

# *Stochastic models for learning language models (Part 1)*

**R. Basili**

*Deep Learning*  
a.a. 2023-24

March 20, 2024

# Outline

## Outline

- 1 *Probability and Language Modeling*
  - Motivations
  - Probability Models for Natural Language
- 2 *Introduction to Markov Models*
  - Hidden Markov Models
  - Advantages
- 3 *HMM and POS tagging*
  - Forward Algorithm and Viterbi
  - About Parameter Estimation for POS
- 4 *References*
- 5 *Exercises*

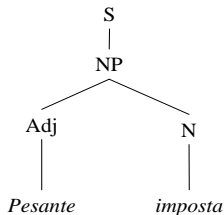
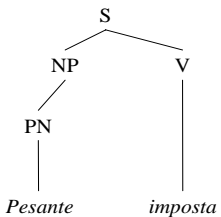
# *Quantitative Models of language structures*

Linguistic structures exhibit syntagmatic information that is crucial for machine learning in Web Mining. The common grammatical modeling framework is the one of (phrase structure) grammars, that can produce often ambiguous readings:

1. S → NP V
2. S → NP
3. NP → PN
4. NP → N
5. NP → Adj N
6. N → "imposta"
7. V → "imposta"
8. Adj → "pesante"
9. PN → "Pesante"
- ...

# The role of Quantitative Approaches

“*Pesante imposta*”



# The role of Quantitative Approaches

Weighted grammars are models of (possibly limited) *degrees of grammaticality*. They are meant to deal with a large range of ambiguity problems:

- 1.     S    -> NP   V                     .7
- 2.     S    -> NP                         .3
- 3.     NP  -> PN                         .1
- 4.     NP  -> N                          .6
- 5.     NP  -> Adj N                      .3
- 6.     N    -> imposta                  .6
- 7.     V    -> imposta                  .4
- 8.     Adj -> Pesante                  .8
- 9.     PN  -> Pesante                  .2



# *Linguistic Ambiguity and weighted grammars*

Weighted grammars allow to compute the degree of grammaticality of different ambiguous derivations, thus supporting disambiguation:

# Linguistic Ambiguity and weighted grammars

Weighted grammars allow to compute the degree of grammaticality of different ambiguous derivations, thus supporting disambiguation:

- 1. S → NP V .7
- 2. S → NP .3
- 3. NP → PN .1
- 4. NP → N .6
- 5. NP → Adj N .3
- 6. N → imposta .6
- 7. V → imposta .4
- 8. Adj → Pesante .8
- 9. PN → Pesante .2
- ...

$$\text{prob}(((\text{Pesante})_{PN} (\text{imposta})_V)_S) = (.7 \cdot .1 \cdot .2 \cdot .4) = 0.0084$$



# Linguistic Ambiguity and weighted grammars

Weighted grammars allow to compute the degree of grammaticality of different ambiguous derivations, thus supporting disambiguation:

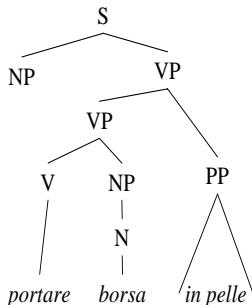
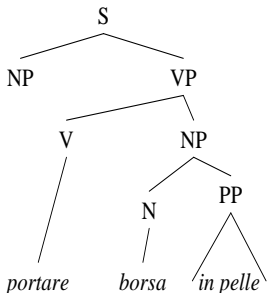
- 1. S → NP V .7
- 2. S → NP .3
- 3. NP → PN .1
- 4. NP → N .6
- 5. NP → Adj N .3
- 6. N → imposta .6
- 7. V → imposta .4
- 8. Adj → Pesante .8
- 9. PN → Pesante .2
- ...

$$\text{prob}(((\text{Pesante})_{PN} (\text{imposta})_V)_S) = (.7 \cdot .1 \cdot .2 \cdot .4) = 0.0084$$

$$\text{prob}(((\text{Pesante})_{Adj} (\text{imposta})_N)_S) = (.3 \cdot .3 \cdot .8 \cdot .6) = 0.0432$$

# Syntactic Disambiguation

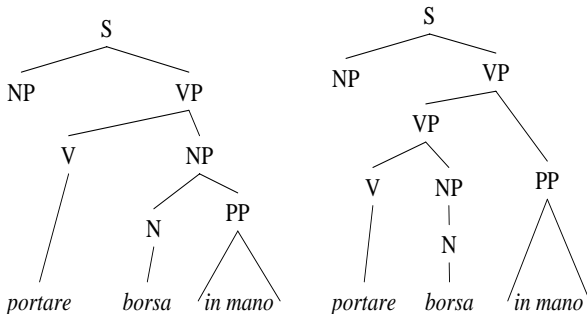
“portare borsa in pelle”



Derivation Trees for a structurally ambiguous sentence

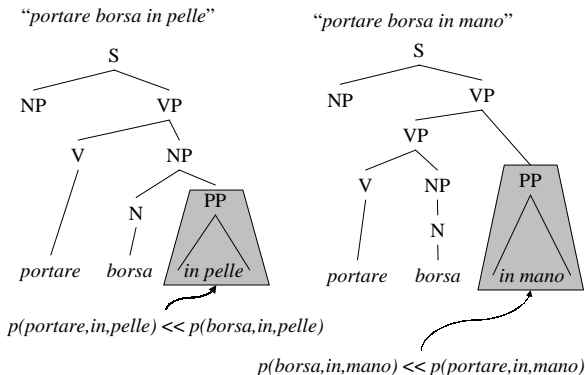
# Syntactic Disambiguation (cont'd)

“portare borsa in mano”



Derivation Trees for a second structurally ambiguous sentence.

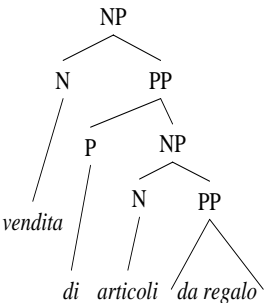
# Structural Disambiguation (cont'd)



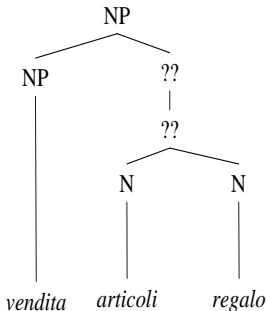
Disambiguation of structural ambiguity.

# Tolerance to errors

“vendita di articoli da regalo”



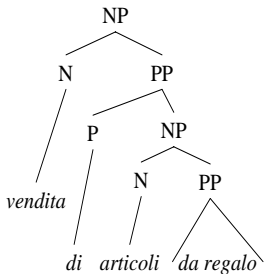
“vendita articoli regalo”



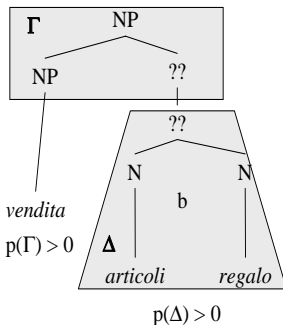
An example of ungrammatical but meaningful sentence

# Error tolerance (cont'd)

“vendita di articoli da regalo”



“vendita articoli regalo”





# Probability and Language Modeling

- Aims
  - to extend grammatical (i.e. rule-based) models with predictive and disambiguation capabilities
  - to offer theoretically well founded *inductive methods*
  - to develop (not merely) quantitative models of linguistic phenomena
- Methods and Resources:
  - Mathematical theories (e.g. Markov models)
  - Systematic testing/evaluation frameworks
  - Extended repositories of examples of *language in use*
  - Traditional linguistic resources (e.g. "models" like dictionaries)



# *Probability and Language Modeling*

- Signals are abstracted via symbols that are not known in advance



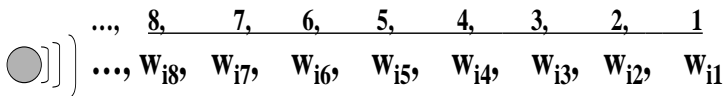
# Probability and Language Modeling

- Signals are abstracted via symbols that are not known in advance
- Emitted signals belong to an alphabet  $A$
- Time is discrete: each time point corresponds to an emitted signal



# Probability and Language Modeling

- Signals are abstracted via symbols that are not known in advance
- Emitted signals belong to an alphabet  $A$
- Time is discrete: each time point corresponds to an emitted signal
- Sequences of symbols  $(w_1, \dots, w_n)$  correspond to sequences of time points  $(1, \dots, n)$



# Probability and Language Modeling

## *A generative language model*

A random variable  $X$  can be introduced so that

# Probability and Language Modeling

## *A generative language model*

A random variable  $X$  can be introduced so that

- It assumes values  $w_i$  in the alphabet  $A$
- Probability is used to describe the uncertainty on the emitted signal

$$p(X = w_i) \quad w_i \in A$$

# Probability and Language Modeling

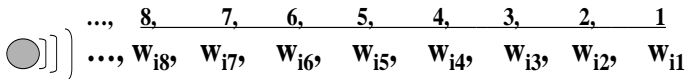
- A random variable  $X$  can be introduced so that
  - $X$  assumes values in  $A$  at each step  $i$ , i.e.  $X_i = w_j$
  - probability is  $p(X_i = w_j)$



# Probability and Language Modeling

- A random variable  $X$  can be introduced so that
  - $X$  assumes values in  $A$  at each step  $i$ , i.e.  $X_i = w_j$
  - probability is  $p(X_i = w_j)$
- Constraints: the total probability is for each step:

$$\sum_j p(X_i = w_j) = 1 \quad \forall i$$

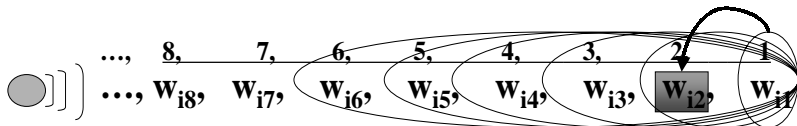


# Probability and Language Modeling

- Notice that time points can be represented as **states** of the emitting source
- An output  $w_i$  can be considered as emitted in a *given state*  $X_i$  by the source, and *given a certain history*

# Probability and Language Modeling

- Notice that time points can be represented as **states** of the emitting source
- An output  $w_i$  can be considered as emitted in a *given state*  $X_i$  by the source, and *given a certain history*



# Probability and Language Modeling

- Formally:

- $P(X_i = w_i, X_{i-1} = w_{i-1}, \dots, X_1 = w_1) =$

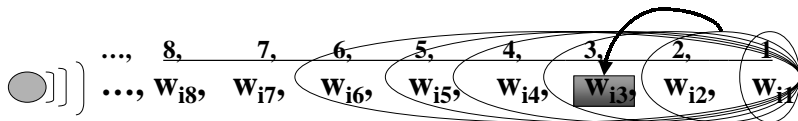
# Probability and Language Modeling

- Formally:

- $$P(X_i = w_i, X_{i-1} = w_{i-1}, \dots, X_1 = w_1) =$$

$$= P(X_i = w_i | X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2}, \dots, X_1 = w_1) \cdot$$

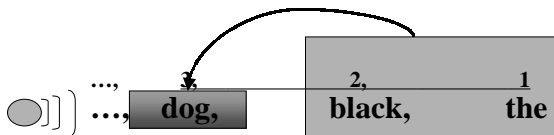
$$P(X_{i-1} = w_{i-1}, X_{i-2} = w_{i-2}, \dots, X_1 = w_1)$$



# Probability and Language Modeling

## What's in a state

$n - 1$  preceding words  $\Rightarrow$   **$n$ -gram language models**

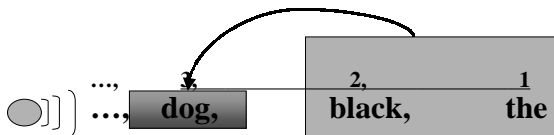


$$p(\text{the, black, dog}) = p(\text{dog} | \text{the, black})$$

# Probability and Language Modeling

## What's in a state

$n - 1$  preceding words  $\Rightarrow$   **$n$ -gram language models**

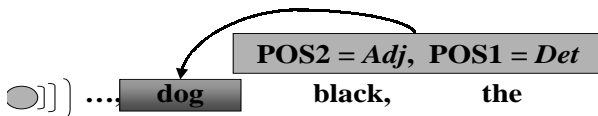


$$p(\text{the, black, dog}) = p(\text{dog}|\text{the, black})p(\text{black}|\text{the})p(\text{the})$$

# Probability and Language Modeling

*What's in a state*

preceding POS tags  $\Rightarrow$  **stochastic taggers**

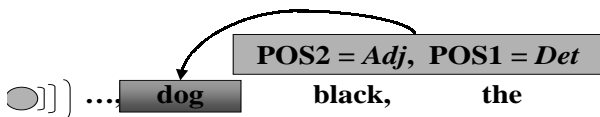




# Probability and Language Modeling

## What's in a state

preceding POS tags  $\Rightarrow$  **stochastic taggers**



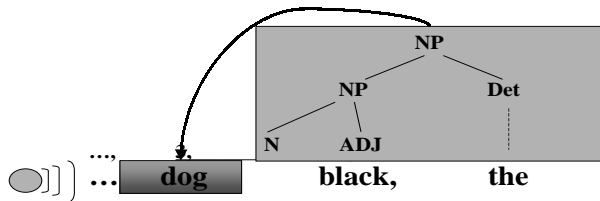
$$p(\text{the}_{DT}, \text{black}_{ADJ}, \text{dog}_N) = p(\text{dog}_N | \text{the}_{DT}, \text{black}_{ADJ}) \dots$$



# Probability and Language Modeling

*What's in a state*

preceding *parses*  $\Rightarrow$  **stochastic grammars**

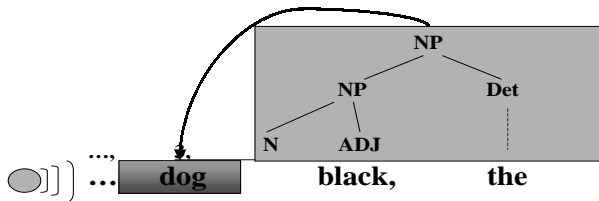


$$p((the_{Det}, (black_{ADJ}, dog_N)_{NP})_{NP}) =$$

# Probability and Language Modeling

*What's in a state*

preceding *parses*  $\Rightarrow$  **stochastic grammars**



$$p((the_{Det}, (black_{ADJ}, dog_N)_{NP})_{NP}) = \\
 p(dog_N | ((the_{Det}), (black_{ADJ}, -))) \dots$$



# Probability and Language Modeling (2)

- Expressivity
  - The predictivity of a statistical grammar can provide a very good explanatory model of the source language (string)
  - Acquiring information from data has a clear definition, with simple and sound induction algorithms
  - Simple but richer descriptions (e.g. grammatical preferences)
  - Optimal Coverage (i.e. better on *more important phenomena*)
- Integrating Linguistic Description
  - Start with poor assumptions and approximate as much as possible *what is known* (early evaluate only performance)
  - *Bias* the statistical model since the beginning and check the results on a *linguistic ground*



# Probability and Language Modeling (3)

## *Advantages: Performances*

- Faster Processing (e.g. through the pruning of the algorithmic search space)
- Faster Design (i.e. **one** probabilistic model for **multiple** tasks)



# Probability and Language Modeling (3)

## *Advantages: Performances*

- Faster Processing (e.g. through the pruning of the algorithmic search space)
- Faster Design (i.e. **one** probabilistic model for **multiple** tasks)
- Linguistic Adequacy
  - Acceptance
  - Psychological Plausibility
  - Explanatory power

# Probability and Language Modeling (3)

## *Advantages: Performances*

- Faster Processing (e.g. through the pruning of the algorithmic search space)
- Faster Design (i.e. **one** probabilistic model for **multiple** tasks)
- Linguistic Adequacy
  - Acceptance
  - Psychological Plausibility
  - Explanatory power
- Tools for further analysis of Linguistic Data

# Markov Models

## Markov Models

Suppose  $X_1, X_2, \dots, X_T$  form a sequence of random variables taking values in a countable set  $W = p_1, p_2, \dots, p_N$  (State space).

- Limited Horizon Property:

$$P(X_{t+1} = p_k | X_1, \dots, X_t) = P(X_{t+1} = k | X_t)$$

# Markov Models

## Markov Models

Suppose  $X_1, X_2, \dots, X_T$  form a sequence of random variables taking values in a countable set  $W = p_1, p_2, \dots, p_N$  (State space).

- Limited Horizon Property:

$$P(X_{t+1} = p_k | X_1, \dots, X_t) = P(X_{t+1} = p_k | X_t)$$

- Time invariant:

$$P(X_{t+1} = p_k | X_t = p_l) = P(X_2 = p_k | X_1 = p_l) \quad \forall t (> 1)$$

# Markov Models

## Markov Models

Suppose  $X_1, X_2, \dots, X_T$  form a sequence of random variables taking values in a countable set  $W = p_1, p_2, \dots, p_N$  (State space).

- Limited Horizon Property:

$$P(X_{t+1} = p_k | X_1, \dots, X_t) = P(X_{t+1} = p_k | X_t)$$

- Time invariant:

$$P(X_{t+1} = p_k | X_t = p_l) = P(X_2 = p_k | X_1 = p_l) \quad \forall t (> 1)$$

It follows that the sequence of  $X_1, X_2, \dots, X_T$  is a **Markov chain**.





# Representation of a Markov Chain

## Graphical Representation (i.e. Automata)

- States as nodes with names
- Transitions from states  $i$ -th and  $j$ -th as arcs labelled by conditional probabilities  $P(X_{t+1} = p_j | X_t = p_i)$

Note that 0 probability arcs are omitted from the graph.

	$S_1$	$S_2$
$S_1$	0.70	0.30
$S_2$	0.50	0.50







# A Simple Example of Hidden Markov Model

## Crazy Coffee Machine

Assume, for example, the following state transition model:

	<i>TP</i>	<i>CP</i>
<i>TP</i>	0.70	0.30
<i>CP</i>	0.50	0.50

and let *CP* be the starting state (i.e.  $\pi_{CP} = 1, \pi_{TP} = 0$ ).

Potential Use:

- 1 What is the probability at time step 3 to be in state *TP*?
- 2 What is the probability at time step *n* to be in state *TP*?
- 3 What is the probability of the following sequence in output: (*Coffee, Tea, Coffee*)?







# Crazy Coffee Machine

*Solution to Problem 3:*

$$\begin{aligned} P(\text{Cof}, \text{Tea}, \text{Cof}) &= \\ &= P(\text{Cof}) \cdot P(\text{Tea}|\text{Cof}) \cdot P(\text{Cof}|\text{Tea}) = 1 \cdot 0.5 \cdot 0.3 = 0.15 \end{aligned}$$

# A Simple Example of Hidden Markov Model (2)

## Crazy Coffee Machine

- **Hidden** Markov model: If the machine output *Tea*, *Coffee* or *Capuccino* **independently** from *CP* and *TP*.

What we need is a description of the random event of output(ting) a drink.



# Crazy Coffee Machine

A description of the random event of output(ing) a drink. Formally we need (for each time step  $n$  and for each kind of output  $O = \{Tea, Cof, Cap\}$ ), the following conditional probabilities:

$$P(O_n = o_k | X_n = p_i, X_{n+1} = p_j)$$

where  $o_k \in \{Tea, Coffee, Capuccino\}$ . This matrix is called the **output matrix** of the machine (or of its Hidden markov Model).

# A Simple Example of Hidden Markov Model (2)

## Crazy Coffee Machine

Given the following output probability for the machine

	Tea	Coffee	Capuccino
TP	0.8	0.2	0.0
CP	0.15	0.65	0.2

and let  $CP$  be the starting state (i.e.  $\pi_{CP} = 1$ ,  $\pi_{TP} = 0$ ).

- Find the following probabilities of output from the machine

# A Simple Example of Hidden Markov Model (2)

## Crazy Coffee Machine

Given the following output probability for the machine

	Tea	Coffee	Capuccino
TP	0.8	0.2	0.0
CP	0.15	0.65	0.2

and let  $CP$  be the starting state (i.e.  $\pi_{CP} = 1, \pi_{TP} = 0$ ).

- Find the following probabilities of output from the machine
  - $(Cappuccino, Coffee)$  given that the state sequence is  $(CP, TP, TP)$

# A Simple Example of Hidden Markov Model (2)

## Crazy Coffee Machine

Given the following output probability for the machine

	Tea	Coffee	Capuccino
TP	0.8	0.2	0.0
CP	0.15	0.65	0.2

and let  $CP$  be the starting state (i.e.  $\pi_{CP} = 1$ ,  $\pi_{TP} = 0$ ).

- Find the following probabilities of output from the machine
  - $(Cappuccino, Coffee)$  given that the state sequence is  $(CP, TP, TP)$
  - $(Tea, Coffee)$  for any state sequence

# A Simple Example of Hidden Markov Model (2)

## Crazy Coffee Machine

Given the following output probability for the machine

	Tea	Coffee	Capuccino
TP	0.8	0.2	0.0
CP	0.15	0.65	0.2

and let  $CP$  be the starting state (i.e.  $\pi_{CP} = 1$ ,  $\pi_{TP} = 0$ ).

- Find the following probabilities of output from the machine
  - 1  $(Cappuccino, Coffee)$  given that the state sequence is  $(CP, TP, TP)$
  - 2  $(Tea, Coffee)$  for any state sequence
  - 3 a generic output  $O = (o_1, \dots, o_n)$  for *any* state sequence















# A Simple Example of Hidden Markov Model (2)

Solutions for the problem 2

In general, for any sequence of three states  $X = (X_1, X_2, X_3)$

$$P(\text{Tea}, \text{Cof} | X_1, X_2, X_3) =$$

$P(\text{Tea}, \text{Cof}) =$  (as sequences are a partition for the sample space)

$$= \sum_{X_1, X_2, X_3} P(\text{Tea}, \text{Cof} | X_1, X_2, X_3) P(X_1, X_2, X_3) \text{ where}$$

$$P(\text{Tea}, \text{Cof} | X_1, X_2, X_3) = P(\text{Tea} | X_1, X_2) P(\text{Cof} | X_2, X_3) =$$

(for the simplified model of the coffee machine )

$$= P(\text{Tea} | X_1) P(\text{Cof} | X_2) \text{ and (for the Markov constraint)}$$

$$P(X_1, X_2, X_3) = P(X_1) P(X_2 | X_1) P(X_3 | X_2)$$

# A Simple Example of Hidden Markov Model (2)

Solutions for the problem 2

In general, for any sequence of three states  $X = (X_1, X_2, X_3)$

$$P(\text{Tea}, \text{Cof} | X_1, X_2, X_3) =$$

$P(\text{Tea}, \text{Cof})$  = (as sequences are a partition for the sample space)

$$= \sum_{X_1, X_2, X_3} P(\text{Tea}, \text{Cof} | X_1, X_2, X_3) P(X_1, X_2, X_3) \text{ where}$$

$$P(\text{Tea}, \text{Cof} | X_1, X_2, X_3) = P(\text{Tea} | X_1, X_2) P(\text{Cof} | X_2, X_3) =$$

(for the simplified model of the coffee machine )

$$= P(\text{Tea} | X_1) P(\text{Cof} | X_2) \text{ and (for the Markov constraint)}$$

$$P(X_1, X_2, X_3) = P(X_1) P(X_2 | X_1) P(X_3 | X_2)$$

The simplified model is concerned with only the following transition chains

$(CP, CP, CP), (CP, TP, CP), (CP, CP, TP)$

$(CP, TP, TP)$





# A Simple Example of Hidden Markov Model (2)

## Solutions for the problem 2

In general, for any sequence of three states  $X = (X_1, X_2, X_3)$

The following probability is given

$$P(\text{Tea}, \text{Cof}) =$$

$$\begin{aligned}
 &P(\text{Tea}|\text{CP})P(\text{Cof}|\text{CP})P(\text{CP})P(\text{CP}|\text{CP})P(\text{CP}|\text{CP}) + && \text{st.: } (\text{CP}, \text{CP}, \text{CP}) \\
 &P(\text{Tea}|\text{CP})P(\text{Cof}|\text{TP})P(\text{CP})P(\text{TP}|\text{CP})P(\text{CP}|\text{TP}) + && \text{st.: } (\text{CP}, \text{TP}, \text{CP}) \\
 &P(\text{Tea}|\text{CP})P(\text{Cof}|\text{CP})P(\text{CP})P(\text{CP}|\text{CP})P(\text{TP}|\text{CP}) + && \text{st.: } (\text{CP}, \text{CP}, \text{TP}) \\
 &P(\text{Tea}|\text{CP})P(\text{Cof}|\text{TP})P(\text{CP})P(\text{TP}|\text{CP})P(\text{TP}|\text{TP}) = && \text{st.: } (\text{CP}, \text{TP}, \text{TP})
 \end{aligned}$$

$$\begin{aligned}
 &= 0.15 \cdot 0.65 \cdot 1 \cdot 0.5 \cdot 0.5 + \\
 &+ 0.15 \cdot 0.2 \cdot 1 \cdot 0.5 \cdot 0.3 + \\
 &+ 0.15 \cdot 0.65 \cdot 1 \cdot 0.5 \cdot 0.5 + \\
 &+ 0.15 \cdot 0.2 \cdot 1.0 \cdot 0.5 \cdot 0.7 =
 \end{aligned}$$

$$\begin{aligned}
 &= 0.024375 + 0.0045 + 0.024375 + 0.0105 = \\
 &= 0.06375
 \end{aligned}$$









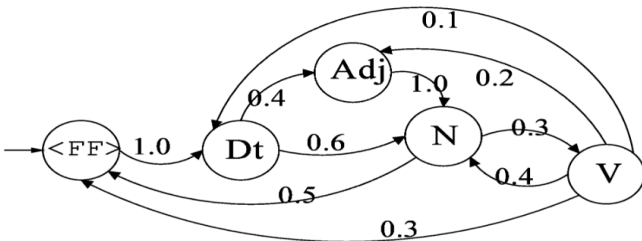






# The task of POS tagging

## An example



Emission

probabilities	.	the	this	cat	kid	eats	runs	fish	fresh	little	big
<FF>	1.0										
Dt		0.6	0.4								
N				0.6	0.1			0.3			
V						0.7	0.3				
Adj									0.3	0.3	0.4

# HMM and POS tagging

Given a sequence of morphemes  $w_1, \dots, w_n$  with ambiguous syntactic descriptions (i.e. part-of-speech tags), derive the sequence of  $n$  POS tags  $t_1, \dots, t_n$  that maximizes the following probability:

$$P(w_1, \dots, w_n, t_1, \dots, t_n)$$

that is

$$(t_1, \dots, t_n) = \operatorname{argmax}_{pos_1, \dots, pos_n} P(w_1, \dots, w_n, pos_1, \dots, pos_n)$$



# HMM and POS tagging

Given a sequence of morphemes  $w_1, \dots, w_n$  with ambiguous syntactic descriptions (i.e. part-of-speech tags), derive the sequence of  $n$  POS tags  $t_1, \dots, t_n$  that maximizes the following probability:

$$P(w_1, \dots, w_n, t_1, \dots, t_n)$$

that is

$$(t_1, \dots, t_n) = \operatorname{argmax}_{pos_1, \dots, pos_n} P(w_1, \dots, w_n, pos_1, \dots, pos_n)$$

Note that this is equivalent to the following:

$$(t_1, \dots, t_n) = \operatorname{argmax}_{pos_1, \dots, pos_n} P(pos_1, \dots, pos_n | w_1, \dots, w_n)$$

as:  $\frac{P(w_1, \dots, w_n, pos_1, \dots, pos_n)}{P(w_1, \dots, w_n)} = P(pos_1, \dots, pos_n | w_1, \dots, w_n)$

and  $P(w_1, \dots, w_n)$  is the same for all the sequences  $(pos_1, \dots, pos_n)$ .

# HMM and POS tagging

## How to map a POS tagging problem into a HMM

The above problem

$$(t_1, \dots, t_n) = \operatorname{argmax}_{pos_1, \dots, pos_n} P(pos_1, \dots, pos_n | w_1, \dots, w_n)$$

can be also written (Bayes law) as:

$$\operatorname{argmax}_{pos_1, \dots, pos_n} P(w_1, \dots, w_n | pos_1, \dots, pos_n) P(pos_1, \dots, pos_n)$$

# HMM and POS tagging

The HMM Model of POS tagging:

- **HMM States are mapped into POS tags** ( $t_i$ ), so that
$$P(t_1, \dots, t_n) = P(t_1)P(t_2|t_1)\dots P(t_n|t_{n-1})$$
- **HMM Output symbols are words**, so that
$$P(w_1, \dots, w_n | t_1, \dots, t_n) = \prod_{i=1}^n P(w_i | t_i)$$
- Transitions represent moves from one word to another

Note that *the Markov assumption is used*

- to model probability of a tag in position  $i$  (i.e.  $t_i$ ) only by means of the preceding part-of-speech (i.e.  $t_{i-1}$ )
- to model probabilities of words (i.e.  $w_i$ ) based only on the tag ( $t_i$ ) appearing in that position ( $i$ ).

# HMM and POS tagging

The final equation is thus:

$$(t_1, \dots, t_n) = \operatorname{argmax}_{t_1, \dots, t_n} P(t_1, \dots, t_n | w_1, \dots, w_n) = \\ \operatorname{argmax}_{t_1, \dots, t_n} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$



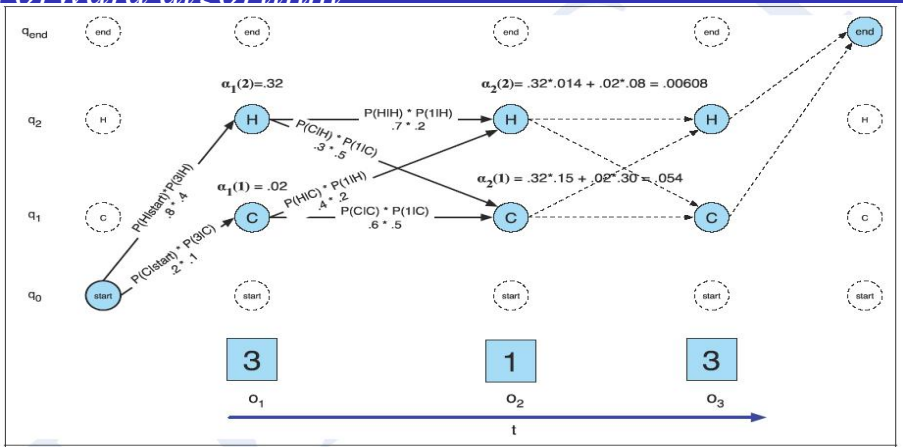








# Forward algorithm



**Figure 6.6** The forward trellis for computing the total observation likelihood for the ice-cream events 3 1 3. Hidden states are in circles, observations in squares. White (unfilled) circles indicate illegal transitions. The figure shows the computation of  $\alpha_t(j)$  for two states at two time steps. The computation in each cell follows Eq. 6.11:  $\alpha_t(j) = \sum_{i=1}^{N-1} \alpha_{t-1}(i) a_{ij} b_j(o_t)$ . The resulting probability expressed in each cell is Eq. 6.10:  $\alpha_t(j) = P(o_1, o_2 \dots o_t, q_t = j | \lambda)$ .

# HMM and POS tagging: Forward Algorithm

```

function FORWARD(observations of len  $T$ , state-graph) returns forward-probability
    num-states  $\leftarrow$  NUM-OF-STATES(state-graph)
    Create a probability matrix forward[num-states+2,  $T$ +2]
    forward[0,0]  $\leftarrow$  1.0
    for each time step  $t$  from 1 to  $T$  do
        for each state  $s$  from 1 to num-states do
            forward[ $s,t$ ]  $\leftarrow$   $\sum_{1 \leq s' \leq \text{num-states}}$  forward[ $s',t-1$ ] *  $a_{s',s}$  *  $b_s(o_t)$ 
    return the sum of the probabilities in the final column of forward
    
```

**Figure 6.8** The forward algorithm; we've used the notation  $\text{forward}[s,t]$  to represent  $\alpha_t(s)$ .

1. Initialization:

$$(6.12) \quad \alpha_1(j) = a_{0j}b_j(o_1) \quad 1 \leq j \leq N$$

2. Recursion (since states 0 and N are non-emitting):

$$(6.13) \quad \alpha_t(j) = \sum_{i=1}^{N-1} \alpha_{t-1}(i)a_{ij}b_j(o_t); \quad 1 < j < N, 1 < t < T$$

3. Termination:

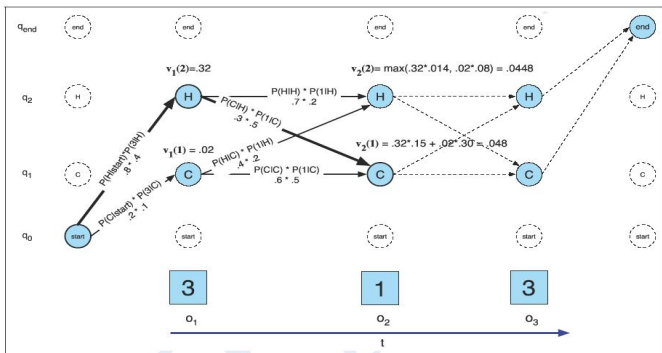
$$(6.14) \quad P(O|\lambda) = \alpha_T(N) = \sum_{i=2}^{N-1} \alpha_T(i) a_{iN}$$



# Viterbi algorithm

In decoding we need to find the most likely state sequence given an observation  $O$ . The Viterbi algorithm follows the same approach (dynamic programming) of the Forward.

Viterbi scores are attached to each possible state in the sequence.



**Figure 6.9** The Viterbi trellis for computing the best path through the hidden state space for the ice-cream eating events 3 1 3. Hidden states are in circles, observations in squares. White (unfilled) circles indicate illegal transitions. The figure shows the computation of  $v_t(j)$  for two states at two time steps. The computation in each cell follows Eq' 6.10:  $v_t(j) = \max_{1 \leq i \leq N-1} v_{t-1}(i) a_{ij} b_j(o_t)$ . The resulting probability expressed in each cell is Eq' 6.16:  $v_t(j) = P(q_0, q_1, \dots, q_{t-1}, o_1, o_2, \dots, o_t, q_t = j | \lambda)$ .



# *HMM and POS tagging: Parameter Estimation*

Supervised methods in tagged data sets:

- Output probs:  $P(w_i|p^j) = \frac{C(w_i,p^j)}{C(p^j)}$
- Transition probs:  $P(p^i|p^j) = \frac{C(p^i \text{ follows } p^j)}{C(p^j)}$
- Smoothing:  $P(w_i|p^j) = \frac{C(w_i,p^j)+1}{C(p^j)+K^i}$   
(see Manning& Schutze, Chapter 6)

# HMM and POS tagging: Parameter Estimation

Unsupervised (few tagged data available):

- With a dictionary:  $P(w_i|p^j)$  are early estimated from  $D$ , while  $P(p^i|p^j)$  are randomly assigned
- With equivalence classes  $u_L$ , (Kupiec92):

$$P(w^i|p^L) = \frac{\frac{1}{|L|} C(u^L)}{\sum_{u_{L'}} \frac{C(u_{L'})}{|L'|}}$$

For example, if  $L = \{\text{noun, verb}\}$  then  $u_L = \{\text{cross, drive, ...}\}$







## Other References

- *"Introduction to Information Retrieval"*, Christopher D. Manning, Prabhakar Raghavan and Hinrich Schütze, Cambridge University Press. 2008. Chapter 12.  
<http://www-csli.stanford.edu/hinrich/information-retrieval-book>.
- Rabiner, Lawrence. "First Hand: The Hidden Markov Model". IEEE Global History Network. Retrieved 2 October 2013. at  
[http://www.ieeeeghn.org/wiki/index.php/First-Hand:The\\_Hidden\\_Markov\\_Model](http://www.ieeeeghn.org/wiki/index.php/First-Hand:The_Hidden_Markov_Model)
- Applet at: <http://www.cs.umb.edu/sreuilak/viterbi/>

## Exercise

Consider a two-bit register. The register has four possible states: 00, 01, 10 and 11. Initially, at time 0, the contents of the register is chosen at random to be one of these four states, each with equal probability. At each time step, beginning at time 1, the register is randomly manipulated as follows: with probability  $1/2$ , the register is left unchanged; with probability  $1/4$ , the two bits of the register are exchanged (e.g., 01 becomes 10); and with probability  $1/4$ , the right bit is flipped (e.g., 01 becomes 00). After the register has been manipulated in this fashion, the left bit is observed. Suppose that on the first three time steps, we observe 0, 0, 1.

- Formulate the register as an HMM. What is the probability of transitioning from every state to every other state? What is the probability of observing output (0 or 1) in each state?
- What is the probability of being in each state at time  $t$  after observing only the first  $t$  bits, for  $t = 1, 2, 3$ .