Sequence analysis for social scientists	Sequence analysis for social scientists Session 1: Background
	Outline
Sequence analysis for social scientists Brendan Halpin, Dept of Sociology, University of Limerick Academica Sinica, Taipei, August 30-31 2016 Sequence analysis for social sciences Sequence Analysis in the social sciences Section 1	Session 1: Background Outline What is sequence analysis? Why it can be worth doing, and how it complements existing approaches How to do it, and how to think about it Practical, hands-on focus, using (<i>inter alia</i>) my SADI add-on for Stata (Halpin, 2014a) Slides available at http://teaching.sociology.ul.ie/taiwan Sequence analysis for social scientess Sequence Analysis in the social sciences Sequence Analysis in the social sciences
Sequence analysis in the social sciences: some background	 What is sequence analysis? Large and active research area From Andrew Abbott in mid-late 1980s, to 2015 special edition of <i>Sociological Methodology</i> Focuses on linear data (such as lifecourse trajectories) as <i>sequences</i>, as wholes Usually proceeds by defining distances between pairs of sequences, creating empirical typologies, etc
Sequence a na lysis for social scientists	Sequence analysis for social scientists
Sequence analysis for social scientists Session 1: Background Sequence Analysis in the social sciences	Sequence analysis for social scientists Session 1: Background Sequence Analysis in the social sciences
Sequence analysis for social scientists Session 1: Background Sequence Analysis in the social sciences A brief history of SA in Sociology	Sequence analysis for social scientists Session 1: Background Sequence Analysis in the social sciences A brief history of SA in Sociology
 Sequence analysis for social scientists Session 1: Background Sequence Analysis in the social sciences A brief history of SA in Sociology Andrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Forrest (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences 	 Sequence analysis for social sciences Sequence Analysis in the social sciences A brief history of SA in Sociology Andrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Hrycak (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences James Coleman: 'No one's gonna pay any attention as long as you write about dead German musicians' (Abbott, 2001, p. 13)
 Sequence Analysis for social sciences A brief history of SA in Socialogy And rew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Hrycak (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences 	 Sequence analysis for social sciences Bedrace Analysis in the social sciences A brief history of SA in Socialogy Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Hrycak (1990) - careers of Baroque musicians Abbott 's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences James Coleman: 'No one's gonna pay any attention as long as you write about dead German musicians' (Abbott, 2001, p. 13)
 Sequence analysis for social sciences Sequence Analysis in the social sciences A brief history of SA in Sociology Andrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Hrycak (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences 	 Sequence analysis for social sciences Sequence Analysis in the social sciences A brief history of SA in Sociology Andrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Hrycak (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences James Coleman: 'No one's gonna pay any attention as long as you write about dead German musicians' (Abbott, 2001, p. 13)
 Sequence analysis for social sciences A brief history of SA in Socialogy And rew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Hrycak (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences 	 Sequence analysis for social sciences Sequence Analysis in the social sciences A brief history of SA in Sociology Andrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Hrycak (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences James Coleman: 'No one's gonna pay any attention as long as you write about dead German musicians' (Abbott, 2001, p. 13)
 Sequence analysis for social scientists Sequence Analysis in the social sciences A brief history of SA in Sociology Andrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Forrest (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences Sequence analysis for social scientists Sequence Analysis in the social sciences	 Sequence analysis for social sciences A brief history of SA in Sociology Andrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Forrest (1990) - careers of Baroque musicians Abbott's main point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences James Coleman: 'No one's gonna pay any attention as long as you write about dead German musicians' (Abbott, 2001, p. 13) Sequence analysis for social sciences Some 1st Backgroud Some 1st wave adopters 2/2
 Sequence analysis for social scienciss A brief history of SA in Sociology A hold to a social scienciss A bott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Forrest (1986) - Morris dancing Abbott and Forrest (1986) - Careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Baroque musicians Abbott and Hycak (1990) - careers of Saroque musicians Abbott and Hycak (1990) - careers of services and a status-based and an achievement-based system, from 1890 to 1970. Wuerker (1996): Treats sequences of services interactions of mental health patients in Los Angeles. A small data set, but of interest because it uses a relatively uncommon form of trajectory. Halpin and Chan (1998): Analyses class careers of British and Irish men to age 35 using retrospective data. 	 Sequence analysis for social sciences A brief history of SA in Sociology A ndrew Abbott's long evangelism Abbott (1984) - earliest, argues for focusing on sequence as well as duration Abbott and Forrest (1986) - Morris dancing Abbott and Forrest (1986) - Morris dancing Abbott and Forrest (1986) - careers of Barcque musicians Abbott smain point: focus on sequences as wholes as an alternative to "variable-based" sociology However, his main practical contribution was to introduce the OM algorithm to the social sciences James Coleman: 'No one's gonna pay any attention as long as you write about dead German musicians' (Abbott, 2001, p. 13)

Sequence analysis for social scientists Session 1: Background	Sequence analysis for social scientists Session 1: Background
Sequence Analysis in the social sciences	Sequence Analysis in the social sciences
2000 debate in SMR	Key developments since
 Position: Abbott and Tsay (2000) Critiques: Levine (2000) and Wu (2000) is it sociologically meaningful? how do we parameterise it? does it have any advantages over conventional approaches? Response: Abbott (2000) 	 Widespread in many fields, especially lifecourse related: transition school to work, labour market, retirement, health outcomes, time use Some focus on multiple domains, dyadic approaches, cohort change in average diversity Much still uses clustering to develop empirical typologies See Aisenbrey and Fasang (2010) and Halpin (2013) for a summary Rather more activity in Europe than in US Two important conferences: LaCOSA1 2012 on Sequence Analysis: Blanchard et al. (2014) (includes historical demographers such as Michel Oris) LaCOSA2 2016 on Sequence Analysis and related methods (Online proceedings: https://lacosa.lives-nccr.ch/online-proceedings)
Sequence analysis for social scientists	Sequence analysis for social scientists
Session 1: Background Someone Analysis in the social sciences	Session 1: Background Why do Seguence Analysis?
Software developments	Why do Sequence Analysis?
Software developments	
 Abbott's optimize program Our own initial work used molecular biology software borrowed from the Oxford Dept of Pathology Götz Rohwer's TDA included an OM module later (mid-late 1990s) Stata: SQ and SADI (mid-late 2000s) R: Traminer (mid-late 2000s) 	 Why would we want to do it Holistic vs analytic? Exploratory vs hypothesis testing? Descriptive, visualisation Complexity of longitudinal processes hard to capture Complementary alternative to stochastic techniques which model data generation process
	(日) (西) (三) (三) (三) (1) (1) (1) (1) (1) (1) (1) (1) (1) (1
10. Sequence analysis for social scientists Session 1: Background	Sequence analysis for social scientists Session 1: Background
Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Seq uences are messy	Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Potentially complex processes
 Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Sequences are messy Lifecourse sequences are epiphenomena of more fundamental underlying processes The processes are potentially complex: difficult to predict distribution of sequences Other techniques (hazard rate models, models of late outcome using history, models of the pattern of transition rates) give a powerful but incomplete view SA clearly allows us visualise complex data; possibly allows us observe features that will otherwise be missed 	 Sequence analysis for social scientists Service Analysis for social scientists Service Analysis? Potentially complex processes individuals bring different characteristics from the beginning individuals bring different characteristics from the beginning history matters, including via duration dependence (individuals accumulate characteristics) time matters: calendar time (e.g. economic cycle), state distribution may change dramatically developmental time (maturation) processes in other lifecourse domains Too many parameters to model, hard to visualise distribution of life courses, also the possibility of <i>emergent</i> features Clear exploratory advantages possibility of detecting things that might not be detected otherwise
 Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Sequences are messys a Lifecourse sequences are epiphenomena of more fundamental underlying processes b The processes are potentially complex: difficult to predict distribution of sequences b Other techniques (hazard rate models, models of late outcome using history, models of the pattern of transition rates) give a powerful but incomplete view SA clearly allows us visualise complex data; possibly allows us observe features that will otherwise be missed 	 Sequence analysis for social scientists Service analysis for social scientists Service Analysis? Potentially complex processes a individuals bring different characteristics from the beginning bistory matters, including via duration dependence (individuals accumulate characteristics) time matters: a calendar time (e.g. economic cycle), state distribution may change dramatically developmental time (maturation) processes in other lifecourse domains Too many parameters to model, hard to visualise distribution of life courses, also the possibility of <i>emergent</i> features Clear exploratory advantages possibility of detecting things that might not be detected otherwise
 Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Sequences are messy Lifecourse sequences are epiphenomena of more fundamental underlying processes The processes are potentially complex: difficult to predict distribution of sequences Other techniques (hazard rate models, models of late outcome using history, models of the pattern of transition rates) give a powerful but incomplete view SA clearly allows us visualise complex data; possibly allows us observe features that will otherwise be missed 	 Sequence analysis for social scientists Servion 1: Background Why do Sequence Analysis? Potentially complex processes individuals bring different characteristics from the beginning history matters, including via duration dependence (individuals accumulate characteristics) time matters: calendar time (e.g. economic cycle), state distribution may change dramatically developmental time (maturation) processes in other lifecourse domains Too many parameters to model, hard to visualise distribution of life courses, also the possibility of <i>emergent</i> features Clear exploratory advantages possibility of detecting things that might not be detected otherwise
 Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Sequences are messy Lifecourse sequences are epiphenomena of more fundamental underlying processes The processes are potentially complex: difficult to predict distribution of sequences Other techniques (hazard rate models, models of late outcome using history, models of the pattern of transition rates) give a powerful but incomplete view SA clearly allows us visualise complex data; possibly allows us observe features that will otherwise be missed 	Sequence analysis for social scientists Servion 1: Background Why do Sequence Analysis? Potentially complex processes • The generating processes are complex: • individuals bring different characteristics from the beginning • history matters, including via duration dependence (individuals accumulate characteristics) • time matters: • calendar time (e.g. economic cycle), state distribution may change dramatically • developmental time (maturation) • processes in other lifecourse domains • Too many parameters to model, hard to visualise distribution of life courses, also the possibility of <i>emergent</i> features • Clear exploratory advantages • possibility of detecting things that might not be detected otherwise
 Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Sequences are messy Lifecourse sequences are epiphenomena of more fundamental underlying processes The processes are potentially complex: difficult to predict distribution of sequences Other techniques (hazard rate models, models of late outcome using history, models of the pattern of transition rates) give a powerful but incomplete view SA clearly allows us visualise complex data; possibly allows us observe features that will otherwise be missed 	Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Potentially complex processes • The generating processes are complex: • individuals bring different characteristics from the beginning • history matters, including via duration dependence (individuals accumulate characteristics) • time matters: • calendar time (e.g. economic cycle), state distribution may change dramatically • developmental time (maturation) • processes in other lifecourse domains • Too many parameters to model, hard to visualise distribution of life courses, also the possibility of <i>emergent</i> features • Clear exploratory advantages • possibility of detecting things that might not be detected otherwise Sequence analysis for social scientists Session 1: Background Non-holistic approaches
 Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Sequences are messy Lifecourse sequences are epiphenomena of more fundamental underlying processes The processes are potentially complex: difficult to predict distribution of sequences Other techniques (hazard rate models, models of late outcome using history, models of the pattern of transition rates) give a powerful but incomplete view SA clearly allows us visualise complex data; possibly allows us observe features that will otherwise be missed Sequence analysis for social scientists Sequence Analysis? Timing, sequence, quantum 	 Sequence analysis for social scientists Session 1: Background Why do Sequence Analysis? Potentially complex processes a individuals bring different characteristics from the beginning a individuals bring different characteristics from the beginning b history matters, including via duration dependence (individuals accumulate characteristics) a time matters:
 Sequence analysis for social sciences Lifecourse sequences are epiphenomena of more fundamental underlying processes The processes are potentially complex: difficult to predict distribution of sequences Other techniques (hazard rate models, models of late outcome using history, models of the pattern of transition rates) give a powerful but incomplete view SA clearly allows us visualise complex data; possibly allows us observe features that will otherwise be missed Sequence analysis for social scientists Sequence analysis for social scientists Sequence analysis for social scientists Sequence in what order do things happen Squantum: how much time is spent in different states (Billari et al., 2006) Many applications in longitudinal social science: annotated bibliography in Halpin (2013) 	 Securice analysis for social scientific Securical Science analysis? Potentially complex processes The generating processes are complex: individuals bring different characteristics from the beginning history matters, including via duration dependence (individuals accumulate characteristics) time matters: calendar time (e.g., economic cycle), state distribution may change dramatically developmental time (maturation)

Sequence analysis for social scientists Session 1: Background	Sequence analysis for social scientists Session 1: Background
Non-holistic approaches	Non-holistic approaches Trancition rate models
 For instance, summarise trajectories in terms of cumulative time in each state Typically use as a predictor (e.g., proportion of time unemployed predicting later ill-health) Or as an outcome: variables measured earlier (e.g., school performance) predicting proportion of time unemployed. 	 Model rates of period-to-period change: e.g., monthly movement between labour market statuses Model origin-destination patterns: e.g., transition between class at entry to labour market, and class at age 35 Markov models Very useful, good overview, can be descriptive or stochastic: tables make categorical data digestible Disadvantage: the focus on the t-1/t or t₀/t_T pattern means a loss of individual continuity Some potential to model longer Markov chains (Gabadinho, 2014)
Sequence analysis for Bockground Session 1: Background Non-holistic approaches	Sequence analysis for social sections: Session 1: Background Non-holistic approaches
Hazard-rate modelling	Latent class analysis
 Hazard-rate modelling is one of the dominant statistical alternative Either in terms of survival tables and curves (essentially descriptive) Or full stochastic models of the determinants of the hazard rate (Cox and/or parametric) Example: what characteristics speed up (or slow down) exit from unemployment? Very nice conceptual model of the temporal process Can test hypotheses Disadvantage: spell orientation, lack of whole-trajectory overview 	 Latent class growth curve models Where theory allows a developmental model of a quantitative outcome Account for the structure of repeated measurement of individuals Not so suitable for categorical variables Latent class models can be applied to careers However, difficult to properly incorporate the longitudinality Examples: Lovaglio and Mezzanzanica (2013); Barban and Billari (2012)
18	19.
Sequence analysis for social scientists Session 1: Background What we do with holizic approaches	sequence analysis for social scientists Sequence In Background What we do with Netricia commonster
Sequence analysis for social scientists Session 1: Background What we do with holistic approaches Holistic approaches	Sequence analysis for social scientists Session 1: Background What we do with holistic approaches Defining similarity
 Sequence analysis for social scientists Sequence analysis for social scientists Session 1: Background What we do with holistic approaches Holistic approaches by definition treat whole trajectories as units Classification of sequences is a typical goal Usually achieved by defining inter-sequence similarity and cluster analysis But other aspects of similarity may be interesting Variation of similarity by grouping variable (cohort, social class) Dyad similarity (couples' time use, mother-daughter fertility etc) Distance to pre-defined ideal types (empirical or theoretical) 	 Sequence analysis for social scientists Session 1: Background What we do with hoistic approaches Defining similarity the key challenge: must be efficient coherent, and sociologically meaningful We will consider a number of methods to do this Hamming distance Optimal Matching distance Time-warping measures Combinatorial subsequence measures
 Sequence analysis for social scientists Session 1: Background What do with holistic approaches Holistic approaches by definition treat whole trajectories as units Classification of sequences is a typical goal Usually achieved by defining inter-sequence similarity and cluster analysis But other aspects of similarity may be interesting Variation of similarity by grouping variable (cohort, social class) Dyad similarity (couples' time use, mother-daughter fertility etc) Distance to pre-defined ideal types (empirical or theoretical) 	 Sequence analysis for social scientists Session 1: Background What we do with holistic approaches Defining similarity the key challenge: must be efficient coherent, and sociologically meaningful We will consider a number of methods to do this Hamming distance Optimal Matching distance Time-warping measures Combinatorial subsequence measures
 Sequence analysis for social scientists Sexuion 1: Background What we do with holistic approaches Holistic approaches by definition treat whole trajectories as units Classification of sequences is a typical goal Usually achieved by defining inter-sequence similarity and cluster analysis But other aspects of similarity may be interesting Variation of similarity by grouping variable (cohort, social class) Dyad similarity (couples' time use, mother-daughter fertility etc) Distance to pre-defined ideal types (empirical or theoretical) 	 Sequence analysis for social scientists Session 1: Background Watwe do with holisist approaches Defining similarity the key challenge: must be efficient coherent, and sociologically meaningful We will consider a number of methods to do this Hamming distance Optimal Matching distance Time-warping measures Combinatorial subsequence measures
 Sequence analysis for social scientists Session 1: Background What we do with holistic approaches Holistic approaches Holistic approaches by definition treat whole trajectories as units Classification of sequences is a typical goal Usually achieved by defining inter-sequence similarity and cluster analysis But other aspects of similarity may be interesting Variation of similarity by grouping variable (cohort, social class) Dyad similarity (couples' time use, mother-daughter fertility etc) Distance to pre-defined ideal types (empirical or theoretical) 	Sequence analysis for social scientists Session 1: Background What we do with holistic approaches Defining similarity the key challenge: must be efficient coherent, and sociologically meaningful Ver will consider a number of methods to do this Hamming distance Optimal Matching distance Optimal Matching distance Time-warping measures Combinatorial subsequence measures Sequence analysis for social scientists Sequence analysis for social scientists Sequence analysis for social scientists
 Sequence analysis for social scientists Session 1: Background What we do with holsitic approaches Holistic approaches by definition treat whole trajectories as units Classification of sequences is a typical goal Usually achieved by defining inter-sequence similarity and cluster analysis But other aspects of similarity may be interesting Variation of similarity by grouping variable (cohort, social class) Dyad similarity (couples' time use, mother-daughter fertility etc) Distance to pre-defined ideal types (empirical or theoretical) Sequence analysis for social scientists Sequence analysis for social scientists Mand Hamming distance and Optimal Matching 	 Sequence analysis for social scientists Session 1: Background Wat we do with holistic approaches Defining similarity the key challenge: must be efficient coherent, and sociologically meaningful We will consider a number of methods to do this Hamming distance Optimal Matching distance Time-warping measures Combinatorial subsequence measures Sequence analysis for social scientists Session 1: Background Mamming distance example
Sequence analysis for social televistics Section 1: Background What we do with heliatic approaches Holistic approaches • Holistic approaches • Classification of sequences is a typical goal • Usually achieved by defining inter-sequence similarity and cluster analysis • But other aspects of similarity may be interesting • Variation of similarity by grouping variable (cohort, social class) • Dyad similarity (couples' time use, mother-daughter fertility etc) • Distance to pre-defined ideal types (empirical or theoretical) • Sequence analysis for social televistics • Steption 1: Background • Mad Hamming Hamming distance and Optimal Matching • The simplest way to compare sequences is element-wise • Given a rule for $d(a, b)$, project it onto $D(A, B)$ as $D(A, B) = \sum_i d(A_i, B_i)$ • Requires sequence of equal length • Hamming distance: recognises match or similarity at same time • Simple but important case of mapping $d(a, b) \rightarrow D(A, B)$	Sequence analysis for social scientists Service 1: Background What we do with holisis approaches Defining similarity the key challenge: must be • efficient • coherent, and • sociologically meaningful • We will consider a number of methods to do this • Hamming distance • Optimal Matching distance • Optimal Matching distance • Dynamic Hamming distance • Time-warping measures • Combinatorial subsequence measures • Combinatorial subsequence measures Sequence analysis for social scientists Sequence analysis for s

Sequence analysis for social scientists	Sequence analysis for social scientists
Session 1: Background OM and Hamming	Session 1: Background OM and Hamming
Hamming distance example	Optimal Matching
Calculate Hamming distance input \$1 \$2 \$3 \$4 \$5 1 2 3 2 3 2 3 2 3 1 4 2 3 2 3 1 1 1 1 1 end // Define the state differences matrix scost = (0,1,2,3 \ /// 1,0,1,2 \ /// 2,1,0,1 \ /// 3,2,1,0) hamming \$1-\$5\$, subs(scost) pwd(ham)	 Hamming recognises similarity at the same time If sequences have similarity that is out of alignment this will not be recognised OM defines similarity like Hamming, but uses insertion and deletion to allow sequences to align I.e., it cuts bits out in order to slide other parts along to match Insertion/deletion also enables comparison of sequences of different lengths Origins in computer science, pattern recognition, extensive use in molecular biology
(ロ)(例)(注)(注) 注) 23	<ロ>(日)(費)(注)(注)を、注 24
Sequence analysis for social scientists	Sequence analysis for social scientists
Session 13 Background OM and Hamming OM example	Session 1: Background OM and Hamming OM example
OMA call . oma s1-s5, subs(scost) indel(1.5) /// pwd(oma) length(5)	OMA call • OM distances . oma s1-s5, subs(scost) indel(1.5) /// pwd(oma) length(5) • Om distances • Hamming distances • Hamming distances symmetric ham[4,4] c1 c2 c3 c4 c1 c2 c3 c4 c1 c2 c3 c4 r1 0 r2 .6 0 r3 .6 .6 0 r4 1.2 1.2 1.8 0 • Hamming distances symmetric ham[4,4] c1 c2 c3 c4 c4 r1 0 r2 1.2 0 r3 .6 .4 0 r4 1.2 1.2 1.8 0 r4 1.2 1.2 1.8 0
2	23
Sequence ana lysis for social scientists Session 1: Background OM and Hamming	Sequence analysis for social scientists Session 1: Background OM and Hamming
OM vs Hamming	A more general example
 For most pairs the OM and Hamming distance is the same For the pairs (1,2) and (2,3), OM distance is less because "alignment" allows a better match 1 vs 2 Seq 1 1 2 3 2 3 - Seq 2 - 2 3 2 3 1 Cost i 0 0 0 0 i 	To convert ABCD into CDAAB the following set of operations gives the cheapest path: <u>Operation</u> Intermediate state Cost Sequence 2 <u>ABCD = 0</u> = = = =
• 2 vs 3 Seq 2 - 2 3 2 3 1 Seq 3 4 2 3 2 3 - Cost i 0 0 0 0 i 	= = Sequence 1 CDAAB =
Snamenen analveis for social scientists	Sequence analysis for social scientists
Session 1: Background OM and Hamming A more general example	Session : Background OM and Hamming A more general example
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 +1.5 = 1.5 = =	To convert ABCD into CDAAB the following set of operations gives the cheapest path: <u>Operation</u> Intermediate state Cost <u>Sequence 2</u> <u>ABCD = 0</u> insert C <u>CABCD +1.5 = 1.5</u> insert D CDABCD +1.5 = 3.0 = =













 Constrained and All and A	Sequence añalysis for social scientists Session 2: Doing SA Cluster ao lusie: e municical translovice frame distances	Sequence analysis for social scientists Session 3 Symmeriding sequences: Duration number of enable concerne
 Dien wich fecours eta, Harming and OM getrate quite air har reata: Hower, where they offer is to with more complex sequences Hower, where they offer is to with more complex sequences Hower, where they offer is to with more complex sequences An off-information of the sequences A complex book of sequences is not offer is and of the sequences of the sequ	Hamming and OM	Complexity of sequences
Image: Standard space (and space (a	 Often with lifecourse data, Hamming and OM generate quite similar results However, where they differ it is with more complex sequences 	 Complexity of sequences is relevant: more complex means less likely to be similar (and perhaps, similarity is more interesting) How to measure? Number of spells is part of it Also distribution of time A single long spell is the simplest sequence Many spells in many different states is very complex
Shannon Entropy In the proportion of membin in planetic distribution of the duration	Sequence a malysis for social scientists	Sequence analysis for social scientists
 Shannon Entropy information theory relates complexity to "entropy" More complex objects are harder to describe, cannot be compressed Shannon Entropy := - Σ_p [kg_p p where p₁ is the ground for man (kg₀ = 0). Takes account of downly of state but ABABAB counts as no more complex than AAABBB Takes account of downly of state but ABABAB counts as no more complex than AAABBB Takes account of downly of state but ABABAB counts as no more complex than AAABBB Takes account of downly of state but ABABAB counts as no more complex than AAABBB Takes account of downly of state but ABABAB counts as no more complex than AAABBB Takes account of downly is proposed that is more appropriate for spid data. It is based on durate organization of their durations to complex than abased of the variance of their durations to complex their space statements and their durations to complex their space statements and their durations to complex their space can we recommend as larg. The other "lobo out" thing to do with pairwes distances s multi dimensional scaling. The other "lobo out" thing to do with happense distances s multi dimensional scaling. The other "lobo out" thing to do with the pairwes distances s multi dimensional scaling. The other "lobo out" thing to do with happense distances	Session 3 Summarising sequences: Duration, number of spells, entropy	Session 3 Summarising sequences: Duration, number of spells, entropy
 Information theory relates complexity to "entropy" More complex objects are harder to describe, cannot be compressed Shone Bitroy: r = - Σ μ² log₂ μ where μ is the proportion of months in state l Takes account of diversity of state but ASABLE counts as no more complex than AAABUB Takes account of diversity of state but ASABLE counts as no more complex than AAABUB Takes account of diversity of state but ASABLE counts as no more complex than AAABUB Takes account of diversity of state but ASABLE counts as no more complex than AAABUB I be account of diversity of state but ASABLE counts as no more complex than AAABUB I be account of diversity of state but ASABLE counts as no more complex than AAABUB I be account of diversity of state but ASABLE counts as no more complex than AAABUB I be account of diversity of state but ASABLE counts as no more complex than AAABUB I be account of diversity of state but ASABLE counts as no more complex than AAABUB I be account of diversity of state but ASABLE counts as no more complex to be accounted diversity as a proportion of more data to the account of diversity is proported that is more appropriate for spell data I to only available in TaM meR I combines a measure based on the number of digiting the specific data to the specif	Shannon Entropy	Example: entropy
 Standard MDS uses principal component analysis Standard MDS uses principal component analysis 	 Information theory relates complexity to "entropy" More complex objects are harder to describe, cannot be compressed Shannon Entropy: ε = -∑p_i log₂ p_i where p_i is the proportion of months in state i Takes account of diversity of state but ABABAB counts as no more complex than AAABBB 	entropy state*, gen(ent) cd(pcd) nstates(4) nspells state*, gen(nsp) gen ent2 = ent*nsp/72 table g8, c(mean ent mean ent2 mean nsp) format(%6.3f)
	Sequence a na lysis for social scientists	Sequence analysis for social scientists
 Elzinga 's turbulence In Elzinga (2010) a measure of complexity is proposed that is more appropriate for spell data It is based on duration weighted spells, and on subsequence counting It combines a measure based on the number of distince subsequences, with a measure of the variance of their durations It is (only) available in TraMineR However, in practice the simpler Shannon entropy correlates highly with it Multi-dimensional scaling (optional) The other "obvious" thing to do with pairwise distances is multi-dimensional scaling The other "obvious" thing to do with pairwise distances is multi-dimensional scaling The network of distances implies a coherent space: can we re-construct if? Preferably with dimensions much less than number of sequences! Standard MDS uses principal component analysis 	Session 3 Summarising sequences: Duration, number of spells, entropy	Session 3 Summarising sequences: Duration, number of spells, entropy
 In Elzinga (2010) a measure of complexity is proposed that is more appropriate for spell data It is based on duration weighted spells, and on subsequence counting It combines a measure based on the number of distince subsequences, with a measure of the variance of their durations It is (only) available in TraMineR However, in practice the simpler Shannon entropy correlates highly with it However, in practice the simpler Shannon entropy correlates highly with it Most add whole durates (seque) Multi-dimensional scaling (optional) The other "obvious" thing to do with pairwise distances is multi-dimensional scaling The network of distances implies a coherent space: can we re-construct if? Preferably with dimensions much less than number of sequences! Standard MDS uses principal component analysis 	Elzinga's turbulence	Regular expressions
Sequence analysis for social scientists Sequence and matrix differ from colum names; ror name used) Obasical metric mitidiancial scaling Instance dimensional scaling The network of distances implies a coherent space: can we re-construct it? Preferably with dimensions much less than number of sequences! Standard MDS uses principal component analysis Standard MDS uses principal component analysis	 In Elzinga (2010) a measure of complexity is proposed that is more appropriate for spell data It is based on duration weighted spells, and on subsequence counting It combines a measure based on the number of distince subsequences, with a measure of the variance of their durations It is (only) available in TraMineR However, in practice the simpler Shannon entropy correlates highly with it 	 If sequences are represented as text, text-processing tools such as "regular expressions" can be used to sort between them Refer to lab notes for more details stripe state*, gen(seqst) list seqst in 1/5, clean count if regexm(seqst, "^A+\$") count if regexm(seqst, "^AAAAAA+.*DDDDDDD.*AAAAAAA.*\$") count if regexm(seqst, "AB.*AB")
Sequence analysis for social scientists Service 3 MDS and pairwise distances (optional) Multi-dimensional scaling (optional) • The other "obvious" thing to do with pairwise distances is multi-dimensional scaling • The network of distances implies a coherent space: can we re-construct it? • Preferably with dimensions much less than number of sequences! • Standard MDS uses principal component analysis	- ロ・ - (ラ・ (三・) - (ロ・ 	(ロ)(母)(そ)(そ)(そ)(日)(日)(日)(日)(日)(日)(日)(日)(日)(日)(日)(日)(日)
 MDS and pairwise distances (optional) Multi-dimensional scaling (optional) The other "obvious" thing to do with pairwise distances is multi-dimensional scaling The network of distances implies a coherent space: can we re-construct it? Preferably with dimensions much less than number of sequences! Standard MDS uses principal component analysis 	Sequence analysis for social scientists	Sequence analysis for social scientists
 Multi-dimensional scaling (optional) The other "obvious" thing to do with pairwise distances is multi-dimensional scaling The network of distances implies a coherent space: can we re-construct it? Preferably with dimensions much less than number of sequences! Standard MDS uses principal component analysis 	Session 3 MDS and pairwise distances (optional)	Session 3 MDS and pairwise distances (optional)
 The other "obvious" thing to do with pairwise distances is multi-dimensional scaling The network of distances implies a coherent space: can we re-construct it? Preferably with dimensions much less than number of sequences! Standard MDS uses principal component analysis At a second secon	Multi-dimensional scaling (optional)	Example
$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	 The other "obvious" thing to do with pairwise distances is multi-dimensional scaling The network of distances implies a coherent space: can we re-construct it? Preferably with dimensions much less than number of sequences! Standard MDS uses principal component analysis 	. mdsmat oma, dim(3) [row names of (dis)similarity matrix differ from column names; row names used) Classical metric multidimensional scaling dissimilarity matrix: oma Eigenvalues > 0 = 188 Number of obs = 940 Eigenvalues > 0 = 188 Nardia fit measure 1 = 0.75566 Retained dimensional = 3 Number of obs = 0.9322



	Substitution costs
Thinking about state spaces and distances	Transitions and substitutions
 Costs can be thought of as distances between states If state space is ℝⁿ, 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) 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, D(ABC, ADC) = f(d(B, D)) (similarity of states) transitions happen within sequences (movement between states)
<ロ> (日)	・ロ・・(費・・ミン・ミン 多 のQの 77
Sequence analysis for social scientists Session 3	Sequence analysis for social scientists Session 3
Substitution costs	Substitution costs Decentive transiton rates
 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 ℝⁿ, 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 usually 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) d(Con, Lab) > d(Con, LibDem) + d(LibDem, Lab) Procedure confuses party state space and voter characteristics Voter polarisation/loyalty is trajectory information, not state information Another type of problem: irrelevant distinctions can cause similar states to have low transition rates
Sequence analysis for social scientists Session 3	Sequence analysis for social scientists Session 3
Substitution costs Take "space" seriously	Substitution costs Looking at state spaces
• Very useful to think in spatial terms	 Two very simple state spaces: Single dimension, equally spaced:
 State space as efficient summary of clustered distribution in ℝⁿ 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 	• All states equidistant $-n-1$ dimensions $ \begin{array}{c ccccccccccccccccccccccccccccccccccc$
 State space as efficient summary of clustered distribution in ℝⁿ 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 	• All states equidistant $-n-1$ dimensions $ \begin{array}{c ccccccccccccccccccccccccccccccccccc$
 State space as efficient summary of clustered distribution in ℝⁿ 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 	• All states equidistant $-n-1$ dimensions $ \begin{array}{c} 0 & 1 & 1 \\ 2 & 2 & 1 & 0 \\ 3 & 2 & 1 & 0 \end{array} $ • All states equidistant $-n-1$ dimensions $ \begin{array}{c} 0 & 1 & 1 & 1 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{array} $
 State space as efficient summary of clustered distribution in Rⁿ 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 	• All states equidistant $-n-1$ dimensions $ \begin{array}{c} \hline 0 & 1 & 1 \\ \hline 2 & 1 & 0 & 1 \\ \hline 3 & 2 & 1 & 0 \end{array} $ • All states equidistant $-n-1$ dimensions $ \begin{array}{c} \hline 0 & 1 & 1 & 1 \\ \hline 1 & 0 & 1 & 1 \\ \hline 1 & 1 & 0 & 1 \\ \hline 1 & 1 & 0 & 1 \\ \hline 1 & 1 & 0 & 1 \end{array} $ Sequence analysis for social scientists Sequence analysis for social scientists Sequence analysis for social scientists
 State space as ethicient summary of clustered distribution in ℝⁿ 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 	• All states equidistant $-n - 1$ dimensions $ \begin{array}{c ccccccccccccccccccccccccccccccccccc$
 State space as efficient summary of clustered distribution in Rⁿ 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 Sequence analysis for social scientists Sequence analysis for social scientists State space structure passes through to trajectory space structure Distances between states clearly affect distances between trajectories containing high proportions of those states If d("A", "B") << d("A", "C") then D("AAAA", "BBB") will tend to be less than D("AAAA", "BBB") will tend to be aligned to match the DD if d("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 	• All states equidistant $-n - 1$ dimensions $\frac{0 1 1 1}{1 0 1}$ $\frac{0 1 1 1}{1 1 0 1}$ $\frac{1}{1 1 0 0}$ $\frac{1}{1 0 0}$ $\frac{1}$

nce analysis for s









Sequence analysis for social scientists Session 4	Sequence analysis for social scientists Session 4
SA and further analysis Beating cumulated duration	SA and further analysis
Logistic regression Logistic regression Log likelihood = -304,3196 working Coef. Std. Err. z P> z [95% Conf. Interval] cdi .0867982 .0303862 1.87 0.0610026585 .116386 cd2 .0448647 .0257586 1.74 0.0610026585 .116386 cd3 .0448647 .0257586 1.74 0.0610026585 .016386 cd4 0 (ontted) state48 Part time., -6551304 .4752833 -1.37 0.170 -1.583082 .2798214 Won-employed -1.42019 .7612029 -1.87 0.062 -2.91212 .0717405 Non-employed -1.92716 0.033 0.351 -1.52235 4.289906 3 .883836 1.427716 0.03 0.351 -1.52235 4.289906 3 .883836 1.427716 0.03 0.351 -1.52235 4.289906 4 1.40097 .515145 2.42 0.015 .28634 0.47583 3.097088 5 1.833173 .746099 2.19 0.029 .167768 2.553106 4 1.40097 .515145 2.42 0.015 .28634 0.47583 3.097088 5 1.833173 .746099 2.19 0.029 .167288 3.097088 5 1.633173 .74609 2.19 0.029 .167288 3.097088 5 1.633173 .74609 2.19 0.029 .167288 3.097088 5 1.633173 .74609 2.19 0.029 .1692883 3.097088 5 1.633173 .74609 2.19 0.029 .169288 3.09708 5 1.660866 .8701866 1.90 0.0560446455 3.36642 5 1.63578 1.23275 2.20 0.026 .4342298 5.55924 	 It may make sense to model with the MDS dimensions set matsize 1000 mdsmat pwd, dim(3) matrix dim=e(Y) svmat dim logit working cd* i.state48 dim* lrtest base
. irtest base Likelihood-ratio test LR chi2(7) = 21.78 (Asymption: base nested in .) Prob > chi2 ビッ (約028 さいくま) ま つので	・ロ・・通・・注・ 注 うくで
Sequence analysis for social scientists Session 4 SA and further analysis MDS dimensions and model	Sequence analysis for social scientists Session 4 SA and further analysis MDS correlated?
Logistic regression Humber of obs = 940 LG chi2() = 660.39 Prob > chi2 = 0.0000 Paeudo R2 = 0.5230 	. corr cd* dim* (obg=940)
Sequence analysis for social scientists Session 4 Discrepancy	Sequence analysis for social scientists Session 4 Discremency
Studer et al's discrepancy	Discrepancy and MVAD
 Studer et al. (2011) propose a method for treating distances matrices analogously to SS in regression and ANOVA The average distance to the centre of the whole matrix is the analogue of total sum of squares With a grouping variable, the distance to the centre for each groups is the residual sum of squares This allows a pseudo-R² and a pseudo-F test Permutation is used to approximate the sampling distribution of pseudo-F 	<pre>use mvad matrix md = (0, 1, 1, 2, 1, 3\ ///</pre>
(ロ) (長) (を) を の(0)	
Sequence a na lysis for social scientists Session 4 Discrepancy Discrepancy results	Sequence analysis for social scientists Session 4 Multichannel SA Multiple domains
. discrepancy funemp, dist(oma) idvar(id) niter(100) dcg(d2c) Discrepancy based R2 and F, 100 permutations for p-value	 Lifecourse analysis recognises the interrelatedness of domains Somewhat hard to handle in many approaches: a potential strength of SA? In practice, not very well developed; most research on single domains Some work (Dijkstra and Taris (1995), Pollock (2007), Gauthier et al. (2010))

Combined distance versus combining distances Combine by cross-stabulation Combined distance versus combining distances Combined problem of the single version Combined problem Combined probl	Sequence analysis for social scientists Session 4	Sequence analysis for social scientists Session 4
Controlled of status and decays and consistency of status and	Multichannel SA Combined distance versus combining distances	Multichannel SA Compline by cross-tabulation
Determining costs Solution of the second 	 How to proceed? Conduct parallel analyses and combine results? Combine domains into a single variable? The former is easy but will be less sensitive to the synchronisation of domains The latter involves a large state space and problem in defining distances However, better sensitivity to cross-domain features makes it attractive 	 The simplest approach is to create a new state space that is the cross-tabulation of the two (or more) domains This yields a large number of states, one for each combination How then to determine costs?
Intermining cols Implementation • S mp set strategy is to sum across the domaine • We take a simple case (four party levels and five employment strategy) is to reale the interaction or crosslabulation of the states • There must be just its can for imposing other patterns, for indexec. • We take a simple case (four party levels and five employment strategy) • We have two substitution out? • We have two aubitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average substitution cost? • equale by average average encrease a substitution cost? • extract state state state state state state statestaverage encresese states	Sequence analysis for social scientists	Sequence analysis for social scientists
Othermining costs Simplest strategy is to sum across the domains In short, $db_{1}^{0} = db_{1}^{0} + db_{1}^{0}$ We take a simple case (four party levels and five employment stratures) In short, $db_{2}^{0} = db_{1}^{0} + db_{2}^{0}$ The major is part from four imposing other patterns, for infrance, import a colling in the dura different substitution cost in the sum of the strate is the different substitution cost? equal be by more substitutis by more substitutis? equal be by more	Session 4 Multichannel SA	Session 4 Multichannel SA
 Simplest strategy is to sum across the domains In both, ddg = df = df, In both, ddg = df, In both, In both	Determining costs	Implementation
Served a whole for excidence and the control is a set of the formation of the control is a set of t	 Simplest strategy is to sum across the domains In short, d^{AB}_{ik,jl} = d^A_{i,j} + d^B_{k,l} There may be justification for imposing other patterns, for instance, imposing a ceiling changing d^A for certain values in domain B weighting the domains differentially Note that with two different substitution matrices it can be difficult to weight equally equalise by max substitution cost? equalise by average substitution cost weighted by occurrence in the data? 	 We take a simple case (four parity levels and five employment statuses) First step is to create the interaction or crosstabulation of the states // Reshape long to work on all months simultaneously reshape long parx emp, i(pid) j(month) // Create a variable that is the interaction of the two gen cross = emp+(parx-1)*5 // Verify the state interaction variable tab cross table parx emp, c(mean cross) // Back to wide, fix the variable order reshape wide parx emp* cross*
<pre>Stread d Middland JA Create the substitution cost matrix • We have two substitution cost matrices, 4x4 and 5x5: matrix spar = (0, 1, 2, 3, /// matrix semp = (0, 1, 2, 3, 3, ///</pre>	Sequence analysis for social scientists	Sequence analysis for social scientists
Create the substitution cost matrix • We have two substitution cost matrices, 4x4 and 5x5: matrix spar = (0,1,2,3) /// matrix semp = (0,1,2,3,3) /// 1,0,1,2,7/// 2,1,0,1,11 /// 3,2,1,0,11 /// 2,1,0,1,11 /// 3,2,1,0	Session 4 Multichannel SA	Session 4 Multichannel SA
 We have two substitution cost matrices, 4x4 and 5x5: matrix spar = (0,1,2,3) /// matrix semp = (0,1,2,3,3) /// 1,0,1,2,1/// 2,1,0,1,1/// 2,1,0,1,1/// 2,1,0,1,1/// 2,1,0,1,1/// 2,1,0,1,1/// 3,2,1,1,0) Both have a max of 3, otherwise perhaps divide each by its max 	Create the substitution cost matrix	Combine into 20x20
Sequence analysis for social scientists Session 4 Multichannel SA Sequence analysis for social scientists Symmetric mcsa[20,20] c1 c2 c3 c4 c5 c6 c7 c8 c9c10c11c12c13c14c15c16c17c18c19c20 r1 0 r3 2 1 0 r4 3 2 1 0 r6 1 2 3 4 4 0 r7 1 2 3 3 1 0 r6 1 2 3 2 1 0 r6 1 2 3 2 1 0 r6 1 2 3 3 1 0 r6 1 2 3 3 1 0 r6 3 2 2 1 2 2 2 1 0 r10 4 3 2 2 1 3 2 1 1 0 r10 4 3 2 2 1 3 2 1 1 0 r10 4 3 2 2 1 3 2 1 1 0 r10 4 3 2 2 1 3 2 1 1 0 r10 2 3 4 4 2 0 r10 2 3 4 4 2 0 r10 2 3 3 4 4 0 r11 2 3 4 4 5 5 1 2 3 4 4 0 r12 2 3 4 4 2 1 0 r13 4 3 2 3 3 3 2 1 2 2 2 1 0 r13 4 3 2 3 3 3 3 2 1 2 2 2 1 0 r13 4 3 2 0 3 4 3 0 r13 4 3 2 0 0 2 4 2 0 0 0 0 r13 4 3 2 0 0 2 4 0 0 0 0 0 0 r14 0 0 0 0 r15 0 0 0 0 0 r15 0 0 0 0 0 0 r15 0 0 0 0 0 r16 0 0 0 r17 0 0 r17 0 0 r17 0 0 r18 0 0 0 r10 0 r10 0 0 r10	 We have two substitution cost matrices, 4x4 and 5x5: matrix spar = (0,1,2,3\ /// matrix semp = (0,1,2,3,3\ /// 1,0,1,2\ /// 1,0,1,2,2\ /// 2,1,0,1\ /// 2,1,0,1,1\ /// 3,2,1,0) 3,2,1,0,1\ /// 3,2,1,1,0) Both have a max of 3, otherwise perhaps divide each by its max 	<pre>// Use Mata to combine the two matrices mata: spar = st_matrix("spar") semp = st_matrix("semp") // each element becomes a 5x5 block sparx = spar # J(1,5,1) # J(5,1,1) // replicate the 5x5 matrix 4x4 times sempx = semp for (i=2; i<=4; i++) { sempx = sempx, semp } sempx = sempx for (i=2; i<=4; i++) { sempxy = sempxysempx } // The combined matrix is the element-wise sum; return it from Mata to Stata st_matrix("mcsa", sempxy :+ sparx) end</pre>
Sequence analysis for social scientists Session 4 Multichannel SA Symmetric mcsa[20,20] c1 c2 c3 c4 c5 c6 c7 c8 c9c10c11c12c13c14c15c16c17c18c19c20 r1 0 r2 1 0 r3 2 1 0 r4 3 2 1 0 r5 3 2 1 1 0 r6 1 2 3 4 4 4 0 r7 2 1 2 3 3 1 0 r6 1 2 3 4 4 4 0 r7 2 1 2 3 2 1 0 r6 3 2 2 1 3 2 1 1 0 r6 3 2 2 1 3 2 1 1 0 r6 3 3 2 1 2 3 2 1 0 r10 4 3 2 2 1 3 2 1 1 0 r11 2 3 4 5 5 1 2 3 3 4 4 0 r11 2 3 4 5 5 1 2 3 3 4 4 0 r12 2 3 3 3 1 0 r13 4 3 2 3 3 3 2 1 2 2 2 1 0 r13 4 3 2 1 3 2 1 1 0 r13 4 3 2 2 1 3 2 1 1 0 r15 5 4 0 2 0 2 6 0 2 0 1 4 0 0 r15 5 5 4 0 2 0 2 0 2 0 2 0 1 0 0 r15 5 5 4 0 2 0 2 0 2 0 2 0 1 0 0 r15 5 5 4 0 2 0 0 2 0 2 0 0 1 0 0 r15 5 5 4 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0	<ロン(費)(注)(そ)(そ)(そ) 128	end 《 다 > 《 라 > 《 군 > 《 문 > 《 E = 》 《
Session 4 Multichannel SA The combined matrix Dyadic sequence analysis symmetric mcsa[20,20] c1 c2 c3 c4 c5 c6 c7 c8 c9c10c11c12c13c14c15c16c17c18c19c20 r1 0 r2 1 0 r2 1 0 r3 2 1 0 r6 1 2 3 4 4 0 cases r7 2 1 2 3 3 1 0 cases r7 2 1 2 3 3 1 0 cases r7 4 3 2 2 1 2 3 2 1 0 cases r7 4 3 2 2 1 2 3 2 1 0 cases r10 4 3 2 2 1 3 2 1 1 0 cases r11 2 3 4 5 5 1 2 3 4 4 0 cases r12 3 2 3 4 4 2 2 1 3 2 1 1 0 couple time-diaries r10 4 3 2 2 1 2 3 2 1 0 dot r11 2 3 4 3 2 3 3 3 2 1 2 2 2 2 1 0 dot r13 4 3 2 3 3 3 2 1 2 2 2 2 1 0 dot r13 4 3 2 3 3 3 2 1 2 2 2 2 1 0 dot r13 4 3 2 3 3 3 2 1 2 2 2 2 1 0 dot r13 4 3 2 3 3 3 2 1 2 2 2 2 1 0 dot r13 4 3 2 2 3 3 3 2 1 2 2 2 2 1 0 dot r13 4 3 2 2 3 3 3 2 1 2 2 2 2 1 0 dot r13 4 5 2 0 2 0 4 0 2 1 4 0 0 dot r14 5 4 0 0 2 3 4 4 0 0 dot r15 4 0 0 2 3 4 0 0 1 0 0 dot r14 5 4 0	Sequence analysis for social scientists	Sequence analysis for social scientists
Symmetric mcsa[20,20] Dyadic SA c1 c2 c3 c4 c5 c6 c7 c8 c9c10c11c12c13c14c15c16c17c18c19c20 0 r1 0 r2 1 0 r3 2 1 0 - r6 1 2 3 4 4 0 - r7 2 1 2 3 3 1 0 - r8 3 2 1 2 2 2 2 1 0 - r9 4 3 2 1 2 3 2 1 0 - r10 4 3 2 2 1 3 2 1 1 0 - r11 2 3 4 5 5 1 2 3 4 4 0 - r12 3 2 3 4 4 2 0 - r13 4 3 2 2 3 3 3 1 0 - r13 4 3 2 3 3 3 2 1 2 2 2 1 0 - r13 4 3 2 2 3 3 3 2 1 2 2 2 1 0 - r13 4 3 2 0 0 3 2 6 0 2 0 0 0 0 0 0 0 - r14 5 6 0 0 3 0 0 2 0 0 0 0 0 0 0 0 0 - r15 5 6 0 0 2 0 0 0 0 0 0 0 0 0 0 -	Sesson 4 Multichannel SA	Session 5 Dyadic sequence a nalysis
symmetric mcsa[20,20] c1 c2 c3 c4 c5 c6 c7 c8 c9c10c11c12c13c14c15c16c17c18c19c20 r1 0 r2 1 0 r4 3 2 1 0 r4 3 2 1 0 r5 3 2 1 1 0 r6 1 2 3 4 4 0 r7 2 1 2 3 3 1 0 r7 2 1 2 3 3 1 0 r8 3 2 1 2 2 2 1 0 r10 4 3 2 2 1 3 2 1 1 0 r10 4 3 2 2 1 3 2 1 1 0 r11 2 3 4 4 2 1 2 3 3 1 0 r12 3 2 3 4 4 2 1 2 3 3 1 0 r13 4 3 2 3 3 3 2 1 2 2 2 1 0 r13 4 3 2 3 3 3 2 1 2 2 2 1 0 r13 4 3 2 3 3 3 2 1 2 2 2 1 0 r13 4 5 2 0 2 0 4 0 0 4 0	The combined matrix	Dyadic SA
r14 0 4 5 2 3 4 5 1 2 3 2 1 0 r16 3 4 5 6 6 2 3 4 5 1 1 0 r17 4 3 4 5 5 1 2 3 3 4 4 0 r17 4 3 4 5 5 1 2 3 3 1 0 r18 5 4 3 4 5 2 3 2 1 0 r18 5 4 3 4 5 2 3 2 1 0 r19 6 5 4 3 4 5 2 3 2 1 2 3 2 1 0 r19 6 5 4 3 4 5 2 3 2 1 2 3 2 1 0 r19 6 5	symmetric mcsa[20,20] c1 c2 c3 c4 c5 c6 c7 c8 c9c10c11c12c13c14c15c16c17c18c19c20 r1 0 r2 1 0 r3 2 1 0 r4 3 2 1 0 r5 3 2 1 1 0 r6 1 2 3 4 4 0 r7 2 1 2 3 3 1 0 r7 2 1 2 3 3 1 0 r7 2 1 2 3 3 1 0 r7 3 2 1 2 2 2 1 1 r8 3 2 1 2 2 2 1 1 r9 4 3 2 1 2 2 2 1 0 r10 4 3 2 2 1 2 3 2 1 0 r11 2 3 4 5 5 1 2 3 4 4 0 r12 3 2 3 4 4 2 1 2 3 3 1 0 r13 4 3 2 3 4 4 2 1 2 3 3 1 0 r14 5 4 3 2 2 4 3 2 1 2 3 2 1 0 r15 5 4 3 3 2 4 3 2 2 1 3 2 1 1 0 r16 3 4 5 6 6 6 2 3 4 5 5 1 2 3 4 4 0 r17 4 3 4 5 5 5 3 2 1 2 3 4 4 0 r17 4 3 2 2 3 4 4 3 2 2 1 0 r14 4 5 4 3 2 3 4 4 3 2 2 1 3 2 1 0 r15 5 4 3 3 2 4 3 2 3 4 4 0 r17 4 3 4 5 6 6 6 2 3 4 5 5 1 2 3 4 4 0 r17 4 3 4 5 5 5 3 2 3 4 4 2 1 2 3 2 1 0 r16 3 4 5 5 5 3 3 2 4 3 2 2 1 3 2 1 0 r17 4 3 4 5 5 5 3 2 3 4 4 2 2 1 2 3 4 4 0 r17 4 3 4 5 5 5 3 3 2 4 3 2 2 1 0 r16 3 4 4 5 6 6 6 2 3 4 4 2 1 2 3 4 4 0 r17 4 3 4 4 5 5 5 3 3 2 4 2 2 1 3 2 1 0 r16 3 4 4 5 5 5 3 3 2 4 3 2 2 1 0 r17 4 3 4 4 5 5 3 3 2 4 3 2 2 1 0 r17 4 3 4 4 5 5 1 2 2 3 4 4 0 r17 4 3 4 4 5 5 3 3 2 4 2 2 1 3 2 1 0 r16 3 4 4 5 6 6 6 2 3 4 4 2 1 2 3 3 4 4 0 r17 4 3 4 4 5 5 5 3 3 2 4 2 2 1 3 2 1 0 r18 5 4 3 4 4 5 6 3 2 3 4 3 2 2 2 1 0 r19 4 3 2 2 1 0 r19 4 3 2 2 1 0 r10 7 10 7 10 r10 7 10 r10 7 10 7 10 r10 7 10 7 10 r10 7 10 7 10 r10 7 10	 SA typically uses all-pair-wise distances, or distance to special cases Dyadic SA is also useful: distance between a specific pair Couple time-diaries Couple labour market histories Mother-daughter fertility histories, etc.
	r19 6 5 4 3 4 5 4 3 2 3 4 3 2 1 2 3 2 1 0	

Sequence analysis for social scientists	Sequence analysis for social scientists
Session 5	Session 5
Dyadic sequence analysis	Dyadic sequence analysis
Research questions	Similarity and difference
 Allows testing hypotheses about dyadic similarity Are couples' time-use patterns or life-course histories aligned Are fertility patterns inherited? Under what conditions are dyadic distances smaller or larger? How do couples arrange joint lifecourses? 	 Couples may coordinate their lives under very different gender constraints Fertility patterns may be similar within the constraints of different cohort patterns of fertility The relationship between sequences may not be one of replication some daughters may completely reject their mother's fertility pattern
<ロト(書) (書) (言) (言) (言) (言) (言) (言) (言) (言) (言) (言	<ロ> (四) (古) (古) (古) (古) (古) (古) (古) (古) (古) (古
Sequence analysis for social scientists	sequence analysis for social scientists
Session 5	Session 5
Dyadic sequence analysis	Dyadic sequence analysis
Literature	Practical issues
 Off-scheduling (Lesnard, 2008) Dyadic in concept but actually creates combined sequences Robette et al. (2015): Mother-daughter labour market careers Fasang and Raab (2014): Intergenerational fertility; notes that focus on similarity ignores heterogeneity Raab et al. (2014): Jun 13 2015 15:18:18 Sibling dyads, fertility 	 We can calculate dyadic distances with standard software For efficiency it might better to just calculate dyads' distances But the cost of calculating all pairs is relatively small, and offers an advantage: Compare dyadic distances with distances to all others
Sequence analysis for social scientists	Sequence analysis for social scientists
Session 5	Section 5
Dyadic sequence analysis	Dyadic sequence analysis
Strategy: Begin with dyad-ordered data	Sort by types
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$
136	337
Sequence analysis for social scientists	Sequence analysis for social scientists
Session 5	Session 5
Dyadic sequence analysis	Dyadic sequence analysis
Submatrices	Submatrices
 Two submatrices, with distances from each mother to each daughter (and transpose) Distance from mother to her own daughter on diagonal (and transpose) Use distance from mother to all daughters to assess whether distance to own daughter is unusual 	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$

Dvadic sequence analysis

Extract diagonals and other information

- The main info is on the diagonals: the dyad distances (repeated across the two submatrices since distance is symmetric)
- Other summaries are also interesting
 - mean distance of each daughter to all mothers (and vice versa)
 - variance, standard deviation of this distance
 - z-score of dyad distance relative to all distances
 - rank of dyad distance compared with all distances

equence analysis for social scientists References

Fasang, A. and Raab, M. (2014). Beyond transmission: Intergenerational patterns of family formation among middle-class american families. *Demography*, 51(5):1703-1728. Gabadinho, A. (2014). Package 'pst'. probabilistic suffix trees and variable length Markov chains Technical report, CRAN.

- Gauthier, J. A., Widmer, E. D., Bucher, P., and Notredame, C. (2010). Multichannel sequence analysis applied to social science data. Sociological Methodology, 40(1):1-38.
- Halpin, B. (2013). Sequence analysis. In Baxter, J., editor, Oxford Bibliographies in Sociology. Oxford University Press, New York.
- Concepts of Less, Herroris, Version Volta, Sequence analysis tools for Stata. Working Paper WP2014-03, Dept of Sociology, University of Limerick, Ireland.
 Halpin, B. (2014b). Three anarctives of sequence analysis. In Blanchard, P., Bühlmann, F., and Gauthier, J.-A., editors, Advances in Sequence Analysis: Theory, Method, Applications, Springer, Berlin.
- Halpin, B. (2016). Cluster analysis stopping rules in stata. Working Paper WP2016-01, Department of Sociology, University of Limerick.
- Halpin, B. and Chan, T. W. (1998). Class careers as sequences: An optimal matching analysis of work-life histories. European Sociological Review, 14(2).
- 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. Lesarad, L. (2006). Optimal matching and social sciences. Document du travail du Centre de Recherche en Economie et Statistique 2006-01, Institut Nationale de la Statistique et des Études Économiques, Paris.
- Pans.
 Lesnard, L. (2008). Off-scheduling within dual-earner couples: An unequal and negative externality for family time. American Journal of Sociology, 114(2):447-90.
 Lesnard, L. (2010). Setting cost in optimal matching to uncover contemporaneous socio-temporal patterns. Sociological Methods and Research, 38(3):389-419.
 Lesnard, L. and de Saint Pol, T. (2009). Patterns of workweek schedules in France. Social Indicators Research, 93:171-176.

equence analysis for social scientists References

- Abbott, A. (1984). Event sequence and event duration: Colligation and measurement. Historical Methods, 17(4):192-204. Abbott, A. (2000). Reply to Levine and Wu. Sociological Methods and Research, 29(1):65-76. Abbott, A. (2001). Time Matters: On Theory and Method. University of Chicago Press, Chicago 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 optimal matching methods in sociology Sociological Methods and Research, 29(1):3-33. Solution and research, 29(1)=>>>. Alsenbrey, S. and Fasang, A. E. (2010). New life for old ideas: The 'second wave' of sequence analysis – bringing the 'course' back into the life course. Sociological Methods and Research, 38(3):420-462. Barban, N. and Bilair, F. (2012). Classifying life course trajectories: A comparison of latent class and sequence analysis. Journal of the Royal Statistical Society Series C, 61(5):765-764. Billari, F. C., Fürnkranz, J., and Prskawetz, A. (2006). Timing, sequencing and quantum of life course events: A machine learning approach. European Journal of Population, 22:37-65. Blair-Loy, M. (1999). Career patterns of executive women in finance: An optimal matching analysis. American Journal of Sociology, 104(5):1346-1397. mersen Journal of Sociology, 104(5):1340–1397.
 Blanchard, P., Bühlmann, F., and Gauthier, J.A., editors (2014). Advances in Sequence Analysis: Theory, Method, Applications. Springer, Berlin. 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. (2010). Complexity of categorical time series. Sociological Methods and Research, 38(3):463-481. ≣ •**ગ** α 14.0 Sequence analysis for social scientists Session 5 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. Lovaglio, P. G. and Mezzanzanica, M. (2013). Classification of longitudinal career paths. Quality and Quantity, 47 (2):989-1008. 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. 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. Raab, M., Fasang, A. E., Karhula, A., and Erola, J. (2014). Sibling similarity in family formation. Demography, 51(6):2127-2154.
 - Robette, N., Bry, X., and Éva Lelièvre (2015). A "global interdependence" approach to multidimensional sequence analysis. Sociological Methodology, Online advance copy.
 - Sankoff, D. and Kruskal, J. B., editors (1983). Time Warps, String Edits and Macromolecules. Addison-Wesley, Reading, MA.
 - Stovel, K. (2001). Local sequential patterns: The structure of lynching in the Deep South, 1882–1930. Social Forces, 79(3):843–880. Stovel, K., Savage, M., and Bearman, P. (1996). Ascription into achievement. American Journal of Sociology, 102(2):358-99.
 - Studer, M. and Ritschard, G. (2014). A comparative review of sequence dissimilarity measures. Working Paper 2014-33, LIVES, Geneva.
 - Studer, M., Ritschard, G., Gabadinho, A., and Müller, N. S. (2011). Discrepancy analysis of state sequences. Sociological Methods and Research, 40(3):471-510.
 - 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. Wuerker, A. (1996). The changing careers of patients which chronic mental illness: A study of sequential patterns in mental health service utilization. The Journal of Behavioral Health Services and Research, 23(4):458–470.