Sequence Analysis in Sociology CREST – GENES OFPR 2008/9

> Brendan Halpin Department of Sociology University of Limerick brendan.halpin@ul.ie

Paris, May 14–28, 2009



- Explore the use of Sequence Analysis (SA) in the social sciences
- Explore the analytical use of "empirical typologies" and other measures derived from SA
- Review some aspects of the SA literature
- Consider alternatives to Optimal Matching (OM)
- Provide enough practical information for participants to conduct their own analyses using software
- How to think about using SA



- Explore the use of Sequence Analysis (SA) in the social sciences
 - Particularly the Optimal Matching Algorithm (OM)
 - Consider issues, problems, solutions, best practices



- Explore the analytical use of "empirical typologies" and other measures derived from SA
 - use of OM to generate classifications for further analysis
 - or to generate other sorts of trajectory level information



- Review some aspects of the SA literature
 - Applications of SA in real research
 - Methodological arguments about SA and OM



• Consider alternatives to Optimal Matching (OM)

- Non-aligning methods
- Combinatorial methods
- Duration-sensitive methods



- Provide enough practical information for participants to conduct their own analyses using software
 - Using SQ add-on for Stata
 - Using my faster but less user-friendly add-on for Stata
 - Using TraMineR for the R statistical language



- How to think about using SA
 - What's it good for?
 - When to use other methods?
 - What method to choose?
 - How to adapt it for your problem



- Is sequence analysis useful?
 - Does it go beyond exploratory and descriptive?
 - If not, is that enough?
- What does it "mean"?
 - How do the results inform us about sociological issues?
 - How can we manipulate the inputs to get better meaning?
 - How should we choose between algorithms for different substantive problems?



- "Sequences" are temporal (or at least linear) trajectories through a state space
- Sequence Analysis treats sequences "holistically"
- Alternative methodologies are more analytical (and often stochastic) and focus on factors such as the generative processes
- Usually at the cost of ignoring some aspect of the information encoded in the sequence
- In contrast, SA tends to be blind to the processes generating the sequence, thus focusing on the epiphenomenal?



- Sequences are ordered "trajectories" through a state space
- Their nature depends on
 - the nature of the state space: multi-dimensional, continuous, real, categorical?
 - the nature of the time dimension: atomic, discrete, continuous, stretchable or rigid?
 - how they start and finish, what they mean substantively



- The state space is important for how we think about distances between points
 - if multi-dimensional in \mathbb{R}^n , point distances are naturally Euclidean or other function of the space
 - if categorical, we often define pairwise point distances *a priori* or empirically (how?)
 - if many categories?



- The nature of time has a large bearing on the adequacy of sequences as a representation of the phenomenon
 - "Atomic" sequences consist of elements that are naturally separate and sequential, such as a series of purchases or votes, or steps or utterances, or CAGT bases in DNA – "time" is naturally discrete
 - Continuous time can be discretised
 - If the state changes very frequently (especially if continuous state) we can consider this as "sampling", for example, digitisation of an audio stream
 - If change is relatively rare, we may wish to represent the sequence as a series of spells (start and end-times of a period in which the state is constant)



- How we define a "whole" sequence is also an important issue where does it start and end?
- For some sequences, it is natural:
 - The steps of a dance, the words of a song
 - The rhetorical structure of a journal article or a folk tale
- For some purposes, fragmentary sequences can be used (e.g., searching for DNA matches)



- For life-course and other sequences, the requirements of the analysis impose structure
 - Usually cannot just match random segments of employment history
 - We impose comparability criteria (e.g., *t*₀ is a specific event, follow until a particular outcome or for a specified duration)
 - Issues of left- and right-censoring become relevant
- The various SA methods tend to be blind to these issues



- For life-course and other sequences, the requirements of the analysis impose structure
- Depending on the nature of the state space and the time dimension different approaches will be required
- A **R**^{*n*} state space simplifies state distance issues compared with categorical states, where we need to find a justification for our distances



- Time that is naturally discrete fits a token-sequence representation naturally
- "Sampling" continuous time raises issues of distortion due to the frequency of sampling
- Whether time has a "ruler" or calendar, or is "developmental" or stretchable, has a bearing on how attractive alignment is



"Conventional" methods for longitudinal data

• Many conventional approaches

- Hazard rate modelling (event history analysis)
- Time-series analysis
- Loglinear or Markov modelling of transition rates
- Start-end tables
- Panel analysis (cross-sectional time-series analysis, multiple-response approaches)
- Use of simple summaries of trajectory to predict future outcomes
- All have analytical strengths allow inference with respect to clear hypotheses
- In what way does SA offer something more than they do?



- Clearly a classification based on SA will do better than one based on summaries such as start/finish state or cumuluated duration in states – respects order
- Relative to hazard rate modelling, SA respects the whole trajectory, rather than looking at time to a single event (note the existence of repeated events hazard models)
- Models based on transition rates have difficulty with transition matrices which change in complex ways through time (due to life course and period effects, for instance)
- Multi-dimensional sequences are even more complex
- None of the conventional methods offer a digestible descriptive overview



- SA is very simple: turn information about the state space into information about the trajectories pairwise distances or similarities
- Permits the use of cluster analysis (CA) to generate a data-driven classification
- Permits comparing all sequences to a set of "typical" sequences
- Permits analysis of the multi-dimensional space implied by the inter-trajectory distances (multi-dimensional scaling, MDS)
- Also permits comparing grouped or paired sequences (e.g., couple's work histories, Han and Moen, 1999)



- Well chosen ideal-typical sequences can make for very interpretable results
- CA may or may not generate a useful classification
- MDS can inform CA, may yield interpretable dimensions itself



What do we do with SA: empirical typology

- Why is data-driven classification an empirical typology attractive?
- Up to $\sum_{i=1}^{m} n^i$ possible sequences, with *n* states and spells up to *m* tokens long (e.g., for 4 states over 20 months, more than a trillion possibilities)
- Observed sequences represent a highly structured subset
- The structure is much more than, say, that summarised by
 - starting state distribution and
 - the transition matrix averaged over the data set (or even changing through time)
- A good classification should (?) pick up this structure



- If we can generate a typology of sequences from theory, we may not need SA
- However, can be difficult to write foolproof rules to assign sequences to groups
- SA linking the observed to the ideal typical sequences will allow us to populate a theoretical typology automatically
- Classification by inspection may be possible, but can be impractical



- Cluster analysis on the pairwise distance matrix creates the classification automatically for us
- But can be difficult to characterise the resulting groups cleanly
- And clustering may be unstable if there are not natural groups



- Key questions: can it be more than exploratory and descriptive?
- Is exploration/description enough justification?



- How do we go from distances within a state space to distances between sequences in that space?
- How do we think about the sociological meaningfulness of inter-sequence distance?
- It is important to have a clear idea of what similarity and difference should mean at an intuitive, substantive level
- Clearly the more a pair of sequences go through the same or similar states, at the same or similar times, the more similar they are
- However, there are multiple ways to approach measuring this relationship



• Counting element-wise matches is intuitive, but only picks the *same* states at the *same* time

$$d = l - \sum_{i=1}^{i=l} (x_{1i} = x_{2i})$$

• The Longest Common Prefix method (e.g., Elzinga, 2009) is even "worse", as it only picks up common states at the beginning



• Hamming distance utilises information about the state-wise differences so it weights the *same and similar* states at the same time

$$d_H = \sum_{i=1}^{i=l} (|x_{1i} - x_{2i}|)$$

or

$$d_H = \sum_{i=1}^{i=l} \delta(x_{1i}, x_{2i})$$

where $\delta(a, b)$ is the distance between state *a* and state *b*

• However, such methods do not recognise similarity of states at *similar but not identical times* – no "alignment" or sliding along



Similarity at same or similar times

- One early method (Degenne et al., 1996) to pick up similarity at similar times has re-emeged recently as "Qualitative Harmonic Analysis" (Robette and Thibauld, 2008)
- This groups the sequence into short intervals and summarises the state distribution within them
- Data reduction (correspondence analysis) is used to simplify these summaries, and the results are clustered
- The same states at the same or different times *within the short intervals* contribute to similarity
- Intuitively appealing, though how to choose interval length and the nature of the summaries may need to be theorised



- Dijkstra and Taris (1995) suggested a method that focused on having the same states in a similar order
 - With a pair of sequences, drop all states that do not occur in both
 - In each sequence, drop duplicate states
 - Count how many moves are needed to reorder one state to match the other
- An intuitively satisfying focus on *same states in similar order*, and a tractable algorithm, but has problems (Abbott, 1995a)
- Not least is throwing away lots of useful information



The Optimal Matching Algorithm

- Going back to computer science work in the 1950s and 1960s by Levenshtein, and described in the 1980s in Sankoff and Kruskal (1983), the Optimal Matching Algorithm was introduced into sociology by Andrew Abbott (e.g., Abbott and Forrest, 1986; Abbott and Hrycak, 1990; Abbott, 1995b; Abbott and Tsay, 2000)
- A very general method for comparing token sequences
- Can recognise same and similar states via its *substitution cost matrix*
- Can recognise similarity that is out of alignment via *insertion and deletion*
- Using the Needleman–Wunsch algorithm to implement it, is quite efficient in computing terms



- OMA is the leading method for SA in social research
- Has received much criticism (some deserved), e.g., from Levine (2000) and Wu (2000)
- Recurrent themes:
 - that while its application in molecular biology is successful, it has poor sociological meaning
 - that there is no good basis for choosing substitution costs or insertion-deletion costs



- The main competition to OM comes from Cees Elzinga (e.g., Elzinga, 2003, 2005, 2009)
- Coming from the Dijkstra–Taris tradition, at least in motivation, his methods focus on the extent to which sequences pass through
 - the same states
 - in the same order (not necessarily consecutively)



- Distinguishes between substrings and subsequences (ABC is both a substring and a subsequence of ABCD, but AD is a subsequence)
- One measure is the longest common subsequence (results equivalent to OM under certain circumstances)
- Main measure is the count of common subsequences
- Also a spell-wise version, with a duration weighted count of common spell subsequences
- Different logic and rather different results to OM
- Beginning to appear in the literature (e.g., Aisenbrey and Fasang, 2007)



- Today I want to look more closely at the Optimal Matching Algorithm
- How it operates
- How to carry out analysis with it



- At base, OM has a very simple logic: the "cost" of "editing" one sequence to match another
- Two sequences which differ at one point have a distance proportional to the difference between the points
 - ABC and ABD have a difference related to $\delta({\tt C}, {\tt D}),$ the cost of substituting C with D
- Two sequences which differ in that one has an extra element differ
 - in relation to the cost of deleting the extra element
 - which is entirely equivalent to inserting it in the sequence where it is absent


- The former operation is called *substitution*, the latter *indel*
- Any sequence can be turned into any other by a concatenation of these operations
- The OM algorithm identifies the "cheapest" concatenation that achieves this
- "Cheap" implies cost: substitutions and *indels* need to have weights applied



• Substitution costs may be determined in many ways

- From *a priori* knowledge about the relations within the state space
- Derived from data about the state space
- Derived from observed transitions in the state space
- Or just give up:

$$c_s = \begin{cases} 0 & (x_i = x_j) \\ 1 & (x_i \neq x_j) \end{cases}$$



- The cost of insertion and deletion affects how prone to "alignment" the algorithm will be
- If *indels* are cheap they will often be chosen in preference to substitutions
- If they are sufficiently expensive they will only be used to equalise unequal sequence lengths
 - In the case of equal-length sequences and high *indel* costs, OM returns the same distance as the Hamming measure (no alignment similar states at same times)



- There are some key limits with respect to the relation between substitution and *indel* costs
- Since substitution is equivalent to an insertion plus a deletion, substitution costs greater than $2 \times indel$ will have no effect
- How high *indel* costs have to be to completely suppress alignment of equal-length sequences seems to depend on the data
 - In my experience, more than about 1.5 to 2 times the largest substitution cost is often enough
- Approx. working range: $0.5 < \frac{indel}{\max(c_s)} < c. 1.5 \rightarrow 2$



- Determining the cheapest set of "elementary operations" is potentially complex a large population of candidates
- However, it can be stated as a recursive problem and programmed very efficiently
- Understanding how it is programmed can help understand the principle of OM



$\Delta_{OM}(A^p, B^q) =$

$$\min \begin{cases} \Delta_{OM}(A^{p-1}, B^q) + indel \\ \Delta_{OM}(A^{p-1}, B^{q-1}) + \delta(a_p, b_q) \\ \Delta_{OM}(A^p, B^{q-1}) + indel \end{cases}$$

(Δ represents distance between sequences, and δ differences within the state space)



Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





Cell value:
$$min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$$

= $min(0 + 2, 2 + 2, 2 + 2) = 2$





ヘロト 人間 とくほとくほとう

44

= 990

Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





Cell value:
$$min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$$

$$= min(2+1, 2+2, 4+2) = 3$$





ヘロト 人間 とくほとくほとう

= 990

Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





Cell value:
$$min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$$

$$= min(4+0,3+2,6+2) = 4$$





ヘロト 人間 とくほとくほとう

= 990

Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





Cell value:
$$min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$$



Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





- 52

Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





Cell value: $min(c_{i-1,j-1} + \omega_{i,j}, c_{i,j-1} + \Delta, c_{i-1,j} + \Delta)$





To convert ABCD into CDAAB the following set of operations gives the cheapest path:

Operation	Intermediate state	Cost
Sequence 2	ABCD	= 0
		=
		=
		=
		=
		=
		=
Sequence 1	CDAAB	=



To convert ABCD into CDAAB the following set of operations gives the cheapest path:

Operation	Intermediate state	Cost
Sequence 2	ABCD	= 0
insert C	CABCD	+2 = 2
		=
		=
		=
		=
		=
Sequence 1	CDAAB	=



To convert ABCD into CDAAB the following set of operations gives the cheapest path:

Operation	Intermediate state	Cost
Sequence 2	ABCD	= 0
insert C	CABCD	+2 = 2
insert D	CDABCD	+2 = 4
		=
		=
		=
		=
Sequence 1	CDAAB	=



э

To convert ABCD into CDAAB the following set of operations gives the cheapest path:

Operation	Intermediate state	Cost
Sequence 2	ABCD	= 0
insert C	CABCD	+2 = 2
insert D	CDABCD	+2 = 4
const A = A	CDABCD	+0 = 4
		=
		=
		=
Sequence 1	CDAAB	=



To convert ABCD into CDAAB the following set of operations gives the cheapest path:

Operation	Intermediate state	Cost
Sequence 2	ABCD	= 0
insert C	CABCD	+2 = 2
insert D	CDABCD	+2 = 4
const A = A	CDABCD	+0 = 4
$\mathrm{subs}\mathtt{B}{ o}\mathtt{A}$	CDAACD	+1 = 5
		=
		=
Sequence 1	CDAAB	=



э

To convert ABCD into CDAAB the following set of operations gives the cheapest path:

Operation	Intermediate state	Cost
Sequence 2	ABCD	= 0
insert C	CABCD	+2 = 2
insert D	CDABCD	+2 = 4
const A = A	CDABCD	+0 = 4
$\mathrm{subs}\mathtt{B}{ o}\mathtt{A}$	CDAACD	+1 = 5
subs $C \rightarrow B$	CDAABD	+1 = 6
		=
Sequence 1	CDAAB	=



To convert ABCD into CDAAB the following set of operations gives the cheapest path:

Operation	Intermediate state	Cost
Sequence 2	ABCD	= 0
insert C	CABCD	+2 = 2
insert D	CDABCD	+2 = 4
const A = A	CDABCD	+0 = 4
$\mathrm{subs}\mathtt{B}{ o}\mathtt{A}$	CDAACD	+1 = 5
subs $C \rightarrow B$	CDAABD	+1 = 6
delete D	CDAAB-	+2 = 8
Sequence 1	CDAAB	= 8



- Where sequence may be of different lengths, the distance is usually divided by the length of the longer sequence
- In this case, 8 units thus become a pairwise distance of 1.6



- It may be interesting to look at the resulting alignments, but this is done much less in social science than in other contexts
- Note that more than one "cheapest" route may exist
- Hence there may be no single "best" alignment



OM and Hamming on example sequences

- In the next slides I present OM and Hamming distances on selected short sequences
- Using two cost configurations:



ヘロト 人間 とくほとくほう

Examples, "flat" substitution costs

AAAAAAA	1	0					OM di	stance	
AAABBBCD	1	5	0						
ABCDDDDD	1	7	6	0					
BAAACCDD	1	5	5	6	0				
BBDDACCC	1	7	7	6	5*	0			
DDDDABCD	1	7	5	6	6*	4	0		
DCCCBBAA	1	6	6	7	8	8	6	0	
DDABCDAA	Ι	5	6	6*	6	6**	4**	5	0
AAAAAAA	I	0					Hammi	ng dis	tance
AAABBBCD	1	5	0						
ABCDDDDD	1	7	6	0					
BAAACCDD	1	5	5	6	0				
BBDDACCC	1	7	7	6	6	0			
DDDDABCD	1	7	5	6	7	4	0		
DCCCBBAA	1	6	6	7	8	8	6	0	
DDABCDAA	1	5	6	7	6	8	6	5	0

Examples, "linear" substitution costs

AAAAAAA	I.	0					OM d	istan	ce	
AAABBBCD	1	8	0							
ABCDDDDD	T	18	10	0						
BAAACCDD	T	11	5	9	0					
BBDDACCC	T	14	10	8	9*	0				
DDDDABCD	T	18	12	10*	13*	6	0			
DCCCBBAA	1	11	13	15	16	11	9	0		
DDABCDAA	Ι	12	14	14*	13	10**	8**	7	0	
AAAAAAA	Ι	0					Hamm	ing d	istan	ce
AAABBBCD	T	8	0					•		
ABCDDDDD	T	18	10	0						
BAAACCDD	T	11	5	9	0					
BBDDACCC	T	14	10	8	11	0				
DDDDABCD	1	18	12	12	15	6	0			
DCCCBBAA	T	11	13	15	16	11	9	0		
DDABCDAA	1	12	14	16	13	16	14	7	0	

- AAAAAAAA with anything: no "order" implies no possible improvement from alignment
- DDABCDAA with ABCDDDDD: alignment reduces cost somewhat (14 vs 16)

D D A B C D A A - -- - A B C D D D D D 2 2 0 0 0 0 3 3 2 2



- BBDDACCC with DDABCDAA: alignment strongly reduces cost (10 vs 16)
 - B B D D A C C C -- - D D A B C D A A 2 2 0 0 0 1 0 1 2 2
- DDDDABCD with ABCDDDDD: alignment reduces cost only for the more diverse substitution cost configuration (10 vs 12)

- - D D D D A B C D A B C D D D - - D D 2 2 1 0 0 0 2 2 1 0



An example: birth and the labour market

- An application of OM: 5 years of labour market history of women who have a birth in month 25
- BHPS data note that even with a quite large sample observing people for c.15 years, only 675 cases fit the criteria
- We classify labour market status thus:

Full-time employed	
Part-time employed	
Unemployed	
Not in labour market	

• *indel* cost is 2



State distribution or "chronogram"



Index plot, unordered


Index plot, lexically sorted



SQ summary data

of observed sequences: 675
overall # of obs. elements: 4
 max sequence length: 73
of producible sequences: 8.920e+43

Cum.	% of observed	Sequences	Dbservations
31.25926	31.25926	211	 1
32	.7407407	5	2
32.14815	.1481481	1	3
32.44445	.2962963	2	5
32.59259	.1481481	1	15
32.74074	.1481481	1	50
32.88889	.1481481	1	127
33.03704	.1481481	1	249
			I
	33.03704	223	Total
S(
UNIV			

イロト イロト イヨト イヨト 二日 二

SQ Tabulating sequences

Sequence-Pattern	1	Freq.	Percent	Cum.
444444444444444444444444444444444444444	-+-	249	36.89	36.89
111111111111111111111111111111111111111	1	127	18.81	55.70
222222222222222222222222222222222222222	1	50	7.41	63.11
333333333333333333333333333333333333333	1	15	2.22	65.33
111111111111111111111114444444444444444	T	5	0.74	66.07
111111111111111111111444444444444444444	T	5	0.74	66.81
111111111111111111144444444444444444444	T	3	0.44	67.26
111111111111111111111111111111111111111	T	2	0.30	67.56
111111111111111111111111111111111111111	T	2	0.30	67.85
1111111111111111111111111144444444444	I	2	0.30	68.15
444444444444444444444444444444444444	T	1	0.15	99.56
444444444444444444444444444444444444444	T	1	0.15	99.70
444444444444444444444444444444444444444	T	1	0.15	99.85
444444444444444444444444444444444444444	I	1	0.15	100.00
Total	-+- 	 675	100.00	



75

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ ▲ □ ● のへで

OM distances



- We proceed by using cluster analysis
- Wards' method reasonably stable, widely available, tends to produce relatively balanced groups
- Hierarchical nested classifications sometimes an advantage
- Hierarchical how to define the appropriate stopping level?



Dendrogram



Indexplot in hierarchical CA order



- 79

Eight-cluster solution



The "interesting" clusters



Chronogram, "boring" clusters (1, 4, 7 & 8)



- 82

Chronogram, interesting clusters (2, 3, 5 & 6)



- We may also wish to characterise the cluster by a "typical" sequence
- Some cluster approaches supply this readily
 - The centroid
 - or medoid
 - or otherwise-defined sequence closest to the cluster centre
- Wards' method with Stata doesn't make this easy
- An alternative is the "modal" sequence



- The modal sequence is composed of the most common token at each time point
- Note this is a synthetic sequence, not drawn from the observed sequences
- While it summarises the time-ordered distribution it does not necessarily reproduce real transistions
- In fact, it can have improbable or impossible features



Cluster Modal sequence 2 3 ppF pppnnnnnnnppp nn nnp 4 5 6 7 8

- On the whole they pick up the main features of the clusters
- Note gaps for cluster 3: no single mode for some months
- However, lots of information discarded



A B >
 A
 B >
 A
 B
 A
 B >
 A
 B
 A
 B
 A
 A
 B
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A
 A

Cluster	Average months in				Av	Ν	
	FT	PT	UE	NonE	N-spells	Turbulence	
1	0.0	0.1	0.3	72.6	1.11	1.17	263
4	0.1	72.3	0.3	0.3	1.11	1.20	54
7	0.4	0.0	67.2	5.4	1.63	2.17	19
8	72.4	0.4	0.0	0.1	1.12	1.24	139
2	14.2	2.0	4.7	52.0	3.21	5.14	67
3	38.1	20.4	1.1	13.4	3.69	6.40	71
5	19.8	30.3	1.2	21.7	4.17	7.73	36
6	11.7	22.5	2.5	36.3	4.31	7.23	26



- Cluster analysis is not the only option for analysis of distance matrices
- Multi-dimensional scaling (MDS) also works on the space implied by the distances
 - Rather than group cases based on closeness,
 - attempt to extract the structure of the space
 - Relatively few "dimensions" rather than many pairwise distances
- The extracted dimensions may be meaningful and useful for further analysis
- Throws light on the cluster analysis



- MDS is primarily a data reduction technique
 - when faced with many variables (*N*), try to reduce them to a "few" dimensions (*n*)
- In normal use, pairwise distances are calculated from the variables; we get distances from OM
- If the distances are metric (see below) they imply a space of unknown dimensions, n << N



- Given pairwise distances for observations, MDS calculates values on *n* dimensions for each observation
- The distance between the points on the calculated dimensions approximates the observed distance
- As *n* → *N* the match between observed and reconstructed distances improves monotonically
- A good match is usually possible for *n* << *N*, depending on the structure of the data



• If |AB| = |BC| = 1 and |AC| = 2 we get a single dimension: A ______ B ____ C

• If |AB| = |BC| = |AC| = 1 we cannot fit it in a single dimension:



In general we may need up to *N* − 1 dimensions but hope to use far fewer



- If the data naturally has few dimensions, this should be picked up
- Noise in the data means that even if there are few dimensions, the observed and reconstructed distances may not match exactly
- Principal Components Analysis is the usual algorithm this extracts the dimensions as eigenvectors of the pairwise distance matrix



- Cluster analysis and MDS deal with space Euclidean space is a special case, and many non-Euclidean spaces are suitable
- However, distances must be *metric*

•
$$d(x,x) = 0$$

•
$$d(x,y) > 0, x \neq y$$

- d(x,y) = d(y,x), symmetry
- $d(x,z) \le d(x,y) + d(y,z)$ the "triangle inequality"
- As is clear, the ordinary space we live in satisfies all these criteria
- But they permit a wide range of types of space and distance (e.g., city-block distance instead of Euclidean)



- If distances are not metric, the implicit space is not coherent for CA or MDS
- Some dissimilarity measures can violate the triangle inequality:
 - A is close to B because of shared characteristic *x*
 - B is close to C because of shared characteristic *y*
 - A is very distant from C because they have no shared characteristic
- Can be useful for comparing sequences against a catalogue of reference sequences e.g., voice recognition
- But non-metric distances are not useful for CA or MDS



<ロト < 四ト < 三ト < 三ト

- Sometimes a transformation of a non-metric similarity measure will yield a metric measure
- Elzinga (2009) asserts that where φ(x, y) measures the amount of a characteristic shared by x and y, and the following holds:

$$\phi(x,y) = \phi(y,x)$$

$$0 \le \phi(x,y) \le \min \left\{ \phi(x,x), \phi(y,y) \right\}$$

then

$$d(x,y) = \phi(x,x) + \phi(y,y) - 2\phi(x,y)$$

is metric



С	lassical metr dissimilarit	ric y n	c multidimensi natrix: omlin	onal scaling			
	Figenwelueg	~ ~	_	05	Number of	obs =	= 675
	Eigenvalues /	- () =	90	Mardia fit	measure 1 -	- 0.7955
		ens			Maruia IIt	measure 2 -	- 0.9937
		Ι		abs(eige	envalue)	(eiger	value)^2
	Dimension	I	Eigenvalue	Percent	Cumul.	Percent	Cumul.
	1	-+- 	964.93863	72.25	72.25	98.79	98.79
	2	Т	67.007783	5.02	77.27	0.48	99.27
	3	I	30.236075	2.26	79.53	0.10	99.37
	4	-+- 	21.574791	1.62	81.15	0.05	99.41
	5	Т	15.707688	1.18	82.33	0.03	99.44
	6	Т	9.2522026	0.69	83.02	0.01	99.45
	7	Т	7.5743302	0.57	83.59	0.01	99.46
	8		5.9310838	0.44	84.03	0.00	99.46
	9	Т	4.5038763	0.34	84.37	0.00	99.46
	10	Ι	4.1936765	0.31	84.68	0.00	99.46

・ロ・・日・・日・・日・ 日・ 今々ぐ

Linear state, first 2 dimensions



First 2 dimensions, with cluster membership



A "string" in cluster 2



- The homogenous clusters (1, 4, 7 and 8) are distinctly located
- Clusters with substantial amounts of more than one state are located between these "vertices" and are quite diffuse
- Some evidence of "strings" adjacent trajectories that differ slightly
- Quite clear structure, but what does it tell us about the meaningfulness of grouping? simple sequences are distinct but complex sequences are more evenly distributed



- Dimension 1 is strongly related to the dimension of the state space
 - Set FT 0, PT 1, UE 2, NonE 3
 - Sum over the time-span to give a weighted cumulative duration
 - Correlation with first dimension is 0.9999
- Dimension 2 has non-employed who become part-timers at the low end, and the transition to non-employment at the top end



"Flat" state, first 2 dimensions



First 2 dimensions, with cluster membership



Main features in the "flat" state space

- Like before, the homogenous clusters are distinct and the others spread between them, with strings etc.
- Overall shape like a distortion of first one
- Dimension 1 is nearly as strongly related to the dimension of the state space
 - Correlation 0.9947
- Dimension 2 runs from 100% part-time to full-timers who exit late to non-employment (100% FT also high on D2)
- Not just the state space but the nature of the trajectories as well
- Trajectory state space is structured by the states, each forming a pole



- As we see MDS throws light on some of the clustering processes
- Distinct tight clusters of simple sequences
- Systematically located intermediate sequences but not in "obvious" groups
- Cluster analysis discriminates but cluster membership is sensitive to small changes
- The main dimensions are often interpretable, and may in some circumstances be useful as variables in further analysis





- Does Optimal Matching make sense for sociological data?
 - Is the algorithm suitable? (see elsewere)
 - How to parameterise it: substitution and *indel* costs
- Repeated claims in the literature:
 - that sociologists don't know how to set substitution costs,
 - that we can't match the effectiveness of molecular biology
- Yes, our analytical goals are often much less well defined than those of the biologists
- No, substitution costs are not an intractable problem



- Wu (2000) treats the choice of substitutions costs as an insurmountable while he has some misunderstandings that made it particularly difficult for him, it is an important stage of the OM process
- Many other writers agonise over the problem
 - Many opt for using transition rates to get data-driven cost
 - Or set all substitution costs equal to 1
- Neither option is really neutral we need to understand substitution cost setting better
- Similarly, how to set *indel* costs?


- The essence of SA is mapping a view of a state space onto a view of a trajectory space: $\delta(s) \rightarrow \Delta(S)$
- We start with *knowledge* or a *view* of how states relate to each other (what states are like each other, what states are dissimilar)
- With a suitable algorithm we map this perspective onto trajectories through the state space: what trajectories are more or less similar
- The nature of the algorithm determines
 - Whether the mapping makes sense
 - Exactly how the structure of the state space affects the structure of the trajectory space



- Can we expect OMA to provide a coherent $\delta(s) \rightarrow \Delta(S)$ mapping?
- Elementary operations are intuitively appealing:
 - $(\textbf{ABC, ADC}) = f(\delta(\textbf{B}, \textbf{D}))$
 - **2** $\Delta(\text{ABCD}, \text{ABD}) = f(indel)$
 - Imminimising concatenation of these two operations to link any pair of trajectories
- If 3 is reasonable, 1 and 2 determine how state space affects trajectory space



- Costs can be thought of as distances between states
- If state space is \mathbb{R}^n , distance is intuitive
- If state space is categorical, how define distance?
 - State space as efficient summary of clustered distribution in Rⁿ: distances are between cluster centroids
 - State space can be mapped onto specific set of quantitative dimensions; each state located at the vector of its mean values; Euclidean or other distances between vectors
 - States can be located relative to each other on theoretical grounds



- Transition rates frequently proposed as basis for substitution costs
- Critics of OMA complain of substitution operations implying impossible transitions (e.g., Wu, 2000)
- Even proponents of OMA are sometimes concerned about "impossible" transitions (e.g., Pollock, 2007),
- But substitutions are not transitions, not even a little bit!
 - substitutions happen across sequences, $\Delta(ABC, ADC) = f(\delta(B, D))$ (similarity of states)
 - transitions happen within sequences (movement between states)



- No logical connection between substitutions and transition rates
- but under certain circumstances transition rates can inform us about state distances
- If state space is a partitioning of an unknown \mathbb{R}^n , movement is random (unstructured), and the probability of a move is inversely related to its length, then
- Distance between states will vary inversely with the transition rates
- However, these conditions are often not met



- Example: using voting intentions as a way of defining inter party distances
- UK: relatively high Con–LibDem two-way flows; ditto Lab–LibDem
- But Con–Lab transitions much lower: implies a potentially incoherent space (non-metric, more below)
 - $\delta(\text{Con, Lab}) > \delta(\text{Con, LibDem}) + \delta(\text{LibDem, Lab})$



- This procedure confuses party state space and voter characteristics
- Voter polarisation/loyalty is trajectory information, not state information
- There is a strong analytical argument for trying to keep the two concepts as separate as possible
- Another type of problem: irrelevant distinctions can cause similar states to have low transition rates



• Very useful to think in spatial terms

- lacksim State space as efficient summary of clustered distribution in \mathbb{R}^n
- State space mapped onto specific set of quantitative dimensions
- State space defined on theoretical grounds
- For 1 and 2, explicitly multidimensional, in case 2 dimensions are explicit
- For 1 and 3, we can attempt to recover the implicit dimensions



Looking at state spaces

- Two very simple state spaces:
 - Single dimension, equally spaced:

0	1	2	3
1	0	1	2
2	1	0	1
3	2	1	0

• All states equidistant – n - 1 dimensions

0	1	1	1
1	0	1	1
1	1	0	1
1	1	1	0



э

- E.g., 2D picture of inter-party distances: location on left–right scale, plus on pro-/anti-EU scale
- Distances are Euclidean or other metric (e.g., L1)
 - Euclidean: $\sqrt{\sum_i (r_i s_i)^2}$
 - L1 (city block): $\sum_i |r_i s_i|$
- Generalises easily to many dimensions
- Problem: how to weight different dimensions?
 - Scale by standard deviation? Substantive importance?







Spatial structure of theoretical spaces

- We can analyse "theoretically-informed" or *ad hoc* state spaces spatially
- Principal components analysis of substitution matrix
- Examples: Halpin and Chan, 1998; McVicar/Anyadike-Danes 2002:

20	02.						
I–II	0	2	2	2	2	3	3
III	2	0	1	1	1	2	2
IVab	2	1	0	1	1	2	2
IVcd	2	1	1	0	1	2	2
V–VI	2	1	1	1	0	2	2
VIIa	3	2	2	2	2	0	1
VIIb	3	2	2	2	2	1	0

E	0	1	1	2	1	3
F	1	0	1	2	1	3
Η	1	1	0	2	1	2
S	2	2	2	0	1	1
Г	1	1	1	1	0	2
U	3	3	2	1	2	0



H&C, 1st two PCA dimensions





э

イロト イポト イヨト イヨト

H&C, dimensions 1 & 3





2

<ロト < 四ト < 三ト < 三ト

MVAD, 1st two dimensions





э

イロト イポト イヨト イヨト

MVAD, dimensions 1 & 3





æ

<ロト < 四ト < 三ト < 三ト

• State space structure passes through to trajectory space structure

- Distances between states clearly affect distances between trajectories containing high proportions of those states
 - If δ("A", "B") << δ("A", "C") then Δ("...AAAA..", "...BBB..") will tend to be less than Δ("...AAAA..", "...CCC..")
- Differential distances promote alignment: AADDAAA and AAADDAA are more likely to be aligned to match the DD if $\delta("A", "D")$ is large
- If the state distances are non-metric, the trajectory distances may also be non-metric (at least between trajectories consisting of near 100% one state)
- Unidimensional states spaces will tend to be reflected strongly in 1st principle component of trajectory space



- How much do different substitution regimes change the results?
- Let's examine a few:
 - All states equidistant
 - 2 All states equally spaced on a single dimension
 - All states on a single dimension but the ends more "extreme"
 - All states on a single dimension but more polarised
 - A 2 dimensional space



Substitution matrices

Equidistant	(0, 3, 3, 3 \	Simple linear	(0,	1,	2,	3	١.
	3, 0, 3, 3 \		1,	0,	1,	2	١.
	3, 3, 0, 3 \		2,	1,	0,	1	١.
	3, 3, 3, 0)		З,	2,	1,	0);
Extreme linear	(0, 1, 2, 6 \	Polarised linear	(0,	1,	5,	6	١
	1, 0, 1, 2 \		1,	0,	1,	5	\
	2, 1, 0, 1 \		5,	1,	0,	1	\
	6, 2, 1, 0);		6,	5,	1,	0);
Two dimensonal	(0.00, 1.00, 2	2.00, 2.24 \					
	1.00, 0.00, 1	.00, 1.41 \					
	2.00, 1.00, 0	.00, 1.00 \					
	2.24, 1.41, 1	.00, 0.00);					

Indel: 2 except for extreme and polarised linear (4)



Comparing effects: scatterplots of distances





<ロト < 四ト < 三ト < 三ト

Equidistant	1.00				
1-D equal	0.85	1.00			
1-D extremes	0.66	0.93	1.00		
1-D polarised	0.79	0.97	0.88	1.00	
2-D	0.87	0.98	0.91	0.94	1.00



- There is a good deal of difference in the distances across the scores
- Systematic agreement at the extremes (very similar or very different)
- Equidistant substitution framework the most different
- Equidistant 1-D very like the 2-D (matrices are very similar)
- Clearly substitution costs matter!



- I now go on to compare the equidistant and 1-D linear frameworks in relative terms
- I consider pairs of sequences where the equidistant distance is relatively greater than the 1-D distance
- Then the converse, i.e., 1-D greater than equidistant, and
- finally sequences where the result is similar.
- What pairs does one see as close and the other different?



Equidistant relatively greater than 1-D





э

<ロト < 四ト < 三ト < 三ト

- The most obvious feature here is the predominance of long spells of adjacent states, particularly 1 with 2 and 3 with 4.
- Relatively little matching of 1 with 3 or 4, or 4 with 1 or 2 these are expensive for 1-D but not for Equdistant



Equidistant relatively less than 1-D





◆ロト < 団 ト < 巨 > < 巨 > < 巨 > 三 のへで 134

- By contrast, here the main feature is 1 alongside 4
- For 1-D these are very different states, so trajectories where they have to be compared are distant
- For Equidistant, they are no more expensive than any other pair
- Relatively long spells, so little room for alignment
- A rather obvious but important observation about how state spaces affect trajectory spaces



Equidistant close to 1-D





э

- The sort of trajectory pairs where the difference is less are composed more of short spells
- There are also relatively good matches (e.g., black with black)
- Alignment is more possible because of the many transitions
- Also, where there is a good level of exact match the substitution costs don't come into play so much



- Be explicit about state spaces and what distances mean
- Think spatially
 - Choose high or low dimensions, but have your reasons
- Simplify state space as far as possible
 - Drop irrelevant distinctions
 - Drop longitudinal information: let the sequence encode the temporal information, make state space cross-sectional



Dropping temporal information

• e.g., Simplify marital status:

	Living a	lone	Living with partner
legally married	Separated		Married
Not legally married	Single,	never	Cohabiting
	married,	post-	
	cohabitatior	٦,	
	divorced		

- The sequence will distinguish adequately between the various "single" states
- Parity sequences: Women's annual fertility history
 - in parity terms: 000112333344444
 - in birth event terms: 000101100010000



イロト イポト イヨト イヨト

- Finally, to adress *indel* costs
- As previously described, there is a specific lower limit and an empirical upper limit:

$$0.5 < \frac{indel}{\max(c_s)} < c. \, 1.5 \rightarrow 2$$

- Substitution costs greater than twice *indel* are ignored
- Raising *indel* costs to as little as twice the highest substitution cost will tend to prevent alignment
- Can we get a clearer idea of the impact?



- The other cost consideration is the *indel*
- It's bottom limit at 1.5 max substitution cost is clear
- What of the top limit? How high does it need to get to reduce OM to Hamming distance?
- With the same data set, I present the effect of varying the *indel* cost from 1.5 to 4.5 with the linear substitution cost



Varying indel and closeness to Hamming result



- The number of cases "misclassified" by CA relative to Hamming is somewhat chaotic up to about 2.5, and then falls sharply
- The correlation between OM and Hamming distances moves more steadily, hitting 1.000 at 3.5
- Note, though, the misclassification of 10% of the cases even though the correlation is 0.9995! CA can be funny.



- Substitution costs make a big difference
 - but largely understandable in operation
 - and an asset more meaningful state space, more meaningful trajectory space
- Think spatially! Use data and geometric models
- Simplify
- Let the sequence do the temporal work




- The main competitor to OM is from Cees H. Elzinga
- In a series of papers (Elzinga, 2003, 2005, 2009; Elzinga and Liefbroer, 2007; Elzinga et al., 2008), he proposes a number of related approaches with a different logical and mathematical underpinning
- In the tradition of Dijkstra and Taris (1995), he focuses on "the same states, in the same order"
- His novelty and power is to bring to bear set theory and combinatorics



- The intuition is that two sequences are more alike, the more they have the same states in the same order
- This explicitly brings the focus on sub-sequences (as distinct from substrings)
 - A subset of elements from the string
 - Not necessarily contiguous
 - But retaining the same order
- AB is both a substring and a subsequence of ABC
- AC is a subsequence but not a substring of ABC



- Enumerating subsequences is the key to this approach
- A number of measures are proposed including
 - The Longest Common Subsequence
 - The Number of Common Subsequences (count how many distinct subsequences both sequences have at least once)
 - The Count of Common Subsequences (for each shared subsequence, the sum of the product on the number in sequence 1 and the number in sequence 2)



- Unlike the Longest Common Prefix (the substring starting at position 1), the LCS can span the whole of the two sequences
- The measure of similarity is the length of the longest subsequence present in both
- Elzinga (2009) shows that under certain cost configurations, OM is equivalent to LCS



- Indentifying the longest common subsequence is intuitively clear, the other measures involve enumerating sequences – combinatorics
- For a sequence of length *l* there are 2^{*l*} combinations of its elements, from length 0 (the "null" sequence) to *l* (the sequence itself)
- Thus ABC has as subsequences
 - . (the null sequence)
 - A, B and C
 - AB, AC and BC, and
 - ABC
- Sequences with repetitions, like

AAC, have repetitions in the subsequences:

- . (the null sequence)
- A, A and C
- AA, AC and AC, and
- AAC



Number/Count of common subsequences

- The "Number of Common Subsequences" measure counts the number of distinct sequences which are subsequences of both sequences
- Elzinga (2009) points out that while it correlates highly with the LCS measure, they are distinctly non-monotone for sequences with moderate similarity
- NCS is not affected by how often matching subsequences occur, and if elements are repeated, subsequences can occur many times
- The "Number of Matching Subsequences" measure takes account of repetition: $\phi(x, y)$ is the sum, for each matching subsequence of $n_{sx} \times n_{sy}$ (i.e., the count of the subsequence in each sequence)



• As mentioned before, Elzinga (2009) asserts that where $\phi(x, y)$ measures the amount of a characteristic shared by *x* and *y*, and the following holds:

$$\phi(x, y) = \phi(y, x)$$
$$0 \le \phi(x, y) \le \min \{\phi(x, x), \phi(y, y)\}$$

then

$$d(x,y) = \phi(x,x) + \phi(y,y) - 2\phi(x,y)$$

is metric



э

- The problem with subsequences is that there are 2^{*l*} of them enumerating them is *O*(2^{*N*})
- Elzinga 2005 outlines an algorithm to enumerate common subsequences in a pair of sequences that is $O(l_1 \times l_2)$
- This is implemented in CHESA, downloadable from http://home.fsw.vu.nl/ch.elzinga/
- I have implemented a "brute-force" algorithm for Stata, which enumerates subsequences in a first pass, and then compares them in a rapid second pass – good for up to about 20 tokens in its current version



- The algorithm is the main difference from OM, but
- The absence of substitution costs is another important difference states are either the same (1) or different (0)
- Elzinga (personal communication) has outlined a measure which takes account of partial as well as complete similarity
- But what does the NMS measure look like in practice?



- First problem: 73 periods is far too many (for my algorithm at least)
- Solution: sample every 4th month to yield 19 tokens per 5 year sequence
- Number of matching sequence measure calculated and distance as

$$d(x,y) = \phi(x,x) + \phi(y,y) - 2\phi(x,y)$$



Cluster analysis

0

• Cluster analysis on the pairwise distances yields the following comparison with OM on the same sequences (8-cluster solutions)

M					NMS					
		1	2	3	4	5	6	7	8	
	+-									-+
1	L	252	0	0	0	0	9	0	0	Ι
2	Τ	0	27	0	0	0	51	0	0	Ι
3	L	0	0	19	0	0	23	0	0	Ι
4	L	0	0	0	50	0	19	0	0	Ι
5	L	0	0	3	0	10	38	0	0	Ι
6	L	0	0	0	0	0	27	0	0	Ι
7	L	0	0	0	0	0	2	15	0	Ι
8	L	0	0	1	0	0	2	0	127	Ι
++										



Eight-cluster solution



The "interesting" clusters



MDS solution: first two dimensions



- Spells, durations and SA
- SA in multiple domains
- TraMineR (including turbulence)



- Lifecourse data is usually spell structured a sequence of periods in a single state, with a given duration
- How to deal with in OM, which works with sequences of tokens?
- Treat spells as tokens, ignore duration?
- Represent time by multiplying tokens by spell duration?



- The latter approach is usual, but it this sociologically optimal?
- For instance, OM says AAAB is as distant from AACB as from AABB (given $\delta(A, B) = \delta(A, C)$)
- Substantively, the first and third are very similar while the second introduces a completely new spell
- Do we need an algorithm that is aware of spells?



- A simple but extreme strategy is to ignore duration
- Makes sequences of spells, ignoring length
- It works the sequences are still sequences, but now are unequal in length
- But duration is important: FFnnnnnnn and FFFnnnnnn are substantively closer to each other than to FFFFFFFnn but all reduce identically to Fn as spell sequences



- Another alternative is to represent events, not spells or monthly tokens,
- For instance, identify the first month of any spell by its state, and subsequent ones by a padding token (meaning no event)
- Optionally mark the final month of the sequence by its state
- FFFnnnnpppUUUFF would thus become F..n...p..U..FF
- If we weight substitutions involving the "padding" state very low, the comparison will focus on the transitions



Limits to cheap costs

• There is a limit to how cheap a transition can be

Non-metric					Metric					
(0,	1,	2,	3,	0.5 \	(0,	1,	2,	3,	1.5 \	
1,	0,	1,	2,	0.5 \	1,	0,	1,	2,	1.5 \	
2,	1,	0,	1,	0.5 \	2,	1,	0,	1,	1.5 \	
З,	2,	1,	0,	0.5 \	З,	2,	1,	0,	1.5 \	
0.5,	0.5,	0.5,	0.5,	0);	1.5,	1.5,	1.5,	1.5,	0);	

- $\delta(A, B)$ cannot be greater than $\delta(A, X) + \delta(X, B)$
- Non-metric substitution matrices may lead to non-metric trajectory spaces
- Therefore difficult to code a truly neutral space in OM
- But perhaps interesting as a general strategy



- Another approach is to represent spells non-linearly
- Represent spell duration as e.g., \sqrt{l}
- Thus the "cost" of deleting a unit is bigger in a short spell (the fixed cost of deleting a unit represents more time in a long spell)
- One way of implementing this is OMv, a relatively simple adaptation of the OM algorithm (Halpin, 2008)



- This approach means that the cost of edits to the token strings are sensitive to the length of the spell the token is in
- It produces distances not very different from OM unless there is a very high level of variation in spell length
- However, as currently implemented it is not a stable solution
 - Sequences composed of few, long spells are judged closer to all other sequences
 - Making for non-metric distances
- Potential solutions in scaling distances according to the most-dissimilar possible sequence



- OMv "warps time" by weighting it differently in different spells
- Harks back to Abbott's use of the term to suggest non-linear time scales (Abbott and Hrycak, 1990)
- In turn informed by Sankoff and Kruskal (1983), *Time Warps, String Edits and Macromolecules*



- Formally, time warping is a family of algorithms that do "continuous time-series to time-series correction" while OM *et al* do "string to string correction" (Marteau, 2007)
- Focus on comparing pairs of continuous-time high-dimensional time-series in \mathbb{R}^n
- Operates by locally compressing or expanding the time scale of one trajectory to minimise the distance to the other
- Distance is usually Euclidean in \mathbb{R}^n or other simple distance



1-dimensional time-warping



- TW used widely: was used for speech recognition, signature verification, other machine learning tasks
- Typically used to match a high-dimensional time-series to a "dictionary" of standard elements
- Conceptually it is a continuous time approach but implementations must be discrete sampling or periodic summaries:
 - e.g., sound at 44 kHz
 - rainfall daily
 - employment history monthly
- Kruskal and Liberman (1983) show that the continuous time logic can be faithfully implemented with discretised series





ABCCCCC



イロト イポト イヨト イヨト

• For cluster analysis and MDS, distance needs to be a "metric":

$$(R,S) = 0 \Rightarrow R = S$$

$$\Delta(R,S) \geq 0$$

$$(R,S) = \Delta(S,R)$$

- (a) $\Delta(R,T) \leq \Delta(R,S) + \Delta(S,T)$ (triangle inequality)
- Conventional TW satisfies 1-3, not 4, thus usually limited to matching against a "dictionary"



- Violation of the triangle inequality is due to TW usually having no cost to expansion or compression, only to the residual point-by-point distance
- Marteau (2007, 2008) proposes a TW algorithm that has a "stiffness" penalty
- Satisfies the triangle inequality
- Can be programmed very similarly to OM (recursive algorithm)
- Stiffness penalty like but not like *indel* cost squeezing/stretching, not inserting/deleting
- Point-to-point distance just like substitution



TW distance, $\delta(A^p, B^q) =$

$$\min \begin{cases} \delta(A^{p-1}, B^{q}) + d_{LP}(a_{p}, a_{p-1}) + \gamma d_{LP}(t_{a_{p}}, t_{a_{p-1}}) + \lambda \\ \delta(A^{p-1}, B^{q-1}) + d_{LP}(a_{p}, b_{q}) + \gamma d_{LP}(t_{a_{p}}, t_{b_{q}}) \\ \delta(A^{p}, B^{q-1}) + d_{LP}(b_{q}, b_{q-1}) + \gamma d_{LP}(t_{b_{q}}, t_{b_{q-1}}) + \lambda \end{cases}$$

(Marteau, 2007)



- Implemented as a Stata plugin
 - alongside similar implementations of OM and OMv
 - fast, comparable to OM plugin
 - but not platfrom independent
- NB experimental implementation, tentative results!
- Will be made available once stable



- Generally very close in pattern to OM correlation of 0.967 for $\gamma = 0.67$
- The higher γ is, the closer to OM...
- Up to a point: at $\gamma = 2$ and above, yields Hamming distance; i.e., "warping" completely suppressed
- Analogously, high indel costs suppress indels in OM...
- But at a rather higher threshold: warping and *indels* are not the same; *indel* costs and *γ* are not on the same scale relative to substitution costs



OM and timewarp distances



TW relatively larger than OM



э

<ロト < 四ト < 三ト < 三ト

TW relatively smaller than OM





4 ロ ト 4 回 ト 4 目 ト 4 目 ト 目 の Q (や)
180
- Trajectories with common spell sequences are closer under TW than OM
- Trajectories with no common spells tend be higher
- Confirmed in pair analysis where it is also clear that TW is better at seeing similarity where spells contain maximally different states, e.g., AAAAAAAADDD and AAADDDDDDDDD



- How to deal with spell data in Elzinga's paradigm?
- As in OM, repeat tokens to indicate duration works but unduly slow
- Treat spells as tokens (ignoring duration) works but ignores duration!
- Elzinga proposes a number of strategies for weighting spell-tokens according to duration:
 - Minimum shared duration: if the subsequences are *A*12/*B*4/*C*5 and *A*8/*B*5/*C*10 the match is weighted by 8 + 4 + 5
 - Duration product: the match is weighted by $12 \times 8 + 4 \times 5 + 5 \times 10$



- I have an experimental implementation of Elzinga's duration-weighted methods (Stata plugin)
- Uses a more naïve algorithm
- Approximates Elzinga's weighting: the contribution of a match is weighted by the product of the total time in each subsequence, thus $(12 + 4 + 5) \times (8 + 5 + 10)$ instead of $12 \times 8 + 4 \times 5 + 5 \times 10$



- How do combinatorial methods stack up against OM, OMv and timewarping?
- First, note that combinatorial methods have no "substitution costs": implicitly states are either the same or different
- For a fair comparison, we need a corresponding substitution matrix:

0	1	1	1
1	0	1	1
1	1	0	1
1	1	1	0



э

Distances including X/t (equidistant state space)



185

OM and X/t distances



م رم 186

OMv and X/t distances



187

- Handling trajectories through multiple domains simultaneously is very attractive
- Quite a few examples in the literature, often combining
 - Labour market
 - Housing
 - Partnership and family formation
- Dijkstra and Taris (1995) use as an example residential, educational and job status
- Pollock (2007) uses a similar trio



- The immediate difficulty is how to deal with multiple state spaces
- One solution is to create a combined space, crosstabulating the others, e.g.,:
 - employed-single
 - employed-partnered
 - not employed-single
 - not employed-partnered
- However, in practice this usually generates a high number of cells
- Practical problem of determining substitutions costs



イロト イポト イヨト イヨト

- It may be possible to set the costs on the cross-tabulated spaces *a priori*, using intuition or theory
- Sometimes it may also be possible to simplify the structure of this space: e.g., collapse certain regions into a single category
- If clear state-space structures exist for the sub-spaces it may be possible to combine them systematically:
 - Euclidean: $\sqrt{\sum_i (x_{ij} x_{ik})^2}$
 - Sum the different distances: $\sum_i |x_{ij} x_{ik}|$
- Weighting subdomains differentially is also possible



- An other option is to conduct parallel analyses in each domain
- This yields multiple distance measures
- Allows analysis of how the different domains cluster independently
- Perhaps less sensitive to coordination issues within lifecourses but should still be interesting



- A team in Switzerland are drawing on newer bioformatics technology: Multi-channel SA
- Explicitly deals with multiple parallel trajectories
- Gauthier et al. (2008) describe their "MCSA" method and claim it is superior to parallel OM analyses, and simpler than handling multiple-state distances
- Software available: http://www.tcoffee.org/saltt
- Not entirely clear from the description how the method works
- Bühlmann (2008) uses their methodology to examine careers of Swiss economists and engineers, using a number of categorical measures of their status



- Sequence complexity is extremely important, and how to classify complex sequences is a real problem
- Simpler sequences cluster well, no so complex ones
- How to measure this complexity?
- Elzinga (2008) proposes a spell-based, combinatorial definition, which he refers to as *turbulence* (see also Elzinga and Liefbroer, 2007)



- Turbulence is higher the more transitions there are
- And the more different states are entered,
- but lower when the variance of spell durations is higher
- Elzinga (2008) offers an efficient algorithm
- Implemented in TraMineR



- TraMiner (http://mephisto.unige.ch/traminer) is an R package for sequence analysis
- R is a free/open-source implementation of the S-Plus language (http://www.r-project.org)
- Available for most platforms (including Windows and Unix)



- Abbott, A. (1995a). A comment on "Measuring the agreement between sequences". *Sociological Methods and Research*, 24(2):232–243.
- Abbott, A. (1995b). Sequence analysis: New methods for old ideas. *Annual Review of Sociology*, 21:93–113.
- Abbott, A. and Forrest, J. (1986). Optimal matching methods for historical sequences. *Journal of Interdisciplinary History*, XVI(3):471–494.
- Abbott, A. and Hrycak, A. (1990). Measuring resemblance in sequence data: An optimal matching analysis of musicians' careers. *American Journal of Sociology*, 96(1):144–85.
- Abbott, A. and Tsay, A. (2000). Sequence analysis and optional matching methods in sociology. *Sociological Methods and Research*, 29(1):3–33.
- Aisenbrey, S. and Fasang, A. (2007). Beyond optimal matching: The 'second wave' of sequence analysis. Working Paper 2007-02, Center for Research on Inequalities and the Life Course, Yale University.
 Bühlmann, F. (2008). The corrosion of career? occupational trajectories of business economists and engineers in swspecified y European Sociological Review, 24(5):601–616.

- Degenne, A., Lebeaux, M.-O., and Mounier, L. (1996). Typologies d'itinéraires comme instrument d'analyse du marché du travail. Troisièmes journées d'études Céreq-Cérétim-Lasmas IdL, Rennes, 23–24 May 1996.
- Dijkstra, W. and Taris, T. (1995). Measuring the agreement between sequences. *Sociological Methods and Research*, 24(2):214–231.
- Elzinga, C. H. (2003). Sequence similarity: A non-aligning technique. *Sociological Methods and Research*, 32(1):3–29.
- Elzinga, C. H. (2005). Combinatorial representations of token sequences. *Journal of Classification*, 22(1):87–118.
- Elzinga, C. H. (2008). Complexity of categorical time series. under review.
- Elzinga, C. H. (2009). Sequence analysis: Metric representations of categorical time series. *Sociological Methods and Research*. (forthcoming).
- Elzinga, C. H. and Liefbroer, A. C. (2007). De-standardization of family-life trajectories of young adults: A cross-nation **Booiners** using sequence analysis. *European Journal of Population*, 23:225–250.

- Elzinga, C. H., Rahmann, S., and Wang, H. (2008). Algorithms for subsequence combinatorics. *Theoretical Computer Science*, 409(3):394–404.
- Gauthier, J.-A., Widmer, E., Bucher, P., and Notredame, C. (2008). Multi-channel sequence analysis applied to social science data. available at http://papers.ssrn.com/.
- Halpin, B. (2008). Optimal Matching Analysis and life course data: The importance of duration. Working paper 2008-01, Dept of Sociology, University of Limerick.
- Han, S.-K. and Moen, P. (1999). Work and family over time: A life course approach. *Annals of the American Academy of Political and Social Science*, 562:98–110.
- Kruskal, J. B. and Liberman, M. (1983). The symmetric time-warping problem. In Sankoff and Kruskal (1983), pages 125–161.
- Levine, J. H. (2000). But what have you done for us lately? Commentary on Abbott and Tsay. *Sociological Methods and Research*, 29(1):34–40.

Marteau, P.-F. (2007). Time Warp Edit Distance with Stiffness

Adjustment for Time Series Matching. *ArXiv Computer Science e-prints*.

Marteau, P.-F. (2008). Time Warp Edit Distance. ArXiv e-prints, 802.

- Pollock, G. (2007). Holistic trajectories: A study of combined employment, housing and family careers by using multiple-sequence analysis. *Journal of the Royal Statistical Society: Series A*, 170(1):167–183.
- Robette, N. and Thibauld, N. (2008). Comparing Qualitative Harmonic Analysis and Optimal Matching. *Population*, 63(4):533–556.
- Sankoff, D. and Kruskal, J. B., editors (1983). *Time Warps, String Edits and Macromolecules*. Addison-Wesley, Reading, MA.
- Wu, L. L. (2000). Some comments on "Sequence analysis and optimal matching methods in sociology: Review and prospect". *Sociological Methods and Research*, 29(1):41–64.

