#### Sequence analysis for social scientists Introduction to sequence analysis

Alexis Gabadinho, Matthias Studer, Gilbert Ritschard, Nicolas S. Müller

Department of Econometrics, University of Geneva http://mephisto.unige.ch/biomining

Summer School on Advanced Methods for the Analysis of Complex Event History Data, Bristol, 28-29 June 2010

1/19

Sequence Sequence analysis 00000 State sequences

#### Definitions

- Alphabet A: finite set of possible states, for example: leaves in the parental home, leaves alone, leaves with partner and no children, leaves with partner and one or more children, leaves with no partner and one or more children
- Sequence of length k: ordered list of k elements taken from A
- State sequences are present in many domains:
  - Text: A = set of letters, words, structure of sociological articles **?**, ...
  - Biology: A = set of nucleotids, proteins, ...
  - Trajectories, biographies: *A* = set of leaving arrangements, professional status, daily activities (time-use diary data, workdays, workweeks), ...

Objectives

Objectives

#### bjectives

DE GENÉVE

#### Basics concepts

- State sequences: ordered lists of states defined on a time axis
- In the social sciences, such state sequences are of interest for studying life trajectories such as
  - occupational histories
  - professional carriers
  - cohabitational life courses
- Data collected for example by panel or retrospective surveys
- The states are values of a categorical variable (hence we talk also of categorical sequences)
- Event sequences: ordered lists of time stamped events

<u> </u>			
1			

Sequence Sequence analysis Objective Sequence analysis Objective Section Secti

- Study from? on transition from school to work in Northern Ireland
- Included in the TraMineR library.
- 712 individuals
- Follow-up starting at the end of the compulsory education (July 1993)
- Time series of 70 status variables: September 1993 to June 1999.
- The alphabet is made of the following statuses: EM (employment), FE (Further Education), HE (Higher Education), JL (Joblessness), SC (School), TR (Training).

#### Sequence

#### Sequence ana

#### Aim of the original study

00000	000	anarysi
he mvad data set -	Variable	list

unique individual identifier

#### Table: List of Variables in the MVAD data set

weight sample weights binary dummy for gender, 1=male male catholic binary dummy for community, 1=Catholic Belfast binary dummies for location of school, one of five Education and Library Board areas in Northern Ireland N.Eastern Southern S.Eastern Western Grammar binary dummy indicating type of secondary education, 1=grammar school funemp binary dummy indicating father's employment status at time of survey, 1=father unemployed binary dummy indicating qualifications gained by the end of compulsory education, 1 = 5 +gcse5eq GCSEs at grades A-C, or equivalent fmpr binary dummy indicating SOC code of father's current or most recent job,1=SOC1 (professional, managerial or related) livboth binary dummy indicating living arrangements at time of first sweep of survey (June 1995), 1=living with both parents jul93 Monthly Activity Variables are coded 1-6, 1=school, 2=FE, 3=employment, 4=training, 5=joblessness, 6=HE

7/

Т

id

jun99



- Holistic approach of the life courses versus "event-oriented" approach [?]
- Approach promoted by ?
- Event history or survival analysis focuses on the occurrence of one specific event in the life course
- Sequence analysis allows a more general overview of the life courses
- The holistic approach is based originally on exploratory analysis (data mining techniques)
- Explanatory methods developed more recently

- The aim of the original study was to
  - Use sequence techniques to characterize young peoples' transitions from school to work into types
  - Distinguish between successful and unsuccessful transitions
  - Link transition type with a collection of static individual, family and school characteristics
  - Identify which young people are more at risk to experience unsuccessful transitions into the labour market

5/19	

Sequence

Sequence analysis

Graphical representation of state sequences

• Here is a graphical representation of the first 10 state sequences



employment
 higher education
 school
 further education
 joblessness
 training



Objectives

#### Sequence analysis 0.00

#### Typical questions

- Some of the typical questions arising in the social sciences:
  - Do life courses obey some social norm? Which are the standard trajectories? What kind of departures do we observe from these standards ? Do we observe a de-standardization of the life trajectories ?

Objectives

- Why are some people more at risk to follow a chaotic trajectory?
- How is the life trajectory related to sex, social origin and other cultural factors?
- How is a given outcome, such as health status or income, related to a trajectory?

11/19		W definit medication
Sequence 000000	Sequence analysis 000	Objectives •00000
Objectives		

During this course, you will learn to:

- Use the R statistical environnement and the TraMineR package.
- Visualize state sequences with various representations



Sep.93 Sep.94 Sep.95 Sep.96 Sep.97 Sep.98

#### Which tools do we need?

- By considering whole sequences we deal with complex objects
- We need special tools for describing and displaying them

000

- Questions regarding the exploration and description of sets of sequences such as
  - Which characteristics of sequences are we interested in?
  - What kind of indicators can we compute for a sequence set?
  - Which plots are suited for rendering sequences?
  - How can we measure similarity between sequences?
- From a more explanatory perspective, how can we measure and assess the link between covariates and whole sequences?

UNIVERSITÉ DE GENÉVE

## **Objectives** - II

• Describe sets of sequences using overall and transversal characteristics.

Sequence analysis



 Describe individual sequences using longitudinal characteristics.

Objectives 00000

Objectives



• Build and visualize a typology of sequences



# Objectives - V

- Focus on event rather than states by using event sequences analysis.
- Find the most frequent event subsequences.
- Find the most discriminating event subsequences.



000000

## Objectives - IV

- Measure and assess the association between sequences and one or several covariates using sequence discrepancy analysis.
- Create sequence regression trees.



17/1



- elefences i
  - Andrew Abbott and Emily Barman. Sequence comparison via alignment and gibbs sampling: A formal analysis of the emergence of the modern sociological article. *Sociological Methodology*, 27:47–87, 1997.
- Andrew Abbott and John Forrest. Optimal matching methods for historical sequences. *Journal of Interdisciplinary History*, 16: 471–494, 1986.
- Francesco C. Billari. Sequence analysis in demographic research. *Canadian Studies in Population*, 28(2):439–458, 2001. Special Issue on Longitudinal Methodology.
- Duncan McVicar and Michael Anyadike-Danes. Predicting successful and unsuccessful transitions from school to work by using sequence methods. *Journal of the Royal Statistical Society. Series A (Statistics in Society)*, 165(2):317–334, 2002. ISSN 09641998. URL http://www.jstor.org/stable/3559930.





Introduction Using R Objects Importing data into R Exploring data frames Working environment References

#### What is R?

## Sequence analysis for social scientists Part II - Introduction to R

Alexis Gabadinho, Matthias Studer, Gilbert Ritschard, Nicolas S. Müller

Department of Econometrics, University of Geneva http://mephisto.unige.ch/biomining

Summer School on Advanced Methods for the Analysis of Complex Event History Data, Bristol, 28-29 June 2010

1/44						Transient Byr ei Friedram Organization di anametrice
Introduction	Using R	Objects	Importing data into R	Exploring data frames	Working environment	References
The R	conso	ole				

• When launching R you get a R console



- R is a free software environment for statistical computing and graphics (can be downloaded at http://www.r-project.org)
- R is also a language derived from the S language
- R is free and open source
- R is available for Linux, MacOS X, Windows

3/44

troduction Using R Objects Importing data into R Exploring data frames Working environment References

## Passing commands to ${\sf R}$

- The prompt '>' indicates that R is waiting commands
- There are several ways of passing commands to R
  - By typing them directly in the console
  - By using the script editor (Windows and MacOS X only)
  - By using the R Commander graphical interface (available from the CRAN)
  - By using an external editor
- R scripts are files containing a series of R commands
- To add comments to your code use '#' before
- The usual extension for such files is '.R'
- Using scripts allows to
  - store, re-use or modify later your statistical analyses
  - send to/share with someone else



- [1] 2.5
- > A <- 8 > a == A
- [1] FALSE
- function is a
- A function is always called by typing its name followed by brackets.
- Inside the brackets you put zero or more arguments. Arguments can be values or existing objects
   b

```
[1] "my object"
> c <- paste(b, "is beautiful")
> c
```

```
[1] "my object is beautiful"
```

better to give explicitly the names of the other optional

• Using argument names allows to pass them in any order

arguments you may have to pass to functions

> seq(from = 1, to = 10, by = 2)

> seq(by = 2, to = 10, from = 1)

[1] 2 3 4 5 6 7 8 9 10

[1] 1 3 5 7 9

[1] 1 3 5 7 9

[1] 1 3 5 7 9

> seq(2, 10, 1)

> seq(1, 10, 2)

![](_page_5_Picture_11.jpeg)

![](_page_6_Figure_0.jpeg)

Higher dimensional objects - Vectors

Vectors are one-dimensional objects containing numeric or character values. The c() function combines values into a vector
v1 <- c(1, 2, 4, 8)</li>
v1
[1] 1 2 4 8
v2 <- c("A", "B", "C", "D")</li>
v2
[1] "A" "B" "C" "D"
Specific elements of vectors can be retrieved with indexes > v1[3]

Objects Importing data into R Exploring data frames Working environment

[1] 4
> v2[1:3]
[1] "A" "B" "C"
> v2[c(1, 4)]
[1] "A" "D"

Indexing with logical expressions

One can use logical expressions to retrieve vectors elements
 v1[v1 >= 4]

Importing data into R Exploring data frames Working environment

[1] 4 8

Objects

- > v2[v2 %in% c("A", "C")]
- [1] "A" "C"

Using R

- Use which() to get the indexes of the elements satisfying one condition
  - > which(v1 >= 4)
    [1] 3 4
    > which(v2 %in% c("A", "C"))
    [1] 1 3

Introduction

Using R

UNIVERSITÉ DE GENÉVE

References

References

Introduction Using R **Objects** Importing data into R Exploring data frames Working environment Reference

#### Higher dimensional objects - Matrices

• Matrices are two dimensional objects containing numeric or character values

```
> m1 <- matrix(1:16. 4. 4)
> m1
     [,1] [,2] [,3] [,4]
[1,]
            5
                 9
                   13
      1
[2,]
       2
            6
                10
                     14
[3.]
       3
            7
                11
                     15
[4,]
                12
                     16
       4
            8
```

- Specific elements of matrices are retrieved with row and column indexes
  - Element in the second row, fourth column
    - > m1[2, 4]
    - [1] 14
  - Whole fourth column (by omitting row index)
     m1[, 4]
    - [1] 13 14 15 16

Introduction	Using R	Objects	Importing data into R	Exploring data frames	Working environment	References
Factor	S					

- Factors are categorical variables, that is variables taking a value that is one of several possible categories
- The Species variable in the iris data frame is a factor > class(iris\$Species)

```
[1] "factor"
```

- Possible categories of a factor are called levels
  - > levels(iris\$Species)
  - [1] "setosa" "versicolor" "virginica"
- You can change the labels for the levels
  - > levels(iris\$Species) <- c("Species 1", "Species 2", "Species 3")
    > head(iris\$Species)
  - [1] Species 1 Species 1 Species 1 Species 1 Species 1
    Levels: Species 1 Species 2 Species 3

## Higher dimensional objects - Data frames

- Data frames combine columns (vectors) of any type: factors, numeric, character strings
  - > data(iris)
  - > iris[1:4, ]

	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3.0	1.4	0.2	setosa
3	4.7	3.2	1.3	0.2	setosa
4	4.6	3.1	1.5	0.2	setosa

- You can access variables in data frames using their name preceded by a '\$' instead of the column number
  - > iris\$Sepal.Width[1:4]
  - [1] 3.5 3.0 3.2 3.1
- > iris[1:4, 2]
- [1] 3.5 3.0 3.2 3.1

17/4

oduction Using R **Objects** Importing data into R Exploring data frames Working environment References

## Objects and methods

- There are many other types of objects in R
- For example contingency tables or outputs of regression models are objects of a specific type
- There are usually specific methods for each type of object like print(), plot() or summary() methods
- Indeed when typing the name of an object to display its content the print() method is called automatically
  - > A
  - [1] 8
  - > print(A)
  - [1] 8

UNIVERSITÉ DE GENÉVE DE GENEVE

![](_page_8_Figure_0.jpeg)

[1] 86

![](_page_9_Picture_0.jpeg)

i- ℝ Ee	liteur de	données											$\mathbf{X}$
ic	id	weight	male	catholic	Belfast	N.Eastern	Southern	S.Eastern	Western	Grammar	funemp	gcse5eq	1^
le 1	1	0.33	no	no	no	no	no	no	yes	no	no	no	1
se 2	2	0.57	no	no	no	no	no	no	yes	no	no	yes	1
3	3	1.59	yes	yes	no	no	no	no	yes	no	no	no	1
ib 4	4	1.59	no	no	no	no	no	no	yes	no	no	no	1
p 5	5	0.57	yes	no	no	no	no	no	yes	no	yes	no	1
6	6	1.59	yes	yes	no	no	no	no	yes	no	no	no	1
7	7	0.57	yes	yes	no	no	no	no	yes	no	no	no	1
8	8	2.75	yes	yes	no	no	no	no	yes	yes	no	no	1
9	9	2	no	no	no	no	no	yes	no	no	no	no	1
AM 10	10	3.6	no	no	no	no	no	yes	no	no	no	no	1
11	11	0.69	yes	no	no	no	no	yes	no	no	no	no	1
12	12	1.1	no	no	no	no	no	no	yes	no	no	yes	1
13	13	1.1	yes	yes	no	no	no	no	yes	no	yes	yes	1
14	14	0.57	no	yes	no	no	no	no	yes	no	no	yes	1
15	15	2	no	yes	no	no	no	no	yes	yes	no	no	11
16	16	0.87	no	yes	no	yes	no	no	no	no	yes	no	1
17	17	4.1	no	yes	yes	no	no	no	no	no	no	no	1
18	18	0.33	yes	yes	yes	no	no	no	no	yes	no	no	1
19	19	0.71	yes	yes	yes	no	no	no	no	no	no	no	1
<	1											8	X

Introduction	Using R	Objects	Importing data into R	Exploring data frames	Working environment	References
Summa	ary					

id	weight	male	catholic	Belfast
Min. : 1.0	Min. :0.130	0 no :342	no :368	no :624
1st Qu.:178.8	1st Qu.:0.450	00 yes:370	yes:344	yes: 88
Median :356.8	Median :0.690	00		
Mean :356.5	Mean :0.999	94		
3rd Qu.:534.2	3rd Qu.:1.070	00		
Max. :712.0	Max. :4.460	00		

[1] "no" "yes"

#### Variable names

To obtain the list of all variables in the 'data frame'
 names(mvad)

E4.1	02.20	Harris and an Islam II	Um a 7 a U	0	IID-14+II	UN Destaurul
LTT	1 <b>d</b>	"weight"	"male"	"catholic"	"Bellast"	"N.Eastern"
[7]	"Southern"	"S.Eastern"	"Western"	"Grammar"	"funemp"	"gcse5eq"
[13]	"fmpr"	"livboth"	"Jul.93"	"Aug.93"	"Sep.93"	"Oct.93"
[19]	"Nov.93"	"Dec.93"	"Jan.94"	"Feb.94"	"Mar.94"	"Apr.94"
[25]	"May.94"	"Jun.94"	"Jul.94"	"Aug.94"	"Sep.94"	"Oct.94"
[31]	"Nov.94"	"Dec.94"	"Jan.95"	"Feb.95"	"Mar.95"	"Apr.95"
[37]	"May.95"	"Jun.95"	"Jul.95"	"Aug.95"	"Sep.95"	"Oct.95"
[43]	"Nov.95"	"Dec.95"	"Jan.96"	"Feb.96"	"Mar.96"	"Apr.96"
[49]	"May.96"	"Jun.96"	"Jul.96"	"Aug.96"	"Sep.96"	"Oct.96"
[55]	"Nov.96"	"Dec.96"	"Jan.97"	"Feb.97"	"Mar.97"	"Apr.97"
[61]	"May.97"	"Jun.97"	"Jul.97"	"Aug.97"	"Sep.97"	"Oct.97"
[67]	"Nov.97"	"Dec.97"	"Jan.98"	"Feb.98"	"Mar.98"	"Apr.98"
[73]	"May.98"	"Jun.98"	"Jul.98"	"Aug.98"	"Sep.98"	"Oct.98"
[79]	"Nov.98"	"Dec.98"	"Jan.99"	"Feb.99"	"Mar.99"	"Apr.99"
[85]	"May.99"	"Jun.99"				-

A description of the data set and variables is available with
 help(mvad)

27/44

Introduction Using R Objects Importing data into R Exploring data frames Working environment References
Frequency tables

 The states we will use to build sequences are in variables Jul.93...Jun.99, that is columns 15 to 86
 mvad[1:4, 15:20]

	Jul.93	Aug.93	Sep.93	Oct.93	Nov.93	Dec.93
1	training	training	employment	employment	employment	employment
2	joblessness	joblessness	FE	FE	FE	FE
3	joblessness	joblessness	training	training	training	training
4	training	training	training	training	training	training

- A frequency table of the gcse5eq variable (qualifications gained by the end of compulsory education)
  - > table(mvad\$gcse5eq)
  - no yes 452 260

26/44

UNIVERSITÉ DE GENÉVE UNIVERSITÉ DE GENÉVE

## Introduction Using R Objects Importing data into R **Exploring data frames** Working environment References Introduction Us

## Contingency tables

- We want a contingency table for variables funemp (father unemployed) and gcse5eq (qualification gained at the end of compulsory school)
- We change first the labels of the two factors (both are dummy variables whose labels are "yes/no" in the original file)

![](_page_10_Picture_4.jpeg)

Introduction	Using R	Objects	Importing data into R	Exploring data frames	Working environment	References
Contin	gency	table	es			

![](_page_10_Figure_6.jpeg)

Row and marginal distributions

	Lower	qual.	Higher	qual.
employed	0.8	008850	0.89	961538
unemployed	0.19	991150	0.10	38462

```
    Margins
```

```
> margin.table(t1, 1)
```

employed unemployed 595 117

> margin.table(t1, 2)

Lower qual. Higher qual. 452 260

31/44

![](_page_10_Picture_14.jpeg)

Histograms are for numerical variables like weight
 hist(mvad\$weight, col = "cyan")

![](_page_10_Figure_16.jpeg)

# We perform a chi-squared contingency table test using the chisq.test() function > chisq.test(t1)

Pearson's Chi-squared test with Yates' continuity correction

```
data: t1
X-squared = 10.2264, df = 1, p-value = 0.001384
```

30/44

UNIVERSITÉ DE GENÉVE

Introduction	Using R	Objects	Importing data into R	Exploring data frames	Working environment	References
Barplo	ts					

• Use the plot() method for factors (categorical variables) > plot(mvad\$gcse5eq, col = c("red", "green"), main = "Variable gcse5eq")

![](_page_11_Figure_2.jpeg)

Introduction	Using R	Objects	Importing data into R	Exploring data frames	Working environment	References
Boxplo	ots					

- The boxplot() function accepts a formula as argument to produce one boxplot for each category of a factor
  - > boxplot(iris\$Sepal.Length ~ iris\$Species, col = "cyan", main = "Sepal length, by specie

![](_page_11_Figure_6.jpeg)

![](_page_11_Figure_7.jpeg)

• Use the plot() method for two numerical variables > plot(iris\$Sepal.Length, iris\$Sepal.Width, col = "red")

![](_page_11_Figure_10.jpeg)

## Saving graphics

- To save graphics in files, depending on the format you can use the pdf(), jpeg() or png() function with the name of the file as argument
- Once you have issued all plotting commands you have to close the file with the dev.off() function
  - > pdf(file = "hist") > plot(mvad\$Sep.93, mvad\$Sep.94) > dev.off() pdf 2

UNIVERSITÉ DE GENÉVE

34/44

UNIVERSITÉ DE GENÉVE

![](_page_12_Figure_0.jpeg)

- library
- Print the variable names
- Add an age variable, by subtracting the birth year from the year of the survey
- What is the min., max., median and mean age in the sample?
- What is the min., max., median and mean age of the woman?
- Add a cohort factor to the biofam data frame grouping the birth years into the following categories: 1900-1929, 1930-1939, 1940-1949, 1950-1959.
- Produce a frequency table of the cohort factor

 When you quit R (either with the menu or with the quit() function, you are asked if you want to save your working environment. If you answer yes, it is saved in a '.RData' file.

" A "

"c"

"m1"

"v2"

"my.data"

"my.seq"

"titanic"

"My.very.first.R.object"

List objects in your R environment

> ls()

[1] "a"

[3] "Ъ"

[13] "t1"

[15] "v1"

[5] "iris"

[7] "mvad"

[9] "my.groups"

[11] "my.table"

## Introduction Using R Objects Importing data into R Exploring data frames Working environment

## The working environment - B

Introduction Using R Objects Importing data into R Exploring data frames Working environment References

#### References I

• You can also save your working environment with the save.image() function or specific objects with the save()
function.

> save.image(file = "myenv.RData")
> save(my.table, file = "data/mytable.RData")

- You get the current working directory with the getwd() function
- You can change the working directory with the setwd() function

- Carlos Alzola and Frank Harrell. An Introduction to S and The Hmisc and Design Libraries, 2006.
- Julian J. Faraway. Practical Regression and Anova using R, July 2002.
- Petra Kuhnert and Bill Venables. An Introduction to R: Software for Statistical Modelling & Computing. CSIRO Mathematical and Information Sciences, Cleveland, Australia, 2005. URL http://www.csiro.au/resources/Rcoursenotes.html.
- Emmanuel Paradis. *R pour les débutants*. Institut des Sciences de l' Evolution Université Montpellier II, F-34095 Montpellier, septembre 2005.

UNIVERSITÉ DE GENÉVE

W. N. Venables, D. M. Smith, and the R Development Core Team. An Introduction to R, 2009.

UNIVERSITÉ DE GENÉVE

References

44/44

#### Formats

#### Sequence analysis for social scientists Part III - Describing and visualizing sequence data sets

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référe

Alexis Gabadinho, Matthias Studer, Gilbert Ritschard, Nicolas S. Müller

Department of Econometrics, University of Geneva http://mephisto.unige.ch/biomining

Summer School on Advanced Methods for the Analysis of Complex Event History Data, Bristol, 28-29 June 2010

- State sequences can be represented in many different ways, depending on the data source and on how the information is organized
  - States observed at the successive time points (panel surveys)

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

- Spells, i.e. of distinct states stamped with their beginning/ending time (retrospective surveys)
- It is then essential for the user to identify the nature and organization of her/his data and to possibly transform it in an input form accepted by the software.
- Data organization and conversion between formats is discussed in details in ?

1/54

## The STS representation

 State sequence of length ℓ : ordered list of ℓ elements successively chosen from a finite set A (the *alphabet*) of size a = |A|.

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

- A natural way of representing a sequence x is x = (x<sub>1</sub>, x<sub>2</sub>, ..., x<sub>ℓ</sub>), with x<sub>j</sub> ∈ A. This is the STS representation.
- Such state sequences have two important properties
  - They are formed by elements that are states, i.e. something that can last as opposed for instance to events that occur at given time points.
  - The position of each element conveys meaningful information in terms of age, date or more generally elapsed time or distance from the beginning of the sequence.

#### Sequence representations

 More compact representations by giving only one of several same successive states (the Distinct Successive State (DSS) sequence) :

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence of

#### A-A-B-B-C-C-C-C-D => A-B-C-D

- To keep time and alignment information successive distinct states are stamped with their duration (State Permanence Sequence (SPS) representation) ...
  - A-B-C-D => (A,3)-(B,2)-(C,3)-(D,1)
- ... or their starting and ending positions (SPELL representation)

Start	End	State
1	3	А
4	5	В
6	9	С
10	10	D

DE GENÉVE

UNIVERSITÉ DE GENÉVE

#### Sequence representations - Example

- Two sequences describing family formation histories of two individuals
- Alphabet : S=single, M=married, MC=married with children, D=divorced

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

Format	Example	e										
	ld	a18	a19	a20	a2	21	a22	a23	a24	a25	a26	a27
STS	101	S	S	S	N	Л	М	MC	MC	MC	MC	D
	102	S	S	S	М	IC	MC	MC	MC	MC	MC	MC
	ld	s1	<i>s</i> 2	<i>s</i> 3	<i>s</i> 4							
DSS	101	S	M	MC	D							
	102	S	MC									
	ld	<i>s</i> 1	s	2	<i>s</i> 3		<i>s</i> 4					
SPS	101	(S,3)	(M	,2)	(MC	,4)	(D,1)					
	102	(S,3)	(MC	2,7)								
	ld	Index	Fro	m	To :	State						
	101	1	18		20 5	S						
	101	2	21	:	22 I	M						
SPELL	101	3	23		26 I	MC						
	101	4	27		27 I	D						
	102	1	18		20 5	S						
	102	2	21	:	27 I	MC						

6/54

State sequence objects

- Converting between formats
  - In this example the sequence s1 is in the State Permanence (SPS) representation and we get its STS representation with the seqformat() function

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

- > s1 <- "(A,2)-(B,10)-(C,3)-(D,2)-(A,3)" > s1
- [1] "(A,2)-(B,10)-(C,3)-(D,2)-(A,3)"
- > seqformat(s1, from = "SPS", to = "STS", compressed = TRUE)

#### Sequence

[1] "A-A-B-B-B-B-B-B-B-B-B-B-C-C-C-D-D-A-A-A"

7/54

The mvad data frame - Variables list

#### • We create a state sequence object with the *mvad* data set

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characterist

#### Table: List of Variables in the MVAD data set

id	unique individual identifier							
weight	sample weights							
male	binary dummy for gender, 1=male							
catholic	binary dummy for community, 1=Catholic							
Belfast binary dummies for location of school, one of five Education and Library Board a								
	Northern Ireland							
N.Eastern	n							
Southern	и 1							
S.Eastern	"							
Western	"							
Grammar	binary dummy indicating type of secondary education, 1=grammar school							
funemp	binary dummy indicating father's employment status at time of survey, 1=father unemployed							
gcse5eq	binary dummy indicating qualifications gained by the end of compulsory education, $1=5+$ GCSEs at grades A-C, or equivalent							
fmpr	binary dummy indicating SOC code of father's current or most recent job,1=SOC1 (profes- sional, managerial or related)							
livboth	binary dummy indicating living arrangements at time of first sweep of survey (June 1995), $1=$ living with both parents							
jul93	Monthly Activity Variables are coded 1-6, 1=school, 2=FE, 3=employment, 4=training, 5=joblessness, 6=HE							
:	"							
jun99	"							

- State sequence objects are required as input by the functions for state sequence analysis
- These objects store the state sequences (as a matrix) together with important attributes

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen

- The alphabet (the list of possible states)
- The (long) labels associated with each state of the alphabet that will be used in the graphics
- The colors for state legends
- Information about missing states in the sequences
- Possibly case weights

UNIVERSITÉ DE GENÉVE

#### Locating the sequence data

• We locate the sequence data in the data frame

> names(mvad)								
[1]	"id"	"weight"	"male"	"catholic"	"Belfast"			
[6]	"N.Eastern"	"Southern"	"S.Eastern"	"Western"	"Grammar"			
[11]	"funemp"	"gcse5eq"	"fmpr"	"livboth"	"Jul.93"			
[16]	"Aug.93"	"Sep.93"	"Oct.93"	"Nov.93"	"Dec.93"			
[21]	"Jan.94"	"Feb.94"	"Mar.94"	"Apr.94"	"May.94"			
[26]	"Jun.94"	"Jul.94"	"Aug.94"	"Sep.94"	"Oct.94"			
[31]	"Nov.94"	"Dec.94"	"Jan.95"	"Feb.95"	"Mar.95"			
[36]	"Apr.95"	"May.95"	"Jun.95"	"Jul.95"	"Aug.95"			
[41]	"Sep.95"	"Oct.95"	"Nov.95"	"Dec.95"	"Jan.96"			
[46]	"Feb.96"	"Mar.96"	"Apr.96"	"May.96"	"Jun.96"			
[51]	"Jul.96"	"Aug.96"	"Sep.96"	"Oct.96"	"Nov.96"			
[56]	"Dec.96"	"Jan.97"	"Feb.97"	"Mar.97"	"Apr.97"			
[61]	"May.97"	"Jun.97"	"Jul.97"	"Aug.97"	"Sep.97"			
[66]	"Oct.97"	"Nov.97"	"Dec.97"	"Jan.98"	"Feb.98"			
[71]	"Mar.98"	"Apr.98"	"May.98"	"Jun.98"	"Jul.98"			
[76]	"Aug.98"	"Sep.98"	"Oct.98"	"Nov.98"	"Dec.98"			
[81]	"Jan.99"	"Feb.99"	"Mar.99"	"Apr.99"	"May.99"			
[86]	"Jun.99"	"region"	"sex"	"religion"				

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

11/54

Creating the state sequence object

- The list of the states found in the data is returned by the seqstatl() function
  - > seqstatl(mvad[, 15:86])

[1] "employment" "FE" "HE" "joblessness" "school"
[6] "training"

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen

- By default, these labels are used for both the state names and their long labels
- To override these default settings we create a vector of state names ...

```
> mvad.shortlab <- c("EM", "FE", "HE", "JL", "SC", "TR")
```

• ... and a vector for the state labels

> mvad.lab <- c("Employment", "Further education", "Higher education",</pre>

+ "Joblessness", "School", "Training")

Overview of the sequences

 The monthly statuses are in variables Jul.93 to Jun.99, that is columns 15 to 86
 mvad[1:4, 15:18]

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

Jul.93 Aug.93 Sep.93 Oct.93 1 training training employment employment 2 joblessness joblessness FE FE 3 joblessness joblessness training training 4 training training training training

- In the mvad data frame, the sequence data are in the STS format, which is the default for the seqdef() function
- If your sequence data is not in the STS format, use the informat argument (or the seqformat() function beforehand to format your data)

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence character

12/54

DE GENÉVE

UNIVERSITÉ DE GENÉVE .

#### Sequence object

• Now we create the mvad.seq sequence object using the seqdef() function

> mvad.seq <- seqdef(mvad, 17:86, states = mvad.shortlab, labels = mvad.lab)</pre>

- We can display the sequences in the 'STS' format (default) ...
  - > mvad.seq[1:4, 1:20]

Sequence

- .. or in the shorter 'SPS' format
  - > print(mvad.seq[1:4, ], format = "SPS")
    - Sequence
  - [1] (EM,4)-(TR,2)-(EM,64)
  - [2] (FE,36)-(HE,34)
  - [3] (TR,24)-(FE,34)-(EM,10)-(JL,2)
  - [4] (TR,47)-(EM,14)-(JL,9)

Sequence representations State sequence objects	000000	000000000	000000000		000000	000000000	000000000
Summary				Retrieving and setting	attributes		
<ul> <li>The summary() function state sequence object</li> <li>&gt; summary(mvad.seq)</li> <li>[&gt;] sequence object</li> <li>[&gt;] 712 sequences in</li> <li>[&gt;] min/max sequence</li> <li>[&gt;] alphabet (state</li> <li>1=EM (Employmention)</li> <li>2=FE (Further edited)</li> <li>4=LL (Joblessnettion)</li> <li>6=TR (Training)</li> <li>[&gt;] dimensionality of [&gt;] colors: 1=#7FC9</li> </ul>	ction returns u ct created with Tr n the data set, e length: 70/70 labels): t) ducation) ucation) ss) of the sequence 7F 2=#BEAED4 3=#	sefull information aMineR version 1.6 490 unique space: 350 FDC086 4=#FFFF99 5=	about a #386CB0 6=#F0027F	<ul> <li>One can also retrie</li> <li>alphabet(mvad.seq)</li> <li>[1] "EM" "FE" "HE"</li> <li>stlab(mvad.seq)</li> <li>[1] "Employment"</li> <li>[4] "Joblessness"</li> <li> or set some attr</li> <li>stlab(mvad.seq) &lt;-</li></ul>	<pre>Ve ) "JL" "SC" "TR"     "Further educ     "School" ibutes with dedie - c("(1)=EMPLOYED" EDUC.", "(4)=JOBL NG") c("green", "blue" ey")</pre>	cation" "Higher educat "Training" cated functions ', "(2)=FURTHER EDUC.", ESSNESS", "(5)=SCHOOL' ', "red", "yellow",	ion" ,
			UNIVERSITÉ DE GENÉVE				UNIVER: DE GENI

DICC

a second a discontraction of the last second s

Subscripts and attribute inheritance

```
• Subsets of sequence objects with indexes like matrices or data
```

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référence

#### frames

```
> seq1 <- mvad.seq[1, ]</pre>
```

```
> print(seq1, format = "SPS")
    Sequence
```

```
[1] (EM,4)-(TR,2)-(EM,64)
```

```
• Subsets inherits attributes
```

```
> alphabet(seq1)
```

```
[1] "EM" "FE" "HE" "JL" "SC" "TR"
```

> stlab(seq1)

[1]	"Employment"	"Further	education"	"Higher	education"
[4]	"Joblessness"	"School"		"Trainin	ng"

#### Missing values - A

- Missing values in STS sequences occur for example when
  - Sequences do not start at a same date and a calendar time axis is used :

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

<mark>ce c</mark>haracteristic

DE GENÉVE

- The follow up time is shorter for some individuals than for others yielding sequences that do not end up at the same position;
- The observation at some positions is missing due to non-response, yielding internal gaps in the sequences.
- We illustrate this with the ex1 data frame
  - > data(ex1)
  - > ex1[, 1:13]

	[P1]	[P2]	[P3]	[P4]	[P5]	[P6]	[P7]	[P8]	[P9]	[P10]	[P11]	[P12]	[P13]
s1	<na></na>	<na></na>	<na></na>	Α	Α	Α	A	A	A	Α	Α	Α	Α
s2	D	D	D	В	В	В	В	В	В	В	<na></na>	<na></na>	<na></na>
s3	<na></na>	D	D	D	D	D	D	D	D	D	D	<na></na>	<na></na>
s4	A	A	<na></na>	<na></na>	В	В	В	В	D	D	<na></na>	<na></na>	<na></na>
s5	A	<na></na>	Α	Α	Α	Α	<na></na>	A	A	Α	<na></na>	<na></na>	<na></na>
s6	<na></na>	<na></na>	<na></na>	С	С	C	C	C	C	С	<na></na>	<na></na>	< <u>NA2</u>

16/54

#### Types of missing values

- Three types of missing values depending on where they appear in the sequence :
  - left-missing-values, that is missing values appearing before the first valid entry in the sequence;

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

- in-between-missing-values or gaps, that is missing values appearing somewhere between the first and last valid entries;
- right-missing values, that is missing values appearing after the last valid entry.
- The seqdef() function provides special options for dealing with each type of missing values (left, gaps and right) with possible value
  - NA : each missing value is left as an explicit missing element
  - DEL : missing elements are deleted
  - a state, belonging to the alphabet or not

#### Missing values - C

• Default values are left=NA, gaps=NA and right="DEL" : the sequence is considered as ending after the last valid state (all missing values encountered after the last (rightmost) valid state are deleted)

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

• In the sequence object remaining missing values are coded with '\*'

```
> seqdef(ex1, 1:13)
```

```
Sequence
s1 *-*-*-A-A-A-A-A-A-A-A-A-A-A
s2 D-D-D-B-B-B-B-B-B-B
s3 *-D-D-D-D-D-D-D-D-D-D-D
s4 A-A-*-*-B-B-B-B-D-D
s5 A-*-A-A-A-A-A-A-A-A
s6 *-*-*-C-C-C-C-C-C-C
```

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

Exercise - State sequence objects

#### Exercise 3.1

Load the biofam data set that comes with the TraMineR library (look at the online help to get more information)

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence character

2 Create a state sequence object named biofam.seq with variables a15 to a30, using the following state names and

labels									
State	Name	Label							
0	Р	Parent							
1	L	Left							
2	M	Married							
3	LM	Left/Married							
4	C	Child							
5	LC	Left/Child							
6	LMC	Left/Married/Child							
7	D	Divorced							

Print the first sequences in biofam.seq, in the STS format and then in the SPS format

## Missing values - D

- By changing to left="DEL", sequences are left aligned > seqdef(ex1, 1:13, left = "DEL")
  - Sequence s1 A-A-A-A-A-A-A-A-A-A s3 D-D-D-D-D-D-D-D-D-D-D
  - s4 A-A-\*-\*-B-B-B-B-D-D s5 A-\*-A-A-A-A-\*-A-A-A
  - s6 C-C-C-C-C-C-C

DE GENÉVE

#### Sequence index plots (A) Sequence index plots (B) • Sequence index plots [??] allow to visualize individual sequences as horizontally stacked boxes coloured according to • When the number of displayed sequences is large these plots the state at the successive positions. are often hard to interpret. • The segiplot() function produces such plots (by default the • For better results we can sort the sequences according to the first 10 sequences, use tlim to override) > seqiplot(mvad.seq, border = NA) values of a covariate. Good choices are for instance the distance to the most frequent sequence or the scores of a multidimensional scaling analysis. • We plot all the sequences with seqlplot() grouped by qualifications gained of compulsory education (gcse5eq covariate) and sorted by distance to the most frequent sequence. > dist.mostfreq <- seqdist(mvad.seq, method = "LCS", refseq = 0)</pre> > seqIplot(mvad.seq, border = NA, group = mvad\$gcse5eq, sortv = dist.mostfreq) \_\_\_\_\_ DE GENÉVE UNIVERSITÉ DE GENÉVE Higher education Schoo Inblessness Trainir

## Sequence index plots by group

Sequence representations State sequence objects

![](_page_19_Figure_2.jpeg)

000000

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen

Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen-

000000

## Frequency table

• A first simple overview of a set of state sequences is the sequence frequency table

000000

• It is returned by the seqtab() function (by default the 10 most frequent sequences are displayed, use tlim to override)

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characterist

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

000000

> seqtab(mvad.seq)

	Freq	Percent	
EM/70	50	7.02	
TR/22-EM/48	18	2.53	
FE/22-EM/48	17	2.39	
SC/24-HE/46	16	2.25	
SC/25-HE/45	13	1.83	
FE/25-HE/45	8	1.12	
FE/34-EM/36	7	0.98	
FE/46-EM/24	7	0.98	
FE/10-EM/60	6	0.84	
FE/24-HE/46	6	0.84	

#### • A graphical view of the sequence frequency table where bar widths are proportional to the frequencies is obtained with the segfplot() function. > seqfplot(mvad.seq, group = mvad\$gcse5eq, border = NA) Exercise 3.2 Moins bons résultats Bons résultats Create a full sequence index plot grouped by values of the cohort variable created in exercise 2.1 and sorted by sex 2 Create a sequence frequency plot grouped by values of the cohort variable and save it as a 'jpeg' file freq. ...... Sep.93 Sep.94 Sep.95 Sep.96 Sep.97 Sep.98 Sep.93 Sep.94 Sep.95 Sep.96 Sep.97 Employment Higher education School Further education Joblessness Training UNIVERSITÉ DE GENÉVE DE GENÉVE Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence character

#### Mean time spent in each state

 A first synthetic information is given by the mean -non necessarily consecutive- time spent in the different states
 > seqmtplot(mvad.seq, group = mvad\$gcse5eq)

•0000000C

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

000000

Visualizing sequence frequency table

![](_page_20_Figure_3.jpeg)

## Transition rates

• The transition rate between each couple of states  $(s_i, s_j)$  is the probability to switch to state  $s_i$  when we are in state  $s_i$ .

00000000

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

00000

Exercise - Visualizing sequence data sets

• The seqtrate() function returns the transition rates matrix.

<pre>&gt; mvad.trate &lt;- seqtrate(mvad</pre>	1.se
--	------

> round(mvad.trate, 2)

	[-> EM]	[-> FE]	[-> HE]	[-> JL]	[-> SC]	[-> TR]
[EM ->]	0.99	0.00	0.00	0.01	0.00	0.00
[FE ->]	0.03	0.95	0.01	0.01	0.00	0.00
[HE ->]	0.01	0.00	0.99	0.00	0.00	0.00
[JL ->]	0.04	0.01	0.00	0.94	0.00	0.01
[SC ->]	0.01	0.01	0.02	0.01	0.95	0.00
[TR ->]	0.04	0.00	0.00	0.01	0.00	0.94

• Transition rates provide information about the most frequent state changes observed in the data together with an assessment of the stability of each state.

#### Transversal statistics Sequence of state distributions • The seqstatd() function returns the series of state distributions together with other information > seqstatd(mvad.seq[, 1:8]) • We can summarize a set of sequences with a series (sequence) [State frequencies] of transversal indicators Sep.93 Oct.93 Nov.93 Dec.93 Jan.94 Feb.94 Mar.94 Apr.94 EM 0.117 0.124 0.133 0.138 0.140 0.140 0.149 0.157 id t<sub>2</sub> tз . . . $t_1$ 0.388 0.382 0.381 0.369 0.364 0.361 FE 0.386 0.353 В В 1 D . . . HE 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 2 А B C . . . JL0.024 0.021 0.020 0.021 0.028 0.038 0.034 0.035 SC 0.251 0.246 0.244 0.242 0.240 0.242 0.240 0.240 3 В B Α . . . TR. 0.222 0.222 0.221 0.219 0.222 0.216 0.216 0.215 • These indicators are based on the state distribution at each [Valid states] position/time point in the sequences Sep.93 Oct.93 Nov.93 Dec.93 Jan.94 Feb.94 Mar.94 Apr.94 712 712 Ν 712 712 712 712 712 712 [Entropy index] Sep.93 Oct.93 Nov.93 Dec.93 Jan.94 Feb.94 Mar.94 Apr.94 0.81 Н 0.77 0.77 0.78 0.78 0.8 0.8 0.8 UNIVERSITÉ DE GENÉVE 34/54

State distribution plots

• The series of transversal state distributions is plotted with the segdplot() function

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen

000000000

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen

000000000

• We plot the series for the two groups defined by the values of the gcse5eq covariate

> seqdplot(mvad.seq, group = mvad\$gcse5eq, border = NA)

![](_page_21_Figure_6.jpeg)

Employment Higher education School
 Further education Joblessness
 Training

DE GENÉVE

Entropy index

• The Shannon entropy of the state distribution at each position is

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characterist

000000000

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

000000000

$$h(p_1,\ldots,p_s)=-\sum_{i=1}^s p_i \log_2(p_i)$$

where  $p_i$  is the frequency of the *i*th state and *s* is the size of the alphabet

- This indicator is called the *entropy index* Billari [2001]
  - It equals 0 when all cases are in the same state (it is thus easy to predict in which state an individual is)
  - It is maximum when the cases are equally distributed between the states of the alphabet (it is thus hard to predict in which state an individual is)

## Plotting the series of transversal entropies

• The series of of transversal entropies can be plotted with the seqHtplot() function

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référe

0000000000

> seqHtplot(mvad.seq, group = mvad\$gcse5eq)

![](_page_22_Figure_3.jpeg)

## Sequence of modal states

• By taking the most frequent (modal) state at each position we get the modal state sequence

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

000000000

State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

- The modal state sequence is returned by the seqmodst() function and can be plotted with the seqmsplot() function
  - > seqmsplot(mvad.seq, group = mvad\$gcse5eq, border = NA)

![](_page_22_Figure_8.jpeg)

37/54

# Exercise - Overall and transversal statistics

#### Exercise 3.3

- $\textcircled{\sc 0}$  Compute the transition rate matrix for the biofam data set
- What is the transition rate between states "Left/Maried" and "Left/Maried/Child"?

esentations State sequence objects Visualizing sequence sets **Overall and transversal statistics** Individual sequence characteristics

- Oisplay the mean times spent in each of the states for each cohort
- Display the sequence of transversal state distributions for each cohort
- O Display the sequence of modal states for each cohort
- Within each cohort, at what age is the diversity of the transversal state distribution at its highest?

## Longitudinal features of individual sequences

• We focus now on the characterization and summarization of longitudinal features of individual sequences

id	$t_1$	$t_2$	t <sub>3</sub>	•••
1	В	В	D	• • •
2	Α	В	С	•••
3	В	В	А	•••

- The aim is to define measures that inform on how each sequence is constituted, i.e. on whether it takes a simple or more complex form.
- TraMineR allows to calculate several indicators of sequence complexity
  - "One-dimensional" indicators : number of transitions, number of subsequences, longitudinal entropy
  - Composite complexity measures : Turbulence, Complexity index

DE GENÉVE

38/54

#### Distinct states and durations

• In SPS form a state sequence is represented as an ordered list of distinct successive states (DSS) with their associated durations, i.e. as a sequence of couples  $(x_j, t_j)$  where  $x_j$  is a state and  $t_j$  its duration

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

> print(mvad.seq[1:3, ], format = "SPS")

#### Sequence

[1] (EM,4)-(TR,2)-(EM,64)
[2] (FE,36)-(HE,34)

```
[3] (TR,24)-(FE,34)-(EM,10)-(JL,2)
```

```
> seqdss(mvad.seq[1:3, ])
```

```
Sequence
```

```
1 EM-TR-EM
2 FE-HE
```

```
3 TR-FE-EM-JL
```

 This suggests to distinguish characteristics of the state sequencing—the distinct successive states (DSS)—from those of the durations.

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen

42/54

## Number of transitions

- A transition in a state sequence is defined as a change of status between positions (time) t and t + 1.
- The seqtransn() function returns the number of transitions contained in each sequence of a state sequence object.
- The number of transitions in a sequence x is  $\ell_d(x) 1$  where  $\ell_d(x)$  is the length of its DSS sequence

![](_page_23_Figure_16.jpeg)

![](_page_23_Figure_17.jpeg)

![](_page_23_Figure_18.jpeg)

## Distinct states and durations

- Complexity measures take into account two main characteristics
  - Distinct state appearing in the sequence : number of transitions, number of distinct subsequences
  - Durations in each of the state : variance of the state durations, entropy of the state distribution

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

• In the two examples below, sequence 2 is more complex than sequence 1

![](_page_23_Figure_24.jpeg)

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics

## Longitudinal entropy (A)

- Total time spent in each state, that is, when there are multiple spells in a same state the sum of the lengths of these spells.
- In sequence 1 of the *mvad.seq* sequence object, (EM,4)-(TR,2)-(EM,64), there are two spells in state EM with respective durations 4 and 64. Hence, the time spent in EM is 68 months.
- The seqistatd() function returns for each sequence the total time spent in each state of the alphabet (the longitudinal state distribution).

> seqistatd(mvad.seq[1:4, ])

 EM
 FE
 HE
 JL
 SC
 TR

 1
 68
 0
 0
 0
 2

 2
 0
 36
 34
 0
 0
 0

 3
 10
 34
 0
 2
 0
 24

 4
 14
 0
 0
 9
 0
 47

DE GENÉVE

Individual sequence characteristic

•00000000

000000000

#### Longitudinal entropy (B)

• The entropy of the state distribution within a sequence is a measure of the diversity of its states.

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

00000000

• The seqient() function returns the vector of the longitudinal Shannon entropies of the sequences, i.e. for each sequence the entropy

$$h(\pi_1,\ldots,\pi_a)=-\sum_{i=1}^a\pi_i\log\pi_i$$

where *a* is the size of the alphabet and  $\pi_i$  the proportion of occurrences of the *i*th state in the considered sequence.

## Longitudinal entropy (C)

- The entropy can be interpreted as the 'uncertainty' in predicting the states in a given sequence :
  - when the state remains the same during the whole sequence, the entropy equals 0

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

• maximum entropy is reached when the time spent in each element of the alphabet is the same.

![](_page_24_Figure_10.jpeg)

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

46/54

Turbulence

- The *Turbulence T*(*x*) of a sequence *x* is a composite measure proposed by Elzinga [Elzinga and Liefbroer, 2007] that accounts for
  - **(**) the number  $\phi(x)$  of distinct subsequences in the DSS sequence

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référen

3 the variance  $s_t^2(x)$  of the consecutive times  $t_j$  spent in the distinct states

• The formula is

$$T(x) = \log_2(\phi(x) \frac{s_{t,max}^2(x) + 1}{s_t^2(x) + 1})$$

where

- s<sup>2</sup><sub>t</sub>(x) is the variance of the state-durations t<sub>j</sub> j = 1,..., l<sub>d</sub>(x) for sequence x
- $s_{t,max}^2(x)$  is the maximum value that this variance can take given the total duration  $\ell(x) = \sum_i t_j$  of that sequence.

UNIVERSITÉ DE GENÉVE

DE GENÉVE

000000000

Turbulence - B

• The turbulence is returned by the seqST() function

- > mvad.turb <- seqST(mvad.seq)</pre>
- > hist(mvad.turb)
- This measure is not standardized
  - > min(mvad.turb)
  - [1] 1
  - > max(mvad.turb)
  - [1] 12.95858

#### Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristic

#### 000000000

#### Complexity index

• The Complexity Index is another composite measure proposed by Gabadinho et al. [2010] that combines

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

- 1 the number of transitions in the sequence
- 2 the longitudinal entropy.
- The Complexity index of a sequence x is defined as

$$C(x) = \sqrt{\frac{\ell_d(x)}{\ell(x)} \frac{h(x)}{h_{max}}}$$

#### where

- $h_{max}$  is the theoretical maximum of the entropy, that is  $h_{max} = \log a$ .
- $\ell_d(x)$  is the length of the Distinct Successive State sequence

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Référer

- Minimum value of C(x) is 0 and can only be reached by a sequence made of a single distinct state, which contains thus 0 transitions and has an entropy of 0.
- Maximum, 1, of C(x) is reached when the two following conditions are fulfilled :
  - Seach of the state in the alphabet is present in the sequence and the same time ℓ/a is spent in each of them;
  - 2 The number of transitions in the sequence is equal to  $\ell 1$ , that is, the length  $\ell_d$  of the DSS is equal to the length of the sequence  $\ell$ .

50/54

#### Complexity measures - Summary

• We compare complexity measures for sequences extracted at regular interval in the distribution of *C*(*s*)

![](_page_25_Figure_18.jpeg)

![](_page_25_Figure_19.jpeg)

000000000

![](_page_25_Figure_20.jpeg)

Trans. H(x) C(x) T(x) 0.0 0.2 0.4 0.6 0.8 1.0 Exercise - Individual sequence characteristics

Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics

#### Exercise 3.4

Complexity index - Values

- Compute the complexity index and turbulence of the sequences in the biofam data set
- ② Display a XY plot of the two indicators
- Oisplay boxplots of turbulence by cohort. Do the same for the complexity index.
- What can you say about the trend in the complexity of family life histories?

Individual sequence character

00000000

#### Sequence representations State sequence objects Visualizing sequence sets Overall and transversal statistics Individual sequence characteristics Références

UNIVERSITÉ DE GENÉVE

#### References I

Francesco C. Billari. The analysis of early life courses : complex descriptions of the transition to adulthood. *Journal of Population Research*, 18(2) : 119(24)–, November 2001. ISSN 1443-2447.

Cees Elzinga and Aart Liefbroer. De-standardization of family-life trajectories of young adults : A cross-national comparison using sequence analysis. *European Journal of Population/Revue européenne de Démographie*, 23(3) : 225–250, October 2007. URL http://dx.doi.org/10.1007/s10680-007-9133-7.

Alexis Gabadinho, Gilbert Ritschard, Matthias Studer, and Nicolas S. Müller. Indice de complexité pour le tri et la comparaison de séquences catégorielles. *Revue des nouvelles technologies de l'information RNTI*, E-19 :61–66, 2010.

![](_page_27_Figure_0.jpeg)

1/64				3/64					
Reminder ○●	Dissimilarity Centrotype	Clustering MDS	Discrepancy analysis References	Reminder 00	Dissimilarity ••••••••••••••••••••••••••••••••••••	Centrotype 000	Clustering 00000000000	MDS 0 0000	Discrepancy analysis
Creatir	ng the state sequ	ence obiect		Dissim	nilarities bet	ween pa	airs of sec	nuenc	es

• Loading TraMineR and the mvad data set

#### R> library(TraMineR)

- R> data(mvad)
- We consider only 70 states from Sept 93, i.e. we skip July and August 93 (i.e. summer holiday).

- R> mvad.shortlab <- c("EM", "FE", "HE", "JL", "SC", "TR")
- R> mvad.seq <- seqdef(mvad, 17:86, states = mvad.shortlab, labels = mvad.lab)

- Distance between sequences
  - Different metrics (LCP, LCS, OM, HAM, DHD, ...)
- With pairwise dissimilarity matrix, we can
  - determine a central sequence (centro-type)
  - measure the discrepancy between sequences
  - build a typology of the sequences using Cluster Analysis
  - represent the sequences using a scatter plot using MDS
  - test and measure the link between one or several covariates and the sequences using discrepancy analysis.

References

## Dissimilarity: Concepts

Dissimilarity

• A dissimilarity is a quantification of how far two objects are.

MDS

Discrepancy analysis

References

References

- For instance, consider two incomes x and y:
  - d(x, y) = |x y|

• 
$$d(x, y) = \log(1 + |x - y|)$$

• 
$$d(x, y) = (x - y)^2$$

- How to do it with categorical sequences?
- Depending on the issue, we want our dissimilarity measure to account for:
  - Order of the states and transitions in each sequence.

Clustering

MDS

Discrepancy analysis

- Temporality of the transitions.
- Duration of stay in each state.

## Optimal matching (optimal alignment)

Centrotype

#### Optimal matching

Dissimilarity

- Based on Levenshtein [1966]'s distance
- Inspired from alignment used in biology (ADN or protein sequences)

Clustering

MDS

Discrepancy analysis

Discrepancy analysis

• Introduced in social sciences by Abbott and Forrest [1986]

Clustering

• Also known as Edit distance.

UNIVERSITÉ De GENEVE		URIVESTIC DEGENERATION
	8/64	

# Optimal matching (OM): principle

Dissimilarity

•	OM distance is the minimal cost needed to transform one	ć
	sequence into the other one.	

- Want to transform one sequence into the other one.
- Using two types of operations
  - Insertion or deletion of an element
  - Substitution of an element
- Each operation has a cost.

indel and substitution costs

Dissimilarity

indel and substitution costs

- indel (insertion-deletion) costs:
  - Same cost for each 'insertion' or 'deletion'.
  - indel cost is a single constant.
- Substitution costs:
  - Each substitution may receive a different cost.
  - Matrix of substitution costs.
  - However: symmetrical cost  $c_{i,j} = c_{j,i}$

![](_page_28_Picture_35.jpeg)

## OM : insertion deletion example.

Consider the two sequences :

Dissimilarity

1	SC	SC	SC	ΕM	ΕM	ΕM	JL
2	SC	SC	SC	ΕM	ΕM	JL	JL

Clustering

MDS

Discrepancy analysis

Using indel cost of 1 and constant substitution cost of 2. Insertion of element 'EM', cost 1.

1	SC	SC	SC	ΕM	ΕM	ΕM	JL	
2	SC	SC	SC	ΕM	ΕM	EM	JL	JL

Deletion of element 'JL', cost 1.

1	SC	SC	SC	ΕM	ΕM	ΕM	JL
2	SC	SC	SC	ΕM	ΕM	ΕM	JL

The two sequences are now identical, total cost: 2.

- Unique cost  $c_{ii} = c$  (user defined).
  - Using c = 2 and *indel* = 1 OM is equivalent to "LCS" metric.
  - Let LLCS be the length of the longest common subsequence and  $|\boldsymbol{x}|$  the length of the sequence  $\boldsymbol{x}$
  - $d_{LCS}(x, y) = |x| + |y| 2 \cdot LLCS(x, y)$
- Based on transition rates (no additional input required)
  - $c_{i,j} = c_{j,i} = 2 p(i_t \mid j_{t-1}) p(j_t \mid i_{t-1})$
- Custom costs

# OM: substitution example

Dissimilarity

	1	SC	SC	SC	ΕM	ΕM	ΕM	JL	
	2	SC	SC	SC	ΕM	ΕM	JL	JL	
Substit	ution of	'JL' Ł	oy eler	nent	'EM',	cost 2			
	1	SC	SC	SC	EM	EM	EM	JL	
	2	SC	SC	SC	ΕM	ΕM	EM	JL	
10/64									

MDS

Discrepancy analysis

## OM: complete exemple

Consider the distance between these two sequences:

I N D U S T R Y I N T E R E S T

With a constant substitution cost of 2 and indel of 1, the optimal matching distance is 8 with the following operations.

1	Ι	Ν	D	-	-	U	S	Т	R	Υ
2	Ι	Ν	Т	Е	R	Е	S	Т	-	-
cost	0	0	2	3	4	6	6	6	7	8

The longest common subsequence is INST and its length is 4.

UNIVERSITÉ DE GENÉVE

## Using Optimal Matching in TraMineR

Dissimilarity

- Create the state sequence object with seqdef()
- Compute matrix of OM distances with seqdist(..., method="OM", indel=..., sm=...)

Clustering

MDS

Discrepancy analysis

References

- Specify the substitution costs using
  - sm="CONSTANT" for constant substitution costs.
  - sm="TRATE" for substitution costs based on transition rates.
  - or a custom substitution cost matrix using the same sm argument.

15/64					OUTWEATE Description		
Reminder 00	Dissimilarity Centrotyp 000000000000000000000000000000000000	e Clustering 000000000	MDS	Discrepancy analysis	References 000000		
Cost Matrix: Custom Costs							

Computing the distances

Dissimilarity

• Using the constant substitution cost, we compute distances

Clustering

R> mvad.dist <- seqdist(mvad.seq, method = "OM", indel = 1, sm = "CONSTANT")</pre>

MDS

Discrepancy analysis

#### R> mvad.dist[1:10, 1:10]

[,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [1,] 0 140 [2.] 140 [3,] [4.] [5,] [6.] [7,] [8,] [9,] [10.] 

-1	~	

![](_page_30_Picture_14.jpeg)

R> subm.custom <- matrix(c(0, 1, 1, 2, 1, 1, 1, 0, 1, 2,

- + 1, 2, 1, 1, 0, 3, 1, 2, 2, 2, 3, 0, 3, 1, 1, 1, 1,
- + 3, 0, 2, 1, 2, 2, 1, 2, 0), nrow = 6, ncol = 6, byrow = TRUE,
- + dimnames = list(mvad.shortlab, mvad.shortlab))

R> subm.custom

- R> mvad.dist.custom <- seqdist(mvad.seq, method = "OM",</pre>
- + indel = 1.5, sm = subm.custom)

Other dissimilarity measures provided by TraMineR:

- Longest Common Prefix (LCP)
- Longest Common Suffix (RLCP)
- Longest Common Subsequence (LCS)
- Hamming distance (HAM)
- Dynamic Hamming Distance (DHD) [Lesnard, 2010]

#### 

• Sequence with minimal sum of distances to other sequences.

![](_page_31_Figure_2.jpeg)

## Plot of representative sequences

![](_page_31_Figure_4.jpeg)

employment indication indica

# Set of representative sequences

Centrotype

• Aim: Find a minimal set of representative sequences such that

Clustering

• they are non redundant (none in the neighborhood of another representative)

MDS

Discrepancy analysis

References

UNIVERSITÉ DE GENÉVE

- Coverage: minimum percentage of all sequences that are in the neighborhood of the representative sequences
- Representatives are found with (tsim: neighborhood diameter, trep: coverage)

R> seqrep(mvad.seq, dist.matrix = mvad.dist, criterion = "density",
+ trep = 0.4, tsim = 0.1)

#### Representatives are plotted with

R> seqrplot(mvad.seq, dist.matrix = mvad.dist, group = mvad\$gcse5eq, + cex.plot = 1, trep = 0.4)

	-	-	
/GA			
04			

Reminder	Dissimilarity	Centrotype	Clustering	MDS	Discrepancy analysis	References
00	000000000000000	000	•00000000000	0000		000
Cluster	analysis					

- Cluster analysis automatically classify different objects in a reduced number of categories.
- It simplifies the large number of distinct sequences in a few different types of trajectories.
- It is used to build a typology of the trajectories.
- It offers a descriptive approach to analyze the sequences.

training

![](_page_31_Picture_22.jpeg)

#### 

- Clustering may be done using a dissimilarity matrix.
- There are several possibilities in R, for instance with the cluster library
  - agnes(): agglomerative nesting, i.e. hierarchical clustering (average, ward, ...).
  - diana(): divisive analysis.
  - pam(): partitioning around medoids.

![](_page_32_Figure_6.jpeg)

UNIVERSITÉ DE GENÉVE

References

UNIVERSITÉ DE GENÉVE

![](_page_32_Figure_7.jpeg)

Clustering

MDS

Discrepancy analysis

## Retrieving cluster membership

- Select the number of clusters,
- Cut the tree at chosen level, and store cluster membership into a vector.

```
R> mvad.cl4 <- cutree(mvad.clusterward, k = 4)
R> mvad.cl4[1:10]
```

```
[1] 1 2 3 3 2 3 1 1 3 4
```

R> clust.labels <- c("Employment", "Higher Education", "Training", + "Joblessness") R> mvad.cl4.factor <- factor(mvad.cl4, levels = 1:4, labels = clust.labels)</pre>

## Hierarchical clustering (Ward)

• Ward is a hierarchical clustering algorithm.

Centrotype

• At each step, it joins together the two less distant groups.

Clustering

MDS

Discrepancy analysis

• Ward aims at minimizing the within cluster discrepancy.

# R> library(cluster) R> myad,clusterward <- agnes(mvad.dist, diss = T, method = "ward") R> plot(mvad.clusterward, ask = F, which.plots = 2) Dendrogram of agnes(x = mvad.dist, diss = T, method = "ward")

![](_page_32_Figure_19.jpeg)

mvad.dist Agglomerative Coefficient = 0.99

MDS

Discrepancy analysis

References

DE GENÊVE

# PAM clustering

- Partitioning Around Medoids.
- Non hierarchical, number of cluster must be set a priori.

Clustering

• Faster, but results may depend on the starting point.

R> mvad.pam4 <- pam(mvad.dist, k = 4, diss = T)

R> plot(mvad.pam4)

![](_page_32_Figure_29.jpeg)

#### Retrieving cluster membership

• With PAM, the cluster membership is stored on the clustering element.

Clustering

MDS

Discrepancy analysis

References

• We may retrieve it using:

#### R> mvad.pam4\$clustering[1:10]

#### [1] 1 2 3 4 2 1 3 3 1 1

Reminder

#### MDS Clustering Discrepancy analysis Warning!!!

- Do not forget to specify the diss = T option.
- Otherwise (i.e. by default) functions agnes(), diana(), pam(), ...
- compute the Euclidean distance matrix between rows of the dissimilarity matrix.

29/64						WINFESTE DECEMBEN Biological Anti- Research of A	30/64						Benny Parkers Provider Provider
Reminder 00	Dissimilarity 000000000000000000000000000000000000	Centrotype	Clustering 00000000000000	MDS 00 0000	Discrepancy analysis	References	Reminder 00	Dissimilarity 000000000000000	Centrotype	Clustering 0000000000000	MDS 0 0000	Discrepancy analysis	References

## Transversal Distributions

![](_page_33_Figure_11.jpeg)

![](_page_33_Figure_12.jpeg)

![](_page_33_Figure_13.jpeg)

Sep.93 Jul.94 May.95 Apr.96 Feb.97 Jan.98 Nov.98

![](_page_33_Figure_15.jpeg)

employment inducation induca

#### • Three types of graphics

Exploring clusters graphically

- Transversal distribution with seqdplot()
- ② Frequency plots with seqfplot()
- Individual index-plots seqiplot()
- Representative sequences seqrplot()
- Required argument: state sequence object.
- Use group = cluster.membership.factor to get plots by cluster.

![](_page_33_Picture_26.jpeg)

References

![](_page_34_Figure_0.jpeg)

![](_page_34_Figure_1.jpeg)

# • Multidimensional Scaling (MDS) seeks numerical factors such

MDS

Discrepancy analysis

MDS

Discrepancy analysis

- that the Euclidean distance in the factor space reproduces at best the dissimilarity matrix.
- Through MDS, we get a scatter plot representation of sequences.
- R > mds2d <- cmdscale(mvad.dist, k = 2)
- R> plot(mds2d, pch = mvad.cl4, col = mvad.cl4)
- R> legend("bottomright", col = 1:4, pch = 1:4, legend = clust.labels)

![](_page_35_Figure_0.jpeg)

![](_page_35_Figure_1.jpeg)

38/64

 Reminder
 Dissimilarity
 Centrotype
 Clustering
 MDS
 Discrepancy analysis
 References

 Sort:
 First factor of MDS analysis
 Analysis

![](_page_35_Figure_4.jpeg)

## Code for scatterplot colored by type of school

Centrotyp

R> plot(mds2d, pch = as.integer(mvad\$Grammar), col = mvad\$Grammar)
R> legend("bottomright", col = 1:2, pch = 1:2, legend = c("Other school",

MDS

Discrepancy analysis

"Grammar school"))

![](_page_35_Figure_8.jpeg)

39/6

# Sequence discrepancy analysis

• How to measure the link between explanatory factors and sequences?

Clustering

- Using cluster and logistic regression is problematic:
  - Clusters are not homogenous.
  - Clustering choice are made on a statistical ground.
  - These choices may artificially hide or create an association.

MDS

Discrepancy analysis

- General principle of sequence discrepancy analysis.
  - Definition of a measure of discrepancy of a set of sequence using any dissimilarity measure.
  - Compute the share of this discrepancy accounted by a given explanatory factor.
  - The method is based on a generalization of ANOVA [Studer et al., 2010].
  - It may be extended to the multi-factor case and regression trees.

DE GENEVE

UNIVERSITÉ DE GENÉVE

![](_page_35_Picture_24.jpeg)

DE GENÉVE

References

#### Discrepancy of the set of sequences: Definition

• In the Euclidean case, the sum of squares SS can be expressed in terms of distances between pairs

Clustering

MDS

Discrepancy analysis

References

$$SS = \sum_{i=1}^{n} (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=i+1}^{n} (y_i - y_j)^2$$
$$= \frac{1}{n} \sum_{i=1}^{n} \sum_{j=i+1}^{n} d_{ij}$$

• Setting d<sub>ij</sub> equal to OM, LCP, LCS ... distance, we get a measure of dispersion.

43/64							
Reminder	Dissimilarity	Centrotype	Clustering	MDS	Discrepancy analysis	References	

## Computing the sequence discrepancy

Computing the sequence discrepancy R> dissvar(mvad.dist)

[1] 42.74502

## Discrepancy of the set of sequences: Interpretation

Clustering

- Interpretation of the discrepancy.
  - The discrepancy measures the between individual variability of the trajectories.

MDS

Discrepancy analysis

- It may reflect a form of precariousness
- r a multiplicity of choices faced by the individuals.
- It is different from the within individual longitudinal state diversity measures

44/64						Factorial and Arginetic Francisco Arginetics of California Organization of California Organization of California California (California)
Reminder	Dissimilarity	Centrotype	Clustering	MDS	Discrepancy analysis Re	ferences
00	00000000000000	000	0000000000000	0000	000000000000000000000000000000000000000	0

## Analysis of sequence discrepancy

- ANOVA like analysis based on pairwise dissimilarities
- We decompose the SS (Sum of squares equivalent)

$$SS_T = SS_B + SS_W$$

• Here, with the formula shown earlier

$$SS_{T} = \frac{1}{n} \sum_{i=1}^{n} \sum_{j=i+1}^{n} d_{ij}$$
  

$$SS_{W} = \sum_{g} \left( \frac{1}{n_{g}} \sum_{i=1}^{n_{g}} \sum_{j=i+1}^{n_{g}} d_{ij,g} \right)$$
  

$$SS_{B} = SS_{T} - SS_{W}$$

Reminder

UNIVERSITÉ DE GENÉVE

## Pseudo R-square and ANOVA Table

• ANOVA table for *m* groups

	Discrepancy	df	Mean Discr.	F
Between	SS <sub>B</sub>	$df_B = m - 1$	$\frac{SS_B}{df_B}$	$rac{SS_B}{SS_W} rac{df_W}{df_B}$
Within	$SS_W$	$df_W = \sum_g n_g - m$	$\frac{SS_W}{df_W}$	
Total	$SS_T$	$df_T = n - 1$		

Clustering

MDS

Discrepancy analysis

• Pseudo R<sup>2</sup>

$$R^2 = \frac{SS_B}{SS_T}$$

• Share of sequence discrepancy explained by the grouping variable.

UNIVERSITÉ DE GENÉVE

References

47/64

Reminder

Reminder	Dissimilarity	Centrotype	Clustering	MDS	Discrepancy analysis	References
00	000000000000000	000	0000000000000	0000		000
Permuta	ation test					

- Estimate the empirical distribution of *F* under independence.
- Compute *F*<sub>perm</sub> the *F* value associated with a random permutation of the profil.
- Permutation: randomly reassign each covariate profile to one of the observed sequence
- Repeat this step *R* times.
- The *p*-value associated with the test is  $p(F_{obs} > F_{perm})$ .
- The confidence interval of the *p*-value is given by  $p \pm 1.96 \sqrt{p(1-p)/R}$
- It is generally admitted that 5000 permutations should be used to assess a significance threshold of 1% and 1000 for a threshold of 5% [Manly, 2007].

#### 

• Pseudo F

$$F = \frac{SS_B/(m-1)}{SS_W/(n-m)}$$

- Normality is not defendable in this setting.
- F cannot be compared with an F distribution.
- The significance is assessed through a permutation test

		c

Clustering

DE GENEVE

References

Discrepancy analysis

## Analysis of sequence discrepancy

• Running an ANOVA like analysis for livboth R> da <- dissassoc(mvad.dist, group = mvad\$livboth, R = 5000) R> print(da) Pseudo ANOVA table: SS df MSE 110.0993 1 110.09927 Exp 30324.3558 710 Res 42.71036 Total 30434.4551 711 42.80514 Test values (p-values based on 4999 permutations): PseudoF PseudoR2 PseudoF\_Pval PseudoT PseudoT\_Pval 2.577812 0.003617586 0.0176 0.333323 0.023 Variance per level: n variance 261 40.88325 no 451 43.57833 yes Total 712 42.74502 

48/64

![](_page_38_Figure_0.jpeg)

Clustering

MDS

Discrepancy analysis

References

#### R> hist(da, col = "cyan")

Reminder

![](_page_38_Figure_2.jpeg)

#### Differences over time

## Homogeneity of the discrepancy

Homogeneity of the discrepancy

• Test of the difference of within group discrepancy.

Clustering

Discrepancy analysis

- Based on a generalization of the Bartlett Test.
- Significance assessed through permutation tests.

52/64						WUNNESSTE Die Generation Report and Provide Antonio	
Reminder 00	Dissimilarity 0000000000000	Centrotype 000	Clustering	MDS 0000	Discrepancy analysis	References	
Differences over time							

- How do differences between groups vary over time?
- At which age do trajectories most differ across birth cohorts?
- Compute  $R^2$  for short sliding windows (length 6)
- We get thus a sequence of  $R^2$ , which can be plotted
- Similarly, we can plot series of
  - total within (residual) discrepancy  $(SS_W)$
  - within discrepancy of each group  $(SS_G)$

#### R> mvad.diff <- seqdiff(mvad.seq, group = mvad\$gcse5eq)</pre>

#### R> mvad.diff\$stat[1:4, ]

 PseudoF
 PseudoR2
 PseudoT

 Sep.93
 29.09196
 0.03936176
 2.313692

 Oct.93
 29.39664
 0.03975760
 2.223468

 Nov.93
 29.76849
 0.04024027
 2.265784

 Dec.93
 30.09793
 0.04066750
 2.304112

#### R> mvad.diff\$variance[1:4, ]

no yes Total Sep.93 0.3688107 0.3113979 0.3620982 Oct.93 0.3691362 0.3127219 0.3629661 Nov.93 0.3704210 0.3133136 0.3642237 Dec.93 0.3725771 0.3146893 0.3663363

![](_page_39_Figure_0.jpeg)

![](_page_39_Figure_1.jpeg)

![](_page_39_Figure_2.jpeg)

- Generalize previous approach for multiple covariates.
- Here, we consider Type III effects
- Measure the additional contribution of each covariate *v* when we accounted for all other covariates.
- The F statistics reads

$$F_{v} = \frac{(SS_{B_{c}} - SS_{B_{v}})/p}{SS_{W_{c}}/(n - m - 1)}$$

where the  $SS_{B_c}$  and  $SS_{W_c}$  are the explained and residual sums of squares of the full model,  $SS_{B_v}$  the explained sum of squares of the model after removing variable v, and p the number of indicators or contrasts used to encode the covariate v.

• Significance is assessed again through permutation tests.

ninder	00000000000000000000000000000000000000	Centrotype 000	00000000000000000000000000000000000000	0000	Occorrepancy analysis
lotting	within di	screpanci	es over t	ime	

R> plot(mvad.diff, lwd = 3, stat = "Variance", legendposition = "bottomleft")

MDS

Discrepancy analysis

![](_page_39_Figure_12.jpeg)

# Running a Multiple factor analysis

R> +	<pre>k&gt; da.mfac &lt;- dissmfac(mvad.dist ~ male + Grammar + funemp +</pre>							
R>	<pre>&gt;&gt; print(da.mfac)</pre>							
	Variable	PseudoF	PseudoR2	p_value				
1	male	3.274802	0.003840223	0.018				
2	Grammar	21.124081	0.024771330	0.000				
3	funemp	4.483016	0.005257046	0.004				
4	gcse5eq	75.725976	0.088800698	0.000				
5	fmpr	2.715988	0.003184926	0.048				
6	livboth	2.314571	0.002714201	0.058				
7	Total	24.829102	0.174448528	0.000				

Clustering

UNIVERSITÉ DE GENÉVE Ren 00

Ρ

Reminde

References

DE GENÉVE

References

#### Tree structured discrepancy analysis

Reminder

• Objective: Find the most important predictors and their interactions.

Clustering

MDS

Discrepancy analysis

References

- Iteratively segment the cases using values of covariates (predictors)
- Such that groups be as homogenous as possible.
- At each step, we select the covariate and split with highest  $R^2$ .
- Significance of split is assessed through a permutation F test.
- Growing stops when the selected split is not significant.

59/64								
Reminder 00	Dissimilarity 0000000000000	Centrotype 000	Clustering	MDS 00000	Discrepancy analysis	References		
Creatin	Creating a Graphviz plot of the tree							

![](_page_40_Picture_8.jpeg)

• The file may be converted to a .jpg image using ImageMagick

R> shell("convert fg\_mvadseqtree.svg fg\_mvadseqtree.jpg")

# Growing the tree

```
R> dt <- disstree(mvad.dist ~ male + Grammar + funemp + gcse5eq +
        fmpr + livboth, data = mvad, R = 5000)
R> print(dt)
Dissimilarity tree
Global R2: 0.113
       |-- Root [ 712 ] var: 42.7
         |-> gcse5eq R2: 0.0821
              |-- no [ 452 ] var: 37.5
                |-> funemp R2: 0.0107
                     |-- no [ 362 ] var: 35.9
                       |-> male R2: 0.0123
                            |-- no [ 146 ] var: 38.7
                            |-- yes [ 216 ] var: 33.3
                     |-- yes [ 90 ] var: 41.8
              |-- yes [ 260 ] var: 42.3
                |-> Grammar R2: 0.0534
                     |-- no [ 183 ] var: 42.2
                     |-- yes [ 77 ] var: 34.9
```

Clustering

MDS

Discrepancy analysis

References

DE GENÉVE

60/64

![](_page_40_Picture_14.jpeg)

![](_page_40_Figure_15.jpeg)

UNIVERSITÉ DE GENÉVE

00	00000000000000000000000000000000000000	000	00000000000000000000000000000000000000	0000	Occorrepancy analysis	Referen
Exercise	s					

#### Exercise 4.2

- Use the previous distance matrix.
- Ompute the association with the cohort covariate using dissassoc.
- Interpret the differences graphically using Index-plot with all sequences sorted according to the first dimension of an MDS.
- Section 2 Explore the evolution of the association using sequiff.
- Fit a regression tree and plot the results.

## References I

Andrew Abbott and John Forrest. Optimal matching methods for historical sequences. *Journal of Interdisciplinary History*, 16:471–494, 1986.

Clustering

MDS

Discrepancy analysis

References

- Laurent Lesnard. Setting cost in optimal matching to uncover contemporaneous socio-temporal patterns. *Sociological Methods Research*, 38:389–419, 2010.
- V. Levenshtein. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet Physics Doklady*, 10:707–710, 1966.
- Bryan F. J. Manly. *Randomization, Bootstrap and Monte Carlo Methods in Biology*. Chapman & Hall, third edition edition, 2007.
- D. McVicar and M. Anyadike-Danes. Predicting successful and unsuccessful transitions from school to work using sequence methods. *Journal of the Royal Statistical Society A*, 165(2):317–334, 2002.
- Matthias Studer, Gilbert Ritschard, Alexis Gabadinho, and Nicolas S. Müller. Discrepancy analysis of complex objects using dissimilarities. In Fabrice Guillet, Gilbert Ritschard, Henri Briand, and Djamel A. Zighed, editors, *Advances in Knowledge Discovery and Management*, Studies in Computational Intelligence. Springer, Berlin, 2010. (forthcoming).

4/64

UNIVERSITÉ DE GENÉVE

	Example data set
Sequence analysis for social scientists	
Part V - Analysis of event sequences	<ul> <li>Data from McVicar and Anyadike-Danes [2002]</li> </ul>
	• 712 individuals
Matthias Studer, Alexis Gabadinho, Gilbert Ritschard, Nicolas	• 14 variables + 72 monthly activity state variables
S. Müller	1 = school
	2 = Further Education (FE)
Department of Econometrics, University of Geneva	3 = employment
http://mephisto.unige.ch/biomining	4 = training
Summer School on Advanced Methods for the Analysis of	5 = joblessness
Complex Event History Data, Bristol 28 20 June 2010	$\mathbf{O}$ $\mathbf{O}$ = Higner Education (HE)
Complex Event history Data, Dristol, 20-29 Julie 2010	

Reminder

References

References

1/24

Reminder

 Reminder
 Mining event sequences

 oo
 Occording the state sequence object

- Loading TraMineR and the mvad data set
- R> library(TraMineR)
  R> data(mvad)
- We consider only 70 states from Sept 93, i.e. we skip July and August 93 (i.e. summer holiday).

Mining event sequences

R> mvad.lab <- c("employment", "further education", "higher education",
+ "joblessness", "school", "training")</pre>

- R> mvad.shortlab <- c("EM", "FE", "HE", "JL", "SC", "TR")
- R> mvad.seq <- seqdef(mvad, 17:86, states = mvad.shortlab, labels = mvad.lab)

3/24

Reminde

6/24

Mining event sequences

Mining event sequences

# Analysis of event sequences Objective

- Focus on events, rather than states.
- Interest in the patterns of events.
  - Pattern of events: events that occur systematically together and in the same order
- Are there typical "patterns" of events?
- Relationship with covariates
  - Which patterns best discriminate specific groups?
  - Typical differences in event sequences between men and women.
- Event sequences analysis may be used for searching typical state sequencing.
- Association rules between event subsequences:
  - Sequence Leaving home → Childbirth generally followed by Marriage → Second Childbirth

References

References

o O	Mining event sequences	References	Reminder 00	Mining event sequences ⊙⊙●⊙○○○○○○○○○○○○	References
Events and tra	ansitions		subsequence		
<ul> <li>Event set</li> <li>Transition</li> <li>Example</li> <li>(LHome, Union)</li> <li>(LHome, Union)</li> <li>* (LHome, Union)</li> </ul>	quence: time ordered transitions. on: set of non-ordered events. on) $\rightarrow$ (Marriage) $\rightarrow$ (Childbirth) , Union) and (Marriage) are transitions. 7, "Union" and "Marriage" are events.		<ul> <li>A subset that</li> <li>eacletee</li> <li>the</li> </ul> Example <ul> <li>A</li> <li>B</li> <li>C</li> </ul> • C is a subset of the subset	quence $B$ of a sequence $A$ is an event seq in event of $B$ is an event of $A$ . events occur in $B$ in the same order as in $A$ . (LHome, Union) $\rightarrow$ (Marriage) $\rightarrow$ (Childli (LHome, Marriage) $\rightarrow$ (Childbirth). (LHome) $\rightarrow$ (Childbirth).	uence such birth).
4 eminder	Mining event sequences	Constant Constant References	respecte B is not occurs b 8/24 Reminder	d. a subsequence of <i>A</i> , since we do not know refore "Marriage". Mining event sequences	w if "LHome" ®
• Frequent and	discriminant subsequences		Data Format	000000000000000000000000000000000000000	

![](_page_44_Picture_0.jpeg)

- id Individual identifier.
- timestamp Time stamp (real valued) of the event.
- event The code (string) of the event.
- One line per event.

1 1 0 PartTime 2 2 0 NoActivity 3 2 4 Start 4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime	1 1 0 PartTime 2 2 0 NoActivity 3 2 4 Start 4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime		id	time	event
2 2 0 NoActivity 3 2 4 Start 4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime	2 2 0 NoActivity 3 2 4 Start 4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime	1	1	0	PartTime
3 2 4 Start 4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime	3 2 4 Start 4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime	2	2	0	NoActivity
4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime	4 2 4 FullTime 5 2 11 Stop 6 3 0 PartTime	3	2	4	Start
5 2 11 Stop 6 3 0 PartTime	5 2 11 Stop 6 3 0 PartTime	4	2	4	FullTime
6 3 0 PartTime	6 3 O PartTime	5	2	11	Stop
		6	3	0	PartTime

Creating an event sequence object From a state sequence object

- Function seqecreate().
- Argument tevent sets the choice for automatic conversion.

Mining event sequences

- Here, we assign an event to the start of each spell spent in a given "state".
- Look for state patterns.
- Other solutions are available.
- R> data(mvad)
- R> mvad.shortlab <- c("EM", "FE", "HE", "JL", "SC", "TR")</pre>
- R> mvad.seq <- seqdef(mvad[, 17:86], labels = mvad.shortlab)</pre>
- R> mvad.seqe <- seqecreate(mvad.seq, tevent = "state")</pre>

Creating an event sequence object  $\mathsf{Using\ the\ TSE\ format}$ 

Reminde

- Function seqecreate().
- With arguments id, timestamp and event we provide the data in the TSE format.

R> actcal.seqe <- seqecreate(id = actcal.tse\$id,</pre>

+ timestamp = actcal.tse\$time, event = actcal.tse\$event)

		UNIVERSITÉ De CENEVE Receive de l'article de
12/24		

 Reminder
 Mining event sequences
 References

 00
 0000000
 00000000

 Event sequence representation

- Each sequence is displayed in the following form (e1,e2,...)-time-(e2,...)-time
- where (e1,e2,...) is the transition defined by the simultaneous occurrences of events e1,e2,....
- time is the time (numerical value) between two transitions (or to the end of the observation time)

R> print(mvad.seqe[2])

[1] (FE)-36.00-(HE)-34.00

Reminder

DE GENÉVE

References

#### Reminder

#### Mining event sequences

#### Finding most frequent subsequences

Function seqefsub(), to which we must provide

- The event sequences (an event sequence object)
- The minimal support (with argument pMinSupport)

```
R> mvad.fsubseq <- seqefsub(mvad.seqe, pMinSupport = 0.05)
R> mvad.fsubseq[1:5]
```

![](_page_45_Picture_7.jpeg)

Reminder 00	Mining event sequences 000000000000000000000000000000000000	Refe
Finding most discrimina	nt subsequences	

- Aim is to identify the frequent sequences that are most strongly related with a given factor.
- Discriminant power is evaluated with *p*-value of a Chi-square independence test.
- Function seqecmpgroup()
- To which we provide the frequent subsequence object and a group factor (gcse5eq).
- A Bonferroni correction is applied when passing argument method="bonferroni"

Reminde

## Graphical display of most frequent subsequences

We can just apply plot() on the object returned by seqefsub()

- Use indexes ([1:15]) for selecting the subsequences to include (subsequences are sorted by decreasing frequencies).
- Other arguments are passed to the function barplot()

#### R> plot(mvad.fsubseq[1:15], col = "cyan", ylab = "Frequency", xlab = "Subsequences", cex = 1.5)

![](_page_45_Picture_22.jpeg)

![](_page_45_Picture_24.jpeg)

Subsequence

```
R> mvad.discr <- seqecmpgroup(mvad.fsubseq, group = mvad$gcse5eq,
+
       method = "bonferroni")
R> mvad.discr[1:5]
  Subsequence
               Support
                             p.value statistic index
                                                        Freq.no Freq.yes
         (HE) 0.2556180 0.000000e+00 193.91346
                                                   8 0.08185841 0.5576923
1
2
    (SC)-(HE) 0.1558989 0.000000e+00 114.94365
                                                  13 0.04424779 0.3500000
3
    (FE)-(HE) 0.1264045 1.865175e-14 65.81448
                                                  16 0.04867257 0.2615385
4
    (TR)-(EM) 0.3174157 1.538769e-13 61.84409
                                                   6 0.42256637 0.1346154
5
         (TR) 0.3497191 1.478151e-12 57.42334
                                                   4 0.45353982 0.1692308
   Resid.no Resid.yes
1 -7.306716 9.633959
2 -6.011885 7.926715
3 -4.648228 6.128722
4
  3.967957 -5.231780
5 3.732446 -4.921258
Computed on 712 event sequences
  Constraint
                       Value
```

countMethod One by sequence

References

UNIVERSITÉ DE GENÉVE

DE GENEVE

#### Reminder Mining event sequences References Graphical display, frequencies

#### R > plot(mvad.discr[1:15], cex = 1.5)

![](_page_46_Figure_2.jpeg)

Sequential association rules

![](_page_46_Figure_4.jpeg)

- Has a minimal support
- When subseq<sub>1</sub> occurs, it is most often followed by subseq<sub>2</sub>
- Extracted from frequent sequences.
- Extraction criteria:
  - Confidence: *p*(subseq<sub>2</sub> | subseq<sub>1</sub>)

• Lift: 
$$\frac{p(\text{subseq}_2 \mid \text{subseq}_1)}{p(\text{subseq}_2)}$$

o ...

Reminde

#### References

## Graphical display, residuals

R> plot(mvad.discr[1:15], ptype = "resid", cex = 1.5)

![](_page_46_Figure_17.jpeg)

#### Reminde Mining event sequences 00 Extracting association rules

• From the mined frequent subsequences, we can extract association rules :

					0	
		H	Rules	Support	Conf	Lift
7	(SC)	=>	(HE)	111	0.5692308	2.226881
50 (SC)	) => (	(HE)-	-(EM)	36	0.1846154	1.961883
23 (JL)	) => (	(EM)-	-(JL)	56	0.2295082	1.201543
9	(TR)	=>	(JL)	96	0.3855422	1.125025
28 (TR)	=> (	(EM)-	-(JL)	53	0.2128514	1.114340
19 (TR)	=> (	(JL)-	-(EM)	63	0.2530120	1.098443
2	(TR)	=>	(EM)	226	0.9076305	1.087934
32 (FE)	-(TR)	=>	(EM)	50	0.8771930	1.051450
36 (EM)	-(TR)	=>	(EM)	46	0.8679245	1.040341
10	(FE)	=>	(HE)	90	0.2578797	1.008848

19/24

References

00	References	00
Exercises		References

#### Exercise 5.1

- Use the biofam data set.
- **2** Create an event sequence from the state sequence object.
- **3** Find the most frequent subsequences.
- Ind the subsequences that discriminate the most the cohort covariable.
- O Plot the results.
- **1** Extract all association rules with a "Lift" greater than 1.

#### UNIVERSITÉ DE GENÉVE

#### References

- R. Agrawal and R. Srikant. Mining sequential patterns. In Philip S. Yu and Arbee L. P. Chen, editors, Proceedings of the International Conference on Data Engeneering (ICDE), Taipei, Taiwan, pages 487-499. IEEE Computer Society, 1995.
- D. McVicar and M. Anyadike-Danes. Predicting successful and unsuccessful transitions from school to work using sequence methods. Journal of the Royal Statistical Society A, 165(2):317-334, 2002.
- Nicolas S. Müller, Matthias Studer, Gilbert Ritschard, and Alexis Gabadinho. Extraction de règles d'association séquentielle à l'aide de modèles semi-paramétriques à risques proportionnels. Revue des nouvelles technologies de l'information RNTI, E-19:25-36, 2010.
- Matthias Studer, Nicolas S. Müller, Gilbert Ritschard, and Alexis Gabadinho. Classer, discriminer et visualiser des séguences d'événements. Revue des nouvelles technologies de l'information RNTI, E-19:37-48, 2010.

24/24

![](_page_47_Picture_18.jpeg)

Importing the mvad data set Description and visualization Computing Optimal Matching distances Clustering References

#### Exporting from R

## Sequence analysis for social scientists Part VI - Overview of sequence analysis with Stata

Alexis Gabadinho, Matthias Studer, Gilbert Ritschard, Nicolas S. Müller

Department of Econometrics, University of Geneva http://mephisto.unige.ch/biomining

Summer School on Advanced Methods for the Analysis of Complex Event History Data, Bristol, 28-29 June 2010

Importing the mvad data set Description and visualization Computing Optimal Matching distances Clustering References

UNIVERSITÉ DE GENÉVE

Importing in Stata and preparing the data

- We import the mvad data set by running the '.do' script produced by the write.foreign() function
  - . do "Z:\...\import-mvad.do"
- We need to reshape the sequence data (data have to be in a kind of "SPELL" format)
  - . reshape long month, i(id) j(order)
- Now we declare the sequence data (in the same way as we do with seqdef() in R)
  - . sqset month id order

- We begin by exporting the mvad data frame
- We do some data preparation in R because it is easier
- We rename some variables because original names contain a
  - '.' that is not supported by Stata
  - > names(mvad)[6] <- "N\_Eastern"</pre>
  - > names(mvad)[8] <- "S\_Eastern"
  - > names(mvad)[15:86] <- paste("month", 1:72, sep = "")</pre>
- We export the mvad data frame in the Stata format
  - > library(foreign)
  - > write.foreign(mvad, datafile = "mvad.dta", codefile = "import-mvad.do",
  - + package = "Stata")

Importing the mvad data set Description and visualization Computing Optimal Matching distances Clustering References

## Sequence frequency tables and index plots

- The sqtab command displays the sequence frequency table . sqtab, ranks(1/10)
- The sqindexplot produces sequence index plots
  - . sqindexplot, by(gcse5eq)

6/12

#### Importing the mvad data set Description and visualization Computing Optimal Matching distances Clustering References

#### Optimal matching - I

- Constant substitution cost matrix (default value is 1 for indels and 2\*indel cost for substitutions)
  - . sqom
- Distances are standardized by dividing each distance by the length of the longest sequence in the data set
- By default, sqom performs Optimal Matching between each sequence and the most frequent sequence in the data
- The distances are stored in the <u>\_SQdist</u> variable

0/10				
8/12				
Importing the mvad data set	Description and visualization	Computing Optimal Matching distances	Clustering	References
Clustering				

- First we have to issue sqclusterdat, which constructs a dataset built from the last sqom command
  - . sqclusterdat
- Now we perform a hierarchical clustering and select the number of clusters
  - . clustermat wardslinkage SQdist, name(wards) add
  - . cluster tree wards, cutnumber(4)
- Finally we merge the cluster results with the original sequence data
  - . sqclusterdat, return

#### Optimal matching - II

- Defining a substitution cost matrix has to be done manually (we use the output of seqsubm(mvad.seq, method="TRATE") to define substitution costs based on observed transition rates)
  - . matrix sub = (
  - 0.000000, 1.967601, 1.987270, 1.951231, 1.984684, 1.959993
  - 1.967601, 0.000000, 1.993341, 1.963360, 1.987531, 1.992045
  - 1.987270,1.993341,0.000000,1.996033,1.982969,1.999488
  - 1.951231, 1.963360, 1.996033, 0.000000, 1.985649, 1.972029
  - 1.984684, 1.987531, 1.982969, 1.985649, 0.000000, 1.994867
  - 1.959993, 1.992045, 1.999488, 1.972029, 1.994867, 0.000000)
- Now we compute the pairwise distances by adding the full option (the distance matrix is stored as \_SQdist)
- . sqom, subcost(sub) full

9/12

Importing the mvad data set Description and visualization Computing Optimal Matching distances Clustering References

References I

Christian Brzinsky-Fay, Ulrich Kohler, and Magdalena Luniak. Sequence analysis with Stata. *The Stata Journal*, 6:435 460, 2006. URL http://www.wz-berlin.de/~brzinsky-fay/publications.de.htm.

![](_page_49_Picture_33.jpeg)

UNIVERSITÉ DE GENÉVE

onclusior	: Sequence	of	analysis	

The TraMineR website

References

UNIVERSITÉ DE GENÉVE Conclusion: Sequence of analysis

The TraMineR website

References

DE GENEVE

Sequence of analysis (state sequences) - I

Sequence analysis for social scientists Part VII - Conclusion

Matthias Studer, Alexis Gabadinho, Gilbert Ritschard, Nicolas S. Müller

Department of Econometrics, University of Geneva http://mephisto.unige.ch/biomining

Summer School on Advanced Methods for the Analysis of Complex Event History Data, Bristol, 28-29 June 2010

1/10

Conclusion: Sequence of analysis The TraMineR website Sequence of analysis (state sequences) - II

- Defining a dissimilarity measure between individual sequences to:
  - find representative sequence (most frequent, centroid, ...)
  - compute the dispersion of the sequences
  - build a typology (cluster analysis)
    - study the relationships between clusters and covariates (sex, cohort, ...), logistic models ...
  - plot sequences in a plan using Multi-dimensional scaling
  - assess the association with a covariable using ANOVA (Analysis of discrepancy: Part of discrepancy explained by one or several factors).
  - segment the sequences in homogenous groups through tree structured approach.

- Exploring sequence distribution (seqdplot, seqfplot, seqiplot)
- Characteristics of the set of sequences
  - Sequence of transversal characteristics (entropies, modal states, ...)
  - Distribution of longitudinal characteristics (entropy, turbulence, time spent in each state, ...)
    - Association between longitudinal characteristics of parallel sequences (family-profession, ego-partner, ...)
- Preceding analyses by groups (sex, birth cohorts, ...), comparisons

Conclusion: Sequence of analysis	The TraMineR website 0	References
Sequence of analysis (event	sequences)	

#### • Seeking frequent subsequences

- Studying the relationship between frequent event subsequences and covariates (sex, cohort, ...), logistic models.
- Studying the effect of experimenting the subsequence on entropy, or other response variable ...
- Finding most discriminant subsequences for a given categorical variable (sex, cluster, ...) Extracting sequential association rules (next release of TraMineR)

#### Conclusion: Sequence of analysis

The TraMineR website

#### The TraMineR website

#### TraMineR website:

• http://mephisto.unige.ch/traminer/

#### There, you may:

- subscribe to the TraMineR user mailing list.
- submit a bug report or a feature request.
- find documentation and articles about TraMineR.
- find a list of articles using TraMineR (maybe yours?)

7/10		
Conclusion: Sequence of analysis	The TraMineR website 0	

## References II

- D. McVicar and M. Anyadike-Danes. Predicting successful and unsuccessful transitions from school to work using sequence methods. *Journal of the Royal Statistical Society A*, 165(2):317–334, 2002.
- Nicolas S. Müller, Alexis Gabadinho, Gilbert Ritschard, and Matthias Studer. Extracting knowledge from life courses: Clustering and visualization. In II-Yeol Song, Johann Eder, and Tho Manh Nguyen, editors, *Data* Warehousing and Knowledge Discovery, 10th International Conference, DAWAK 2008, Turin, Italy, September 2-5, volume LNCS 5182 of Lectures Notes in Computer Science, pages 176–185. Springer, Berlin Heidelberg, 2008a. doi: 10.1007/978-3-540-85836-2\\_17.
- Nicolas S. Müller, Sylvain Lespinats, Gilbert Ritschard, Matthias Studer, and Alexis Gabadinho. Visualisation et classification des parcours de vie. *Revue des nouvelles technologies de l'information RNTI*, E-11, II:499–510, 2008b.
- Nicolas S. Müller, Matthias Studer, Gilbert Ritschard, and Alexis Gabadinho. Extraction de règles d'association séquentielle à l'aide de modèles semi-paramétriques à risques proportionnels. *Revue des nouvelles technologies de l'information RNTI*, E-19:25–36, 2010.

#### References

#### References I

Andrew Abbott and John Forrest. Optimal matching methods for historical sequences. *Journal of Interdisciplinary History*, 16:471–494, 1986.

- Andrew Abbott and Alexandra Hrycak. Measuring resemblance in sequence data: An optimal matching analaysis of musician's carrers. *American Journal of Sociolgy*, 96(1):144–185, 1990.
- Alexis Gabadinho, Gilbert Ritschard, Matthias Studer, and Nicolas S. Müller. Summarizing sets of categorical sequences. In *International Conference on Knowledge Discovery and Information Retrieval, Madeira, 6-8 October, 2009*, pages 62–69. INSTICC, 2009a. (Received the Best Paper Award).
- Alexis Gabadinho, Gilbert Ritschard, Matthias Studer, and Nicolas S. Müller. Mining sequence data in R with the TraMineR package: A user's guide. Technical report, Department of Econometrics and Laboratory of Demography, University of Geneva, Geneva, 2009b. URL http://mephisto.unige.ch/traminer/.
- Alexis Gabadinho, Gilbert Ritschard, Matthias Studer, and Nicolas S. Müller. Indice de complexité pour le tri et la comparaison de séquences catégorielles. *Revue des nouvelles technologies de l'information RNTI*, E-19:61–66, 2010.

8/10

Conclusion: Sequence of analysis The TraMineR website o
References III

Matthias Studer, Gilbert Ritschard, Alexis Gabadinho, and Nicolas S. Müller. Analyse de dissimilarités par arbre d'induction. *Revue des nouvelles technologies de l'information RNTI*, E-15:7–18, 2009.

- Matthias Studer, Nicolas S. Müller, Gilbert Ritschard, and Alexis Gabadinho. Classer, discriminer et visualiser des séquences d'événements. *Revue des nouvelles technologies de l'information RNTI*, E-19:37–48, 2010a.
- Matthias Studer, Gilbert Ritschard, Alexis Gabadinho, and Nicolas S. Müller. Discrepancy analysis of complex objects using dissimilarities. In Fabrice Guillet, Gilbert Ritschard, Henri Briand, and Djamel A. Zighed, editors, *Advances in Knowledge Discovery and Management*, Studies in Computational Intelligence. Springer, Berlin, 2010b. (forthcoming).

References

UNIVERSITÉ DE GENÉVE

References