

Understanding state preferences with text as data: Introducing the UN General Debate corpus

Research and Politics
 April-June 2017: 1–9
 © The Author(s) 2017
 DOI: 10.1177/2053168017712821
journals.sagepub.com/home/rap


Alexander Baturo¹, Niheer Dasandi² and Slava J. Mikhaylov³

Abstract

Every year at the United Nations (UN), member states deliver statements during the General Debate (GD) discussing major issues in world politics. These speeches provide invaluable information on governments' perspectives and preferences on a wide range of issues, but have largely been overlooked in the study of international politics. This paper introduces a new dataset consisting of over 7300 country statements from 1970–2014. We demonstrate how the UN GD corpus (UNGDC) can be used as a resource from which country positions on different policy dimensions can be derived using text analytic methods. The article provides applications of these estimates, demonstrating the contribution the UNGDC can make to the study of international politics.

Keywords

Policy preferences, foreign policy, United Nations, text as data

Introduction

Every September, the heads of state and other high-level country representatives gather in New York at the start of a new session of the United Nations General Assembly (UNGA) and address the Assembly in the General Debate (GD). The GD provides the governments of the almost 200 UN member states with an opportunity to present their views on international conflict and cooperation, terrorism, development, climate change and other key issues in international politics. As such, the statements made during the GD are an invaluable and largely untapped source of information on governments' policy preferences across a wide range of issues over time.

Government preferences are central to the study of international relations and comparative politics. As preferences cannot be directly observed, they must be inferred from states' observed behaviour. One approach has been to use military alliances as an indicator of preference similarity (e.g. Bueno de Mesquita, 1983). This approach, however, provides little information about preferences when states do not have alliances. Scholars have instead overwhelmingly relied on UNGA voting records to estimate foreign policy preferences (see Bailey et al., 2015; Voeten, 2013). However, UNGA voting-based methods – like all measures of preference – rely on certain assumptions and, as such,

have both strengths and limitations (see Voeten, 2013). For example, one shortcoming is that estimates of state preference are derived from the limited number of issues that are voted on in the UNGA in a given year.¹ Therefore, it is essential that researchers can draw on additional data and measures to avoid producing findings about government preferences that are based on one type of observed state behaviour.

We argue that the application of text analytic methods to GD statements can provide much-needed additional measures and tools that can broaden our understanding of government preferences and their effects. The use of text analytic methods is rapidly gaining ground in comparative politics and legislative studies (see Herzog and Benoit, 2015; Laver et al., 2003; Proksch and Slapin, 2010). To date, however, there has been little effort to use speeches to

¹Dublin City University, Ireland

²University of Birmingham, UK

³University of Essex, UK

Corresponding author:

Slava J. Mikhaylov, University of Essex, Wivenhoe Park, Colchester, CO4 3SQ, UK.

Email: s.mikhaylov@essex.ac.uk



estimate policy preferences in international relations. The formal and institutionalised setting of the GD, its inclusion of all UN member states, which are provided with equal opportunity to address the Assembly, and the fact that it takes place every year, makes it an ideal resource from which to derive, using text analysis, estimates of state preferences that can be applied to systematic analyses of international politics.

This paper introduces a new dataset, the UN GD corpus (UNGDC), consisting of 7314 GD statements from 1970–2014, that we have preprocessed, categorised and prepared for empirical applications. In the next section, we discuss the characteristics, content and purpose of the UN GD. Second, we explain the process of collecting and pre-processing the statements, and provide an overview of the UNGDC. We then use the text as data approach to show how the UNGDC can be used as a resource from which estimates of government preferences can be derived, providing applications of these estimates. We conclude by outlining potential uses of the UNGDC in future research.

The UN GD and world politics

The GD marks the start of the UNGA regular session each year. By tradition, the opening speech is made by Brazil, with the US also scheduled to speak on the first day. Typically, the heads of state and governments speak during the first days of the GD, followed by vice-presidents, deputy prime ministers and foreign ministers, and concluding with the heads of delegation to the UN (Bailey, 1960; Luard and Heater, 1994; Nicholas, 1959; Smith, 2006).

The GD provides governments with an opportunity to declare, and to have on public record, their official position on various major international events of the past year (Smith, 2006). In addition, country representatives use the GD venue to present their governments' perspectives on broader underlying issues in international politics. Their speeches frequently deal with issues of mutual concern such as terrorism, nuclear non-proliferation, development and aid, and climate change, often appealing to the international community to do more to tackle these issues. For example, in 1995, the US discussed UN reform, non-proliferation, terrorism, money laundering and the narcotics trade in its GD statement. In turn, the UK and France both drew attention to the challenges of UN peacekeeping, while India discussed terrorism, disarmament, human rights and concerns about the inability of global institutions, such as the WTO, to address the needs of the Global South.

There are several important characteristics of GD speeches that have implications for their use in deriving estimates of state preferences from them. In contrast to UNGA roll-call votes, which are directly linked to the adoption of UN resolutions, GD speeches are not institutionally connected to decision-making within the UN. As a consequence, states face lower external constraints and

pressures when delivering GD statements than when voting in the UNGA. Indeed, studies that use UNGA voting highlight the various constraints countries face when voting as a result of, among other things, aid relationships and strategic voting blocs (see Alesina and Dollar, 2000; Kim and Russett, 1996; Voeten, 2000). The lack of external constraints means that when delivering their GD statements, governments have more leverage with the positions they take and the issues they emphasise. Hence, GD statements provide more information on key national priorities than the limited number of votes in the UNGA.

This view is supported by interviews conducted by the authors with members of the diplomatic community. The Deputy Representative of the Finnish Mission to the UN, for example, explained, 'speeches at the General Debate are interesting because they flesh out national policies – what states think ... the General Debate is the one place where states can speak their mind; it reflects the issues that states consider important'. Similarly, a spokesperson for the German Mission to the UN stated that the absence of external pressures when delivering GD statements means 'these speeches are the most sovereign thing that a country does as a member of the UN'.² It is clear that non-democratic regimes also attach great importance to GD statements. For example, members of Russia's inner political circle not only viewed the 2015 GD statement as a key summary of the country's foreign policy concerns, they were also apparently aware of its content weeks in advance.³

A significant consequence of the relative lack of external constraints in the GD is that member states can more freely express their government's perspectives on issues deemed important – including more contentious issues. As Smith (2006: 155) argues, a key function of the GD is that 'it provides members with the opportunity to blow off steam on contentious issues without causing undue damage'. This is particularly relevant for smaller nations who can use the GD to raise more disagreeable political issues (see Nicholas, 1959). For example, in 2014, Antigua and Barbuda's statement emphasised the failure of the US government to adhere to a ruling from the WTO's Dispute Settlement Body that stated that the US should pay compensation to Antigua and Barbuda. In making this complaint, the Antiguan representative highlighted the importance of the GD for smaller nations, stating 'my small nation has no military might, no economic clout. All that we have is membership of the international system as our shield and our voice in this body as our sword.'

The fewer external constraints on representatives when delivering GD statements does not, however, imply that these speeches are not strategic. Scholars have long recognised that 'member states present themselves exclusively in the guise in which they wish to be known' during these annual debates (Nicholas, 1959: 98). In fact, a key purpose of the GD is that it provides governments with the opportunity to 'influence international perceptions of their state,

aiming to position their states favorably, as well as to influence the perception of other states' (Hecht, 2016: 10). Therefore, governments use GD speeches strategically to signal their preferences among the community of states. This use of strategic signalling in the GD can be seen when we compare references to Iran in the US statements of 2012 and 2013. In the 2012 address, President Obama⁴ was highly critical of Iran:

In Iran we see where the path of a violent and unaccountable ideology leads [...] Time and again, it has failed to take the opportunity to demonstrate its nuclear program is peaceful [...] Make no mistake: a nuclear-armed Iran is not a challenge that can be contained. It would threaten the elimination of Israel, the security of Gulf nations and the stability of the global economy [...] and that is why the United States will do what we must to prevent Iran from obtaining a nuclear weapon.

In contrast, speaking a year later⁴, the US president was more reconciliatory, offering to give diplomacy one last chance in relation to Iran's nuclear programme:

if we can resolve the issue of Iran's nuclear program, that can serve as a major step down a long road towards a different relationship, one based on mutual interests and mutual respect [...] America prefers to resolve its concerns over Iran's nuclear program peacefully [...] We are not seeking regime change, and we respect the right of the Iranian people to access peaceful nuclear energy.

A few hours later during the same session, President Rouhani in his address⁴ also emphasised diplomacy and the hope of reaching a compromise. The world has subsequently learned that the US and Iran held secret talks in the background, which eventually led to the breakthrough and signing of the intermediate deal (Borger and Kamali, 2013). As such, the change in rhetoric between 2012 and 2013 demonstrates the strategic nature of GD speeches. A further example of both the importance placed by governments on the GD address and its strategic purpose is provided by the Chilcot Inquiry into the UK's role in the Iraq War. The report contains a memo sent by Prime Minister Tony Blair to President George W. Bush, complimenting the US president on the speech delivered at the 2002 GD that set out the case for war, 'It was a brilliant speech ... it puts us on exactly the right strategy to get the job done.'⁵ Hence, the US speech was seen as part of the US and UK strategy to build support for intervention in Iraq.

The lack of external constraints on member states in delivering GD statements means that they can use their address to indicate the issues considered most important by devoting more attention to these topics. As governments can choose what issues to discuss or ignore, and how strongly to emphasise certain issues, the GD provides detailed information about a government's *position*

on a policy issue and, also, the *importance* – or *salience* – of an issue for a government. As Smith (2006: 155) notes, the GD acts 'as a barometer of international opinion on important issues, even those not on the agenda for that particular session'. The focus on position and salience means that GD speeches can be used to uncover the most important topics that emerge in international politics over time.

UNGDC: The UN General Debate corpus

The speeches made in the GD are subsequently deposited at the UN Dag Hammarskjöld Library. However, statements made before 1992 are stored as image copies of typewritten documents. These are of very poor quality and require additional preprocessing using optical character recognition software. We collected speeches through the dedicated webpages of the individual UNGA GDs and the UN Bibliographic Information System (UNBIS).

Speeches are typically delivered in the native language. Based on the rules of the Assembly, all statements are then translated by UN staff into the six official languages of the UN. If a speech was delivered in a language other than English, we use the official English version provided by the UN. Therefore, all of the speeches in the UNGDC are in English.

The annual sessions are assigned numbers, starting with the first session in 1946 up to the most recent seventieth session in 2015. We collected all GD speeches from 1970 (Session 25) to 2014 (Session 69). In total, there are 7314 country statements delivered between 1970–2014. The number of countries participating in the GD increased from 70 in 1970 to 193 in 2014 in line with the increase in UN membership. Non-member states may also participate in the GD (e.g. the Holy See and Palestine). Several states that previously participated in the GD have ceased to exist. Where possible, we linked such states to their legal successor states (e.g. USSR and the Russian Federation). If this was not possible we kept speeches in the data under the country's last known name (e.g. German Democratic Republic). Overall, the corpus contains the GD contributions from 198 countries. On average, speeches contain 123 sentences and 945 unique words.⁶

Table 1 provides an overview of the UNGDC. It shows average frequency of types (unique form of a word), tokens (individual words) and sentences for each individual speech in the text corpus. In terms of who delivered the statement, for sessions with identifiable speakers and their posts, 1909 (44.3%) were delivered by heads of state or government (e.g. presidents, prime ministers, kings), 2126 (49.3%) by vice-presidents, deputy prime ministers and foreign ministers and 276 (6.4%) by country representative at the UN.⁷

Table 1. UN GD corpus.

Year	UN membership	GD statements	Types (mean freq)	Tokens (mean freq)	Sentences (mean freq)
1970	127	70	1569	8230	257
1971	132	116	1336	5927	230
1972	132	125	1157	4895	180
1973	135	120	1291	5923	230
1974	138	129	1093	4248	191
1975	144	126	1041	4280	165
1976	147	134	951	3720	151
1977	149	140	965	3452	135
1978	151	141	1159	4169	163
1979	152	144	1220	4804	200
1980	154	149	1173	4663	183
1981	157	145	1159	4357	183
1982	157	147	1134	3986	151
1983	158	149	1078	3669	157
1984	159	150	1160	3951	172
1985	159	137	1142	3605	113
1986	159	149	895	2715	85
1987	159	152	922	3010	102
1988	159	154	985	3463	124
1989	159	153	1036	3365	117
1990	159	156	1076	3606	125
1991	166	162	1086	3519	127
1992	179	167	932	2962	103
1993	184	175	1062	3433	135
1994	185	178	1142	4040	140
1995	185	172	1255	4306	168
1996	185	181	1220	4149	157
1997	185	176	915	2659	122
1998	185	181	892	2749	115
1999	188	181	857	2567	91
2000	189	178	937	2677	88
2001	189	189	681	1925	78
2002	191	188	588	1465	58
2003	191	189	666	1761	72
2004	191	192	557	1400	61
2005	191	185	505	1311	51
2006	192	193	554	1393	63
2007	192	191	573	1392	52
2008	192	192	609	1498	59
2009	192	193	662	1754	65
2010	192	189	631	1668	58
2011	193	193	709	2097	79
2012	193	194	626	1671	66
2013	193	192	776	2306	71
2014	193	193	555	1451	50

Note: Descriptive statistics for the UNGDC containing 7314 statements delivered by heads of state or their representative from 1970–2014. From 2011, the president of the European Commission made a separate statement on behalf of the EU. UN: United Nations; GD: General Debate.

Empirical application: Preferences on single-issue dimensions

The UNGDC can be used by scholars who require easy access to the statements and may want to read a particular

text, or compare selected statements. Primarily, however, we envision the UNGDC to be used in quantitative applications looking at the nature, formation and effects of state preferences in world politics. Treating text as data has a long tradition in political science (for a review, see Laver,

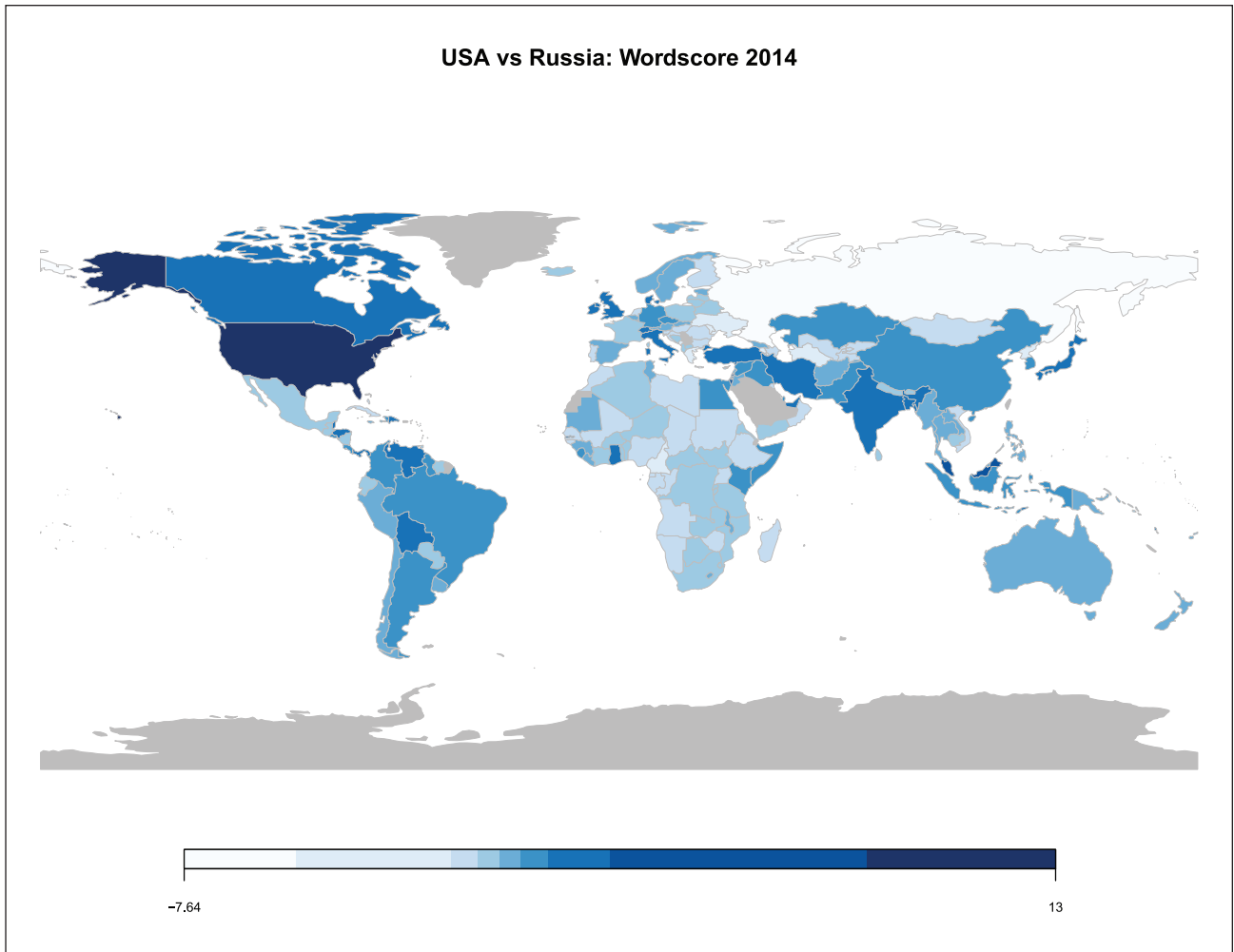


Figure 1. Wordscores map 2014.

Note: The scores are estimated in `quanteda` package (version 0.9.9-3) in R (Benoit et al., n.d.). We follow standard preprocessing during the tokenisation stage, remove English stopwords and perform stemming. We also trim the document-feature matrix to have features that appear at least five times in three documents. The US is given reference score (+1) and Russia (-1). Results are rescaled using classical LBG rescale, hence predicted scores may be beyond the (-1;+1) range.

2014). Since the earlier introduction of text scaling methods, such as Wordscores (Laver et al., 2003) and Wordfish (Slapin and Proksch, 2008), to estimate policy positions on dimensions of interest, the availability and complexity of methods has increased dramatically (Grimmer and Stewart, 2013; Herzog and Benoit, 2015). The majority of such methods are either derived directly from, or can be traced to, the natural language processing literature in computer science and computational linguistics (e.g. Benoit and Nulty, 2013; Lowe, 2008). Wordscores is by far the most popular text scaling method in political science based on a Google Scholar citation count. It is related to the Naive Bayes classifier deployed for text categorisation problems (Benoit and Nulty, 2013).

Working with text as data generally involves using the bag-of-words approach, whereby each document can be represented by a multiset (bag) of its words that disregards

grammar and word order. Word frequencies in the document are then used to classify the document into one of two categories. In Wordscores, the learning is supervised by providing training documents that are a priori known to belong to either category, so that the chosen dimension is substantively defined by the choice of training documents.

As an illustration of this approach, we derive from our resource estimates of preferences on the very specific issue of US–Russia rivalry in world politics. Figure 1 maps Wordscores estimates for the 2014 UN GD. We use statements by the US and Russia as reference texts. We therefore a priori define the policy dimension as Russia vs US. We do not use the resulting scores as an explanatory variable in an empirical application here due to limited space. However, such an application would clearly be of value for research on international relations. Here, we simply demonstrate how it is possible to derive estimates of

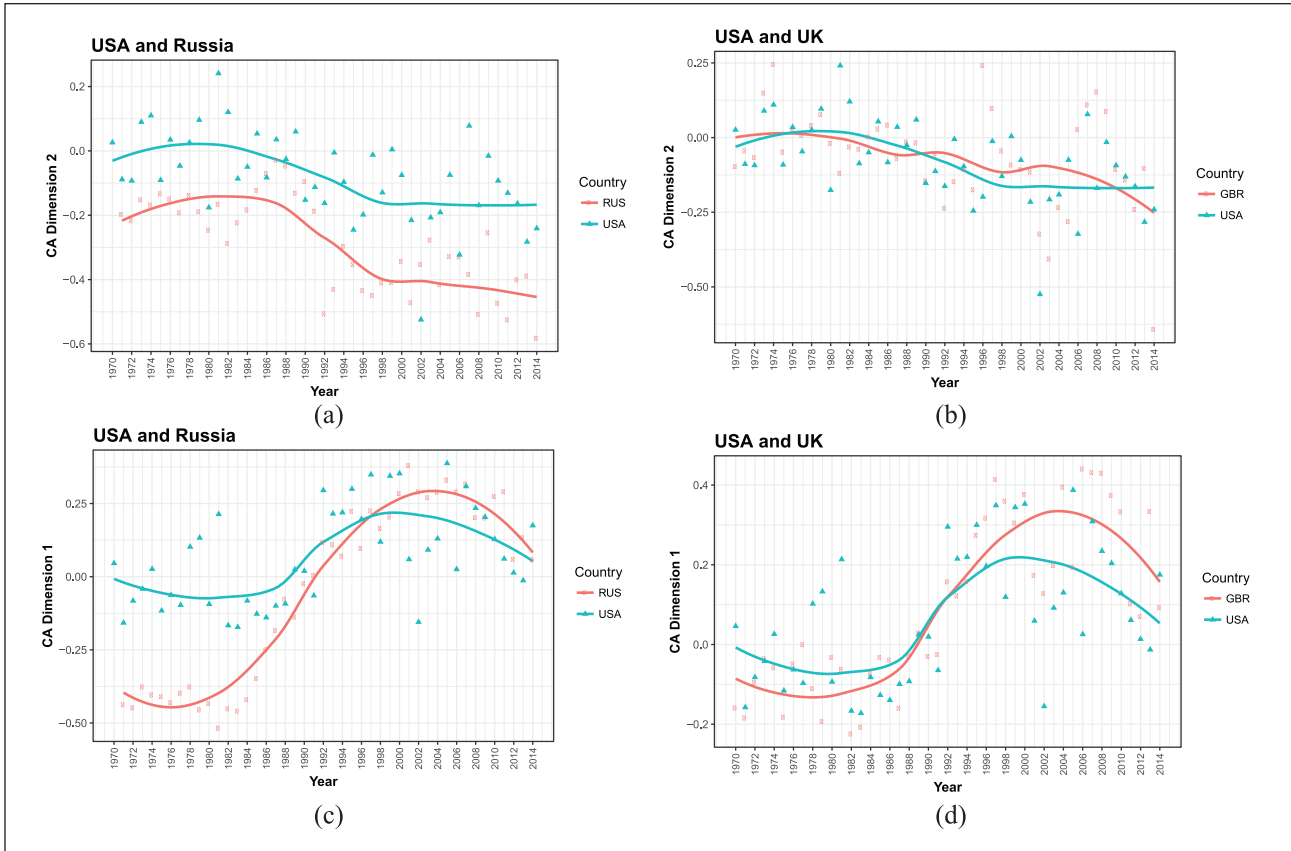


Figure 2. CA1 and CA2 of allies and opponents.

Note: The four subplots show CA estimates for the US, Russia and the UK on the first and second dimensions, as discussed in the text. Overlaid lines are loess smoothers.

differences between UN member states from our resource using the text as data approach.

Empirical application: Preferences on multiple dimensions

While estimating state preferences on single-issue dimensions has many benefits, countries routinely express preferences on multiple dimensions of foreign policy. We therefore turn to correspondence analysis (CA) – a dimensionality reduction technique (e.g. Bonica, 2013). In CA, the first dimension is fitted to explain maximal variation in the data, while subsequent dimensions explain maximal residual variation (which means dimensions are orthogonal to each other). Unlike Wordscores, the definition of the dimensions produced by CA must be discerned inductively, a posteriori (Laver, 2014). This also implies that the dimensions produced by CA may correspond to single, multiple or meta issues. Figure 2 presents the positions over time of USA and Russia (opponents) and USA and the UK (allies) on the first and second dimensions (CA1 and CA2) uncovered using CA.

Lowe (2016) suggests that position estimated by such models is a low dimensional summary of the relative

emphasis of one topic over another, compared to what would be expected by chance. This is consistent with a key assumption of the saliency theory of party competition (Budge et al., 2001), which posits that the policy differences between parties are determined by their contrasting emphases on different issues. In the context of GD statements, the CA model fitted to the count data of unique words captures countries' relative emphasis of different issues – and therefore the differences in their policy preferences.

A benefit of using CA is that it allows us to easily estimate positions on multiple dimensions. We demonstrate the ease of using multidimensional text scaling by including the new CA measures in an existing analysis of the International Criminal Court (ICC) and US nonsurrender agreements (Kelley, 2007). The format of this article prevents us from covering issues in detail; therefore, the following is intended merely as an illustrative example. In brief, the US sought to pressure other states to sign bilateral agreements not to surrender US citizens to the ICC. This attempt to seek exceptional treatment was widely criticised for inconsistency with international norms, and many countries (but not all) turned it down. Kelley (2007: 573) argues that, for these states, normative preferences trumped strategic concerns. Overall, the views on the

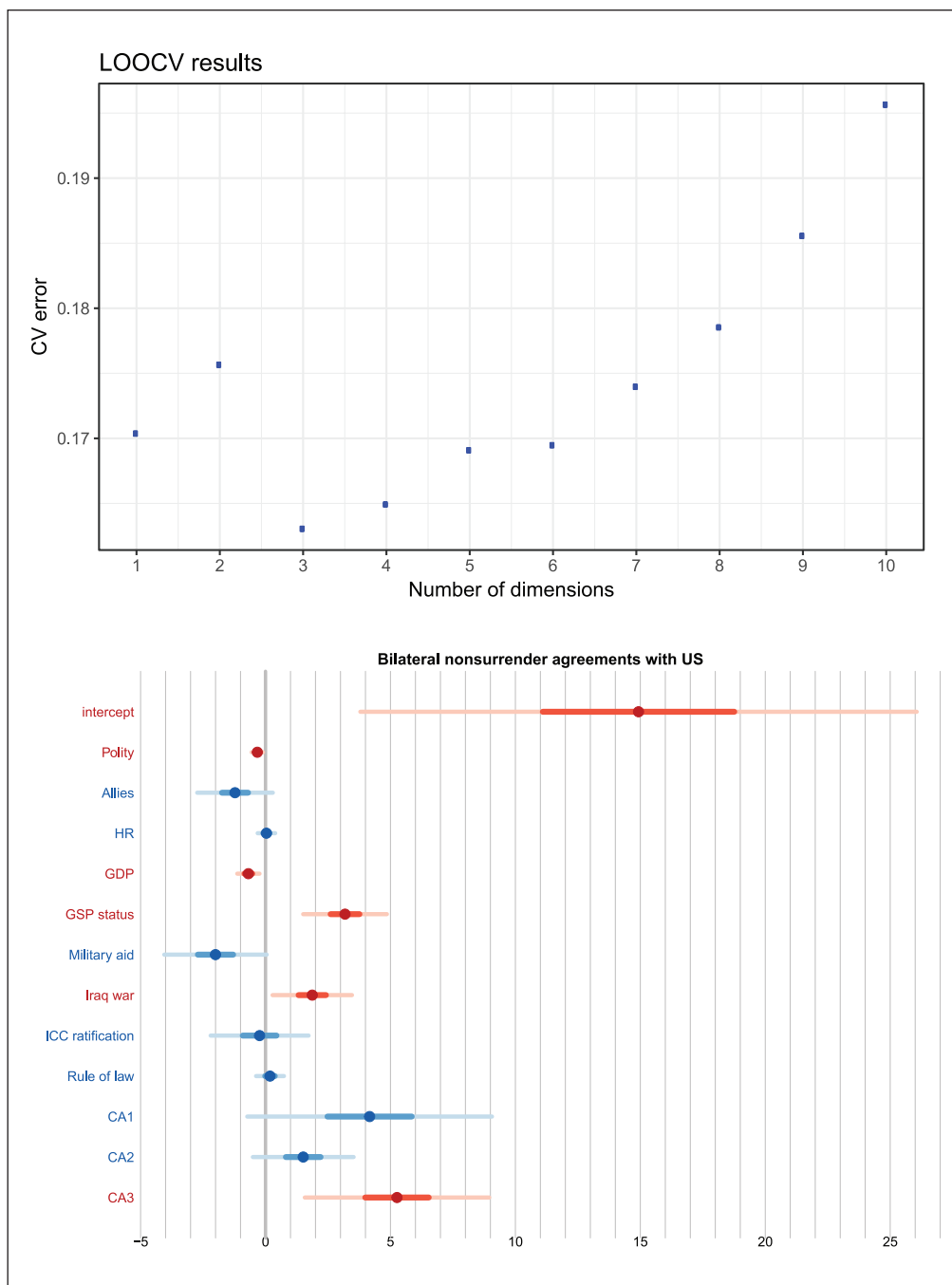


Figure 3. Choosing optimal number of CA dimensions and the estimated model results. *Note:* The upper subplot displays the results of the LOOCV analysis to choose the optimal number of dimensions for the estimated model, as discussed in the text. The bottom subplot displays the coefficients of the estimated model, as discussed in the text, with 50% and 95% CIs. Coefficients in red are statistically significant (at 95% level). The specification has fewer observations than the original analysis as Kelley (2007) includes non-UN members or states that did not participate in the GD that year, i.e. with absent text data.

nonsurrender agreements were complex and unlikely to be reduced to an easily identifiable single-issue dimension.

To determine the optimal number of dimensional estimates to include in the estimate we rely on the leave-one-out cross-validation (LOOCV) method (James et al., 2013: 178). Given the sample size, we considered alternative specifications with up to 10 CA dimensions, as presented in

Figure 3.⁸ For each alternative model we calculate the cross-validation error. As the LOOCV indicates that the optimal number of CA dimensions is three, we include three dimensions to the original specification that predicts whether countries signed nonsurrender agreements (Kelley, 2007).

The results presented in the second subplot in Figure 3 indicate that the CA3 coefficient is statistically significant.

5. Section 3.4 of the Iraq Inquiry, p.187, see http://www.iraqinquiry.org.uk/media/248175/the-report-of-the-iraq-inquiry_section-34.pdf (accessed 25 January 2017).
6. We have developed a browsing and visualisation tool that allows users to explore individual documents and the topics covered, including the top words that characterise topics, the evolution of topics over time, word distributions across topics, the underlying digitised texts of speeches, and the source PDFs. The tool can be accessed here: <http://www.smikhaylov.net/ungdc/>. The UNGDC can be downloaded from the Harvard Dataverse “United Nations General Debate Corpus” here: <http://dx.doi.org/10.7910/DVN/0TJX8Y>
7. Detailed information is available for sessions 49–69; transcripts from earlier sessions do not provide the same degree of detail regarding the post of the speaker. In the rare cases where the post of the speaker was missing in the transcript for sessions 49–69, we added this information.
8. We implement the simplest specification search using additive models. Users can implement more extensive searches using a similar approach, e.g. including interaction terms.

Carnegie Corporation of New York grant

This publication was made possible (in part) by a grant from Carnegie Corporation of New York. The statements made and views expressed are solely the responsibility of the author.

References

- Alesina A and Dollar D (2000) Who gives foreign aid to whom and why? *Journal of Economic Growth* 5(1): 33–63.
- Bailey M, Strezhnev A and Voeten E (2015) Estimating dynamic state preferences from United Nations voting data. *Journal of Conflict Resolution* 61(2): 430–456.
- Benoit K and Nulty P (2013) Classification methods for scaling latent political traits. Paper prepared for presentation at the *annual meeting of the Midwest Political Science Association 2013*, Chicago, 11–14 April 2013.
- Bailey S (1960) *The General Assembly of the United Nations: A Study of Procedure and Practice*. London: Stevens and Sons.
- Benoit K, Watanabe K, Nulty P, et al. (n.d.) Quanteda: Quantitative analysis of textual data. R package version 0.9.9–60. Available at: <http://quanteda.io>
- Bonica A (2013) Ideology and interests in the political marketplace. *American Journal of Political Science* 57(2): 294–311.
- Borger J and Kamali S (2013) Secret talks helped forge Iran nuclear deal. Available at: <http://www.theguardian.com/world/2013/nov/24/secret-usa-iran-talks-nuclear-deal> (accessed 28 November 2016).
- Budge I, Klingemann H-D, Volkens A, et al. (eds) (2001) *Mapping Policy Preferences: Estimates for Parties, Electors, and Governments, 1945–1998*. Oxford: Oxford University Press.
- Bueno de Mesquita B (1983) *The War Trap*. New Haven, CT: Yale University Press.
- Grimmer J and Stewart B (2013) Text as data: The promise and pitfalls of automatic content analysis methods for political texts. *Political Analysis* 21(3): 267–297.
- Häge F and Hug S (2016) Consensus decisions and similarity measures in international organizations. *International Interactions* 42(3): 503–529.
- Hecht C (2016) The shifting salience of democratic governance: Evidence from the United Nations General Assembly General Debates. *Review of International Studies* 42(5): 915–938.
- Herzog A and Benoit K (2015) The most unkindest cuts: Speaker selection and expressed government dissent during economic crisis. *Journal of Politics* 77(4): 1157–1175.
- James G, Witten D, Hastie, et al. (2013) *An Introduction to Statistical Learning*. New York: Springer.
- Kelley J (2007) Who keeps international commitments and why? The International Criminal Court and bilateral nonsurrender agreements. *American Political Science Review* 101(03): 573–589.
- Kim SY and Russett B (1996) The new politics of voting alignments in the United Nations General Assembly. *International Organization* 50(4): 629–652.
- Laver M (2014) Measuring policy positions in political space. *Annual Review of Political Science* 17: 207–223.
- Laver M, Benoit K and Garry J (2003) Extracting policy positions from political texts using words as data. *American Political Science Review* 97: 311–331.
- Lowe W (2016) Scaling things we can count. Paper presented at the *annual meeting of the American Political Science Association 2013*, Chicago, 29 August–1 September 2013.
- Lowe W (2008) Understanding Wordscores. *Political Analysis* 16(4): 356–371.
- Luard E and Heater D (1994) *The United Nations: How it Works and What it Does*. New York: St. Martin’s Press.
- Nicholas HG (1959) *The United Nations as a Political Institution*. Oxford: Oxford University Press.
- Proksch S-O and Slapin JB (2010) Position-taking in European Parliament speeches. *British Journal of Political Science* 40(03): 587–611.
- Slapin J and Proksch S-O (2008) A scaling model for estimating time-series party positions from texts. *American Journal of Political Science* 52(3): 705–722.
- Smith C (2006) *Politics and Process at the United Nations: The Global Dance*. Boulder, CO: Lynne Rienner.
- Voeten E (2000) Clashes in the assembly. *International Organization* 54(2): 185–215.
- Voeten E (2013) Data and analyses of voting in the United Nations General Assembly. In: Reinalda B (ed.) *Routledge Handbook of International Organization*. London: Routledge, pp.54–66.