

Appendix:
Using Sequence Analysis to Understand Career
Progression: an Application to the UK House of
Commons

August 14, 2018

Optimal Matching

Optimal matching relies on the Needleman-Wunsch algorithm (Needleman and Wunsch 1970). By taking into consideration the costs set up by the researcher, this algorithm calculates the minimum number of substitutions and indels needed to align two sequences, namely to make them identical. This algorithm was originally introduced to calculate the largest number of amino acids of one protein than can be matched with those of another protein allowing for all possible substitutions and deletions in either sequence (Needleman and Wunsch 1970). In this Appendix we briefly illustrate how this algorithm works. Take as example the sequence number 1 and the sequence number 2 mentioned in the main text and take also the following substitution cost matrix:

$$\begin{array}{cc} & \begin{array}{ccc} BB & FB & GO \end{array} \\ \begin{array}{c} BB \\ FB \\ GO \end{array} & \left(\begin{array}{ccc} 0 & 1 & 2 \\ 1 & 0 & 2 \\ 2 & 2 & 0 \end{array} \right) \end{array}$$

Assume that the cost of substituting the element FB with the element BB is 1 and with the element GO the cost is 2. Also, assume that the substitution cost of an element with itself is equal to 0 and that substituting FB with BB is equally costly to substituting BB with FB , for instance. By considering these substitution costs, we calculate the distance between the sequences $BB - BB - BB - FB - FB - GO$ and $BB - BB - BB - GO - GO - GO$ with the following formula:

$$\begin{aligned} S(BB, BB) + S(BB, BB) + S(BB, BB) + S(FB, GO) + S(FB, GO) + S(GO, GO) = \\ 0 + 0 + 0 + 2 + 2 + 0 = 4 \end{aligned}$$

This is rather straightforward when we compare two sequences of identical length and

consider only substitution costs. In this case, only one alignment is possible. Yet, when we compare sequences of different lengths, different alignments are possible. To solve this issue Needleman and Wunsch (1970) introduce a way to find the optimal alignment between two sequences by using a series of pairwise comparisons between their single elements. For the sake of simplicity, we illustrate how to compare two sub-sequences of the one mentioned above, namely $FB - FB - FB - BB$ and $FB - FB - FB$, characterized by different length. In this vein, four different alignments are possible:

$$FB - FB - FB - BB$$

1. $FB - FB - FB - XX$
2. $FB - FB - XX - FB$
3. $FB - XX - FB - FB$
4. $XX - FB - FB - FB$

If we compare those possible alignments, we conclude that the first one is the optimal one being the distance equal to the cost of inserting a gap. The others are associated with distance indexes equal to the cost of inserting a gap plus the cost of substituting BB with FB . What follows is the illustration of how the algorithm works, as explained by Needleman and Wunsch (1970).¹ A matrix $S(i, j)$ is created recursively such that:

$$\min \begin{cases} S(i-1, j-1) + s(x_i, y_j) \\ S(i-1, j) + g \\ S(i, j-1) + g \end{cases}$$

The element $s(x_i, y_j)$ represents the substitution cost between the elements in the row i and in the column j and g is the cost of inserting a gap. We assume that the cost of

1. Differently from the original discussion we do not use similarities but distances between elements, based on substitution costs. In this vein, some parts of the discussion as well as of the formalization are modified accordingly, such as that the calculations at the basis of the algorithm in this work use the minimum and not maximum function. Indeed, in the original discussion the aim of the algorithm is to maximise the similarity between sequences while in this case it is to minimize their distance.

inserting a gap is 0.5 in this example. For the sake of this example it is logical that the cost of inserting a gap is half of the lowest substitution cost (which equals 1, in this example), in that a substitution is equivalent to a deletion of an element and the insertion of a gap. We start by compiling the matrix from the top-left cell which is associated the value 0 and from the values on the first row and the first column to which the values of insertion of gaps are associated. As clarified below, indeed, if the sequences are aligned the cell takes value 0 whereas along the first row and first column a movement rightwards or downwards refers to the insertion of a gap. For instance, $S(1, 5)$ is four steps on the right from $S(0, 0)$: this means it will assume the value of the insertion of four gaps. Once we set those cells up, all other cells are calculated through the algorithm.

Take the example of cell $S(2, 2)$, its value will be:

$$\min \begin{cases} S(1, 1) + s(FB, FB) \\ S(1, 2) + 0.5 \\ S(2, 1) + 0.5 \end{cases}$$

$$\min \begin{cases} 0 \\ 1 \\ 1 \end{cases}$$

By applying the same reasoning to all the other cells we obtain the following matrix:

		<i>FB</i>	<i>FB</i>	<i>FB</i>	<i>BB</i>
	$S(1, 1)0$	$S(1, 2)0.5$	$S(1, 3)1$	$S(1, 4)1.5$	$S(1, 5)2$
<i>FB</i>	$S(2, 1)0.5$	$S(2, 2)0$	$S(2, 3)0.5$	$S(2, 4)1$	$S(2, 5)1.5$
<i>FB</i>	$S(3, 1)1$	$S(3, 2)0.5$	$S(3, 3)0$	$S(3, 4)0.5$	$S(3, 5)1$
<i>FB</i>	$S(4, 1)1.5$	$S(4, 2)1$	$S(4, 3)0.5$	$S(4, 4)0$	$S(4, 5)0.5$

We now show how the algorithm works in an intuitive manner. A movement on the diagonal means that the elements are aligned whereas a movement rightwards or downwards

represents either an insertion of a gap or a substitution. Starting from $S(1,1)$ we move diagonally into $S(2,2)$: the latter contains the minimum value with respect to the other options, namely $S(1,2)$ and $S(2,1)$. Inserting a gap respectively in the column or in the row sequence at this stage is not an optimal move. We do the same until we get to cell $S(3,3)$ where we have two choices: either we proceed diagonally and we add a gap thus moving rightwards or we insert a gap and then we substitute the element BB with the element FB (the dash arrows). As we can see the second strategy is more costly: the distance index is 1.5. Contrariwise, the first strategy represents the optimal one in that it minimizes the distance index, namely 0.5.

Finally, we show how to compare two sequences which differ due to the different length of episodes, namely ‘chunks’ of sequences (and only the sequence length, as done above). These sequences will be used below as well, to compare sequence analysis and cluster analysis. We take the career paths of observations $id = 676$ and $id = 682$, respectively

$BB - BB - BB - BB - FB - FB - FB - BB - BB - BB - BB - BB - BB - BB$
 $BB - BB - FB - FB - FB - FB - FB - FB - BB - BB - BB - BB$

Remember that the cost of substituting the element FB with the element BB (and the other way round) is 1 and that the cost of inserting a gap or deleting an element is 0.5 (hence the indel costs, namely the combined costs of insertion and deletion, equals 1). First, try to replace all the elements FB in the second sequence and to add the necessary gaps to make the two sequences of the same length.

$BB - BB - BB - BB - FB - FB - FB - BB - BB - BB - BB - BB - BB - BB$
 $BB - BB - BB\cancel{FB} - BB\cancel{FB} - FB - FB - FB - BB\cancel{FB} - BB - BB - BB - BB -$
 $XX - XX$

The cost of this option is 4, namely the cost of substituting the element FB with BB three times, which is 3, plus the cost of adding two gaps, which is 1.

Instead, if we consider using only the insertion of gaps and the deletion of elements, our option will be

$BB - BB - BB - BB - FB - FB - FB - BB - BB - BB - BB - BB - BB - BB$
 $BB - BB - XX - XX - FB - FB - FB - \cancel{FB} - \cancel{FB} - \cancel{FB} - BB - BB - BB -$
 $BB - XX - XX - XX$

The total cost of this option is 4, namely the cost of five gap insertions and the cost of three element deletions. In this case, the two options are the same in terms of costs. Yet, it should be noted that in this example we set the cost of both the insertion of gaps and deletion of elements to 0.5 (making the indel cost equal to 1), which is half of the cost of substituting the element FB with the element BB . We would find different results, if the substitution costs were higher, which could be the case in this example if the second sequence contained some GO elements (as shown in the matrix above, substituting BB with GO is more costly than substituting BB with FB). In conclusion, we have shown that in case of two sequences which differ due to the length of episodes, the optimal alignment depends on the cost of substitutions and indels.

TraMineR

For the analysis in this paper, we use the TraMineR R package. This package provides a series of functions which allow mining and visualizing sequences, calculating distance matrices, deriving sequence clusters and running basic inferential statistics, such as the discrepancy analysis, which is used in the paper.²

In the analysis in the paper, we use the ‘TRATE’ (transition rates) method to calculate the substitution matrix, where the substitution costs between two elements depend on the transition rates between the two elements, namely the probabilities of transition from one element to the other in the dataset (Studer et al. 2011). This is common practice where there is no theoretical motivation to assign specific substitution costs. As for missing values in sequences, we choose not to treat missing values as separate elements. Treating missing

2. For more information, see <http://traminer.unige.ch/index.shtml>

values as separate elements would be useful if missing values were present inside the sequences and if the absolute meaning of time points (in this case, the specific years when the elements are measured) was of interest for the analysis. This is not the case in our analysis. Moreover, to compute pairwise distances between sequences we use Optimal Matching, by relying on the substitution matrix. In computing the distances, we leave indel costs at the default setting of 1 (Studer et al. 2011). Also, we cluster sequences based on the distance matrix using Ward (hierarchical) clustering.³ Finally, as a robustness test, we run a multi-factor analysis of variance using the distance matrix with a specific function provided by the TraMineR package.

Cluster Analysis v. Sequence Analysis

In this section we provide a simple comparison between cluster analysis (used by itself) and sequence analysis (in conjunction with cluster analysis). As it is clear from the paper, one of the main steps in sequence analysis is the clustering of the sequences based on the dissimilarity/distance matrix obtained with optimal matching. In this vein, in order to show the validity of sequence analysis, we compare the results of the clustering using different dissimilarity matrices, obtained respectively with optimal matching (sequence analysis) and the Euclidean distance (the default option in cluster analysis). To be consistent, in both cases we use Ward (hierarchical) clustering.

Optimal matching is explained in detail above. The Euclidean distance is

$$\sqrt{\sum (x_i - y_i)^2}$$

It should be noted that, in order to perform the Euclidean distance on the variables measuring the different stages of the sequences, we assume that these variables are interval

3. Ward’s method starts with each point (sequence, in this case) being in its own cluster and keeps merging cluster by minimizing the ‘merging costs’, which in turn depend on the distance between (the center of the) clusters and the number of points in those clusters (Mirkin 2013).

level (i.e. $BB = 0$, $FB = 1$ and $GO = 2$) and we assign a distinct value to missing data (i.e. they are included in the analysis).⁴

We calculate the Euclidean distance across the variables measuring the different stages of the sequences, we obtain a dissimilarity/distance matrix and then we cluster the sequences using this matrix. To be consistent with what we do in the paper, we choose eight clusters and then we compare the results. The variables resulting from the two clustering procedures are highly correlated, with p value lower than 0.01. Yet, by looking at the data, we find that sequence analysis provides a more fine-grained measure of sequences, which take into consideration the duration of career paths and the order of stages of a career within a path.

Take the career paths of observations $id = 508$ and $id = 281$ (where $BB = 0$, $FB = 1$ and $GO = 2$), respectively

0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 1 – 1 – 1 – 1 – 1 – 1 – 1 – 1
1 – 1 – 1 – 1 – 1 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0 – 0

The hierarchical clustering with the default option of the Euclidean distance groups these two observations into the same cluster ('Cluster 3'), whereas the clustering with the dissimilarity matrix calculated with optimal matching groups them into different clusters (respectively, 'Career 4' and 'Career 6'). The traditional clustering procedures does not take into consideration the order of the elements in a sequence, whereas sequence analysis does. Especially when studying career paths, the order of elements is a rather important piece of information. The first MP above started in the backbench and moved his/her way up to the top, whereas the second MP started in the frontbench and then, for some reasons, he/she moved to the backbench. By relying on this information, it might be concluded that, for instance, the first MP has a brilliant career ahead, whereas the second MP has somehow burnt out and his/her political career is dead.

Also, take the career paths of observations $id = 676$ and $id = 682$ (where $BB = 0$, $FB = 1$

4. We assign a numerical value to missing values only in the cluster analysis used by itself. In the sequence analysis in conjunction with cluster analysis used in the main text, we drop missing values, as explained above.

and $GO = 2$; these sequences are reported below without missing values for illustrative purposes), respectively

0 – 0 – 0 – 0 – 1 – 1 – 1 – 0 – 0 – 0 – 0 – 0 – 0 – 0

0 – 0 – 1 – 1 – 1 – 1 – 1 – 1 – 0 – 0 – 0 – 0

The hierarchical clustering with the default option of the Euclidean distance groups these two observations into the same cluster (‘Cluster 6’), whereas the clustering with the dissimilarity matrix calculated with optimal matching groups them into different clusters (respectively, ‘Career 6’ and ‘Career 3’). The traditional clustering misses an important piece of information, at least for the study of career paths: the duration of episodes (i.e. chunks of the sequence). Both MPs started in the backbench, made their way up to the top and then got back to the backbench. Yet, the second MP spent more time in the frontbench than the first MP. This is an important distinction in the study of career paths. This might suggest, for instance, that as the second MP has more experience in the cabinet, maybe he/she is more likely to get back to power soon.

In conclusion, sequence analysis (by which we mean here the use of optimal matching to create a dissimilarity matrix to be used in cluster analysis) provides a more fine-grained measure of sequences. Differently from cluster analysis used alone, sequence analysis takes into consideration the order of elements in a sequence, namely whether the MP’s career is upward or downward, in this case, and the duration of episodes, namely whether the MP has spent more or less time in the frontbench, in this case. These two pieces of information are rather important in the study of political careers.

Descriptive Statistics

Table A1: University Education

University Degree	Freq.	Percent	Cum.
No University	83	11.42	11.42
Undergraduate	407	55.98	67.40
Postgraduate	237	32.60	100.00
Total	727	100.00	

Table A2: Legislature of Entry

Legislature	Freq.	Percent	Cum.
1997-2001	280	38.51	38.51
2001-2005	90	12.38	50.89
2005-2010	123	16.92	67.81
2010-2015	234	32.19	100.00
Total	727	100.00	

Table A3: Gender, State School, Job and Age at Entry

Variable	Obs	Mean	Std. Dev.	Min	Max
Gender	727	.2654746	.4418896	0	1
State School	697	.6140603	.4871661	0	1
Job	727	.2984869	.4579093	0	1
Age at Entry	725	43.27862	8.099836	25	69

Robustness Checks

Table A4: Clustering

Cluster Numbers	R Squared
Eight	0.81368903
Seven	0.80837233
Six	0.76708005
Five	0.68583622
Four	0.62082634

Figure A1: State Distribution Plot - Seven Clusters

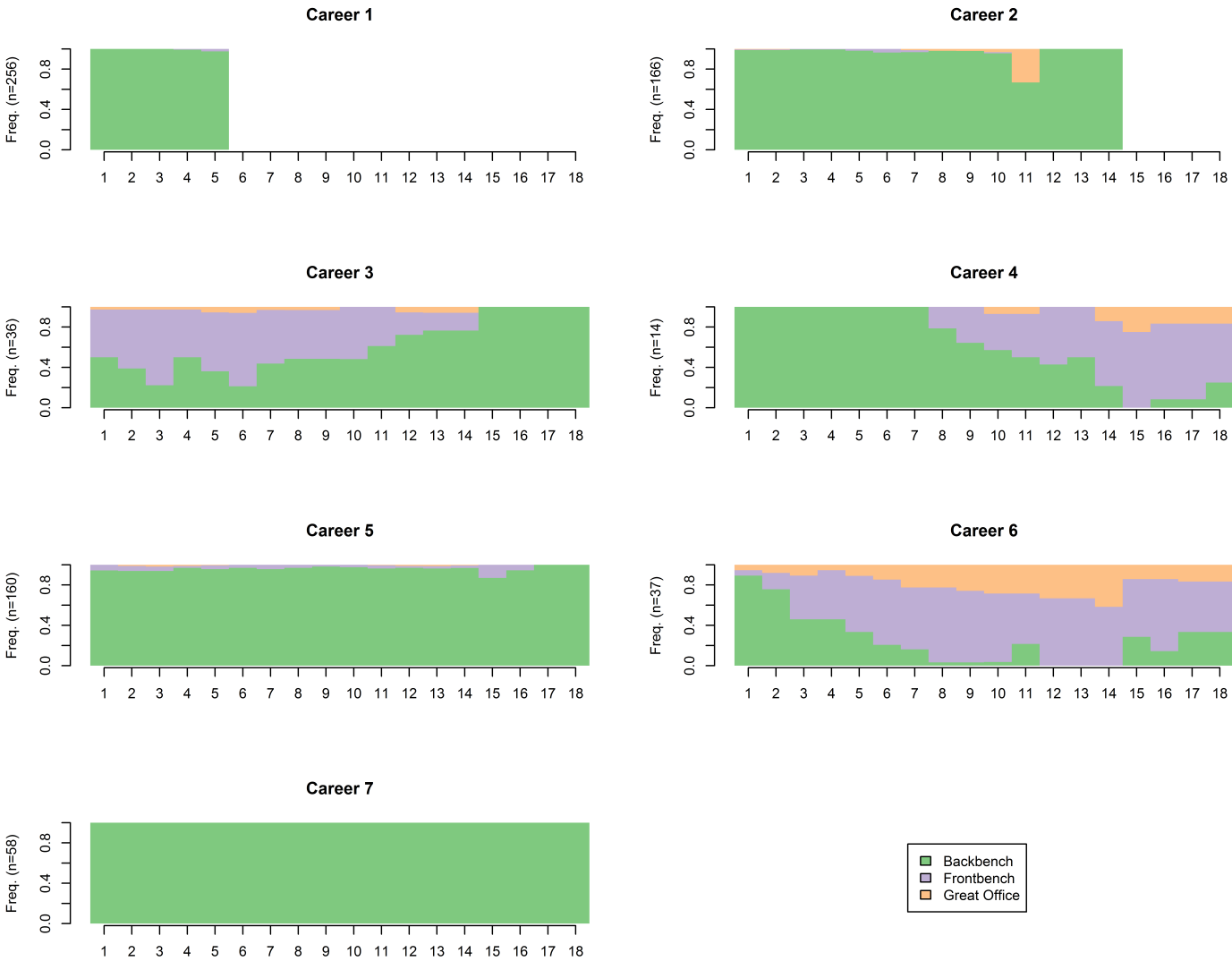


Figure A2: State Distribution Plot - Six Clusters

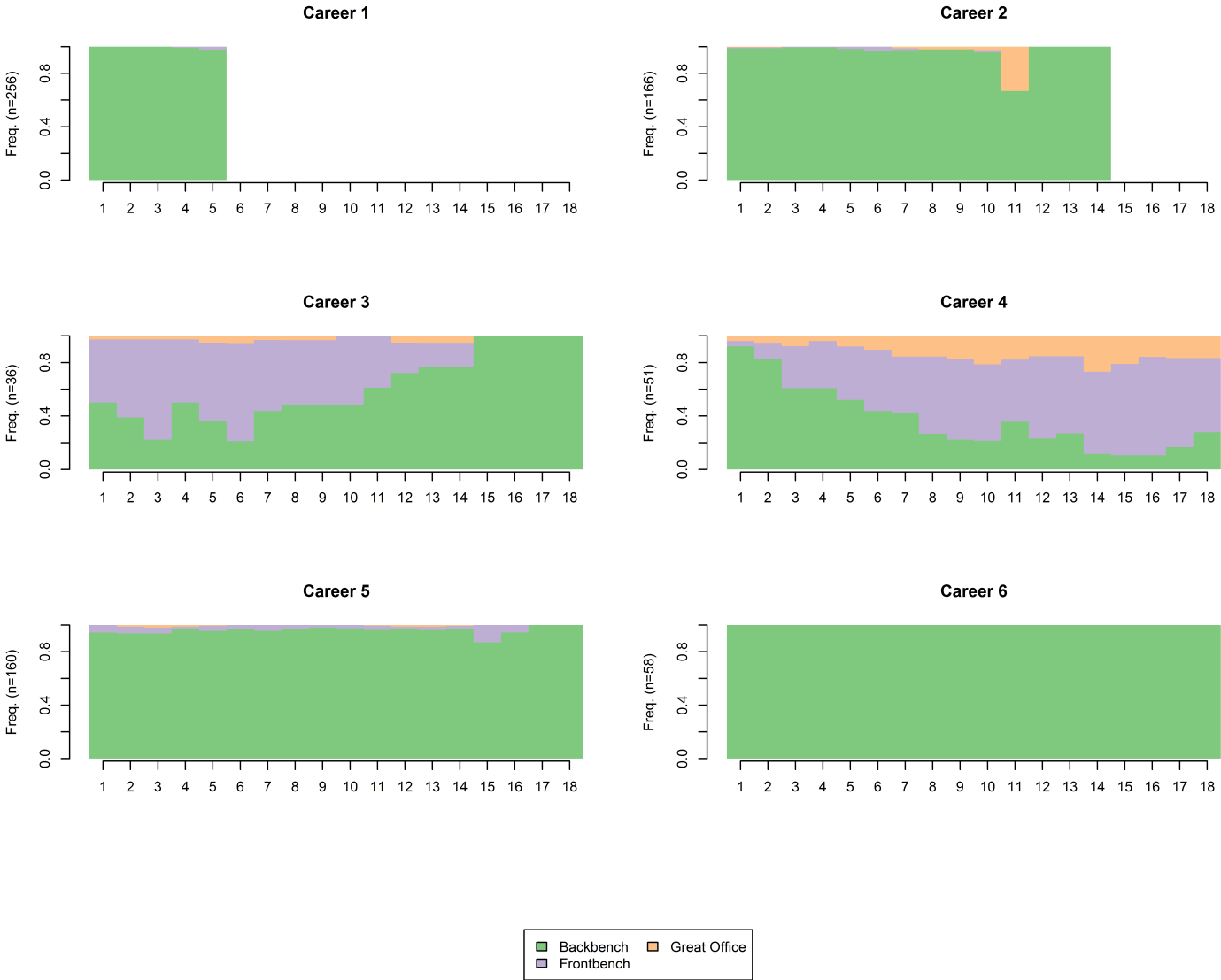


Figure A3: State Distribution Plot - Five Clusters

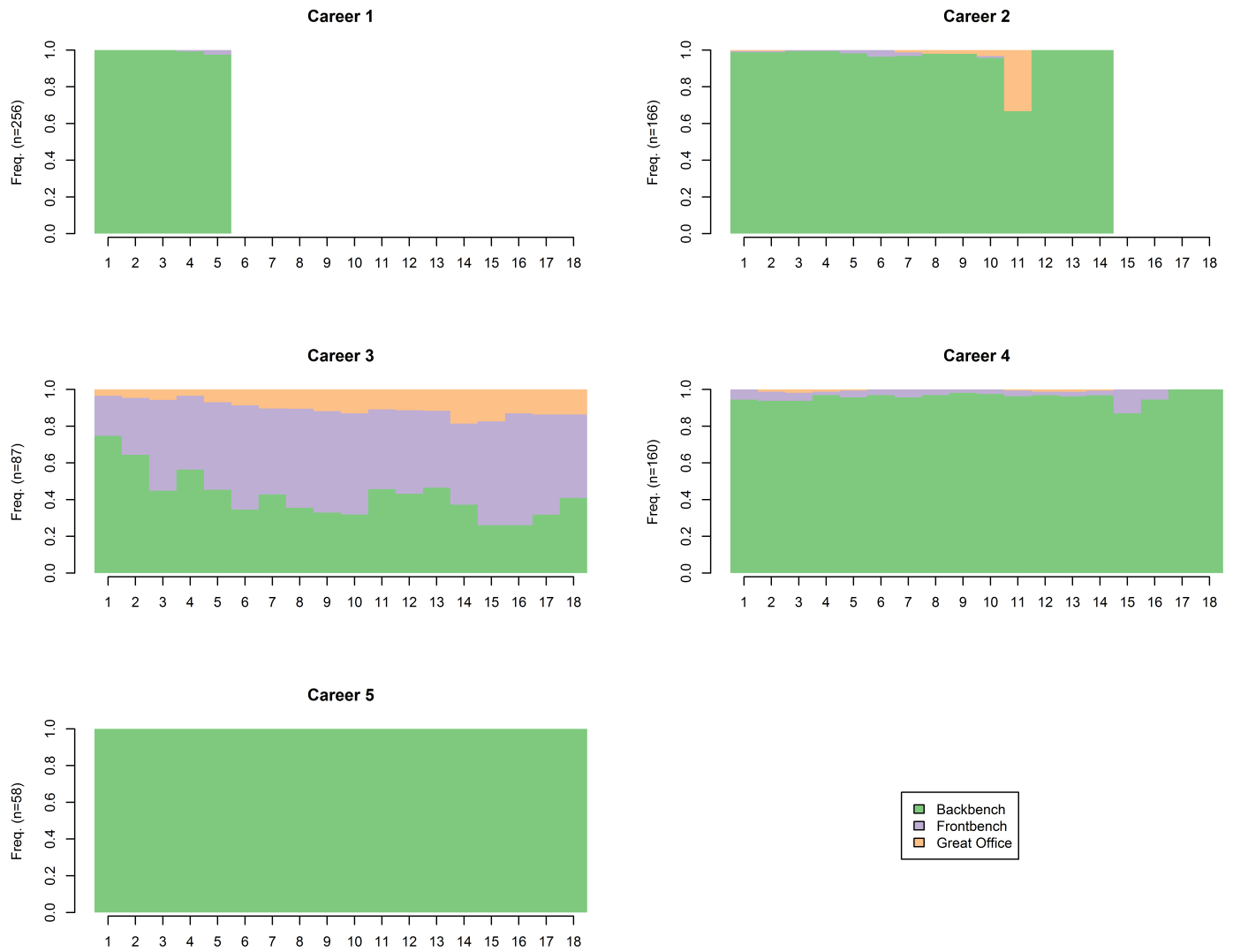


Figure A4: State Distribution Plot - Four Clusters

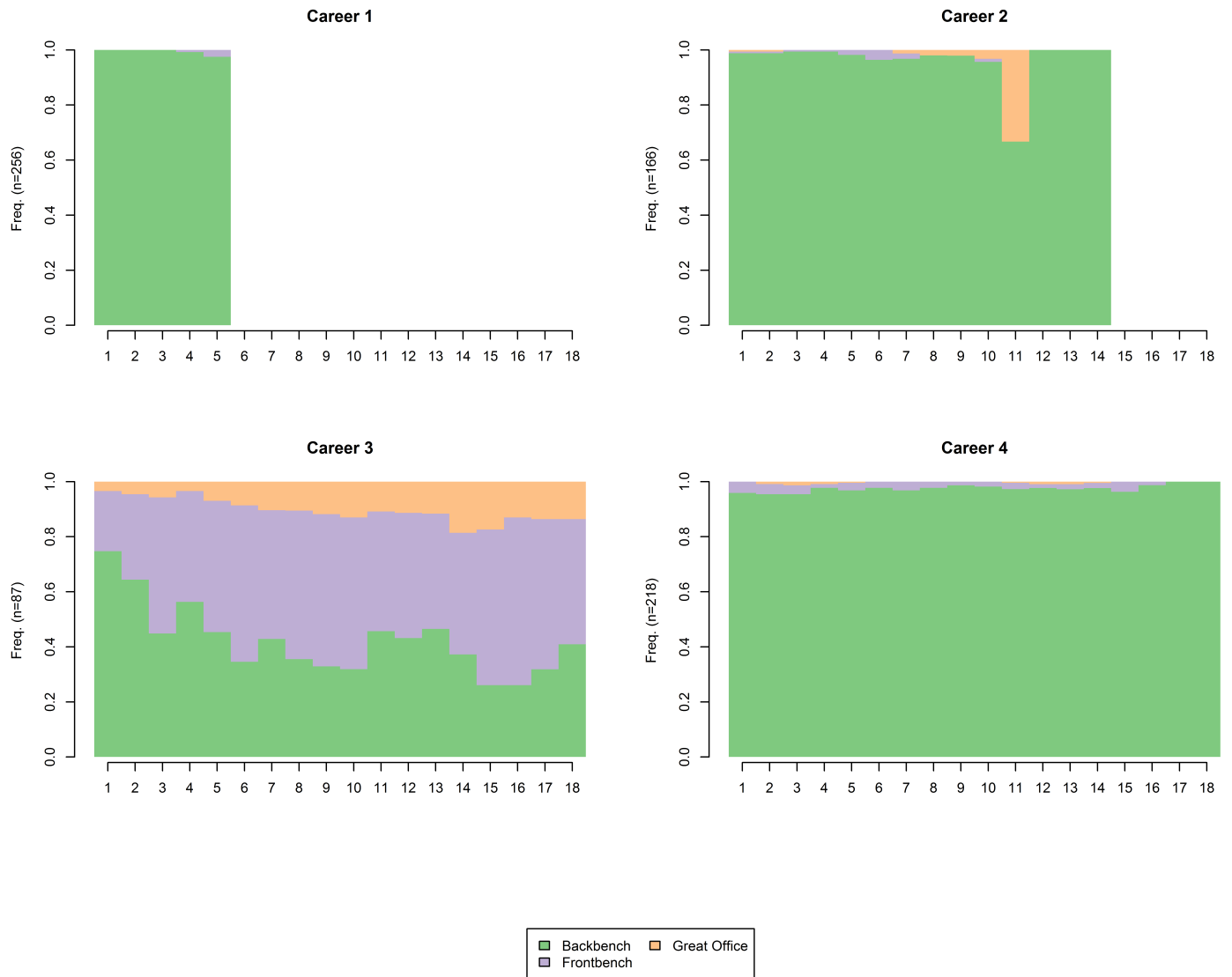


Table A5: Multinomial Analysis - Six Clusters

	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6
Gender	-0.225 (0.239)	-0.420 (0.440)	-0.0326 (0.360)	-0.0326 (0.360)	-0.486** (0.245)	-0.535 (0.361)	-0.523** (0.254)	-0.323** (0.254)	0.0434 (0.498)	-0.118 (0.406)	-0.941*** (0.267)	-1.012*** (0.380)	-0.367 (0.544)	-0.367 (0.544)	0.432 (0.682)	0.216 (0.579)	-0.225 (0.514)	-0.612 (0.598)
State School	0.108 (0.215)	-0.0941 (0.383)	-0.0884 (0.329)	-0.0884 (0.329)	0.467** (0.221)	-0.0459 (0.308)	-0.287 (0.237)	-0.287 (0.237)	-0.393 (0.461)	-0.477 (0.356)	-0.139 (0.241)	-0.508 (0.331)	-0.721 (0.475)	-0.721 (0.475)	-0.821 (0.603)	-0.791 (0.512)	-0.179 (0.447)	-0.549 (0.523)
Job	0.0206 (0.227)	0.198 (0.369)	-0.153 (0.348)	-0.153 (0.348)	0.0356 (0.232)	0.0907 (0.328)	-0.117 (0.231)	-0.117 (0.231)	0.234 (0.487)	-0.277 (0.371)	-0.174 (0.255)	-0.126 (0.331)	0.734 (0.553)	0.734 (0.553)	0.915 (0.706)	0.416 (0.598)	0.552 (0.513)	0.738 (0.506)
Age at Entry	0.00505 (0.0138)	0.0182 (0.0256)	-0.0835*** (0.0232)	-0.0835*** (0.0232)	0.0190 (0.0128)	-0.00591 (0.0163)	0.00265 (0.0148)	0.00265 (0.0148)	0.00290 (0.0256)	-0.0622*** (0.0224)	0.0157 (0.0144)	-0.00973 (0.0177)	-0.0420 (0.0382)	-0.0420 (0.0382)	-0.0478 (0.0422)	-0.137*** (0.0419)	-0.0377 (0.0359)	-0.0859** (0.0382)
University Undergraduate	0.133 (0.316)	0.712 (0.633)	0.555 (0.648)	0.555 (0.648)	0.513 (0.353)	-0.0472 (0.475)	0.166 (0.328)	0.166 (0.328)	0.949 (0.739)	0.578 (0.646)	0.549 (0.366)	-0.0326 (0.487)	-0.458 (0.758)	-0.458 (0.758)	0.404 (1.001)	0.0717 (0.894)	-0.0438 (0.773)	-0.582 (0.830)
University Postgraduate	0.0121 (0.351)	0.634 (0.671)	0.743 (0.671)	0.743 (0.671)	0.991*** (0.374)	0.586 (0.481)	-0.0727 (0.370)	-0.0727 (0.370)	0.875 (0.809)	0.659 (0.680)	0.859** (0.391)	0.488 (0.497)	-1.076 (0.813)	-1.076 (0.813)	0.0444 (1.063)	-0.241 (0.953)	-0.250 (0.803)	-0.674 (0.868)
Party Conservative							-0.607 (0.450)	-0.607 (0.450)	5.120*** (1.050)	1.266** (0.495)	-1.642*** (0.564)	-1.498* (0.802)	-2.051** (0.817)	-2.051** (0.817)	4.305*** (1.150)	0.255 (0.780)	-2.241*** (0.741)	-1.887** (0.914)
Party LibDem							-1.440*** (0.253)	-1.440*** (0.253)	0.785 (1.126)	-1.254*** (0.427)	-2.109*** (0.266)	-1.772*** (0.345)	-1.251** (0.556)	-1.251** (0.556)	1.741 (1.175)	-0.589 (0.621)	-0.741 (0.498)	0.134 (0.576)
By-elections							0.0385 (0.582)	0.0385 (0.582)	0.162 (1.294)	0.615 (0.640)	-0.289 (0.584)	-13.29*** (0.472)	-0.123 (1.248)	-0.123 (1.248)	-0.164 (1.697)	0.448 (1.382)	-0.148 (1.206)	-16.60*** (1.191)
Legislature 2001-2005																		
Legislature 2005-2010																		
Legislature 2010-2015																		
Constant	-0.685 (0.735)	-3.180** (1.262)	1.405 (1.190)	1.405 (1.190)	-2.034*** (0.698)	-1.188 (0.897)	0.541 (0.793)	0.541 (0.793)	-5.059*** (1.745)	2.032 (1.315)	-0.287 (0.798)	0.450 (0.964)	4.717** (1.968)	4.717** (1.968)	-0.217 (2.387)	6.981*** (2.180)	4.351** (1.825)	6.399*** (1.967)
Observations	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table A6: Multinomial Analysis - Seven Clusters

	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6	Career 7	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6	Career 7	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6	Career 7
Gender	-0.225 (0.239)	-0.419 (0.440)	-0.114 (0.670)	-0.485** (0.245)	-0.0100 (0.407)	-0.535 (0.361)	-0.535 (0.361)	-0.329** (0.255)	0.0921 (0.501)	-0.639 (0.672)	-0.950*** (0.268)	0.0884 (0.473)	-1.022*** (0.381)	-0.392 (0.547)	0.469 (0.686)	-0.392 (0.844)	0.469 (0.686)	-0.392 (0.844)	-0.272 (0.520)	0.368 (0.636)	-0.674 (0.606)
State School	0.107 (0.215)	-0.0968 (0.383)	0.974 (0.682)	0.465** (0.221)	-0.427 (0.370)	-0.0452 (0.308)	-0.0452 (0.308)	-0.282 (0.237)	-0.437 (0.462)	0.330 (0.683)	-0.130 (0.240)	-0.731* (0.411)	-0.409 (0.330)	-0.725 (0.476)	-0.878 (0.605)	0.0680 (0.794)	-0.878 (0.605)	0.0680 (0.794)	-0.126 (0.449)	-1.054* (0.547)	-0.484 (0.525)
Job	0.0207 (0.228)	0.199 (0.369)	-0.416 (0.625)	0.0358 (0.232)	-0.0562 (0.392)	0.0904 (0.328)	0.0904 (0.328)	-0.121 (0.231)	0.257 (0.486)	-0.536 (0.621)	-0.180 (0.255)	-0.139 (0.424)	-0.134 (0.332)	0.729 (0.557)	0.937 (0.709)	0.255 (0.826)	0.937 (0.709)	0.255 (0.826)	0.538 (0.517)	0.513 (0.604)	0.711 (0.604)
Age at Entry	0.00507 (0.0139)	0.0185 (0.0257)	-0.180*** (0.0491)	0.0190 (0.0129)	-0.0528** (0.0238)	-0.00591 (0.0164)	-0.00591 (0.0164)	0.00208 (0.0149)	0.00751 (0.0252)	-0.178*** (0.058)	0.0146 (0.0146)	-0.0497** (0.0225)	-0.0111 (0.0179)	-0.0439 (0.0374)	-0.0410 (0.0408)	-0.262*** (0.0632)	-0.0410 (0.0408)	-0.262*** (0.0632)	-0.0450 (0.0353)	-0.0963** (0.0381)	-0.0966** (0.0381)
University Undergraduate	0.133 (0.317)	0.715 (0.633)	-0.660 (0.905)	0.512 (0.354)	1.428 (1.046)	-0.0482 (0.475)	-0.0482 (0.475)	0.157 (0.329)	1.036 (0.744)	-0.576 (0.865)	0.531 (0.366)	1.539 (1.006)	-0.0542 (0.487)	-0.499 (0.783)	0.458 (1.011)	-0.540 (1.245)	0.458 (1.011)	-0.540 (1.245)	-0.112 (0.799)	0.844 (1.196)	-0.672 (0.855)
University Postgraduate	0.0118 (0.352)	0.636 (0.671)	-0.412 (0.968)	0.990*** (0.374)	1.587 (1.066)	0.585 (0.481)	0.585 (0.481)	-0.0838 (0.370)	0.960 (0.814)	-0.581 (0.969)	0.839** (0.391)	1.635 (1.034)	0.464 (0.496)	-1.101 (0.839)	0.126 (1.077)	-1.082 (1.292)	0.126 (1.077)	-1.082 (1.292)	-0.321 (0.831)	0.622 (1.261)	-0.767 (0.892)
Party Conservative								-0.610 (0.451)	5.153*** (1.050)	-0.887 (1.093)	-1.642*** (0.564)	1.989*** (0.556)	-1.520* (0.801)	-1.987** (0.800)	4.261*** (1.144)	-1.867 (1.174)	4.261*** (1.144)	-1.867 (1.174)	-2.314*** (0.736)	0.799 (0.832)	-2.135** (0.915)
Party LibDem								-1.443*** (0.254)	0.804 (1.127)	-2.507*** (0.792)	-2.110*** (0.267)	-0.650 (0.521)	-1.779*** (0.345)	-1.281** (0.565)	1.654 (1.173)	-0.981 (1.071)	1.654 (1.173)	-0.981 (1.071)	-0.764 (0.504)	-0.289 (0.717)	0.0852 (0.583)
By-elections								0.0191 (0.585)	0.347 (1.318)	-0.430 (1.047)	-0.317 (0.584)	0.999 (0.728)	-14.97*** (0.473)	-0.0975 (1.245)	-0.0123 (1.727)	-0.987 (1.782)	-0.0975 (1.727)	-0.987 (1.782)	-0.225 (1.197)	0.727 (1.482)	-17.81*** (1.189)
Legislature 2001-2005																					
Legislature 2005-2010																					
Legislature 2010-2015																					
Constant	-0.685 (0.740)	-3.192** (1.268)	4.097* (2.171)	-2.035*** (0.702)	-0.813 (1.434)	-1.187 (0.904)	-1.187 (0.904)	0.576 (0.805)	-5.363*** (1.750)	5.740*** (2.133)	-0.224 (0.812)	-0.037 (1.509)	0.534 (0.980)	4.879** (1.979)	-0.594 (2.379)	11.77*** (2.815)	-0.594 (2.379)	11.77*** (2.815)	4.768*** (1.842)	3.618 (2.332)	6.999*** (1.999)
Observations	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695	695

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table A7: Multinomial Analysis with Oxbridge - Eight Clusters

	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6	Career 7	Career 8	Career 1	Career 2	Career 3	Career 4	Career 5	Career 6	Career 7	Career 8
Gender	-0.363 (0.312)	-0.321 (0.456)	-0.217 (0.688)	-0.0927 (0.328)	-0.451* (0.253)	0.00728 (0.412)	-0.341 (0.380)	-0.341 (0.380)	-0.506 (0.743)	-0.506 (0.743)	0.826 (0.739)	-0.461 (0.936)	-0.154 (0.603)	-0.179 (0.569)	0.512 (0.678)	-0.489 (0.676)
State School	-0.0115 (0.304)	0.0250 (0.414)	1.074** (0.819)	0.0355 (0.330)	0.264 (0.246)	-0.141 (0.387)	-0.450 (0.354)	-0.450 (0.354)	-1.015* (0.615)	-1.015* (0.615)	-0.991 (0.648)	0.586 (0.944)	-1.353** (0.544)	-0.870* (0.487)	-0.996* (0.589)	-1.780*** (0.632)
Job	-0.415 (0.321)	0.157 (0.386)	-0.615 (0.686)	0.332 (0.304)	0.00806 (0.244)	-0.0774 (0.399)	-0.170 (0.364)	-0.170 (0.364)	-0.0491 (0.667)	-0.0491 (0.667)	0.643 (0.724)	0.271 (0.832)	0.832 (0.571)	0.530 (0.530)	0.189 (0.642)	0.425 (0.651)
Age at Entry	-0.0182 (0.0167)	0.0140 (0.0208)	-0.187*** (0.0594)	0.0421** (0.0203)	0.0184 (0.0242)	-0.0490* (0.0258)	-0.00968 (0.0185)	-0.00968 (0.0185)	-0.0456 (0.0470)	-0.0456 (0.0470)	-0.0559 (0.0420)	-0.301*** (0.0696)	-0.0345 (0.0394)	-0.056 (0.0365)	-0.104** (0.0435)	-0.119*** (0.0422)
Oxbridge	-0.441 (0.351)	0.399 (0.409)	1.291** (0.612)	0.212 (0.341)	-0.346 (0.271)	0.704* (0.386)	-1.438*** (0.492)	-1.438*** (0.492)	-1.551*** (0.576)	-1.551*** (0.576)	-0.344 (0.621)	0.448 (0.809)	-1.045* (0.568)	-1.703*** (0.516)	-0.101 (0.564)	-3.331*** (0.712)
Party Conservative									-1.274 (0.945)	-1.274 (0.945)	4.349*** (1.160)	-1.618 (1.337)	-2.223** (0.930)	-2.383*** (0.838)	0.903 (0.870)	-2.832** (1.105)
Party LibDem									-0.880 (0.721)	-0.880 (0.721)	1.824 (1.182)	-0.467 (0.989)	-1.208* (0.645)	-0.433 (0.560)	-0.140 (0.758)	0.547 (0.673)
By-elections									-18.51*** (1.508)	-18.51*** (1.508)	-0.388 (1.950)	-0.171 (1.734)	-0.388 (1.431)	-0.150 (1.417)	0.437 (1.752)	-17.83*** (1.400)
Legislature 2001-2005									18.06*** (2.188)	18.06*** (2.188)	-0.201 (1.063)	-2.221** (0.894)	-2.282** (0.894)	-0.170 (0.567)	0.216 (0.721)	-20.34*** (0.672)
Legislature 2005-2010									21.88*** (0.966)	21.88*** (0.966)	1.503 (1.183)	-18.42*** (1.161)	2.183*** (1.008)	-18.60*** (0.985)	2.786** (1.103)	-19.54*** (1.072)
Legislature 2010-2015									-3.719*** (0.519)	-3.719*** (0.519)	-20.53*** (0.661)	-22.49*** (0.834)	-22.16*** (0.632)	-22.92*** (0.613)	-4.716*** (1.139)	-24.40*** (0.863)
Constant	0.0815 (0.814)	-2.524** (1.215)	2.718 (2.207)	-3.111*** (0.966)	-1.098 (0.671)	0.0753 (1.184)	-0.270 (0.838)	-0.270 (0.838)	-13.60*** (2.207)	-13.60*** (2.207)	0.487 (2.311)	11.48*** (2.791)	4.271** (1.982)	5.932*** (1.831)	4.558** (2.176)	8.835*** (2.134)
Observations	613	613	613	613	613	613	613	613	613	613	613	613	613	613	613	613

Robust standard errors in parentheses
*** p<0.01, ** p<0.05, * p<0.1

Table A8: Discrepancy Analysis

Variable	PseudoF	PseudoR2	p_value
Gender	1.5406378	0.0013683279	0.165
State School	1.0299627	0.0009147683	0.385
Job	0.8401325	0.0007461693	0.535
Age at Entry	3.1711299	0.0028164604	0.005
By-elections	6.0032943	0.0053318663	0.005
Legistature 2001-2005	10.2285072	0.0090845176	0.005
Legistature 2005-2010	60.3632637	0.0536120395	0.005
Legistature 2010-2015	46.5763286	0.2189984284	0.005
University Undergraduate	0.4084218	0.0003627426	0.930
University Postgraduate	0.6823863	0.0006060660	0.655
Party Conservative	69.4499886	0.0616824755	0.005
Party LibDem	3.2946912	0.0029262022	0.005
Total	36.9939491	0.3942771028	0.005

References

- Mirkin, B. 2013. *Mathematical Classification and Clustering*. London: Springer.
- Needleman, S., and C. Wunsch. 1970. “A General Method Applicable to the Search for Similarities in the Amino Acid Sequence of Two Proteins.” *Journal of Molecular Biology* 48 (3): 443–453.
- Studer, M., A. Gabadinho, N. S. Muller, and R. G. 2011. *An Introduction to Sequence Analysis Using the TraMineR R Package*. Workshop on Trajectories Paris Institute for Demographic and Life Course Studies, University of Geneva, October.