

Stochastic models for learning language models

R. Basili

Corso di *Web Mining e Retrieval*
a.a. 2015-16

March 23, 2016



Outline

Outline

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

Quantitative Models of language structures

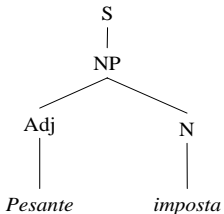
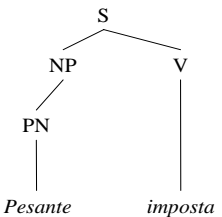
Linguistic structures are example of structures where syntagmatic information is crucial for machine learning. The most used modeling here are grammars:

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*”

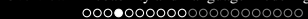




The role of Quantitative Approaches

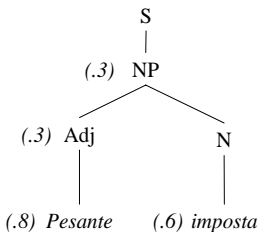
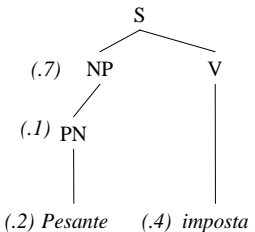
Weighted grammars are models of (possibly limited) *degrees of grammaticality*. They are meant to deal with a large range of ambiguity problems:

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

“*Pesante imposta*”





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:

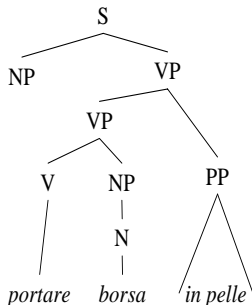
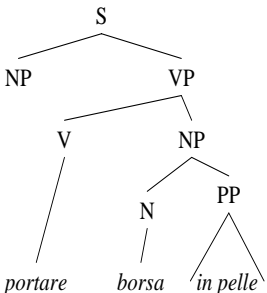
- 1. S → NP V .7
- 2. S → NP .3
- 3. NP → PN .1
- 4. NP → N .6
- 5. NP → Adj N .3
- 6. N → imposta .6
- 7. V → imposta .4
- 8. Adj → Pesante .8
- 9. PN → Pesante .2
- ...

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

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

Syntactic Disambiguation

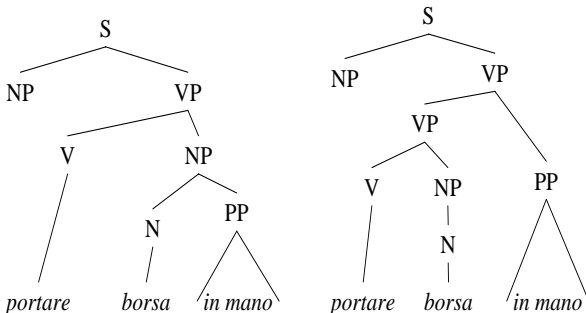
“portare borsa in pelle”



Derivation Trees for a structurally ambiguous sentence

Syntactic Disambiguation (cont'd)

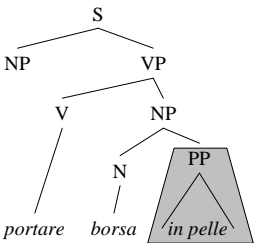
“portare borsa in mano”



Derivation Trees for a second structurally ambiguous sentence.

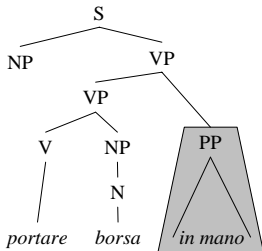
Structural Disambiguation (cont'd)

“portare borsa in pelle”



$p(\text{portare, in, pelle}) \ll p(\text{borsa, in, pelle})$

“portare borsa in mano”



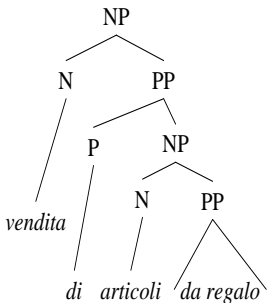
$p(\text{borsa, in, mano}) \ll p(\text{portare, in, mano})$

Disambiguation of structural ambiguity.

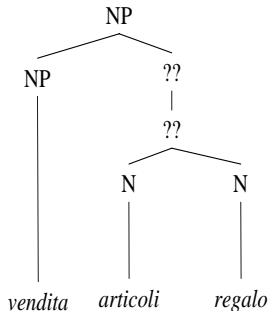


Tolerance to errors

“vendita di articoli da regalo”



“vendita articoli regalo”



An example of ungrammatical but meaningful sentence



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



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



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

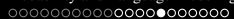
- It assumes values w_i in the alphabet A
- Probability is used to describe the uncertainty on the emitted signal

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



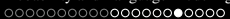
Probability and Language Modeling

- 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

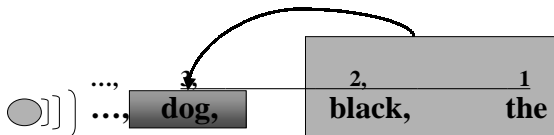
- Notice that time points can be represented as **states** of the emitting source
- An output w_i can be considered as emitted in a *given state* X_i by the source, and *given a certain history*



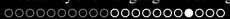
Probability and Language Modeling

What's in a state

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



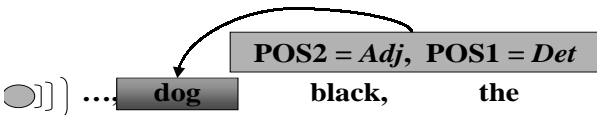
$$p(\text{the}, \text{black}, \text{dog}) = p(\text{dog} | \text{the}, \text{black})$$



Probability and Language Modeling

What's in a state

preceding POS tags \Rightarrow **stochastic taggers**





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*)



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*



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:

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

- 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

Markov Models: Matrix Representation

- A (transition) matrix A :

$$a_{ij} = P(X_{t+1} = p_j | X_t = p_i)$$

Note that $\forall i, j \quad a_{ij} \geq 0$ and $\forall i \quad \sum_j a_{ij} = 1$

Representation of a Markov Chain

Markov Models: Matrix Representation

- A (transition) matrix A :

$$a_{ij} = P(X_{t+1} = p_j | X_t = p_i)$$

Note that $\forall i, j \quad a_{ij} \geq 0$ and $\forall i \quad \sum_j a_{ij} = 1$

- Initial State description (i.e. probabilities of initial states):

$$\pi_i = P(X_1 = p_i)$$

Note that $\sum_{j=1}^n \pi_{ij} = 1$.

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

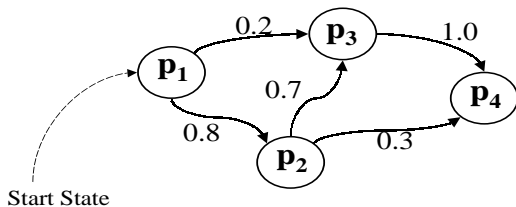
Representation of a Markov Chain

Graphical Representation

$$P(X_1 = p_1) = 1 \quad \leftarrow \text{StartState}$$

$$P(X_k = p_3 | X_{k-1} = p_2) = 0.7 \quad \forall k$$

$$P(X_k = p_4 | X_{k-1} = p_1) = 0 \quad \forall k$$



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 2

$$\begin{aligned}
 P(X_n = TP) &= \\
 \sum_{CP, p_2, p_3, \dots, TP} &P(X_1 = CP)P(X_2 = p_2|X_1 = CP)P(X_3 = p_3|X_1 = \\
 &CP, X_2 = p_2) \cdot \dots \cdot P(X_n = TP|X_1 = CP, X_2 = p_2, \dots, X_{n-1} = \\
 &p_{n-1}) = \\
 &= \sum_{CP, p_2, p_3, \dots, TP} P(CP)P(p_2|CP)P(p_3|p_2) \cdot \dots \cdot P(TP|p_{n-1}) = \\
 &= \sum_{CP, p_2, p_3, \dots, TP} P(CP) \cdot \prod_{t=1}^{n-2} P(p_{t+1}|p_t) \cdot P(p_n = TP|p_{n-1}) \\
 & (= \sum_{p_1, \dots, p_n} P(p_1) \cdot \prod_{t=1}^{n-1} P(p_{t+1}|p_t))
 \end{aligned}$$



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)

Solution for the problem 1 For the given state sequence

$$X = (CP, TP, TP)$$

$$P(O_1 = Cap, O_2 = Cof, X_1 = CP, X_2 = TP, X_3 = TP) =$$

$$P(O_1 = Cap, O_2 = Cof | X_1 = CP, X_2 = TP, X_3 = TP) P(X_1 = CP, X_2 = TP, X_3 = TP) =$$

$$P(Cap, Cof | CP, TP, TP) P(CP, TP, TP) \text{ Now:}$$

$P(Cap, Cof | CP, TP, TP)$ is the probability of output *Cap, Cof* during transitions from *CP* to *TP* and *TP* to *TP*

and $P(CP, TP, TP)$ is the probability of the transition chain.

Therefore,

$$= P(Cap | CP, TP) P(Cof | TP, TP) = (\text{in our simplified model})$$

$$= P(Cap | CP) P(Cof | TP) = 0.2 \cdot 0.2 = 0.04$$



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)

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



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)

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



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

$$\begin{aligned}
 P(\text{Tea}, \text{Cof}) = & \\
 & P(\text{Tea}|\text{CP})P(\text{Cof}|\text{CP})P(\text{CP})P(\text{CP}|\text{CP})P(\text{CP}|\text{CP}) + \text{ st.}: (\text{CP}, \text{CP}, \text{CP}) \\
 & P(\text{Tea}|\text{CP})P(\text{Cof}|\text{TP})P(\text{CP})P(\text{TP}|\text{CP})P(\text{CP}|\text{TP}) + \text{ st.}: (\text{CP}, \text{TP}, \text{CP}) \\
 & P(\text{Tea}|\text{CP})P(\text{Cof}|\text{CP})P(\text{CP})P(\text{CP}|\text{CP})P(\text{TP}|\text{CP}) + \text{ st.}: (\text{CP}, \text{CP}, \text{TP}) \\
 & P(\text{Tea}|\text{CP})P(\text{Cof}|\text{TP})P(\text{CP})P(\text{TP}|\text{CP})P(\text{TP}|\text{TP}) = \text{ st.}: (\text{CP}, \text{TP}, \text{TP}) \\
 & \\
 = & 0.15 \cdot 0.65 \cdot 1 \cdot 0.5 \cdot 0.5 + \\
 + & 0.15 \cdot 0.2 \cdot 1 \cdot 0.5 \cdot 0.3 + \\
 + & 0.15 \cdot 0.65 \cdot 1 \cdot 0.5 \cdot 0.5 + \\
 + & 0.15 \cdot 0.2 \cdot 1.0 \cdot 0.5 \cdot 0.7 =
 \end{aligned}$$



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

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

$$P(\text{Tea}|\text{CP})P(\text{Cof}|\text{CP})P(\text{CP})P(\text{CP}|\text{CP})P(\text{CP}|\text{CP}) + \text{st.: } (\text{CP}, \text{CP}, \text{CP})$$

$$P(\text{Tea}|\text{CP})P(\text{Cof}|\text{TP})P(\text{CP})P(\text{TP}|\text{CP})P(\text{CP}|\text{TP}) + \text{st.: } (\text{CP}, \text{TP}, \text{CP})$$

$$P(\text{Tea}|\text{CP})P(\text{Cof}|\text{CP})P(\text{CP})P(\text{CP}|\text{CP})P(\text{TP}|\text{CP}) + \text{st.: } (\text{CP}, \text{CP}, \text{TP})$$

$$P(\text{Tea}|\text{CP})P(\text{Cof}|\text{TP})P(\text{CP})P(\text{TP}|\text{CP})P(\text{TP}|\text{TP}) = \text{st.: } (\text{CP}, \text{TP}, \text{TP})$$

$$= 0.15 \cdot 0.65 \cdot 1 \cdot 0.5 \cdot 0.5 +$$

$$+ 0.15 \cdot 0.2 \cdot 1 \cdot 0.5 \cdot 0.3 +$$

$$+ 0.15 \cdot 0.65 \cdot 1 \cdot 0.5 \cdot 0.5 +$$

$$+ 0.15 \cdot 0.2 \cdot 1.0 \cdot 0.5 \cdot 0.7 =$$

$$= 0.024375 + 0.0045 + 0.024375 + 0.0105 =$$

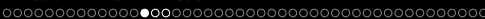
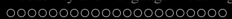
$$= 0.06375$$

A Simple Example of Hidden Markov Model (2)

Solution to the problem 3 (*Likelihood*)

In the general case, a sequence of n symbols $O = (o_1, \dots, o_n)$ out from any sequence of $n + 1$ transitions $X = (p_1, \dots, p_{n+1})$ can be predicted by the following probability:

$$\begin{aligned} P(O) &= \sum_{p_1, \dots, p_{n+1}} P(O|X)P(X) = \\ &= \sum_{p_1, \dots, p_{n+1}} P(CP) \prod_{t=1}^n P(O_t|p_t, p_{t+1})P(p_{t+1}|p_t) \end{aligned}$$



Modeling linguistic tasks as Stochastic Processes

Advantages

There are several advantages to model a linguistic problem as an HMM

- It is a powerful mathematical framework for modeling
- It provides clear problems settings for different applications: **estimation**, **decoding** and **model induction**
- HMM-based models provides sound solutions for the above applications

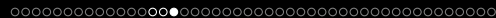
We will see an example as the HMM modeling of POS tagging

Fundamental problems for HMM

Fundamental Questions for HMM

The complexity of training and decoding can be limited by the use of optimization techniques

- Given the observation sequence $O = O_1, \dots, O_n$ and a model $\lambda = (E, T, \pi)$, how to efficiently compute $P(O|\lambda)$? (*Language Modeling*)
- Given the observation sequence $O = O_1, \dots, O_n$ and a model $\lambda = (E, T, \pi)$, how do we choose the optimal state sequence $Q = q_1, \dots, q_n$ responsible of generating O ? (*Tagging/Decoding*)
- How to adjust model parameters $\lambda = (E, T, \pi)$ so to maximize $P(O|\lambda)$? (*Model Induction*)



HMM: Mathematical Methods

All the above problems can be approached by several optimization techniques able to limit the complexity.

- Language Modeling via *dynamic programming* (**Forward algorithms**) ($O(n)$)
- Tagging/Decoding via *dynamic programming* ($O(n)$) (**Viterbi**)
- Parameter estimation via *entropy minimization* (EM)

A relevant issue is the availability of source data: supervised training cannot be applied always



The task of POS tagging

POS tagging

Given a sequence of morphemes w_1, \dots, w_n with ambiguous syntactic descriptions (i.e. part-of-speech tags) t_j , compute the sequence of n POS tags t_{j_1}, \dots, t_{j_n} that characterize correspondingly all the words w_i .

The task of POS tagging

POS tagging

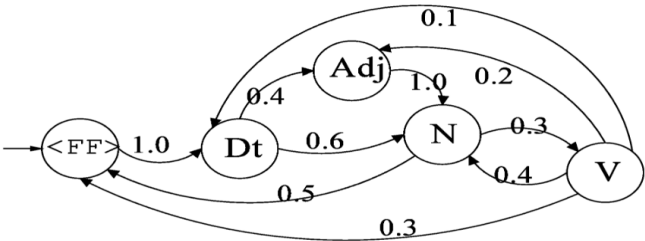
Given a sequence of morphemes w_1, \dots, w_n with ambiguous syntactic descriptions (i.e. part-of-speech tags) t_j , compute the sequence of n POS tags t_{j_1}, \dots, t_{j_n} that characterize correspondingly all the words w_i .

Examples:

- *Secretariat is expected to race tomorrow*
- \Rightarrow NNP VBZ VBN TO VB NR
- \Rightarrow NNP VBZ VBN TO NN NR

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



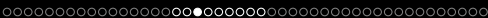
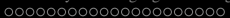
HMM and POS tagging

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:

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

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

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:

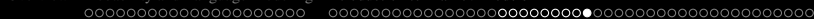
$$(t_1, \dots, t_n) = \underset{t_1, \dots, t_n}{\operatorname{argmax}} P(t_1, \dots, t_n | w_1, \dots, w_n) = \underset{t_1, \dots, t_n}{\operatorname{argmax}} \prod_{i=1}^n P(w_i | t_i) P(t_i | t_{i-1})$$

Fundamental Questions for HMM in POS tagging

- 1 Given a model what is the probability of an output sequence, O :
Computing Likelihood.
- 2 Given a model and an observable output sequence O (i.e. words), how to determine the sequence of states (t_1, \dots, t_n) such that it is the best explanation of the observation O :
Decoding Problem
- 3 Given a sample of the output sequences and a space of possible models how to find out the best model, that is the model that best explains the data:
how to estimate parameters?

Fundamental Questions for HMM in POS tagging

- 1. Not much relevant for POS tagging, where (w_1, \dots, w_n) are always known.
Trellis and dynamic programming technique.
- 2. (Decoding) Viterbi Algorithm for evaluating $P(W|O)$.
Linear in the sequence length.
- 3. Baum-Welch (or Forward-Backward algorithm), that is a special case of Expectation Maximization estimation.
Weakly supervised or even unsupervised.
Problems: Local minima can be reached when initial data are poor.



HMM and POS tagging

Advantages for adopting HMM in POS tagging

- An elegant and sound theory
- Training algorithms:
 - Estimation via EM (Baum-Welch)
 - Unsupervised (or possibly weakly supervised)
- Fast Inference algorithms: Viterbi algorithm
Linear wrt the sequence length ($O(n)$)
- Sound methods for comparing different models and estimations
(e.g. cross-entropy)



Forward algorithm

In computing the likelihood $P(O)$ of an observation we need to sum up the probability of all paths in a Markov model. Brute force computation is not applicable in most cases. The forward algorithm is an application of dynamic programming.

Forward algorithm

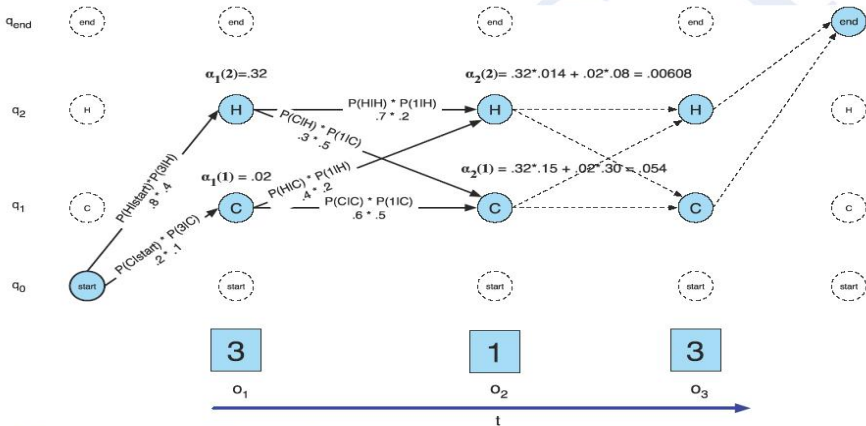


Figure 6.6 The forward trellis for computing the total observation likelihood for the ice-cream events 3 1 3. Hidden states are in circles, observations in squares. White (unfilled) circles indicate illegal transitions. The figure shows the computation of $\alpha_t(j)$ for two states at two time steps. 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 ← NUM-OF-STATES(*state-graph*)

Create a probability matrix *forward*[*num-states*+2, *T*+2]

forward[0,0] ← 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*] ← $\sum_{1 \leq s' \leq \text{num-states}}$ *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 *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.

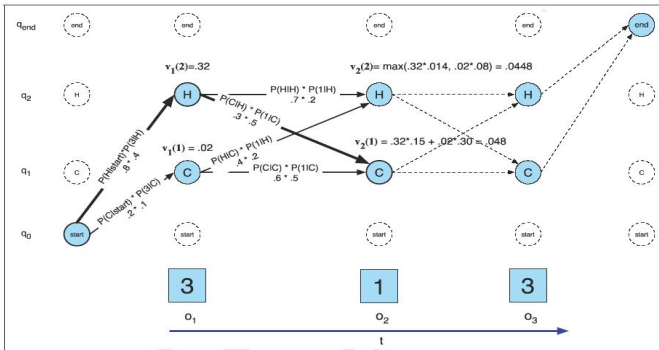


Figure 6.9 The Viterbi trellis for computing the best path through the hidden state space for the ice-cream eating events 3 1 3. Hidden states are in circles, observations in squares. White (unfilled) circles indicate illegal transitions. The figure shows the computation of $v_t(j)$ for two states at two time steps. The computation in each cell follows Eq 6.10: $v_t(j) = \max_{i \leq N-1} v_{t-1}(i) a_{ij} b_j(o_t)$ 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: the Viterbi Algorithm

function VITERBI(*observations* of len T , *state-graph*) **returns** *best-path*

$num\text{-}states \leftarrow \text{NUM-OF-STATES}(state\text{-}graph)$

Create a path probability matrix $viterbi[num\text{-}states+2, T+2]$

$viterbi[0,0] \leftarrow 1.0$

for each time step t **from** 1 **to** T **do**

for each state s **from** 1 **to** $num\text{-}states$ **do**

$viterbi[s,t] \leftarrow \max_{1 \leq s' \leq num\text{-}states} viterbi[s',t-1] * a_{s',s} * b_s(o_t)$

$backpointer[s,t] \leftarrow \operatorname{argmax}_{1 \leq s' \leq num\text{-}states} viterbi[s',t-1] * a_{s',s}$

Backtrace from highest probability state in final column of $viterbi[[]]$ and return path

Figure 6.10 Viterbi algorithm for finding optimal sequence of tags. Given an observation sequence and an HMM $\lambda = (A, B)$, the algorithm returns the state-path through the HMM which assigns maximum likelihood to the observation sequence. Note that states 0 and $N+1$ are non-emitting *start* and *end* states.



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

$$u_L = \{\text{cross, drive, \dots}\}$$



A survey of the Baum-Welch method

The learning Problem

Given a HMM $\lambda = (E, T, \pi)$ and an observation history $Z = (z_1, z_2, \dots, z_t)$, and a new HMM $\lambda' = (E', T', \pi')$ that explains the observations at least as well, or possibly better, i.e., such that $Pr[Z|\lambda'] \geq Pr[Z|\lambda]$.

- Ideally, we would like to find the model that **maximizes** $Pr[Z|\lambda]$; however, this is in general an intractable problem.
- We will be satisfied with an algorithm that converges to local maxima of such probability.
- Notice that in order for learning to be effective, we need **lots of data**, i.e., many, long observation histories!

The Baum-Welch estimation as a EM process

Baum-Welch re-estimation: the idea

Baum-Welch reestimation is also called the **Forward-Backward** algorithm It is special case of the **Expectation Maximization (EM)** algorithm

- 1 Start with initial probability estimates
- 2 Compute expectations of how often each transition/emission is used
- 3 Re-estimate the probabilities based on those expectations
...and repeat until convergence

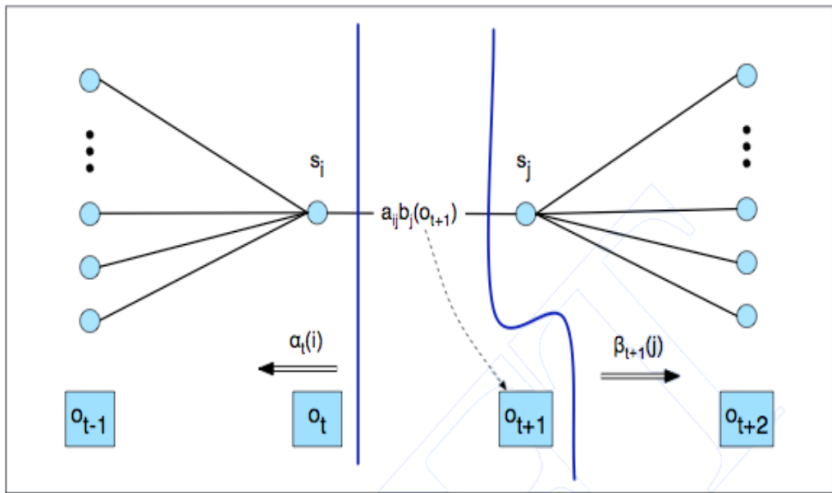


Baum-Welch: Expectation of (state) counts

- Let us define: $\gamma_k(s) = Pr[X_k = s|Z, \lambda]$
i.e., $\gamma_k(s)$ is the probability that the system is at state s at the k -th time step, given the observation sequence Z and the model λ .
- We already know how to compute this, e.g., using smoothing:

$$\gamma_k(s) = \frac{\alpha_k(s)\beta_k(s)}{Pr[X_k|Z, \lambda]} = \frac{\alpha_k(s)\beta_k(s)}{\sum_{q \in S} \alpha_k(q)}$$

The forward backward probabilities





Baum-Welch: Forward and Backward probabilities

- Forward probabilities (DEF):

$$\alpha_k(s) = Pr[o_1, \dots, o_k, x_k = s | \lambda]$$

Recursively

$$\alpha_{k+1}(q) = \sum_{s \in S} \alpha_k(s) a_{sq} b_q(o_{k+1}) \quad (\text{with } \alpha_1(q)) = \pi_q)$$



Baum-Welch: Forward and Backward probabilities

- Forward probabilities (DEF):

$$\alpha_k(s) = Pr[o_1, \dots, o_k, x_k = s | \lambda]$$

Recursively

$$\alpha_{k+1}(q) = \sum_{s \in S} \alpha_k(s) a_{sq} b_q(o_{k+1}) \quad (\text{with } \alpha_1(q)) = \pi_q)$$

- Backward probabilities (DEF):

$$\beta_k(s) = Pr[o_k, \dots, o_t | x_k = s, \lambda]$$

Recursively:

$$\beta_k(s) = \sum_{q \in S} a_{sq} b_q(o_{k+1}) \beta_{k+1}(q)$$



Baum-Welch: Expectation of (state) counts

- Let us define: $\gamma_k(s) = Pr[X_k = s|Z, \lambda]$
- We already know how to compute this, e.g., using smoothing:

$$\gamma_k(s) = \frac{\alpha_k(s)\beta_k(s)}{Pr[X_k|Z, \lambda]} = \frac{\alpha_k(s)\beta_k(s)}{\sum_{q \in S} \alpha_k(q)}$$

- **New concept:** how many times is the state trajectory expected to transition from state s ?

$$E[\# \text{ of transitions from } s] = \sum_{k=1}^{t-1} \gamma_k(s)$$



Baum-Welch: Expectation of (transitions) counts

- In much the same vein, let us define $\xi_k(q, s) = Pr[X_k = q, X_{k+1} = s | Z, \lambda]$ (i.e., $\xi_k(q, s)$ is the probability of being at state q at time k , and at state s at time $k + 1$, given the observations and the current HMM model)
- We have that $\xi_k(q, s) = \eta_k \alpha_k(q) T_{q,s} E_{s, o_{k+1}} \beta_{k+1}(s)$ where η_k is a normalization factor, such that $\sum_{q,s} \xi_k(q, s) = 1$.
- **New concept:** how many times it the state trajectory expected to transition *from* state q to state s ?
 $E[\# \text{ of transitions from } q \text{ to } s] = \sum_{k=1}^{t-1} \xi_k(q, s)$



Baum-Welch algorithm

- Based on the probability estimates and expectations computed so far, using the original HMM model $\lambda = (E, T, \pi)$, we can construct a new model $\hat{\lambda} = (\hat{E}, \hat{T}, \hat{\pi})$ (notice that the two models share the states and observations):
- The new initial condition distribution is the one obtained by smoothing: $\hat{\pi}_s = \gamma_1(s)$
- The entries of the new transition matrix can be obtained as follows:
$$\hat{T}_{q,s} = \frac{E[\text{\# of transitions from } q \text{ to } s]}{E[\text{\# of transitions from } q]} = \frac{\sum_{k=1}^{t-1} \xi_k(q,s)}{\sum_{k=1}^{t-1} \gamma_k(s)}$$



Baum-Welch algorithm

- The entries of the new emission matrix can be obtained as follows:

$$\begin{aligned}\hat{E}_{s,o} (= \hat{b}_s(o)) &= \frac{E[\text{\# of times in state } s, \text{ when the observation was } o]}{E[\text{\# of times in state } s]} = \\ &= \frac{\sum_{k=1}^t \gamma_k(s) \mathbf{1}(z_k=o)}{\sum_{k=1}^t \gamma_k(s)}\end{aligned}$$

- In this way, new estimated version for \hat{E} , \hat{T} and $\hat{\pi}$ are available:

They correspond to a new model $\hat{\lambda} = (\hat{E}, \hat{T}, \hat{\pi})$



Baum-Welch as an EM iterative model refinement

E-step (expectaton)

$\sum_{k=1}^t \gamma_k(i) =$ expected number of transitions involving q_i

$\sum_{k=1}^{t-1} \xi_k(i,j) =$ expected number of transitions from q_i to q_j

M-step (Likelyhood Maximimization)

We can re-estimate parameters by ratio of expected counts

$$\hat{a}_{i,j} = \frac{\sum_{k=1}^{t-1} \xi_k(i,j)}{\sum_{k=1}^{t-1} \gamma_k(j)}$$
$$\hat{b}_i(s) = \frac{\sum_{k=1}^{t-1} \gamma_k(i) \cdot \mathbf{1}(o_k=s)}{\sum_{k=1}^{t-1} \gamma_k(j)}$$

Baum-Welch: an example on the soft drink machine

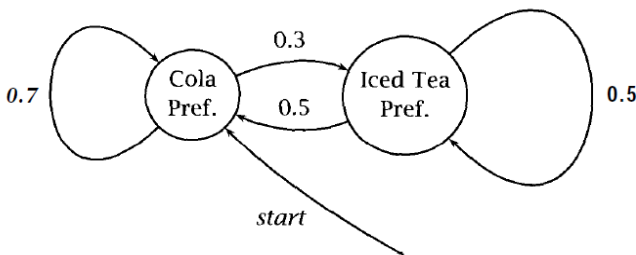


Figure 9.2 The crazy soft drink machine, showing the states of the machine and the state transition probabilities.

Output probability given From state

	cola	iced tea (ice_t)	lemonade (lem)
CP	0.6	0.1	0.3
IP	0.1	0.7	0.2



About Parameter Estimation for POS

Baum-Welch re-estimation on the soft drink machine

training on the observation sequence (lem, ice_t, cola)
 values for $p_t(i, j)$:

		Time (and j)								
		1		2			3			y_3
i		CP	IP	y_1	CP	IP	y_2	CP	IP	
	CP		0.3	0.7	1.0	0.28	0.02	0.3	0.616	0.264
IP		0.0	0.0	0.0	0.6	0.1	0.7	0.06	0.06	0.12

and so the parameters will be reestimated as follows:

		Original			Reestimated		
II	CP	1.0			1.0		
	IP	0.0			0.0		
A	CP	CP	IP	CP	IP		
	IP	0.7	0.3	0.5486	0.4514		
B	CP	0.5	0.5	0.8049	0.1951		
	IP	cola	ice_t	lem	cola	ice_t	lem
	CP	0.6	0.1	0.3	0.4037	0.1376	0.4587
	IP	0.1	0.7	0.2	0.1363	0.8537	0.0

Baum-Welch algorithm: convergence

- It can be shown [Baum et al., 1970] that the new model $\hat{\lambda}$ is such that
 - $Pr[Z|\hat{\lambda}] \geq Pr[Z|\lambda]$, as desired.
 - $Pr[Z|\hat{\lambda}] = Pr[Z|\lambda]$ only if λ is a critical point of the likelihood function

$$f(\lambda) = Pr[Z|\lambda]$$

Other Approaches to POS tagging

- Church (1988):

$$\prod_{i=n}^3 P(w_i | t_i) P(t_{i-2} | t_{i-1}, t_i) \text{ (backward)}$$

Estimation from tagged corpus (Brown)

No HMM training

Performances: > 95%

- De Rose (1988):

$$\prod_{i=1}^n P(w_i | t_i) P(t_{i-1} | t_i) \text{ (forward)}$$

Estimation from tagged corpus (Brown)

No HMM training Performance: 95%

- Merialdo et al., (1992), ML estimation vs. Viterbi training
Propose an incremental approach: small tagging and then Viterbi training
- $\prod_{i=1}^n P(w_i | t_i) P(t_{i+1} | t_i, w_i) ???$



POS tagging: References

- F. Jelinek, Statistical methods for speech recognition, Cambridge, Mass.: MIT Press, 1997.
- Manning & Schutze, Foundations of Statistical Natural Language Processing, MIT Press, Chapter 6.
- Church (1988), A Stochastic Parts Program and Noun Phrase Parser for Unrestricted Text, <http://acl.ldc.upenn.edu/A/A88/A88-1019.pdf>
- Rabiner, L. R. (1989). A tutorial on Hidden Markov Models and selected applications in speech recognition. Proceedings of the IEEE, 77(2), 257-286.
- Viterbi, A. J. (1967). Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. IEEE Transactions on Information Theory, IT-13(2), 260-269.
- Parameter Estimation (slides):
<http://jan.stanford.edu/fsnlp/statest/henke-ch6.ppt>



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.ieeeahn.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.

- Show how the register can be formulated as an HMM. What is the probability of transitioning from every state to every other state? What is the probability of observing each 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 $t = 1, 2, 3$.