# Algorithmic Data Analysis

Esther Galbrun

Spring 2024
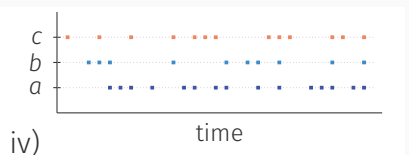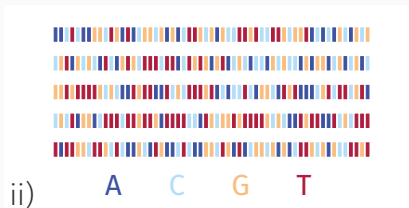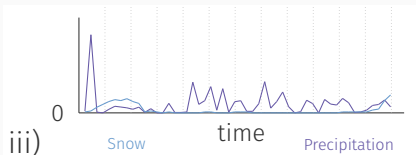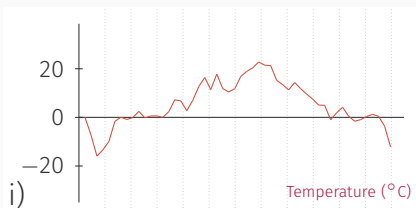
UNIVERSITY OF
EASTERN FINLAND

Associate characteristics to the datasets

univariate vs. multivariate     sequential data vs. time-series
regularly vs. irregularly sampled     real vs. symbolic values



i) Temperature (°C)

ii) A   C   G   T

iii) Snow   time   Precipitation

iv) time

## Q4.2: Distances (i)

We consider discrete sequences of items.

The distance between two sequence elements is defined as
$d(s, s') = 0$ if and only if $s = s'$, 1 otherwise.

Given two sequences $S$ and $S'$, we denote

$\text{DTW}(S, S')$ the *Dynamic Time Warping* (DTW) distance,

$\text{DTW}_w(S, S')$ the DTW distance with window constraint $w$,
   i.e. matching elements no further than $w$ positions apart

$\text{E}(S, S')$ the *Edit* distance with insertion and deletion operations
   such that $c_{\text{ins}} = 1$ and $c_{\text{del}} = 0.5$,

$\text{H}(S, S')$ the *Hamming distance*,

$\text{L}(S, S')$ the length of
   the *Longest contiguous common subsequence.*

## Q4.2: Distances (i)

$S_A$ and $S_B$ are two sequences of length 10 such that
$L(S_A, S_B) = 5$ and $S_C$ is a third sequence obtained by deleting
the first and the last elements of $S_A$.

What can you say about the following statements?

 i) $H(S_A, S_B) \leq H(S_C, S_B) + 2$
 ii) $DTW_3(S_A, S_B) < DTW(S_A, S_B)$
iii) $DTW(S_A, S_B) \leq H(S_A, S_B)$
 iv) $E(S_A, S_B) \leq H(S_A, S_B)$
 v) $E(S_A, S_B) < 8$
 vi) $E(S_A, S_C) \leq 1$
vii) E is a metric
viii) H is a metric
 ix) L is a metric

## Q4.3: Frequent sequences

We want to mine frequent sequences from a long sequence of itemsets, with minimum support threshold set to 4

The frequent sequences of length 3, with their support, are

| | | | |
|---|---|---|---|
| $\{a, b, c\} : 7$ | $\{a, b\}\{a\} : 4$ | $\{a, b\}\{c\} : 5$ | $\{a\}\{a, b\} : 6$ |
| $\{a\}\{a, c\} : 8$ | $\{a\}\{b, c\} : 6$ | $\{b\}\{a, c\} : 5$ | $\{b\}\{b, c\} : 4$ |

Provide the tightest possible upper bound on the support of each of the following sequences of length 4

$\text{supp}(\{a, b\}\{a, c\}) \leq ?$
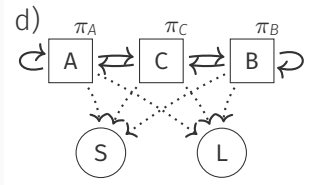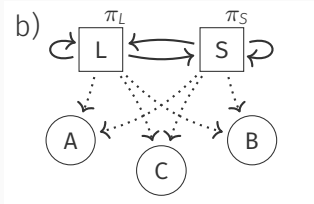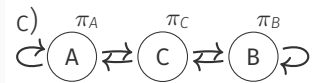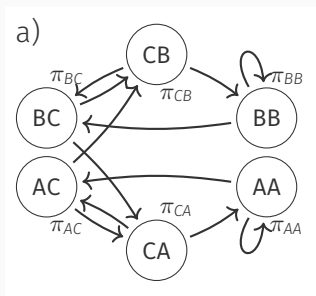$\text{supp}(\{a, b\}\{b, c\}) \leq ?$
$\text{supp}(\{a\}\{a, b, c\}) \leq ?$
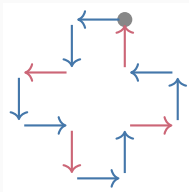$\text{supp}(\{a\}\{a, c\}\{a\}) \leq ?$

Associate each model to its name and to the size of the corresponding transition and/or emission matrices

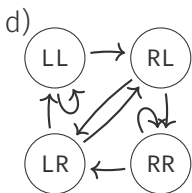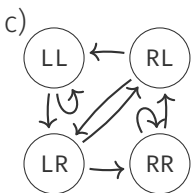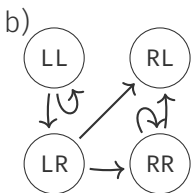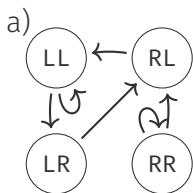Markov Chain    Hidden Markov Model    first/second order



a)

b)

c)

d)

## Q4.5: Markov unchained

At each time step, $\mathcal{R}$ turns either left (L) or right (R) then rolls forward by one unit of distance. The path followed during one run, starting and ending at the gray the dot, facing up, is shown on the right.
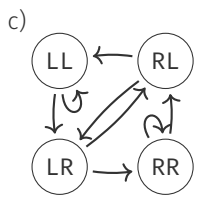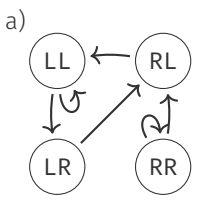


From which of the Markov chains below did it possibly arise?



a)
b)
c)
d)

The path can be represented as sequence $S = $ LLRLLRLLRLLR

From which of the Markov chains below did it most likely arise?

a)



c)



$$\pi_a = \begin{pmatrix} \text{LL} & \text{LR} & \text{RL} & \text{RR} \\ 0.25 & 0.25 & 0.25 & 0.25 \end{pmatrix}$$

$$A_a = \begin{pmatrix} 0.90 & 0.10 & 0.00 & 0.00 \\ 0.00 & 0.00 & 1.00 & 0.00 \\ 1.00 & 0.00 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.90 & 0.10 \end{pmatrix}$$

$$\pi_c = \begin{pmatrix} \text{LL} & \text{LR} & \text{RL} & \text{RR} \\ 0.35 & 0.15 & 0.15 & 0.35 \end{pmatrix}$$

$$A_c = \begin{pmatrix} 0.25 & 0.75 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.40 & 0.60 \\ 0.45 & 0.55 & 0.00 & 0.00 \\ 0.00 & 0.00 & 0.25 & 0.75 \end{pmatrix}$$

## Q4.7: HMM problems

Match tasks, solutions and algorithms

### Tasks

Evaluation
Explanation
Training

### Solutions

$P_{\mathcal{M}}(O)$
$\arg \max_{X \in \mathcal{X}} P_{\mathcal{M}}(X, O)$
$\arg \max_{\mathcal{M} \in \mathcal{H}} P_{\mathcal{M}}(O)$

### Algorithms

Backward algorithm
Baum–Welch algorithm
Forward algorithm
Viterbi algorithm