximum probability of path to each cell
  - Store backpointers to recover path

#### Conclusion

#### **Overall Summary**

- HMM relate a sequence of observation to a sequence of hidden states.
- The process of discovering the sequence of hidden states given a sequence of observation is called decoding
- Viterbi algorithm is an efficient way to perform decoding.

Daniel Jurafsky and James H Martin. Speech and Language Processing : An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition, volume 163 of Prentice Hall Series in Artificial Intelligence. Prentice Hall, 2009.

The blue slides :

http://courses.washington.edu/ling570/gina\_fall11/
slides/ling570\_class12\_viterbi.pdf
(Slides designed by Fei Xia and Gina-Anne Levow used with their

kind authorisation.)