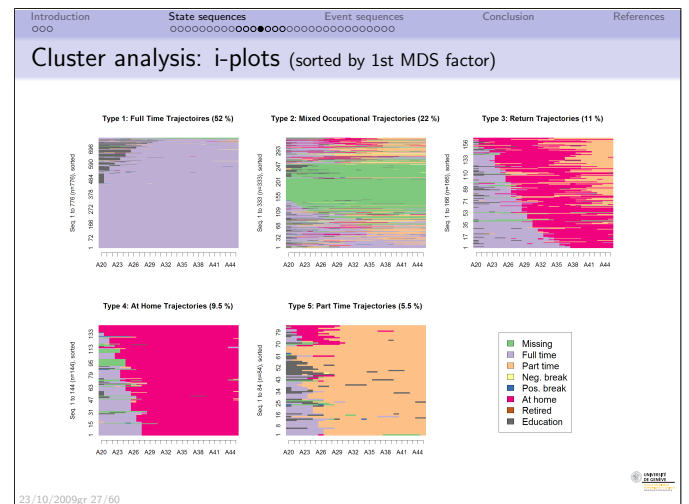
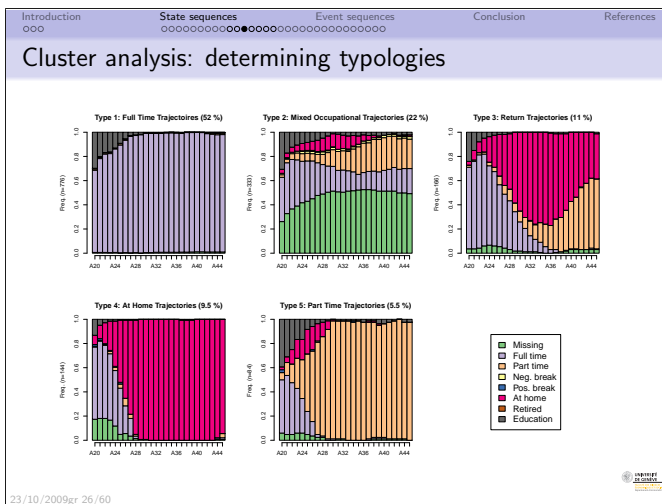


- Introduction State sequences Event sequences Conclusion References
- ## Pairwise dissimilarities between sequences
- Distance between sequences
 - Different metrics (LCP, LCS, OM, HAM, DHD, ...)
 - Once we have pairwise dissimilarities, we can
 - Determine a **central sequence** (centro-type)
 - Measure the **discrepancy between sequences**
 - Cluster a set of sequences
 - MDS scatterplot representation of sequences
 - Discrepancy analysis of a set of sequences (ANOVA)
 - Tree Structured Discrepancy Analysis (Induction trees)
- 23/10/2009gr 24/60

- Introduction State sequences Event sequences Conclusion References
- ## Deriving clusters from pairwise dissimilarities
- For each of the two sets of sequences: cohabitational and occupational
 - Compute **Pairwise dissimilarities** (a 1503×1503 matrix)
 - Here, we used **Optimal Matching (OM)**
 - For each pair $\{x, y\}$ of sequences, OM is the minimal cost of transforming one sequence into the other
 - insert/deletion (indel) cost = 1
 - substitution cost $c_{i,j} = c_{j,i} = 2 - p(i | j_{i-1}) - p(j | i_{i-1})$
 - Cluster by plugging obtained dissimilarity matrix in any cluster algorithm
 - We used an agglomerative hierarchical method with **Ward's** criteria
 - and retained partition into 5 clusters
- 23/10/2009gr 25/60



Introduction State sequences Event sequences Conclusion References

Pseudo F

- Pseudo F

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

- Normality is not defensible in this setting.
- F cannot be compared with an F distribution.
- The significance is assessed through a **permutation test**
- Permutation test: iteratively randomly reassign each covariate profile to one of the observed sequence and recompute the F .
- Empirical distribution** of F under independence.

23/10/2009gr 35/60

Introduction State sequences Event sequences Conclusion References

Analysis of sequence discrepancy

- Running an ANOVA like analysis for cohort3b

Pseudo ANOVA table:

	SS	df	MSE
Exp	106.4437	2	53.22183
Res	15645.8712	1500	10.43058
Total	15752.3148	1502	10.48756

Test values (p-values based on 999 permutation):

PseudoF	PseudoR2	PseudoF_Pval	PseudoT	PseudoT_Pval
5.10248	0.006757335	0	7.361347	0

Variance per level:

	n	variance
1910-1924	71	7.713761
1925-1945	659	9.651546
1946-1957	773	11.303784
Total	1503	10.480582

23/10/2009gr 36/60

Introduction State sequences Event sequences Conclusion References

Distribution of pseudo F

23/10/2009gr 37/60

Introduction State sequences Event sequences Conclusion References

Multiple factor analysis

- 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.

23/10/2009gr 38/60

Introduction State sequences Event sequences Conclusion References

Running a Multiple factor analysis

Variable	PseudoF	PseudoR2	p_value
1 sex	486.157573	0.222836269	0.000000000
2 cohort3b	5.297978	0.004856786	0.000999001
3 edu_lev	33.998319	0.046750636	0.000000000
4 Total	114.523325	0.314748465	0.000000000

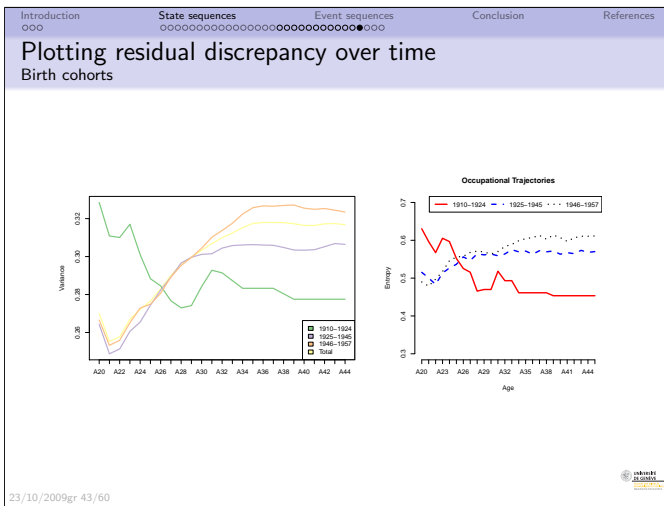
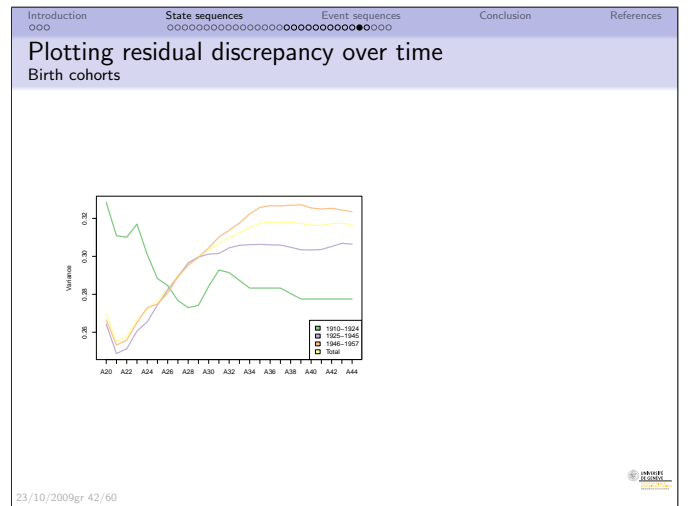
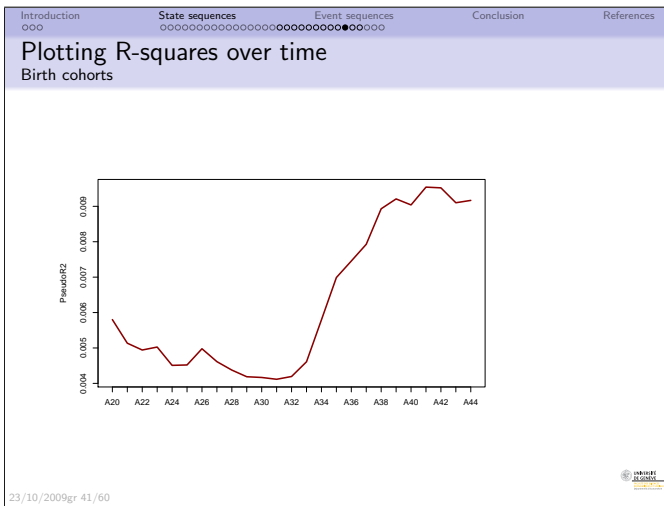
23/10/2009gr 39/60

Introduction State sequences Event sequences Conclusion 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 2)
- We get thus a **sequence of R^2** , which can be plotted
- Similarly, we can plot series of
 - total residual discrepancy (SS_W)
 - residual discrepancy of each group (SS_G)

23/10/2009gr 40/60



- Introduction State sequences Event sequences Conclusion References
- Tree structured discrepancy analysis
- Objective: Find the most important predictors and their interactions.
 - 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.
- 23/10/2009gr 45/60

Introduction State sequences Event sequences Conclusion References

Growing the tree

```

Dissimilarity tree
Global R2: 0.229
|-- Root [ 1503 ] var: 10.5
|-> sex R2: 0.179
|   |-- man [ 752 ] var: 4.37
|   |   |--> edu_lev R2: 0.143
|   |   |   |-- University [ 157 ] var: 6.28
|   |   |   |-- Compulsory/College+Prof/Prof.HS [ 595 ] var: 3.08
|   |   |-- woman [ 751 ] var: 12.8
|   |   |   |--> edu_lev R2: 0.0206
|   |   |   |-- Compulsory/College+Prof [ 632 ] var: 12.5
|   |   |   |   |--> edu_lev R2: 0.00905
|   |   |   |   |-- Compulsory [ 116 ] var: 12.0
|   |   |   |   |-- College+Prof [ 516 ] var: 12.5
|   |   |   |   |   |--> cohort3b R2: 0.00714
|   |   |   |   |   |   |-- 1946-1957 [ 280 ] var: 12.5
|   |   |   |   |   |   |-- 1910-1924/1925-1945 [ 236 ] var: 12.2
|   |   |   |   |   |   |-- Prof.HS/University [ 119 ] var: 13.1

```

23/10/2009gr 46/60

