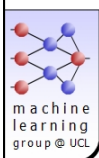


Sequence modelling



Marco Saerens (UCL)



Université catholique de Louvain

Université partenaire de l'Académie universitaire 'Louvain'

Slides references

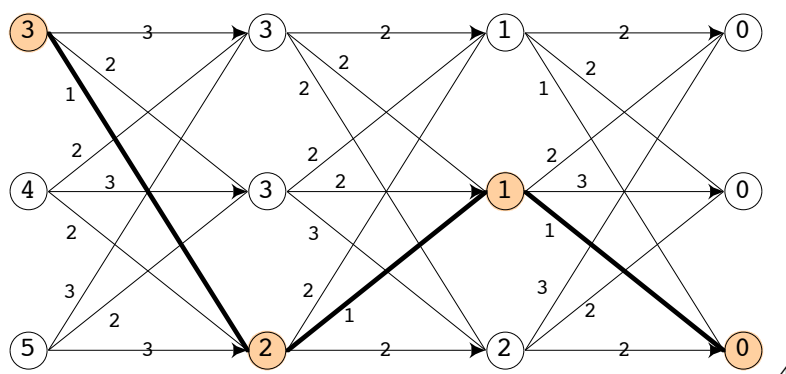
- Many slides and figures have been adapted from the slides associated to the following books:
 - Alpaydin (2004), Introduction to machine learning. The MIT Press.
 - Duda, Hart & Stork (2000), Pattern classification, John Wiley & Sons.
 - Han & Kamber (2006), Data mining: concepts and techniques, 2nd ed. Morgan Kaufmann.
 - Tan, Steinbach & Kumar (2006), Introduction to data mining. Addison Wesley.
 - Grinstead and Snell's Introduction to Probability (2006; GNU Free Documentation License)
- As well as from Wikipedia, the free online encyclopedia
- I really thank these authors for their contribution to machine learning and data mining teaching

Dynamic programming and edit-distance

Marco Saerens (UCL)

Dynamic programming

- Suppose we have a lattice with N levels:



Dynamic programming

- The problem is to reach level N
- From level 1
- With minimum cost
- = shortest-path problem

5

Dynamic programming

- Some definitions
 - s_k = variable containing the state at level k
- The local cost associated to the decision to jump to the state $s_k = j$ at level k

$$d(s_k = j \mid s_{k-1} = i)$$

- Given that we were in state $s_{k-1} = i$ at level $k-1$,

Dynamic programming

- The **total cost** of a path (s_0, s_1, \dots, s_N)

is

$$D(s_0, s_1, \dots, s_N) = \sum_{i=1}^N d(s_i | s_{i-1})$$

- The **optimal cost** when starting from state s_0 is

$$\begin{aligned} D^*(s_0) &= \min_{(s_1, \dots, s_N)} \{D(s_0, s_1, \dots, s_N)\} \\ &= \min_{(s_1, \dots, s_N)} \left\{ \sum_{i=1}^N d(s_i | s_{i-1}) \right\} \end{aligned}$$

7

Dynamic programming

- The **optimal cost**, whatever the initial state, is

$$D^* = \min_{s_0} \{D^*(s_0)\}$$

- The **optimal cost** when starting from some intermediate state s_k

$$D^*(s_k) = \min_{(s_{k+1}, \dots, s_N)} \left\{ \sum_{i=k+1}^N d(s_i | s_{i-1}) \right\}$$

8

Dynamic programming

- Here are the **recurrence relations** allowing to obtain the optimal cost:

$$D^*(s_N) = 0$$

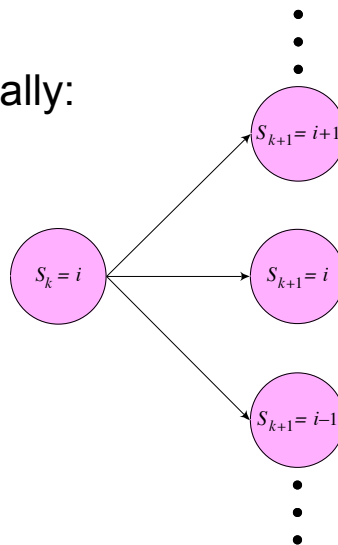
$$D^*(s_k = i) = \min_{s_{k+1}} \{d(s_{k+1} | s_k = i) + D^*(s_{k+1})\}$$

$$D^* = \min_{s_0} \{D^*(s_0)\}$$

9

Dynamic programming

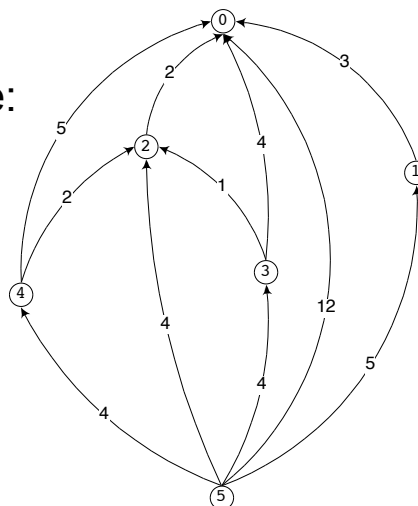
- Graphically:



10

Dynamic programming

■ Example:



$$D^*(s_{N-1} = i) = \min_{s_N} \{d(s_N | s_{N-1} = i)\}$$

11

Dynamic programming

■ Proof:

$$\begin{aligned}
 D^*(s_k) &= \min_{(s_{k+1}, \dots, s_N)} \left\{ \sum_{i=k+1}^N d(s_i | s_{i-1}) \right\} \\
 &= \min_{(s_{k+1}, \dots, s_N)} \left\{ d(s_{k+1} | s_k) + \sum_{i=k+2}^N d(s_i | s_{i-1}) \right\} \\
 &= \min_{s_{k+1}} \left\{ \min_{(s_{k+2}, \dots, s_N)} \left\{ d(s_{k+1} | s_k) + \sum_{i=k+2}^N d(s_i | s_{i-1}) \right\} \right\} \\
 &= \min_{s_{k+1}} \left\{ d(s_{k+1} | s_k) + \underbrace{\min_{(s_{k+2}, \dots, s_N)} \left[\sum_{i=k+2}^N d(s_i | s_{i-1}) \right]}_{D^*(s_{k+1})} \right\} \\
 &= \min_{s_{k+1}} \{d(s_{k+1} | s_k) + D^*(s_{k+1})\}
 \end{aligned}$$

Dynamic programming

- In a symmetric way:

$$D^*(s_0) = 0$$

$$D^*(s_k = i) = \min_{s_{k-1}} \{d(s_k = i \mid s_{k-1}) + D^*(s_{k-1})\}$$

$$D^* = \min_{s_N} \{D^*(s_N)\}$$

13

Dynamic programming

- Now, if there are jumps bypassing some levels,

$$D^*(s_N) = 0$$

$$D^*(s_k = i) = \min_{\{s \mid s_k = i \rightarrow s\}} \{d(s \mid s_k = i) + D^*(s)\}$$

$$D^* = \min_{s_0} \{D^*(s_0)\}$$

- where $\{s \mid s_k = i \rightarrow s\}$ is the set of states to which there is a direct jump from $s_k = i$



Dynamic programming

- To find the **optimal path**, we need to keep track of
 - the previous state in each node (a pointer from the previous state to the current state)
- And use backtracking from the last node
 - in order to retrieve the optimal path

15



Application to edit-distance

- Computation of a distance between two strings
 - It computes the **minimal** number of **insertions**, **deletions** and **substitutions**
 - That is, the minimum number of **editions**
- For transforming one character string x into another character string y

16

Application to edit-distance

- Also called the “Levenstein distance”
- We thus have two character strings

$$\begin{cases} \mathbf{y} = y_1 y_2 \dots y_{|\mathbf{y}|} \\ \mathbf{x} = x_1 x_2 \dots x_{|\mathbf{x}|} \end{cases}$$

- x_i being the character i of string \mathbf{x}

17

Application to edit-distance

- The length of the string \mathbf{x} is denoted by $|\mathbf{x}|$

- In general, we have

$$|\mathbf{x}| \neq |\mathbf{y}|$$

- The **substring** of \mathbf{x} , beginning at character i and ending at j is defined by

$$\mathbf{x}_i^j = x_i x_{i+1} \dots x_j$$

18

Application to edit-distance

- We will now read the characters of x one by one
- In order to construct the string y
- To this aim, we define the three editing operations:
 - **Insertion** of a character in y (without taking any character from x)
 - **Deletion** of a character from x (without concatenating it to y)
 - **Substitution** of a character from y by one of x

19

Application to edit-distance

- First convention:

$$x_i^{|x|}$$

- Means that we have read (extracted) the $i-1$ first characters of x
- Thus, x has been cutted from its $i-1$ first characters
- They have been taken in order to construct y

20

Application to edit-distance

- Second convention:

$$\mathbf{y}_0^j$$

- Means that the j first characters of \mathbf{y} have been transcribed

- We progressively read the first characters of \mathbf{x} in order to build \mathbf{y}

21

Application to edit-distance

- It corresponds to a process with **levels** (steps) and **states**

- We will now apply dynamic programming

$$\left\{ \begin{array}{l} \text{One state corresponds to } (\mathbf{x}_i^{|\mathbf{x}|}, \mathbf{y}_0^j) \\ \text{The level } k \text{ corresponds to } i + j = \text{const} = k \end{array} \right.$$

- One state is characterized by a couple (i, j)

22

Application to edit-distance

- Here is the formal definition of the three operations

$$\left\{ \begin{array}{l} \text{Insertion with respect to } \mathbf{x}: (\mathbf{x}_i^{|\mathbf{x}|}, \mathbf{y}_0^{j-1}) \rightarrow (\mathbf{x}_i^{|\mathbf{x}|}, \mathbf{y}_0^j) \\ \text{Deletion with respect to } \mathbf{x}: (\mathbf{x}_{i-1}^{|\mathbf{x}|}, \mathbf{y}_0^j) \rightarrow (\mathbf{x}_i^{|\mathbf{x}|}, \mathbf{y}_0^j) \\ \text{Substitution with respect to } \mathbf{x}: (\mathbf{x}_{i-1}^{|\mathbf{x}|}, \mathbf{y}_0^{j-1}) \rightarrow (\mathbf{x}_i^{|\mathbf{x}|}, \mathbf{y}_0^j) \end{array} \right.$$

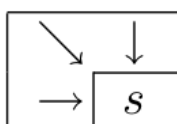
- For the two first operations, we are jumping from level k to level $k+1$
- For the last one, we directly jump to level $k+2$ (one level is passed)

23

Application to edit-distance

- This situation can be represented in a **two-dimensional table**

- One level is represented by $(i + j) = \text{constant}$
- One state is represented by (i, j)
- One operation corresponds to a valid transition in this table



24

Application to edit-distance

- Exemple of a table computing the distance between “livre” and “lire”

		y					
		∅	l	i	v	r	e
x	∅	0	1	2	3	4	5
	l	1	0	1	2	3	4
	i	2	1	0	1	2	3
	r	3	2	1	1	1	2
	e	4	3	2	2	2	1

25

Application to edit-distance

- Each level corresponds to a **diagonal** of the table:

$$-(i + j) = \text{constant}$$

		y					
		∅	l	i	v	r	e
x	∅	0	1	2	3	4	5
	l	1	0	1	2	3	4
	i	2	1	0	1	2	3
	r	3	2	1	1	1	2
	e	4	3	2	2	2	1

26

Application to edit-distance

- A **cost** (or **penalty**) is associated to each operation (insertion, deletion, substitution); for instance

$$\begin{cases} \delta_{\text{ins}}(y_i) = 1 \\ \delta_{\text{del}}(x_i) = 1 \\ \delta_{\text{sub}}(x_i, y_j) = \delta_{ij} \end{cases}$$

- with

$$\begin{cases} \delta_{ij} = 1 \text{ si } x_i \neq y_j \\ \delta_{ij} = 0 \text{ si } x_i = y_j \end{cases}$$

27

Application to edit-distance

- The dynamic programming formula can be applied to this problem:

– Initialization:

$$D^*(\mathbf{x}_0^{|\mathbf{x}|}, \mathbf{y}_0^0) = 0$$

– Then:

$$D^*(\mathbf{x}_i^{|\mathbf{x}|}, \mathbf{y}_0^j) = \min \begin{cases} D^*(\mathbf{x}_i^{|\mathbf{x}|}, \mathbf{y}_0^{j-1}) + 1 \\ D^*(\mathbf{x}_{i-1}^{|\mathbf{x}|}, \mathbf{y}_0^j) + 1 \\ D^*(\mathbf{x}_{i-1}^{|\mathbf{x}|}, \mathbf{y}_0^{j-1}) + \delta_{ij} \end{cases}$$

28

Application to edit-distance

- And finally:

$$\text{dist}(\mathbf{x}, \mathbf{y}) = D^*(\mathbf{x}_{|\mathbf{x}|}, \mathbf{y}_0^{|\mathbf{y}|})$$

- This value is the **edit-distance** or **Levenstein distance**

29

Application to edit-distance

- One example:

		y					
		∅	l	i	v	r	e
x	∅	0	1	2	3	4	5
	l	1	0	1	2	3	4
	i	2	1	0	1	2	3
	r	3	2	1	1	1	2
	e	4	3	2	2	2	1

- $\text{dist}(\text{lire}, \text{livre}) = 1$

30

Application to edit-distance

- In order to find the optimal path (optimal sequence of operations)
 - A pointer to the previous state has to be maintained for each cell of the table
 - The full path can be found by **backtracking** from the final state
- Numerous extensions and generalizations of this basic algorithm have been developed

31

Application to edit-distance

- Notice that the related quantity, the **longest common subsequence**, $\text{lcs}(\mathbf{x}, \mathbf{y})$, can be obtained by (no proof provided)

$$\begin{aligned} \text{dist}(\mathbf{x}, \mathbf{y}) &= \text{lcs}(\mathbf{x}, \mathbf{x}) + \text{lcs}(\mathbf{y}, \mathbf{y}) - 2 \text{lcs}(\mathbf{x}, \mathbf{y}) \\ &= |\mathbf{x}| + |\mathbf{y}| - 2 \text{lcs}(\mathbf{x}, \mathbf{y}) \end{aligned}$$

- and thus

$$\text{lcs}(\mathbf{x}, \mathbf{y}) = \frac{1}{2}(|\mathbf{x}| + |\mathbf{y}| - \text{dist}(\mathbf{x}, \mathbf{y}))$$

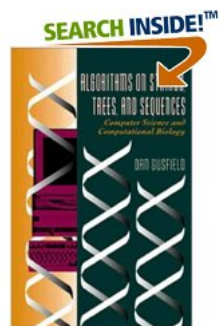
32

Application to edit-distance

- Thus, since
 - $\text{dist}(\text{lire}, \text{livre}) = 1$
- We have
 - $\text{lcs}(\text{lire}, \text{livre}) = 0.5 (4 + 5 - 1) = 4$

33

Application to edit-distance



34

A brief introduction to Markov chains

Introduction to Markov chains

- We have a **set of states**, $S = \{1, 2, \dots, n\}$
 - $s_t = k$ means that the **process**, or **system**, is in state k at time t
- **Example:**
 - We take as states the kind of weather R (rain), N (nice), and S (snow)

$$\mathbf{P} = \begin{array}{c} \text{R} \\ \text{N} \\ \text{S} \end{array} \begin{array}{ccc} \text{R} & \text{N} & \text{S} \\ \left(\begin{array}{ccc} 1/2 & 1/4 & 1/4 \\ 1/2 & 0 & 1/2 \\ 1/4 & 1/4 & 1/2 \end{array} \right) \end{array}$$

36



Introduction to Markov chains

- Markov chains are models of sequential discrete-time and discrete-state stochastic processes

37



Introduction to Markov chains

- The entries of this matrix \mathbf{P} are p_{ij} with
$$p_{ij} = \mathbb{P}(s_{t+1} = j | s_t = i)$$
- We assume that the probability of jumping to a state only depends on the current state
 - And not on the past, before this state
 - This is the **Markov property**

38

Introduction to Markov chains

- The matrix \mathbf{P} is called the one-step transition probabilities matrix
 - It is a stochastic matrix
 - The row sums are equal to 1
- We also assume that these transition probabilities are stationary
 - That is, independent of time
- Let us now compute $P(s_{t+2} = j | s_t = i)$:

39

Introduction to Markov chains

$$\begin{aligned}
 P(s_{t+2} = j | s_t = i) &= \sum_{k=1}^n P(s_{t+2} = j, s_{t+1} = k | s_t = i) \\
 &= \sum_{k=1}^n P(s_{t+2} = j | s_t = i, s_{t+1} = k) P(s_{t+1} = k | s_t = i) \\
 &= \sum_{k=1}^n P(s_{t+2} = j | s_{t+1} = k) P(s_{t+1} = k | s_t = i) \\
 &= \sum_{k=1}^n p_{kj} p_{ik} \\
 &= [\mathbf{P}^2]_{ij}
 \end{aligned}$$

40

Introduction to Markov chains

- The matrix \mathbf{P}^2 is the two-steps transition probabilities matrix
- By induction, \mathbf{P}^τ is the τ -steps transition probabilities matrix containing elements

$$p_{ij}^{(\tau)} = P(s_{t+\tau} = j | s_t = i)$$

41

Introduction to Markov chains

- If $\mathbf{x}(t)$ is the column vector containing the **probability distribution** of finding the process in each state of the Markov chain at time step t , we have

$$\begin{aligned} x_i(t) &= P(s_t = i) \\ &= \sum_{k=1}^n P(s_t = i, s_{t-1} = k) \\ &= \sum_{k=1}^n P(s_t = i | s_{t-1} = k) P(s_{t-1} = k) \\ &= \sum_{k=1}^n p_{ki} x_k(t-1) \end{aligned}$$

42

Introduction to Markov chains

- In matrix form, $\mathbf{x}(t) = \mathbf{P}^T \mathbf{x}(t-1)$
- Or, in function of the initial distribution, $\mathbf{x}(t) = (\mathbf{P}^T)^t \mathbf{x}(0)$
- Now, $x_j(t) = \mathbf{x}(t)^T \mathbf{e}_j = \mathbf{x}(0)^T \mathbf{P}^t \mathbf{e}_j$

43

Introduction to Markov chains

- Thus, when starting from state i , $\mathbf{x}(0) = \mathbf{e}_i$, and
- $$x_j(t | s_0 = i) = x_{ji}(t) = (\mathbf{e}_i)^T \mathbf{P}^t \mathbf{e}_j$$
- It is the probability of observing the process in state j at time t when starting from state i at time $t = 0$

44

Introduction to Markov chains

$$P^1 = \begin{array}{c} \text{Rain} \quad \text{Nice} \quad \text{Snow} \\ \text{Rain} \begin{pmatrix} .500 & .250 & .250 \\ .500 & .000 & .500 \\ .250 & .250 & .500 \end{pmatrix} \\ \text{Nice} \\ \text{Snow} \end{array}$$

$$P^2 = \begin{array}{c} \text{Rain} \quad \text{Nice} \quad \text{Snow} \\ \text{Rain} \begin{pmatrix} .438 & .188 & .375 \\ .375 & .250 & .375 \\ .375 & .188 & .438 \end{pmatrix} \\ \text{Nice} \\ \text{Snow} \end{array}$$

$$P^3 = \begin{array}{c} \text{Rain} \quad \text{Nice} \quad \text{Snow} \\ \text{Rain} \begin{pmatrix} .406 & .203 & .391 \\ .406 & .188 & .406 \\ .391 & .203 & .406 \end{pmatrix} \\ \text{Nice} \\ \text{Snow} \end{array}$$

$$P^4 = \begin{array}{c} \text{Rain} \quad \text{Nice} \quad \text{Snow} \\ \text{Rain} \begin{pmatrix} .402 & .199 & .398 \\ .398 & .203 & .398 \\ .398 & .199 & .402 \end{pmatrix} \\ \text{Nice} \\ \text{Snow} \end{array}$$

$$P^5 = \begin{array}{c} \text{Rain} \quad \text{Nice} \quad \text{Snow} \\ \text{Rain} \begin{pmatrix} .400 & .200 & .399 \\ .400 & .199 & .400 \\ .399 & .200 & .400 \end{pmatrix} \\ \text{Nice} \\ \text{Snow} \end{array}$$

$$P^6 = \begin{array}{c} \text{Rain} \quad \text{Nice} \quad \text{Snow} \\ \text{Rain} \begin{pmatrix} .400 & .200 & .400 \\ .400 & .200 & .400 \\ .400 & .200 & .400 \end{pmatrix} \\ \text{Nice} \\ \text{Snow} \end{array}$$

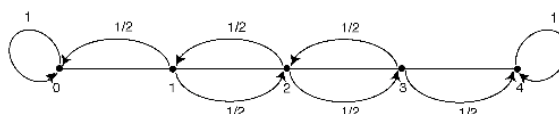
45

Introduction to Markov chains

- Let us now introduce **absorbing Markov chains**
- A state i of a Markov chain is called **absorbing** if it is impossible to leave it, $p_{ii} = 1$
- An **absorbing Markov chain** is a Markov chain containing absorbing states
 - The other states being called **transient** (TR^{ns})

Introduction to Markov chains

- Let us take an example: the drunkard's walk (from Grinstead and Snell)



$$\mathbf{P} = \begin{matrix} & \begin{matrix} 0 & 1 & 2 & 3 & 4 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 1/2 & 0 & 1/2 & 0 & 0 \\ 0 & 1/2 & 0 & 1/2 & 0 \\ 0 & 0 & 1/2 & 0 & 1/2 \\ 0 & 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}$$

47

Introduction to Markov chains

- The transition matrix can be put in canonical form:

$$\mathbf{P} = \begin{matrix} & \begin{matrix} \text{TR.} & \text{ABS.} \end{matrix} \\ \begin{matrix} \text{TR.} \\ \text{ABS.} \end{matrix} & \left(\begin{array}{c|c} \mathbf{Q} & \mathbf{R} \\ \hline \mathbf{0} & \mathbf{I} \end{array} \right) \end{matrix}$$

- \mathbf{Q} is the transition matrix between transient states

48

Introduction to Markov chains

- **R** is the transition matrix between transient and absorbing states
- Both **Q** and **R** are sub-stochastic
 - Their row sums are ≤ 1 and at least one row sum is < 1

49

Introduction to Markov chains

- Now, \mathbf{P}^t can be computed as

$$\mathbf{P}^t = \begin{array}{c} \text{TR.} \\ \text{ABS.} \end{array} \left(\begin{array}{c|c} \text{TR.} & \text{ABS.} \\ \hline \mathbf{Q}^t & \dots \\ \hline \mathbf{0} & \mathbf{I} \end{array} \right)$$

- Since **Q** is sub-stochastic, it can be shown that: $\mathbf{Q}^t \rightarrow \mathbf{0}$ as $t \rightarrow \infty$

50

Introduction to Markov chains

- The matrix

$$\mathbf{N} = \mathbf{I} + \mathbf{Q} + \mathbf{Q}^2 + \mathbf{Q}^3 + \dots = (\mathbf{I} - \mathbf{Q})^{-1}$$

– is called the **fundamental matrix** of the absorbing Markov chain

- Let us interpret the elements $n_{ij} = [\mathbf{N}]_{ij}$ of the fundamental matrix, where i, j are transient states

51

Introduction to Markov chains

- Recall that since i, j are transient states,

– we have $x_j(t | s_0 = i) = x_{ji}(t) = (\mathbf{e}_i)^T \mathbf{P}^t \mathbf{e}_j = (\mathbf{e}_i)^T \mathbf{Q}^t \mathbf{e}_j$

- Thus, entry i, j of the matrix \mathbf{N} for transient states, n_{ij} , is

$$\begin{aligned} n_{ij} &= \mathbf{e}_i^T \mathbf{N} \mathbf{e}_j \\ &= \mathbf{e}_i^T \left(\sum_{t=0}^{\infty} \mathbf{Q}^t \right) \mathbf{e}_j \\ &= \sum_{t=0}^{\infty} (\mathbf{e}_i^T \mathbf{Q}^t \mathbf{e}_j) \\ &= \sum_{t=0}^{\infty} x_{ji}(t) \end{aligned}$$

52

Introduction to Markov chains

- Thus, element n_{ij} contains the **expected number of passages** through transient state j when starting from transient state i
- The expected number of visits (and therefore steps) before being absorbed when starting from each state is

$$\mathbf{n} = \mathbf{N}\mathbf{e}$$

53

Introduction to Markov chains

- For the drunkard's walk,

$$\mathbf{P} = \begin{array}{c} \begin{array}{ccccc} & 1 & 2 & 3 & 0 & 4 \\ 1 & 0 & 1/2 & 0 & 1/2 & 0 \\ 2 & 1/2 & 0 & 1/2 & 0 & 0 \\ 3 & 0 & 1/2 & 0 & 0 & 1/2 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 4 & 0 & 0 & 0 & 0 & 1 \end{array} \end{array} .$$

54

Introduction to Markov chains

$$\mathbf{Q} = \begin{pmatrix} 0 & 1/2 & 0 \\ 1/2 & 0 & 1/2 \\ 0 & 1/2 & 0 \end{pmatrix}$$

$$\mathbf{I} - \mathbf{Q} = \begin{pmatrix} 1 & -1/2 & 0 \\ -1/2 & 1 & -1/2 \\ 0 & -1/2 & 1 \end{pmatrix}$$

$$\mathbf{N} = (\mathbf{I} - \mathbf{Q})^{-1} = \begin{matrix} & \begin{matrix} 1 & 2 & 3 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{pmatrix} 3/2 & 1 & 1/2 \\ 1 & 2 & 1 \\ 1/2 & 1 & 3/2 \end{pmatrix} \end{matrix}$$

55

Introduction to Markov chains

$$\begin{aligned} \mathbf{n} = \mathbf{N}\mathbf{e} &= \begin{pmatrix} 3/2 & 1 & 1/2 \\ 1 & 2 & 1 \\ 1/2 & 1 & 3/2 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \\ &= \begin{pmatrix} 3 \\ 4 \\ 3 \end{pmatrix} \end{aligned}$$

56

Introduction to Markov chains

- We can also compute **absorption probabilities** from each starting state
- We can compute the probability of being absorbed by absorbing state j given that we started in transient state i by

$$b_{ij} = \sum_{t=0}^{\infty} \sum_{\substack{k=1 \\ k \in \text{TR}}}^{n_{tr}} x_{k|i}(t) r_{kj}$$

where n_{tr} is the number of transient states and the sum over k is taken on the set of transient states (TR) only

57

Introduction to Markov chains

- The formula states that the probability of reaching absorbing node j at time $(t+1)$ is given by
 - the probability of passing through any state k at t and then jumping to state j from k at $(t+1)$
- The absorption probability is then given by taking the sum over all possible time steps

58

Introduction to Markov chains

- Let us compute this quantity

$$\begin{aligned}
 b_{ij} &= \sum_{t=0}^{\infty} \sum_{\substack{k=1 \\ k \in \text{TR}}}^{n_{tr}} x_{k|i}(t) r_{kj} \\
 &= \sum_{\substack{k=1 \\ k \in \text{TR}}}^{n_{tr}} \left[\sum_{t=0}^{\infty} \mathbf{e}_i^T \mathbf{Q}^t \mathbf{e}_k \right] r_{kj} \\
 &= \sum_{\substack{k=1 \\ k \in \text{TR}}}^{n_{tr}} n_{ik} r_{kj} \\
 &= [\mathbf{NR}]_{ij}
 \end{aligned}$$

59

Introduction to Markov chains

- The absorption probabilities are put in the **B** matrix
- Let us reconsider the drunkard's example

$$\mathbf{R} = \begin{matrix} & & 0 & 4 \\ \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} & \begin{pmatrix} 1/2 & 0 \\ 0 & 0 \\ 0 & 1/2 \end{pmatrix} \end{matrix}$$

60

Introduction to Markov chains

$$\mathbf{B} = \mathbf{NR} = \begin{pmatrix} 3/2 & 1 & 1/2 \\ 1 & 2 & 1 \\ 1/2 & 1 & 3/2 \end{pmatrix} \begin{pmatrix} 1/2 & 0 \\ 0 & 0 \\ 0 & 1/2 \end{pmatrix}$$

$$= \begin{matrix} & 0 & 4 \\ 1 & \begin{pmatrix} 3/4 & 1/4 \end{pmatrix} \\ 2 & \begin{pmatrix} 1/2 & 1/2 \end{pmatrix} \\ 3 & \begin{pmatrix} 1/4 & 3/4 \end{pmatrix} \end{matrix}$$

61

Introduction to Markov chains

■ Some additional definitions

- If, in a Markov chain, it is possible to go to every state from each state, the Markov chain is called **irreducible**
- Moreover, the Markov chain is called **regular** if some power of the transition matrix has only positive (non-0) elements

62



Introduction to Markov chains

- It can further be shown that the powers \mathbf{P}^t of a regular transition matrix tend to a matrix with all rows the same

$$\lim_{t \rightarrow \infty} \mathbf{P}^t = \begin{pmatrix} \boldsymbol{\pi}^T \\ \boldsymbol{\pi}^T \\ \vdots \\ \boldsymbol{\pi}^T \end{pmatrix}$$

63



Introduction to Markov chains

- Moreover, the limiting probability distribution of states is independent of the initial state:

$$\lim_{t \rightarrow \infty} \mathbf{x}(t) = \boldsymbol{\pi}$$

64

Introduction to Markov chains

- The stationary vector π is the **left eigenvector** of \mathbf{P} , corresponding to eigenvalue 1 and normalized to a probability vector:

$$\begin{aligned}
 \pi &= \lim_{t \rightarrow \infty} \mathbf{x}(t) \\
 &= \lim_{t \rightarrow \infty} \mathbf{P}^T \mathbf{x}(t-1) \\
 &= \mathbf{P}^T \lim_{t \rightarrow \infty} \mathbf{x}(t-1) \\
 &= \mathbf{P}^T \pi
 \end{aligned}$$

65

Introduction to Markov chains

- It provides the probability of finding the process in each state on the long run
- One can prove that this vector is unique for a regular Markov chain

66



Introduction to Markov chains

- Notice that the **fundamental matrix** for absorbing chains can be generalized
- To regular chains
 - It, for instance, allows to compute the **average first-passage times** in matrix form
 - See, for instance, Grinstead and Snell

67



Application to marketing

Application to marketing

- Suppose we have the following model
 - We have a number n of customer clusters or segments
 - Based, for instance, on RFM (Recency, Frequency, Monetary value)
- Each cluster is a state of a Markov chain
- The last (n th) cluster corresponds to lost customers
 - It is absorbing and generates no benefit

69

Application to marketing

- Each month, we observe the movements from cluster to cluster
- Transition probabilities are estimated
 - by counting the observed frequencies of jumping from one state to another in the past
 - This provides the entries of the transition probabilities matrix

70

Application to marketing

- Suppose also there is an **average profit**, m_i per month, associated to each customer in state i
 - which could be negative
- There is also a discounting factor:

$$0 < \gamma < 1$$
- The expected profit on an infinite time horizon can be computed

71

Application to marketing

- It is given by

$$\bar{m} = \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^n x_i(t) m_i$$

- It provides the **expected profit** on a infinite horizon

72

Application to marketing

- Which finally provides

$$\begin{aligned}
 \bar{m} &= \sum_{t=0}^{\infty} \gamma^t \sum_{i=1}^n x_i(t) m_i \\
 &= \sum_{t=0}^{\infty} \gamma^t \mathbf{m}^T (\mathbf{P}^T)^t \mathbf{x}(0) \\
 &= \sum_{t=0}^{\infty} \gamma^t \mathbf{x}^T(0) \mathbf{P}^t \mathbf{m} \\
 &= \mathbf{x}^T(0) \left(\sum_{t=0}^{\infty} \gamma^t \mathbf{P}^t \right) \mathbf{m} \\
 &= \mathbf{x}^T(0) (\mathbf{I} - \gamma \mathbf{P})^{-1} \mathbf{m}
 \end{aligned}$$

73

Application to marketing

- This is an example of the computation of the **lifetime value** of a customer
- Which is the expected profit provided by the customer until it leaves the company

74

A brief introduction to dynamic time warping

Marco Saerens, UCL
Alain Soquet, ULB

Dynamic time warping

- **Context:** word recognition
- Suppose we have a database of **word templates** or **references**
 - Someone pronounces a word which is recorded and analysed
 - We want to recover the **nearest** template word in the database = nearest neighbor
 - Which will be the recognized word

76

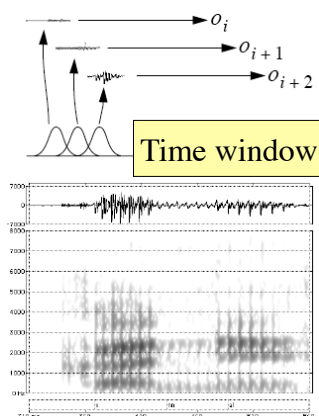
Dynamic time warping

- However, the timing in which the word has been pronounced can differ greatly
- Hence, we have to account for **distortions** or **warping** of the signal
=> Dynamic time warping

77

Dynamic time warping

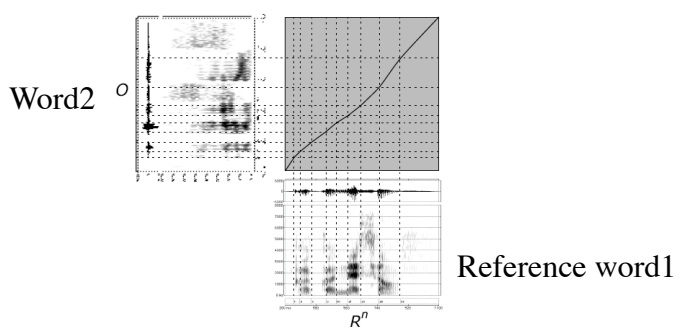
- Here is an example: a **spectrogram** is produced



78

Dynamic time warping

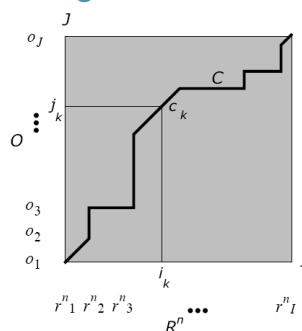
- We have to compare word1 and word2
 - Thus, align the two signals



79

Dynamic time warping

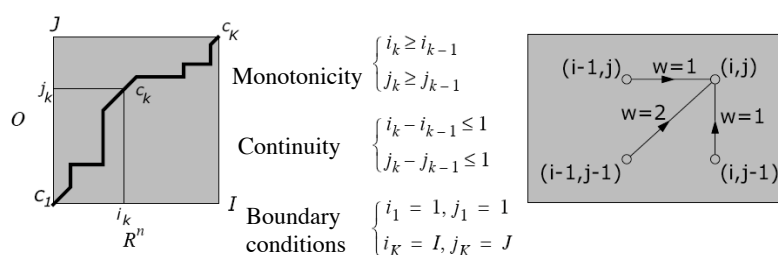
- The two signals are aligned by
 - Defining a **distance** $d(i,j)$ between two frames, for instance the Euclidean distance
 - Defining a **time alignment** that allows for warpings



80

Dynamic time warping

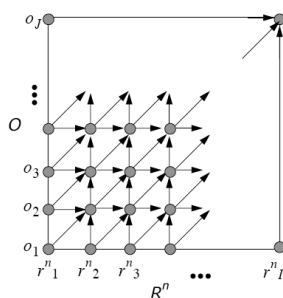
- We have to add **monotonicity constraints** in order to obtain meaningful alignments



81

Dynamic time warping

- The problem can be solved by **dynamic programming**
- By considering only the valid transitions



82

Dynamic time warping

- Here are the dynamic programming recurrence relations

$$g(1, 1) = d(r_1^n, o_1)$$

$$g(i, j) = \min \left\{ \begin{array}{l} g(i-1, j) + d(r_i^n, o_j) \\ g(i-1, j-1) + 2d(r_i^n, o_j) \\ g(i, j-1) + d(r_i^n, o_j) \end{array} \right\} \quad \text{for } \begin{cases} i = 1, 2, \dots, I \\ j = 1, 2, \dots, J \end{cases}$$

$$D(R^n, O) = \frac{1}{I+J} g(I, J)$$

83

A brief introduction to hidden Markov models

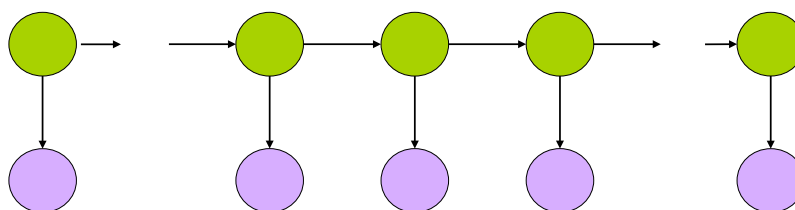
Marco Saerens, UCL

Hidden Markov models

- Most slides are adapted from David Meir Blei's slides, SRI International
- Thanks to David Meir Blei

85

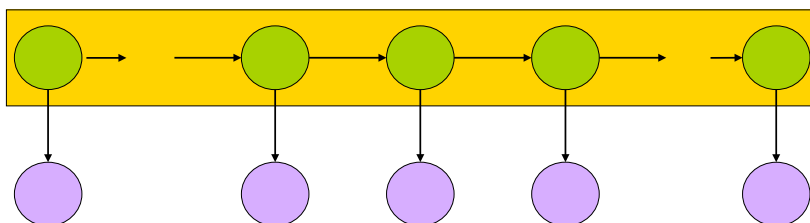
What is an HMM ?



- Graphical model
- Green circles indicate **states**
- Purple circles indicate **observations**
- Arrows indicate probabilistic dependencies between states

86

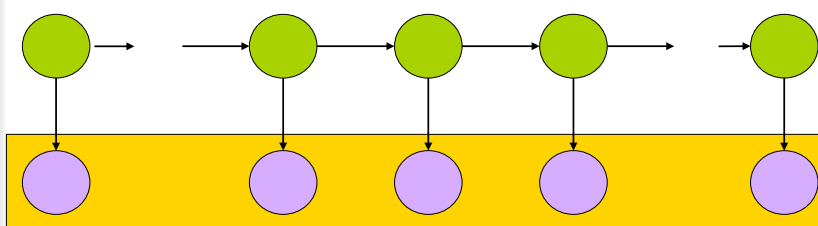
What is an HMM?



- Green circles are **hidden**, unobserved, **states** (n in total)
- Dependent only on **the previous state** (arrow)
- “The past is independent of the future given the present”: Markov property

87

What is an HMM ?

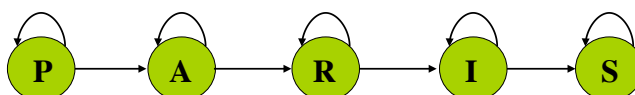


- Purple nodes are **discrete observations** (p in total)
- They depend only on their corresponding **hidden state** in which they were generated

88

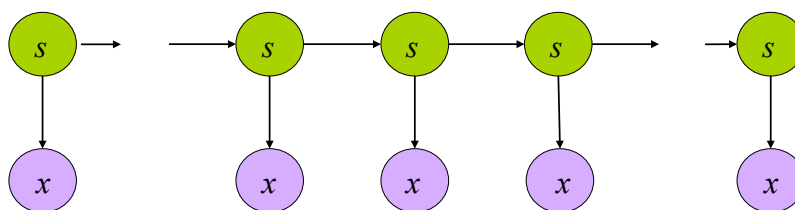
What is an HMM ?

- Example of a word model



89

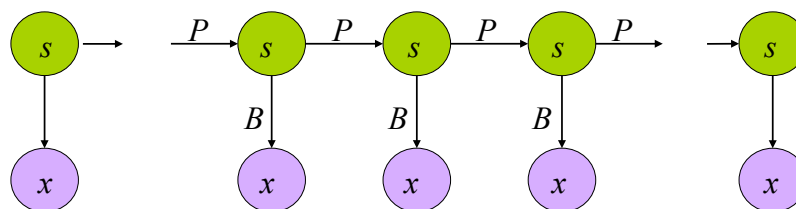
HMM formalism



- $\{s, x, \Pi, P, B\}$
- $s : \{1, \dots, n\}$ is the random variable for the hidden states taking its values from 1 to n
- $x : \{o_1, \dots, o_p\}$ is the random variable for the observations whose values are denoted o_i

90

HMM formalism

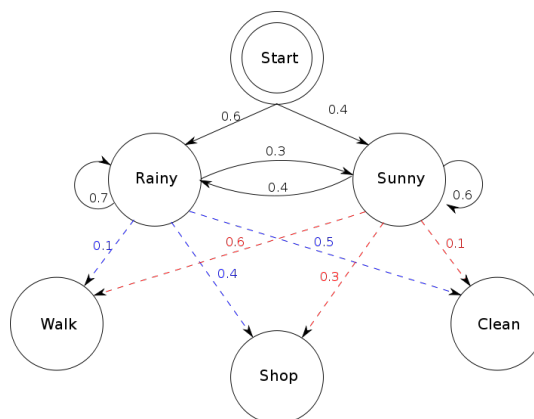


- $\{s, o, \Pi, P, B\}$
- $\Pi = \{\pi_i\}$ are the initial state probabilities, $P(s_1=i)$
- $P = \{p_{ij}\}$ are the **state transition probabilities**, $P(s_{t+1}=j | s_t=i)$
- $B = \{b_i(o_k)\}$ are the **observation or emission probabilities**, $P(x_t=o_k | s_t=i)$

91

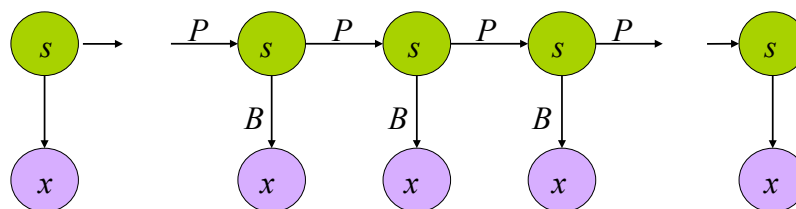
A HMM example

Taken from wikipedia



92

Inference in an HMM



■ Three basic problems:

1. Compute the **probability/likelihood** of a given observation sequence (classification) ?
2. Given an observation sequence, compute the **most likely hidden state sequence** (decoding) ?
3. Given an observation sequence, which **model parameters most closely fit the data** (parameters estimation) ?

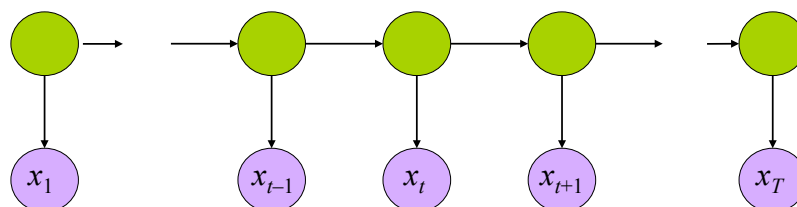
93

Likelihood computation



94

Likelihood computation



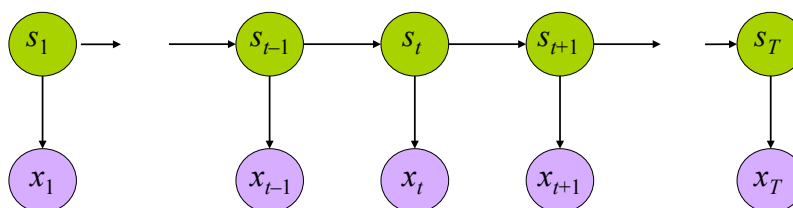
Given an observation sequence and a model, compute the **likelihood** of the observation sequence

$$\mathbf{x} = [x_1, \dots, x_T]^T, \theta = \{\Pi, P, B\}$$

Compute $P(\mathbf{x} | \theta)$

95

Likelihood computation



$$P(\mathbf{x} | \mathbf{s}, \theta) = b_{s_1}(x_1) b_{s_2}(x_2) \dots b_{s_T}(x_T)$$

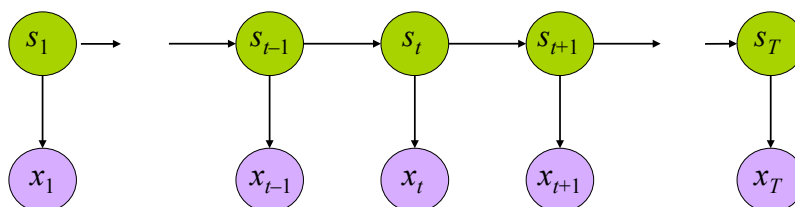
$$P(\mathbf{s} | \theta) = \pi_{s_1} p_{s_1 s_2} p_{s_2 s_3} \dots p_{s_{T-1} s_T}$$

$$P(\mathbf{x}, \mathbf{s} | \theta) = P(\mathbf{x} | \mathbf{s}, \theta) P(\mathbf{s} | \theta)$$

$$P(\mathbf{x} | \theta) = \sum_{\mathbf{s}} P(\mathbf{x} | \mathbf{s}, \theta) P(\mathbf{s} | \theta)$$

96

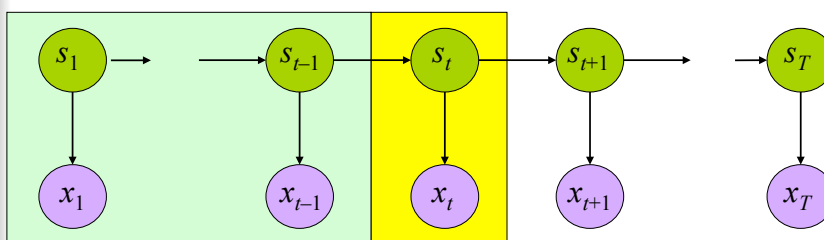
Likelihood computation



$$P(\mathbf{x} | \theta) = \sum_{(s_1 \dots s_T)} \pi_{s_1} b_{s_1}(x_1) \prod_{t=1}^{T-1} p_{s_t s_{t+1}} b_{s_{t+1}}(x_{t+1})$$

It involves a sum over **all possible sequences of states** ! 97

Forward procedure



- The special structure gives us an efficient solution using some **recurrence formulas**
- **Define:** $\alpha_i(t) = P(x_1 \dots x_t, s_t = i | \theta)$
- We will omit the dependency on θ

98

Forward procedure

$$\alpha_j(1) = P(x_1, s_1 = j)$$

$$= P(x_1 | s_1 = j)P(s_1 = j)$$

$$= b_j(x_1)\pi_j$$

99

Forward procedure

$$\alpha_j(t+1) = P(x_1 \dots x_{t+1}, s_{t+1} = j)$$

$$= \sum_{i=1}^n P(x_1 \dots x_t, s_t = i, x_{t+1}, s_{t+1} = j)$$

$$= \sum_{i=1}^n P(x_{t+1}, s_{t+1} = j | x_1 \dots x_t, s_t = i) P(x_1 \dots x_t, s_t = i)$$

$$= \sum_{i=1}^n P(x_{t+1}, s_{t+1} = j | s_t = i) \alpha_i(t)$$

$$= \sum_{i=1}^n P(x_{t+1} | s_{t+1} = j) P(s_{t+1} = j | s_t = i) \alpha_i(t)$$

$$= \left(\sum_{i=1}^n p_{ij} \alpha_i(t) \right) b_j(x_{t+1})$$

100

Forward procedure

$$\alpha_j(1) = b_j(x_1)\pi_j$$

$$\alpha_j(t+1) = \left(\sum_{i=1}^n p_{ij}\alpha_i(t) \right) b_j(x_{t+1})$$

101

Forward procedure

From Rabiner et al.:

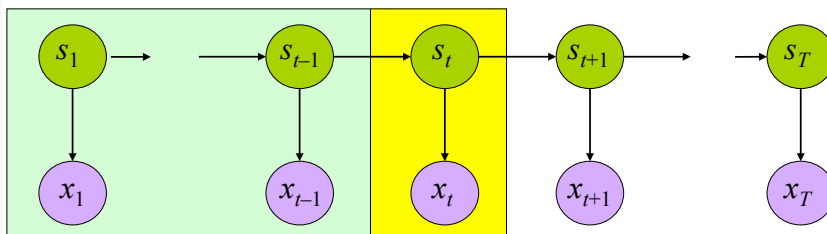
(a)

(b)

Fig. 4. (a) Illustration of the sequence of operations required for the computation of the forward variable $\alpha_{t+1}(j)$.
 (b) Implementation of the computation of $\alpha_t(i)$ in terms of a lattice of observations t , and states i .

102

Likelihood computation: solutions

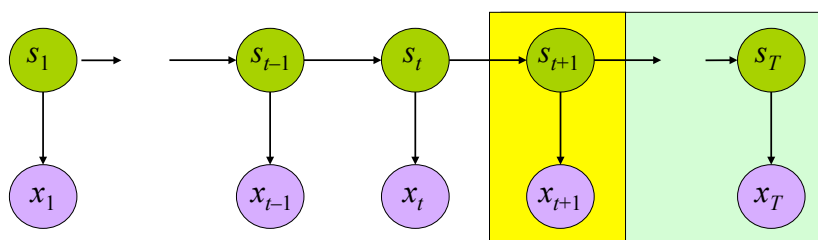


The likelihood is:

$$P(\mathbf{x} | \theta) = \sum_{j=1}^n P(x_1 \dots x_T, s_T = j | \theta) = \sum_{j=1}^n \alpha_j(T)$$

103

Backward procedure



- We now introduce the **backward variable**

- **Define:** $\beta_i(t) = P(x_{t+1} \dots x_T | s_t = i, \theta)$

- We will omit the dependency on θ

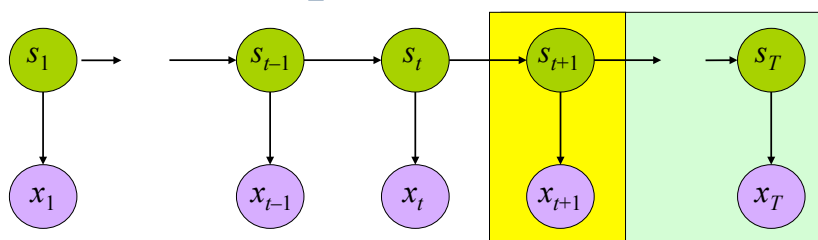
104

Backward procedure

$$\begin{aligned}
 \beta_i(t) &= P(x_{t+1} \dots x_T \mid s_t = i) \\
 &= \sum_{j=1}^n P(x_{t+1} x_{t+2} \dots x_T, s_{t+1} = j \mid s_t = i) \\
 &= \sum_{j=1}^n P(x_{t+1} x_{t+2} \dots x_T \mid s_t = i, s_{t+1} = j) P(s_{t+1} = j \mid s_t = i) \\
 &= \sum_{j=1}^n p_{ij} P(x_{t+1} x_{t+2} \dots x_T \mid s_{t+1} = j) \\
 &= \sum_{j=1}^n p_{ij} P(x_{t+1} \mid x_{t+2} \dots x_T, s_{t+1} = j) P(x_{t+2} \dots x_T \mid s_{t+1} = j) \\
 &= \sum_{j=1}^n p_{ij} b_j(x_{t+1}) \beta_j(t+1)
 \end{aligned}$$

105

Backward procedure

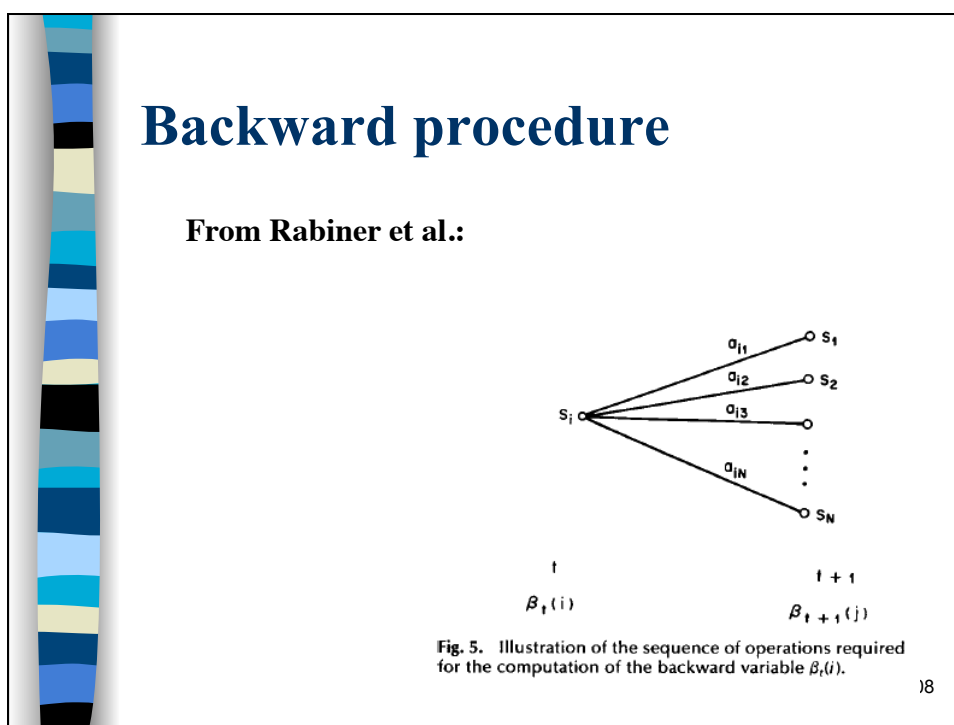
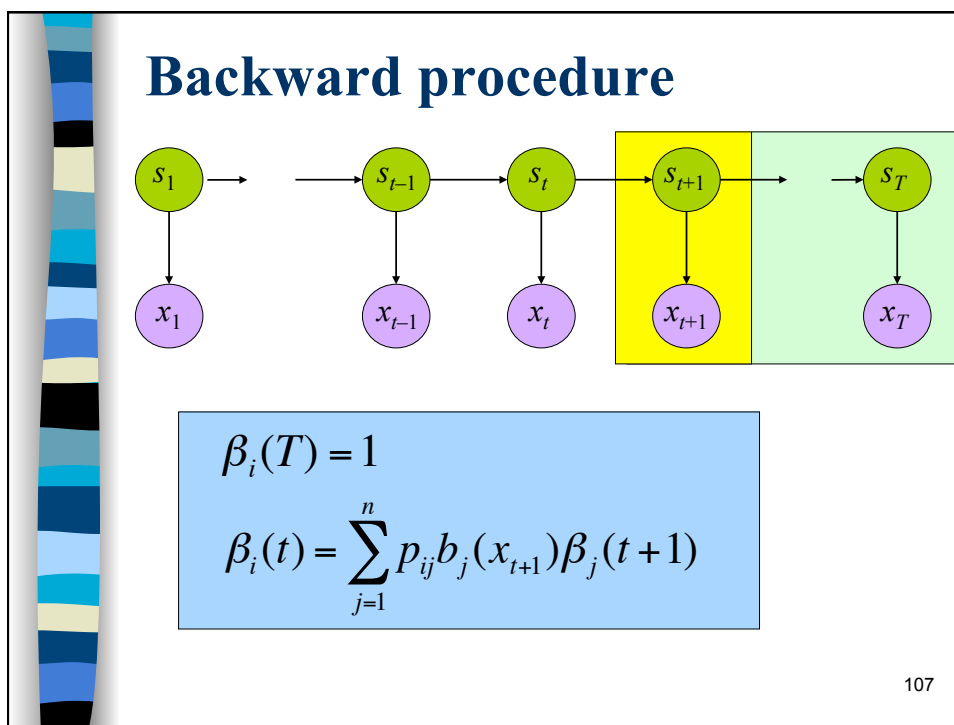


- When $t = T - 1$, we have

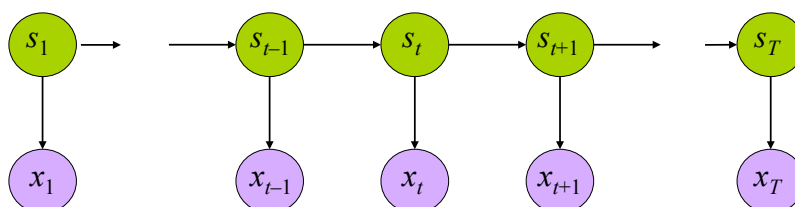
$$\begin{aligned}
 \beta_i(T-1) &= P(x_T \mid s_{T-1} = i) \\
 &= \sum_{j=1}^n p_{ij} b_j(x_T) \\
 \beta_i(T-1) &= \sum_{j=1}^n p_{ij} b_j(x_T) \beta_j(T)
 \end{aligned}$$

$$\Rightarrow \beta_j(T) = 1$$

106



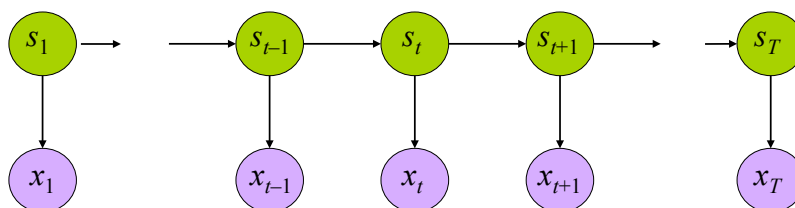
Likelihood computation: solutions



$$\begin{aligned}
 P(\mathbf{x}, s_t = i | \theta) &= P(x_1 \dots x_t, s_t = i, x_{t+1} \dots x_T | \theta) \\
 &= P(x_{t+1} \dots x_T | x_1 \dots x_t, s_t = i, \theta) P(x_1 \dots x_t, s_t = i | \theta) \\
 &= P(x_{t+1} \dots x_T | s_t = i, \theta) P(x_1 \dots x_t, s_t = i | \theta) \\
 &= \alpha_i(t) \beta_i(t)
 \end{aligned}$$

109

Likelihood computation: solutions



$$P(\mathbf{x} | \theta) = \sum_{i=1}^n \alpha_i(T) \quad \text{Forward procedure}$$

$$P(\mathbf{x} | \theta) = \sum_{i=1}^n \pi_i \beta_i(1) \quad \text{Backward procedure}$$

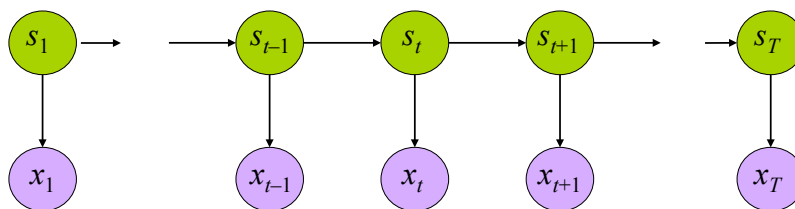
$$P(\mathbf{x} | \theta) = \sum_{i=1}^n \alpha_i(t) \beta_i(t) \quad \text{Combination}$$

Optimal state sequence (decoding)



111

Best state



- Find the **most probable state** at time t given the observations

$$\begin{aligned} \gamma_i(t) &= P(s_t = i | \mathbf{x}, \theta) = \frac{P(\mathbf{x}, s_t = i | \theta)}{P(\mathbf{x} | \theta)} \\ &= \frac{\alpha_i(t) \beta_i(t)}{\sum_{j=1}^n \alpha_j(t) \beta_j(t)} \end{aligned}$$

112

Best state sequence

- Find the **state sequence** that best explains the observations
- **Viterbi algorithm** = **dynamic programming algorithm**

$$\arg \max_s P(s | \mathbf{x}, \theta) = \arg \max_s P(s, \mathbf{x} | \theta)$$

113

Viterbi algorithm

$$\delta_j(t) = \max_{s_1 \dots s_{t-1}} P(s_1 \dots s_{t-1}, x_1 \dots x_{t-1}, s_t = j, x_t | \theta)$$

The state sequence which maximizes the probability of generating the observations up to time $t-1$, landing in state j , and emitting the observation at time t

114

Viterbi algorithm

$\delta_j(t) = \max_{s_1 \dots s_{t-1}} \text{P}(s_1 \dots s_{t-1}, x_1 \dots x_{t-1}, s_t = j, x_t \mid \theta)$

$$\delta_j(1) = \pi_j b_j(x_1)$$

$$\delta_j(t+1) = \max_i \{ \delta_i(t) p_{ij} b_j(x_{t+1}) \}$$

$$\psi_j(1) = 0$$

$$\psi_j(t+1) = \arg \max_i \{ \delta_i(t) p_{ij} b_j(x_{t+1}) \}$$

Recursive computation

115

Viterbi algorithm

- Indeed, $\delta_j(t+1)$ is equal to

$$\begin{aligned} \max_{s_1 \dots s_t} \text{P}(s_1 \dots s_{t+1} = j, x_1 \dots x_{t+1} \mid \theta) &= \max_{s_1 \dots s_t} \text{P}(s_{t+1} = j, x_{t+1} \mid s_1 \dots s_t, x_1 \dots x_t, \theta) \text{P}(s_1 \dots s_t, x_1 \dots x_t \mid \theta) \\ &= \max_{s_1 \dots s_t} \{ \text{P}(x_{t+1} \mid s_{t+1} = j, \theta) \text{P}(s_{t+1} = j \mid s_t, \theta) \text{P}(s_1 \dots s_t, x_1 \dots x_t \mid \theta) \} \\ &= \max_{s_t} \max_{s_1 \dots s_{t-1}} \{ b_{s_{t+1}}(x_{t+1}) p_{s_t s_{t+1}} \text{P}(s_1 \dots s_t, x_1 \dots x_t \mid \theta) \} \\ &= \max_{s_t} \left\{ b_j(x_{t+1}) p_{s_t j} \max_{s_1 \dots s_{t-1}} \text{P}(s_1 \dots s_t, x_1 \dots x_t \mid \theta) \right\} \\ &= \max_{s_t} \{ b_j(x_{t+1}) p_{s_t j} \delta_{s_t}(t) \} \end{aligned}$$

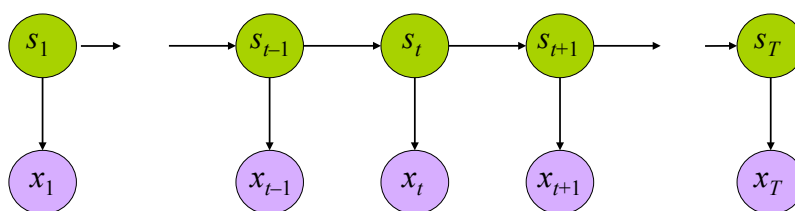
116

Viterbi algorithm

- And we can apply **dynamic programming** to the log-likelihood
 - The cost is then additive

117

Viterbi algorithm



$$s^*(T) = \arg \max_i \delta_i(T)$$

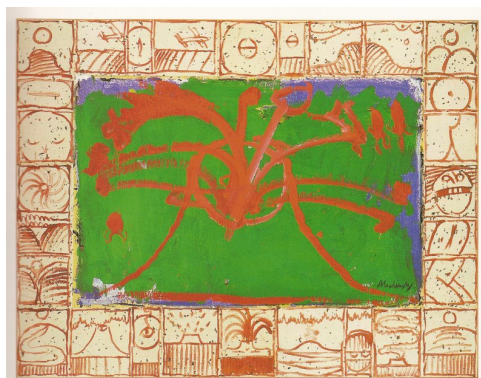
$$s^*(t-1) = \psi_{s^*(t)}(t)$$

$$P(\mathbf{x}^*) = \max_i \delta_i(T)$$

Compute the most likely state sequence by working backwards

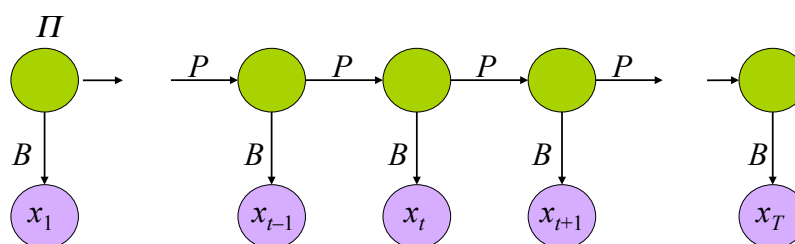
118

Parameter estimation



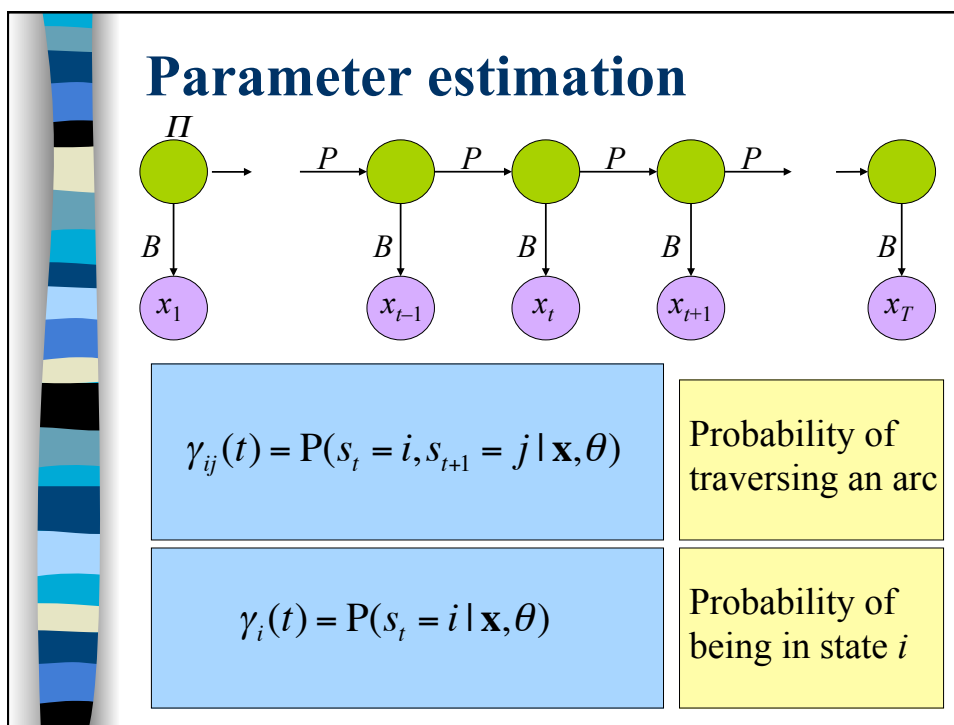
119

Parameter estimation



- Given an observation sequence, find the model parameters Π, P, B , that most likely produce that sequence (maximum likelihood)
- No closed-form solution
- Instance of the iterative [EM algorithm](#)

120



Parameter estimation

Recall that we already computed

$$\gamma_i(t) = P(s_t = i \mid \mathbf{x}, \theta) = \frac{P(\mathbf{x}, s_t = i \mid \theta)}{P(\mathbf{x} \mid \theta)}$$

$$= \frac{\alpha_i(t)\beta_i(t)}{\sum_{j=1}^n \alpha_j(t)\beta_j(t)}$$

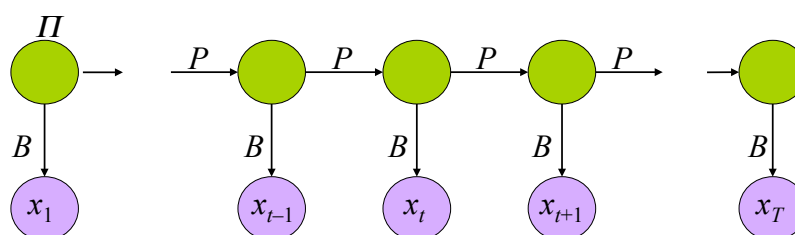
122

Parameter estimation

$$\begin{aligned}
 \gamma_{ij}(t) &= P(s_t = i, s_{t+1} = j | \mathbf{x}, \theta) = \frac{P(x_1 \dots x_t, s_t = i, s_{t+1} = j, x_{t+1} \dots x_T | \theta)}{P(\mathbf{x} | \theta)} \\
 &= \frac{P(s_{t+1} = j, x_{t+1} \dots x_T | x_1 \dots x_t, s_t = i, \theta) P(x_1 \dots x_t, s_t = i | \theta)}{P(\mathbf{x} | \theta)} \\
 &= \frac{P(s_{t+1} = j, x_{t+1} \dots x_T | s_t = i, \theta) \alpha_i(t)}{P(\mathbf{x} | \theta)} \\
 &= \frac{P(x_{t+1} \dots x_T | s_t = i, s_{t+1} = j, \theta) P(s_{t+1} = j | s_t = i, \theta) \alpha_i(t)}{P(\mathbf{x} | \theta)} \\
 &= \frac{P(x_{t+1} \dots x_T | s_{t+1} = j, \theta) p_{ij} \alpha_i(t)}{P(\mathbf{x} | \theta)} \\
 &= \frac{P(x_{t+1} | x_{t+2} \dots x_T, s_{t+1} = j, \theta) P(x_{t+2} \dots x_T | s_{t+1} = j, \theta) p_{ij} \alpha_i(t)}{P(\mathbf{x} | \theta)} \\
 &= \frac{\alpha_i(t) b_j(x_{t+1}) p_{ij} \beta_j(t+1)}{\sum_{k=1}^N \alpha_k(t) \beta_k(t)}
 \end{aligned}$$

123

Parameter estimation



$\hat{\pi}_i$ = probability of starting from i

$\hat{\pi}_i = P(s_1 = i | \mathbf{x}, \hat{\theta}) = \gamma_i(1)$

Now we can compute the new estimates of the model parameters.

124

Parameter estimation

$\hat{p}_{ij} = \frac{\text{expected number of transitions from } i \text{ to } j}{\text{expected number of transitions out of state } i}$

125

Parameter estimation

$$\hat{p}_{ij} = \frac{\sum_{t=1}^{T-1} P(s_t = i, s_{t+1} = j | \mathbf{x}, \hat{\theta})}{\sum_{t=1}^{T-1} P(s_t = i | \mathbf{x}, \hat{\theta})} = \frac{\sum_{t=1}^{T-1} \gamma_{ij}(t)}{\sum_{t=1}^{T-1} \gamma_i(t)}$$

126

Parameter estimation

$$\hat{b}_i(o_k) = \frac{\text{expected number of emissions of } o_k \text{ in state } i}{\text{total number of emissions in state } i}$$

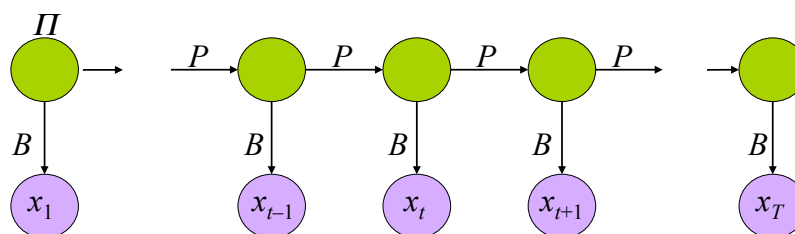
127

Parameter estimation

$$\hat{b}_i(o_k) = \frac{\sum_{t=1}^T P(s_t = i, x_t = o_k | \mathbf{x}, \hat{\theta})}{\sum_{t=1}^T P(s_t = i | \mathbf{x}, \hat{\theta})} = \frac{\sum_{t=1}^T P(s_t = i | \mathbf{x}, \hat{\theta}) \delta(x_t = o_k)}{\sum_{t=1}^T P(s_t = i | \mathbf{x}, \hat{\theta})}$$

$$= \frac{\sum_{\{t: x_t = o_k\}} P(s_t = i | \mathbf{x}, \hat{\theta})}{\sum_{t=1}^T P(s_t = i | \mathbf{x}, \hat{\theta})} = \frac{\sum_{\{t: x_t = o_k\}} \gamma_t(i)}{\sum_{t=1}^T \gamma_t(i)}$$

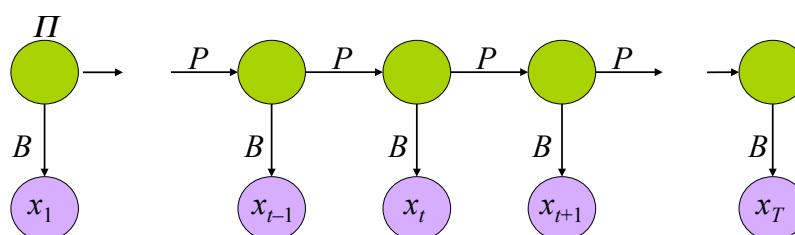
Parameter estimation



- The two following two steps are iterated until convergence:
 - Recompute the forward and backward variables α and β
 - Recompute the parameter estimates for the Π , P , B

129

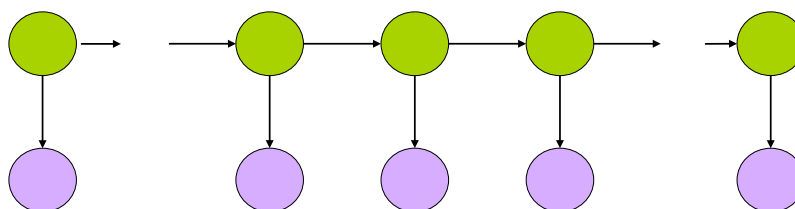
Parameter estimation



- It can be shown that this iterative algorithm increases the likelihood

130

HMM applications



- Generating parameters for n-gram models
- Tagging speech
- Speech recognition
- Bioinformatics sequence modeling

131

HMM applications

- Part-of-speech tagging
 - The representative put chairs on the table
 - AT NN VBD NNS IN AT NN
 - AT JJ NN VBZ IN AT NN
- Some tags :
 - AT: article, NN: singular or mass noun, VBD: verb, past tense, NNS: plural noun, IN: preposition, JJ: adjective

132



HMM applications

- Bioinformatics
 - Durbin et al. Biological Sequence Analysis, Cambridge University Press.
- Several applications, e.g. proteins
 - From primary structure `ATCPLELLLD`
 - Infer secondary structure `HHHBBBBBC..`