Group and Topic Discovery from Relations and Text

Xuerui Wang, Natasha Mohanty, Andrew McCallum Department of Computer Science University of Massachusetts Amherst, MA 01003 {xuerui,nmohanty,mccallum}@cs.umass.edu

ABSTRACT

We present a probabilistic generative model of entity relationships and textual attributes that simultaneously discovers groups among the entities and topics among the corresponding text. Block-models of relationship data have been studied in social network analysis for some time. Here we simultaneously cluster in several modalities at once, incorporating the words associated with certain relationships. Significantly, joint inference allows the discovery of groups to be guided by the emerging topics, and vice-versa. We present experimental results on two large data sets: sixteen years of bills put before the U.S. Senate, comprising their corresponding text and voting records, and 43 years of similar data from the United Nations. We show that in comparison with traditional, separate latent-variable models for words or Blockstructures for votes, the Group-Topic model's joint inference improves both the groups and topics discovered.

Categories and Subject Descriptors

I.2.6 [Artificial Intelligence]: Learning; H.2.8 [Database Management]: Database Applications—data mining

General Terms

Algorithms, experimentation

Keywords

Graphical models, text modeling, relational learning

1. INTRODUCTION

Research in the field of social network analysis (SNA) has led to the development of mathematical models that discover patterns in interaction between entities [21, 5, 14]. One of the objectives of SNA is to detect salient groups of entities. Group discovery has many applications, such as understanding the social structure of organizations [6] or native

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

LinkKDD-2005 August 21, 2005, Chicago, Illinois, USA Copyright 2005 ACM 1-59593-215-1 ...\$5.00.

tribes [8], uncovering criminal organizations [19], and modeling large-scale social networks in Internet services such as Friendster.com or LinkedIn.com.

Social scientists have conducted extensive research on group detection, especially in fields such as anthropology [8] and political science [11, 7]. Recently, statisticians and computer scientists have begun to develop models that specifically discover group memberships [15, 3, 17, 13]. One such model is the stochastic Blockstructures model [17], which discovers the latent structure, groups or classes based on pair-wise relation data. A particular relation holds between a pair of entities (people, countries, organizations, etc.) with some probability that depends only on the class (group) assignments of the entities. The relations between all the entities can be represented with a directed or undirected graph. The class assignments can be inferred from a graph of observed relations or link data using Gibbs sampling [17]. This model is extended in [13] to automatically select an arbitrary number of groups by using a Chinese Restaurant Process prior.

The aforementioned models discover latent groups only by examining whether one or more relations exist between a pair of entities. The Group-Topic (GT) model presented in this paper, on the other hand, considers not only the relations between objects but also the attributes of the relations (for example, the text associated with the relations) when assigning group membership.

The GT model can be viewed as an extension of the stochastic Blockstructures model [17, 13] with the key addition that group membership is conditioned on a latent variable associated with the attributes of the relation. In our experiments, the attributes of relations are words, and the latent variable represents the topic responsible for generating those words. Unlike previous methods, our model captures the (language) attributes associated with interactions between entities, and uses distinctions based on these attributes to better assign group memberships.

Consider a legislative body and imagine its members forging alliances (forming groups), and voting accordingly. However, different alliances arise depending on the topic of the resolution up for a vote. For example, one grouping of the legislators may arise on the issue of taxation, while a quite different grouping may occur for votes on foreign trade. Similar patterns of topic-based affiliations would arise in other types of entities as well, e.g., research paper co-authorship relations between people and citation relations between papers, with words as attributes on these relations.

In the GT model, the discovery of groups is guided by the emerging topics, and the discovery of topics is guided by emerging groups. Both modalities are driven by the common goal of increasing data likelihood. Consider the voting example again; resolutions that would have been assigned the same topic in a model using words alone may be assigned to different topics if they exhibit distinct voting patterns. Distinct word-based topics may be merged if the entities vote very similarly on them. Likewise, multiple different divisions of entities into groups are made possible by conditioning them on the topics.

The importance of modeling the language associated with interactions between people has recently been demonstrated in the Author-Recipient-Topic (ART) model [16]. In ART the words in a message between people in a network are generated conditioned on the author, recipients and a set of topics that describes the message. The model thus captures both the network structure within which the people interact as well as the language associated with the interactions. In experiments with Enron and academic email, the ART model is able to discover role similarity of people better than SNA models that consider network connectivity alone. However, the ART model does not explicitly capture groups formed by entities in the network.

The GT model simultaneously clusters entities to groups and clusters words into topics, unlike models that generate topics solely based on word distributions such as Latent Dirichlet Allocation [4]. In this way the GT model discovers salient topics relevant to relationships between entities in the social network—topics which the models that only examine words are unable to detect.

We demonstrate the capabilities of the GT model by applying it to two large sets of voting data: one from US Senate and the other from the General Assembly of the UN. The model clusters voting entities into coalitions and simultaneously discovers topics for word attributes describing the relations (bills or resolutions) between entities. We find that the groups obtained from the GT model are significantly more cohesive (p-value < .01) than those obtained from the Blockstructures model. The GT model also discovers new and more salient topics in both the Senate and UN datasets—in comparison with topics discovered by only examining the words of the resolutions, the GT topics are either split or joined together as influenced by the voters' patterns of behavior.

2. GROUP-TOPIC MODEL

The Group-Topic Model is a directed graphical model that clusters entities with relations between them, as well as attributes of those relations. The relations may be either directed or undirected and have multiple attributes. In this paper, we focus on undirected relations and have words as the attributes on relations.

In the generative process for each event (an interaction between entities), the model first picks the topic t of the event and then generates all the words describing the event where each word is generated independently according to a multinomial (discrete) distribution ϕ_t , specific to the topic t. To generate the relational structure of the network, first the group assignment, g_{st} for each entity s is chosen conditionally from a particular multinomial (discrete) distribution θ_t over groups for each topic t. Given the group assignments on an event b, the matrix $V^{(b)}$ is generated where each cell $V_{ij}^{(b)}$ represents if the groups of two entities (i and j) behaved the same or not during the event b, (e.g., voted the same or not

SYMBOL	DESCRIPTION
g_{it}	entity i 's group assignment in topic t
t_b	topic of an event b
$w_k^{(b)} \\ V_{ij}^{(b)}$	the k th token in the event b
$V_{ij}^{(b)}$	entity i and j 's groups behaved same (1)
,	or differently (2) on the event b
S	number of entities
T	number of topics
G	number of groups
B	number of events
V	number of unique words
N_b	number of word tokens in the event b

number of entities who participated in the event b

DECODIDATON

 S_b

Table 1: Notation used in this paper

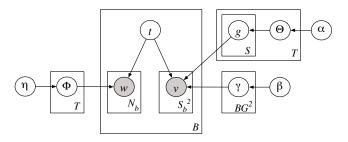


Figure 1: The Group-Topic model

on a bill). Each element of V is sampled from a binomial (Bernoulli) distribution $\gamma_{g_ig_j}^{(b)}$. Our notation is summarized in Table 1, and the graphical model representation of the model is shown in Figure 1.

Without considering the topic of an event, or by treating all events in a corpus as reflecting a single topic, the simplified model (only the right part of Figure 1) becomes equivalent to the stochastic Blockstructures model [17]. To match the Blockstructures model, each event defines a relationship, e.g., whether in the event two entities' groups behave the same or not. On the other hand, in our model a relation may have multiple attributes (which in our experiments are the words describing the event, generated by a per-topic multinomial (discrete) distribution).

When we consider the complete model, the dataset is dynamically divided into T sub-blocks each of which corresponds to a topic. The complete GT model is as follows,

$$t_b \sim \operatorname{Uniform}(\frac{1}{T})$$

$$w_{it}|\phi_t \sim \operatorname{Multinomial}(\phi_t)$$

$$\phi_t|\eta \sim \operatorname{Dirichlet}(\eta)$$

$$g_{it}|\theta_t \sim \operatorname{Multinomial}(\theta_t)$$

$$\theta_t|\alpha \sim \operatorname{Dirichlet}(\alpha)$$

$$V_{ij}^{(b)}|\gamma_{g_ig_j}^{(b)} \sim \operatorname{Binomial}(\gamma_{g_ig_j}^{(b)})$$

$$\gamma_{ab}^{(b)}|\beta \sim \operatorname{Beta}(\beta).$$

We want to perform joint inference on (text) attributes and relations to obtain topic-wise group memberships. Since inference can not be done exactly on such complicated probabilistic graphical models, we employ Gibbs sampling to conduct inference. Note that we adopt conjugate priors in our setting, and thus we can easily integrate out θ , ϕ and γ to decrease the uncertainty associated with them. This simplifies the sampling since we do not need to sample θ , ϕ and γ at all, unlike in [17]. In our case we need to compute the conditional distribution $P(g_{st}|\mathbf{w}, \mathbf{V}, \mathbf{g}_{-st}, \mathbf{t}, \alpha, \beta, \eta)$ and $P(t_b|\mathbf{w}, \mathbf{V}, \mathbf{g}, \mathbf{t}_{-b}, \alpha, \beta, \eta)$, where \mathbf{g}_{-st} denotes the group assignments for all entities except entity s in topic t, and \mathbf{t}_{-b} represents the topic assignments for all events except event b. Beginning with the joint probability of a dataset, and using the chain rule, we can obtain the conditional probabilities conveniently. The derivations are provided in detail in Appendix A. In our setting, the relationship we are investigating is always symmetric, so we do not distinguish R_{ij} and R_{ji} in our derivations (only $R_{ij}(i \leq j)$ remain). Thus

$$P(g_{st}|\mathbf{V}, \mathbf{g}_{-st}, \mathbf{w}, \mathbf{t}, \alpha, \beta, \eta)$$

$$\propto \frac{\alpha_{g_{st}} + n_{tg_{st}} - 1}{\sum_{g=1}^{G} (\alpha_g + n_{tg}) - 1} \prod_{b=1}^{B} \left(I(t_b = t) \right)$$

$$\times \prod_{h=1}^{G} \frac{\prod_{k=1}^{2} \prod_{x=1}^{d_{g_{st}hk}^{(b)}} \left(\beta_k + m_{g_{st}hk}^{(b)} - x \right)}{\prod_{x=1}^{\sum_{k=1}^{2} d_{g_{st}hk}^{(b)}} \left(\left(\sum_{k=1}^{2} (\beta_k + m_{g_{st}hk}^{(b)}) - x \right) \right),$$

where n_{tg} represents how many entities are assigned into group g in topic t, c_{tv} represents how many tokens of word v are assigned to topic t, $m_{ghk}^{(b)}$ represents how many times group g and h vote same (k = 1) and differently (k = 2) on event b, $I(t_b = t)$ is an indicator function, and $d_{gsthk}^{(b)}$ is the increase in $m_{gsthk}^{(b)}$ if entity s were assigned to group g_{st} than without considering s at all (if $I(t_b = t) = 0$, we ignore the increase in event b). Furthermore,

$$P(t_{b}|\mathbf{V},\mathbf{g},\mathbf{w},\mathbf{t}_{-b},\alpha,\beta,\eta)$$

$$\propto \frac{\prod_{v=1}^{V}\prod_{x=1}^{e_{v}^{(b)}}(\eta_{v}+c_{t_{b}v}-x)}{\prod_{x=1}^{\sum_{v=1}^{V}e_{v}^{(b)}}\left(\sum_{v=1}^{V}(\eta_{v}+c_{t_{b}v})-x\right)}$$

$$\times \prod_{g=1}^{G}\prod_{h=g}^{G}\frac{\prod_{k=1}^{2}\Gamma(\beta_{k}+m_{ghk}^{(b)})}{\Gamma(\sum_{k=1}^{2}(\beta_{k}+m_{ghk}^{(b)}))},$$

where $e_v^{(b)}$ is the number of tokens of word v in event b. Note that $m_{ghk}^{(b)}$ is not a constant and changes with the assignment of t_b since it influences the group assignments of all entities that vote on event b.

The GT model uses information from two different modalities. In general, the likelihood of the two modalities is not directly comparable, since the number of occurrences of each type may vary greatly (e.g., there may be far more pairs of voting entities than word occurrences). Thus we use a weighting parameter to rescale the likelihoods from different modalities, as is also common in speech recognition when the acoustic and language models are combined.

3. RELATED WORK

There has been a surge of interest in models that describe relational data, or relations between entities viewed as links in a network, including recent work in group discovery. One such algorithm, presented by Bhattacharya and Getoor [3], is a bottom-up agglomerative clustering algorithm that partitions links in a network into clusters by considering the change in likelihood that would occur if two clusters were

merged. Once the links have been grouped, the entities connected by the links are assigned to groups.

Another model due to Kubica et al. [15] considers both link evidence and attributes on entities to discover groups. The Group Detection Algorithm (GDA) uses a Bayesian network to group entities from two datasets, demographic data describing the entities and link data. Unlike our model, neither of these models [3, 15] consider attributes associated with the links between the entities. The model presented in [15] considers attributes of an entity rather than attributes of relations between entities.

The central theme of GT is that it simultaneously clusters entities and attributes on relations (words). There has been prior work in clustering different entities simultaneously, such as information theoretic co-clustering [9], and multi-way distributional clustering using pair-wise interactions [2]. However, these models do not also cluster attributes based on interactions between entities in a network.

In our model, group membership defines pair-wise relations between nodes. The GT model is an enhancement of the stochastic Blockstructures model [17] and the extended model of Kemp et al. [13] as it takes advantage of information from different modalities by conditioning group membership on topics. In this sense, the GT model draws inspiration from the Role-Author-Recipient-Topic (RART) model [16]. As an extension of ART model, RART clusters together entities with similar roles. In contrast, the GT model presented here clusters entities into groups based on their relations to other entities.

Exploring the notion that the behavior of an entity can be explained by its (hidden) group membership, Jakulin and Buntine [12] develop a discrete PCA model for discovering groups. In the model each entity can belong to each of the kgroups with a certain probability, and each group has its own specific pattern of behaviors. Therefore, an entity's behavior depends on the probability of belonging to a group and the probability that the group has that behavior. They apply this model to voting data in the 108th US Senate where the behavior of an entity is its vote on a resolution. A similar model is developed in [18] that examines group cohesion and voting similarity in the Finnish Parliament. We apply our GT model also to voting data. However, unlike [12, 18], since our goal is to cluster entities based on the similarity of their voting patterns, we are only interested in whether a pair of entities voted the same or differently, not their actual yes/no votes. Two resolutions on the same topic may differ only in their goal (e.g., increasing vs. decreasing budget), thus the actual votes on one could be the converse of votes on the other. However, pairs of entities who vote the same on one resolution would tend to vote same on the other resolution. To capture this, we model relations as agreement between entities, not the yes/no vote itself. This kind of "contentignorant" feature is similarly found in some work on web log clustering [1].

There has been a considerable amount of previous work in understanding voting patterns [10, 11, 7], including research on voting cohesion of countries in the EU parliament [11] and partisanship in roll call voting [7]. In these models roll call data are used to estimate *ideal points* of a legislator (which refers to a legislator's preferred policy in the Euclidean space of possible policies). The models assume that each vote in the roll call data is independent of the remaining votes, i.e., each individual is not connected to anyone else who is voting.

Datasets	Avg. Al for GT	Avg. Al for Baseline	p-value
Senate	0.8294	0.8198	< .01
UN	0.8664	0.8548	< .01

Table 2: Average AI for GT and Baseline for both Senate and UN datasets. The group cohesion in GT is significantly better than in baseline.

However, in reality, legislation is shaped by the coalitions formed by like-minded legislators. The GT model attempts to capture this interaction.

4. EXPERIMENTAL RESULTS

We present experiments applying the GT model to the voting records of members of two legislative bodies: the US Senate and the UN General Assembly.

For comparison, we present the results of a baseline method that first uses a mixture of unigrams to discover topics and associate a topic with each resolution, and then runs the Blockstructures model [17] separately on the resolutions assigned to each topic. This baseline approach is similar to the GT model in that it discovers both groups and topics, and has different group assignments on different topics. However, whereas the GT model performs joint inference simultaneously, the baseline performs inference serially. Note that our baseline is still more powerful than the Blockstructures models, since it models the topic associated with each event, and allows the creation of distinct groupings dependent on different topics.

In this paper, we are interested in the quality of both the groups and the topics. In the political science literature, group cohesion is quantified by the Agreement Index (AI) [12, 18], which measures the similarity of votes cast by members of a group during a particular roll call. The AI for a particular group on a given roll call i is based on the number of group members that vote $Yes(y_i)$, $No(n_i)$ or $Abstain(a_i)$ in the roll call i. Higher AI index means better cohesion

$$AI_i = \frac{\max\{y_i, n_i, a_i\} - \frac{y_i + n_i + a_i - \max\{y_i, n_i, a_i\}}{2}}{y_i + n_i + a_i}$$

The Blockstructures model assumes that members of a legislative body have the same group affiliations irrespective of the topic of the resolution on vote. However, it is likely that members form their groups based on the topic of the resolution being voted on. We quantify the extent to which a member s switches groups with a *Group Switch Index* (GSI).

$$GSI_s = \sum_{i,j}^{T} \frac{\text{abs}(\vec{s_i} - \vec{s_j})}{|G(s,i)| - 1 + |G(s,j)| - 1}$$

where $\vec{s_i}$ and $\vec{s_j}$ are bit vectors of the length of the size of the legislative body. The k_{th} bit of $\vec{s_i}$ is set if k is in the same group as s on topic i and similarly $\vec{s_j}$ corresponds to topic j. G(s,i) is the group of s on topic i which has a size of |G(s,i)| and G(s,j) is the group of s on topic j. We present entities that frequently change their group alliance according to the topics of resolutions.

The group cohesion using the GT model is found to be significantly greater than the baseline group cohesion under pairwise t-test, as shown in Table 2 for both the Senate

Economic	Education	Military	Energy
		Misc.	
federal	education	government	energy
labor	school	military	power
insurance	aid	foreign	water
aid	children	tax	nuclear
tax	drug	congress	gas
business	students	aid	petrol
employee	elementary	law	research
care	prevention	policy	pollution

Table 3: Top words for topics generated with the mixture of unigrams model on the Senate dataset. The headers are our own summary of the topics.

Economic Education		Foreign	Social Security	
	+ Domestic		+ Medicare	
labor	education	foreign	social	
insurance	school	trade	security	
tax federal		chemicals	insurance	
congress	aid	tariff	medical	
income	government	congress	care	
minimum	tax	drugs	medicare	
wage	energy	communicable	disability	
business	research	diseases	assistance	

Table 4: Top words for topics generated with the GT model on the Senate dataset. The topics are influenced by both the words and votes on the bills.

and the UN datasets, which indicates that the GT model is better able to capture cohesive groups. We find that nearly every document has a higher Agreement Index across groups using the GT model as compared to the baseline.

4.1 The US Senate Dataset

Our Senate dataset consists of the voting records of Senators in the 101st-109th US Senate (1989-2005) obtained from the Library of Congress THOMAS database. During a roll call for a particular bill, a Senator may respond Yea or Nay to the question that has been put to vote, else the vote will be recorded as Not Voting. We do not consider Not Voting as a unique vote since most of the time it is a result of a Senator being absent from the session of the US Senate. The text associated with each resolution is composed of its index terms provided in the database. There are 3423 resolutions in our experiments (we excluded roll calls that were not associated with resolutions). Each bill may come up for vote many times in the U.S. Senate, each time with an attached amendment, and thus many relations may have the same attributes (index terms). Since there are far fewer words than pairs of votes, we adjust the text likelihood to the 5th power (weighting factor 5) in the experiments with this dataset so as to balance its influence during inference.

We cluster the data into 4 topics and 4 groups (cluster sizes are chosen somewhat arbitrarily) and compare the results of GT with the baseline. The most likely words for each topic from the traditional mixture of unigrams model is shown in Table 3, whereas the topics obtained using GT are shown in Table 4. The GT model collapses the topics Education and Energy together into Education and Domestic,

Group 1	Group 3	Group 4		
73 Republicans	Cohen(R-ME)	Armstrong(R-CO)		
Krueger(D-TX)	Danforth(R-MO)	Garn(R-UT)		
Group 2	Durenberger(R-MN)	Humphrey(R-NH)		
90 Democrats	Hatfield(R-OR)	McCain(R-AZ)		
Chafee(R-RI)	Heinz(R-PA)	McClure(R-ID)		
Jeffords(I-VT)	Kassebaum(R-KS)	Roth(R-DE)		
	Packwood(R-OR)	Symms(R-ID)		
	Specter(R-PA)	Wallop(R-WY)		
	Snowe(R-ME)	Brown(R-CO)		
	Collins(R-ME)	DeWine(R-OH)		
		Thompson(R-TN)		
		Fitzgerald(R-IL)		
		Voinovich(R-OH)		
	Miller(D-G.			
		Coleman(R-MN)		

Table 5: Senators in the four groups corresponding to Topic Education + Domestic in Table 4.

Senator	Group Switch Index
Shelby(D-AL)	0.6182
Heflin(D-AL)	0.6049
Voinovich(R-OH)	0.6012
Johnston(D-LA)	0.5878
Armstrong(R-CO)	0.5747

Table 6: Senators that switch groups the most across topics for the 101st-109th Senates

since the voting patterns on those topics are quite similar. The new topic Social Security + Medicare did not have strong enough word coherence to appear in the baseline model, but it has a very distinct voting pattern, and thus is clearly found by the GT model. Thus GT discovers topics that are salient in that they correlate with people's behavior and relations, not simply word co-occurrences.

Examining the group distribution across topics in the GT model, we find that on the topic Economic the Republicans form a single group whereas the Democrats split into $3\,$ groups indicating that Democrats have been somewhat divided on this topic. With regard to Education + Domestic and Social Security + Medicare, Democrats are more unified whereas the Republicans split into 3 groups. The group membership of Senators on Education + Domestic issues is shown in Table 5. We see that the first group of Republicans include a Democratic Senator from Texas, a state that usually votes Republican. Group 2 (majority Democrats) includes Sen. Chafee who is known to be pro-environment and is involved in initiatives to improve education, as well as Sen. Jeffords who left the Republican Party to become an Independent and has championed legislation to strengthen education and environmental protection.

Nearly all the Senators in Group 4 (in Table 5) are advocates for education and many of them have been awarded for their efforts (e.g., Sen. Fitzgerald has been honored by the NACCP for his active role in Early Care and Education, and Sen. McCain has been added to the ASEE list as a *True Hero* in American Education). Sen. Armstrong was a member of the Education committee; Sen. Voinovich and Sen. Symms are strong supporters of early education

Everything Nuclear	Human Rights	Security	
		in Middle East	
nuclear	rights	occupied	
weapons	human	israel	
use	palestine	syria	
implementation	situation	security	
countries	israel	calls	

Table 7: Top words for topics generated from mixture of unigrams model with the UN dataset (1990-2003). Only text information is utilized to form the topics, as opposed to Table 8 where our GT model takes advantage of both text and voting information.

and vocational education, respectively; and Sen. Roth has constantly voted for tax deductions for education. It is also interesting to see that Sen. Miller (D-GA) appears in a Republican group; although he is in favor of educational reforms, he is a conservative Democrat and frequently criticizes his own party—even backing Republican George W. Bush over Democrat John Kerry in the 2004 Presidential election.

Many of the Senators in Group 3 have also focused on education and other domestic issues such as energy, however, they often have a more liberal stance than those in Group 4, and come from states that are historically less conservative. Senators Hatfield, Heinz, Snowe, Collins, Cohen and others have constantly promoted pro-environment energy options with a focus on renewable energy, while Sen. Danforth has presented bills for a more fair distribution of energy resources. Sen. Kassebaum is known to be uncomfortable with many Republican views on domestic issues such as education, and has voted against voluntary prayer in school. Thus, both Groups 3 and 4 differ from the Republican core (Group 2) on domestic issues, and also differ from each other.

The Senators that switch groups the most across topics in the GT model are shown in Table 6 based on their GSIs. Sen. Shelby(D-AL) votes with the Republicans on Economic, with the Democrats on Education + Domestic and with a small group of maverick Republicans on Foreign and Social Security + Medicare. Both Sen. Shelby and Sen. Heflin are Democrats from a fairly conservative state (Alabama) and are found to side with the Republicans on many issues.

4.2 The United Nations Dataset

The second dataset involves the voting record of the UN General Assembly [20]. We focus first on the resolutions discussed from 1990-2003, which contain votes of 192 countries on 931 resolutions. If a country is present during the roll call, it may choose to vote Yes, No or Abstain. Unlike the Senate dataset, a country's vote can have one of three possible values instead of two. Because we parameterize agreement and not the votes themselves, this 3-value setting does not require any change to our model. In experiments with this dataset, we use a weighting factor 500 for text (adjusting the likelihood of text by a power of 500 so as to make it comparable with the likelihood of pairs of votes for each resolution). We cluster this dataset into 3 topics and 5 groups (numbers are chosen somewhat arbitrarily).

The most probable words in each topic from the mixture of unigrams model is shown in Table 7. For example, Every-

G	Nuclear Arsenal	Human Rights	Nuclear Arms Race
\mathbf{R}	nuclear	rights	nuclear
О	states	human	arms
U	united	palestine	prevention
Р	weapons	occupied	race
\downarrow	nations	israel	space
	Brazil	Brazil	UK
	Columbia	Mexico	France
1	Chile	Columbia	Spain
	Peru	Chile	Monaco
	Venezuela	Peru	East-Timor
	USA	Nicaragua	India
	Japan	Papua	Russia
2	Germany	Rwanda	Micronesia
	UK	Swaziland	
	Russia	Fiji	
	China	USA	Japan
	India	Japan	Germany
3	Mexico	Germany	Italy
	Iran	UK	Poland
	Pakistan	Russia	Hungary
	Kazakhstan	China	China
	Belarus	India	Brazil
4	Yugoslavia	Indonesia	Mexico
	Azerbaijan	Thailand	Indonesia
	Cyprus	Philippines	Iran
	Thailand	Belarus	USA
	Philippines	Turkmenistan	Israel
5	Malaysia	Azerbaijan	Palau
	Nigeria	Uruguay	
	Tunisia	Kyrgyzstan	

Table 8: Top words for topics generated from the GT model with the UN dataset (1990-2003) as well as the corresponding groups for each topic (column). The countries listed for each group are ordered by their 2005 GDP (PPP) and only the top 5 countries are shown in groups that have more than 5 members.

thing Nuclear constitutes all resolutions that have anything to do with the use of nuclear technology, including nuclear weapons. Comparing these with topics generated from the GT model shown in Table 8, we see that the GT model splits the discussion about nuclear technology into two separate topics, Nuclear Arsenal which is generally about countries obtaining nuclear weapons and management of nuclear waste, and Nuclear Arms Race which focuses on the arms race between Russia and the US and preventing a nuclear arms race in outer space. These two issues had drastically different voting patterns in the U.N., as can be seen in the contrasting group structure for those topics in Table 8. The countries in Table 8 are ranked by their GDP in 2005. Thus, again the GT model is able to discover salient topics—topics that reflect the voting patterns and coalitions, not simply word co-occurrence alone.

As seen in Table 8, groups formed in Nuclear Arms Race are unlike the groups formed in the remaining topics. These

groups map well to the global political situation of that time when, despite the end of the Cold War, there was mutual distrust between Russia and the US with regard to the continued manufacture of nuclear weapons. For missions to outer space and nuclear arms, India was a staunch ally of Russia, while Israel was an ally of the US.

4.2.1 Overlapping Time Intervals

In order to understand changes and trends in topics and groups over time, we run the GT model on resolutions that were discussed during overlapping time windows of 15 years, from 1960-2000, each shifted by a period of 5 years. We consider 3823 unique resolutions in this way. The topics as well as the group distribution for the most dominant topic during each time period are shown in Table 9.

Over the years there is a shift in the topics discussed in the UN, which corresponds well to the events and issues in history. During 1960-1975 the resolutions focused on countries having the right to self-determination, especially countries in Africa which started to gain their freedom during this time. Although this topic continued to be discussed in the subsequent time period, the focus of the resolutions shifted to the role of the UN in controlling nuclear weapons as the Cold War conflict gained momentum in the late 70s. While there were few resolutions condemning the racist regime in South Africa between 1965-1980, this was the topic of many resolutions during 1970-1985—culminating in the UN censure of South Africa for its discriminatory practices.

Other topics discussed during the 70s and early 80s were Israel's occupation of neighboring countries and nuclear issues. The reduction of arms was primarily discussed during 1975-1990, the time period during which the US and Soviet Union had talks about disarmament. During 1980-1995 the central topic of discussion was the Israeli-Palestinian conflict; this time period includes the beginning of the *Intifada* revolt in Palestine and the Gulf War. This topic continued to be important in the next time period (1985-2000), but in the most recent slice (1990-2003, Table 8) it has become a part of a broader topic on human rights by combining other human rights related resolutions that appear as a separate topic during 1985-2000. The human rights issue continues to be the primary topic of discussion during 1990-2003.

Throughout the history of the UN, the US is usually in the same group as Europe and Japan. However, as we can see in Table 9 during 1985-2000, when the Israeli-Palestinian conflict was the most dominant topic, US and Israel form a group of their own separating themselves from Europe. In other topics discussed during 1985-2000, US and Israel are found to be in the same group as Europe and Japan.

Another interesting result of considering the groups formed over the years is that, except for the last time period (1990-2003), countries in eastern Europe such as Poland, Hungary, Bulgaria, etc., form a group along with USSR (Russia). However, in the last time window on most topics they become a part of the group that consists of the western Europe, Japan and the US. This shift corresponds to the end of the communist regimes in these countries that were supported by the Soviet Union. It is also worth mentioning that before 1990, our model assigned East Germany to the same group as other eastern European countries and USSR (Russia), while it assigned West Germany to the same group as western European countries.²

¹http://en.wikipedia.org/wiki/List_of_countries_by_GDP_% 28PPP%29. In Table 8, we omit some countries (represented by ...) in order to incorporate other interesting but relatively low ranked countries (for example, Russia) in the GDP list.

 $^{^2\}mathrm{Not}$ shown in Table 9 because missing from 2005 GDP data.

Time				Group distributions for Topic 3				
Period	Topic 1	Topic 2	Topic 3	Group 1	Group2	Group3	Group4	Group5
	Nuclear	Procedure	Africa Indep.	India	USA	Argentina	USSR	Turkey
	operative	committee	calling	Indonesia	Japan	Colombia	Poland	
60-75	general	amendment	right	Iran	UK	Chile	Hungary	
	nuclear	assembly	africa	Thailand	France	Venezuela	Bulgaria	
	power	deciding	self	Philippines	Italy	Dominican	Belarus	
	Independence	Finance	Weapons	Cuba	India	Algeria	USSR	USA
	territories	budget	nuclear	Albania	Indonesia	Iraq	Poland	Japan
65-80	independence	appropriation	UN		Pakistan	Syria	Hungary	UK
	self	contribution	international		Saudi	Libya	Bulgaria	France
	colonial	income	weapons		Egypt	Afghanistan	Belarus	Italy
	N. Weapons	Israel	Rights	Mexico	China	USA	Brazil	India
	nuclear	israel	africa	Indonesia		Japan	Turkey	USSR
70-85	international	measures	territories	Iran		UK	Argentina	Poland
	UN	hebron	south	Thailand		France	Colombia	Vietnam
	human	expelling	right	Philippines		Italy	Chile	Hungary
	Rights	Israel/Pal.	Disarmament	Mexico	USA	Algeria	China	India
	south	israel	UN	Indonesia	Japan	Vietnam	Brazil	
75-90	africa	arab	international	Iran	UK	Iraq	Argentina	
	israel	occupied	nuclear	Thailand	France	Syria	Colombia	
	rights	palestine	disarmament	Philippines	USSR	Libya	Chile	
	Disarmament	Conflict	Pal. Rights	USA	China	Japan	Guatemala	Malawi
	nuclear	need	rights	Israel	India	UK	St Vincent	
80-95	US	israel	palestine		Russia	France	Dominican	
	disarmament	palestine	israel		Spain	Italy		
	international	secretary	occupied		Hungary	Canada		
	Weapons	Rights	Israel/Pal.	Poland	China	USA	Russia	Cameroon
	nuclear	rights	israeli	Czech R.	India	Japan	Argentina	Congo
85-00	weapons	human	palestine	Hungary	Brazil	UK	Ukraine	Ivory C.
	use	fundamental	occupied	Bulgaria	Mexico	France	Belarus	Liberia
	international	freedoms	disarmament	Albania	Indonesia	Italy	Malta	

Table 9: Results for 15-year-span slices of the UN dataset (1960-2000). The top probable words are listed for all topics, but only the groups corresponding the most dominant topic are shown (Topic 3). We list the countries for each group ordered by their 2005 GDP (PPP) and only show the top 5 countries in groups that have more than 5 members. We do not repeat the results in Table 8 for the most recent window (1990-2003).

5. CONCLUSIONS

We present the Group-Topic model that jointly discovers latent groups in a network as well as clusters of attributes (or topics) of events that influence the interaction between entities in the network. The model extends prior work on latent group discovery by capturing not only pair-wise relations between entities but also multiple attributes of the relations (in particular, the model considers words describing the relations). In this way the GT model obtains more cohesive groups as well as fresh topics that influence the interaction between groups. The model could be applied to variables of other data types in addition to voting data. We are now using the model to analyze the citations in academic papers to capture the topics of research papers and discover research groups. It would also apply to a much larger network of entities (people, organizations, etc.) that frequently appear in newswire articles.

The model can be altered suitably to consider other attributes characterizing relations between entities in a network. In ongoing work we are extending the Group-Topic model to capture a richer notion of topic, where the attributes describing the relations between entities are represented by a mixture of topics.

6. ACKNOWLEDGMENTS

This work was supported in part by NSF grant #IIS-0326249, and by the Defense Advanced Research Projects Agency, through the Department of the Interior, NBC, Acquisition Services Division, under contract #NBCHD030010. We would also like to greatly thank Prof. Vincent Moscardelli, Chris Pal and Aron Culotta for helpful discussion.

7. REFERENCES

- [1] D. Beeferman and A. Berger. Agglomerative clustering of a search engine query log. In SIGKDD, 2000.
- [2] R. Bekkerman, R. E. Yaniv, and A. McCallum. Multi-way distributional clustering via pairwise interactions. In *ICML*, 2005.
- [3] I. Bhattacharya and L. Getoor. Deduplication and group detection using links. In *LinkKDD*, 2004.
- [4] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation. JMLR, 3:993–1022, 2003.
- [5] K. Carley. A theory of group stability. American Sociological Review, 56(3):331–354, 1991.
- [6] K. Carley. A comparison of artificial and human organizations. *Journal of Economic Behavior and Organization*, 56:175–191, 1996.

- [7] G. Cox and K. Poole. On measuring the partisanship in roll-call voting: The U.S. House of Representatives, 1887-1999. American Journal of Political Science, 46(1):477-489, 2002.
- [8] W. W. Denham, C. K. McDaniel, and J. R. Atkins. Aranda and Alyawarra kinship: A quantitative argument for a double helix model. *American Ethnologist*, 6(1):1–24, 1979.
- [9] I. S. Dhillon, S. Mallela, and D. S. Modha. Information-theoretic co-clustering. In SIGKDD, 2003.
- [10] D. Fenn, O. Suleman, J. Efstathiou, and N. Johnson. How does Europe make its mind up? Connections, cliques, and compatibility between countries in the Eurovision song contest. arXiv:physics/0505071, 2005.
- [11] S. Hix, A. Noury, and G. Roland. Power to the parties: Cohesion and competition in the European Parliament, 1979-2001. British Journal of Political Science, 35(2):209-234, 2005.
- [12] A. Jakulin and W. Buntine. Analyzing the US Senate in 2003: Similarities, networks, clusters and blocs, 2004.
- [13] C. Kemp, T. L. Griffiths, and J. Tenenbaum. Discovering latent classes in relational data. Technical report, MIT CSAIL, 2004.
- [14] D. Krackhardt and K. M. Carley. A PCANS model of structure in organization. In *Int. Sym. on Command* and Control Research and Technology, June 1998.
- [15] J. Kubica, A. Moore, J. Schneider, and Y. Yang. Stochastic link and group detection. In AAAI, 2002.
- [16] A. McCallum, A. Corrada-Emanuel, and X. Wang. Topic and role discovery in social networks. In *IJCAI*, 2005.
- [17] K. Nowicki and T. A. Snijders. Estimation and prediction for stochastic blockstructures. *Journal of the American Statistical Association*, 96(455), 2001.
- [18] A. Pajala, A. Jakulin, and W. Buntine. Parliamentary group and individual voting behavior in Finnish Parliamentin year 2003: A group cohesion and voting similarity analysis, 2004.
- [19] M. Sparrow. The application of network analysis to criminal intelligence: an assessment of prospects. *Social Networks*, 13:251–274, 1991.
- [20] E. Voeten. Documenting votes in the UN General Assembly. http://home.gwu.edu/~voeten/UNVoting.htm.
- [21] S. Wasserman and K. Faust. Social Network Analysis: Methods and Applications. Cambridge University Press, 1994.

APPENDIX

A. GIBBS SAMPLING DERIVATIONS

Begin with the joint distribution $P(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t} | \alpha, \beta, \eta)$, we can take the advantages of conjugate priors to simplify the formulae. All symbols are defined in Sec. 2.

$$P(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t} | \alpha, \beta, \eta)$$

$$= \iiint p(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t}, \theta, \gamma, \phi | \alpha, \beta, \eta) d\theta d\gamma d\phi$$

$$= \iiint \prod_{b=1}^{B} P(t_b) \prod_{t=1}^{T} \left(p(\theta_t | \alpha) \prod_{s=1}^{S} P(g_{st} | \theta_t) p(\phi_t | \eta) \right)$$

$$\times \prod_{b=1}^{B} \prod_{g=1}^{G} \prod_{h=g}^{G} p(\gamma_{gh}^{(b)}|\beta) \prod_{b=1}^{B} \prod_{i=1}^{N_{b}} P(w_{i}^{(b)}|\phi_{t_{b}})$$

$$\times \prod_{b=1}^{B} \prod_{i=1}^{S} \prod_{j=i+1}^{S} P(V_{ij}^{(b)}|\gamma_{g_{i}g_{j}}^{(b)})) d\theta d\gamma d\phi$$

$$= \iiint \left(\frac{1}{T}\right)^{B} \prod_{t=1}^{T} \left(\frac{\Gamma(\sum_{g=1}^{G} \alpha_{g})}{\prod_{g=1}^{G} \Gamma(\alpha_{g})} \prod_{g=1}^{G} \theta_{tg}^{\alpha_{g}-1} \prod_{g=1}^{G} \theta_{tg}^{n_{tg}}\right)$$

$$\times \prod_{t=1}^{T} \left(\frac{\Gamma(\sum_{v=1}^{V} \eta_{v})}{\prod_{v=1}^{V} \Gamma(\eta_{v})} \prod_{v=1}^{V} \phi_{tv}^{\eta_{v}-1}\right)$$

$$\times \prod_{b=1}^{B} \prod_{g=1}^{G} \prod_{h=g}^{G} \left(\frac{\Gamma(\sum_{k=1}^{2} \beta_{k})}{\prod_{k=1}^{2} \Gamma(\beta_{k})} \prod_{k=1}^{2} (\gamma_{ghk}^{(b)})^{\beta_{k}-1}\right)$$

$$\times \prod_{t=1}^{T} \prod_{v=1}^{V} \phi_{tv}^{c_{tv}} \prod_{b=1}^{B} \prod_{g=1}^{G} \prod_{h=g}^{G} \prod_{k=1}^{2} (\gamma_{ghk}^{(b)})^{m_{ghk}^{(b)}} d\theta d\gamma d\phi$$

$$\times \iiint_{t=1}^{T} \prod_{g=1}^{G} \theta_{tg}^{\alpha_{g}+n_{tg}-1} \prod_{t=1}^{T} \prod_{v=1}^{V} \phi_{tv}^{\eta_{v}+c_{tv}-1}$$

$$\times \prod_{b=1}^{B} \prod_{g=1}^{G} \prod_{h=g}^{G} \prod_{k=1}^{2} (\gamma_{ghk}^{(b)})^{\beta_{k}+m_{ghk}^{(b)}-1} d\theta d\gamma d\phi$$

$$\propto \prod_{t=1}^{T} \left(\frac{\prod_{g=1}^{G} \Gamma(\alpha_{g}+n_{tg})}{\Gamma(\sum_{g=1}^{G} (\alpha_{g}+n_{tg}))} \frac{\prod_{v=1}^{V} \Gamma(\eta_{v}+c_{tv})}{\Gamma(\sum_{v=1}^{V} (\eta_{v}+c_{tv}))} \right)$$

$$\times \prod_{b=1}^{B} \prod_{g=1}^{G} \prod_{h=g}^{G} \frac{\prod_{k=1}^{C} \Gamma(\beta_{k}+m_{ghk}^{(b)})}{\Gamma(\sum_{k=1}^{C} (\beta_{k}+m_{ghk}^{(b)})}$$

Using the chain rule, we can get the conditional probability conveniently,

$$P(g_{st}|\mathbf{V}, \mathbf{g}_{-st}, \mathbf{w}, \mathbf{t}, \alpha, \beta, \eta)$$

$$= \frac{P(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t}|\alpha, \beta, \eta)}{P(\mathbf{g}_{-st}, \mathbf{V}, \mathbf{w}, \mathbf{t}|\alpha, \beta, \eta)}$$

$$\propto \frac{P(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t}|\alpha, \beta, \eta)}{P(\mathbf{g}_{-st}, \mathbf{V}_{-st}, \mathbf{w}, \mathbf{t}|\alpha, \beta, \eta)}$$

$$\propto \frac{\alpha_{g_{st}} + n_{tg_{st}} - 1}{\sum_{g=1}^{G} (\alpha_g + n_{tg}) - 1} \prod_{b=1}^{B} \left(I(t_b = t) \right)$$

$$\times \prod_{h=1}^{G} \frac{\prod_{k=1}^{2} \prod_{x=1}^{d_{g_{st}hk}^{(b)}} \left(\beta_k + m_{g_{st}hk}^{(b)} - x \right)}{\prod_{r=1}^{2} \prod_{x=1}^{d_{g_{st}hk}^{(b)}} \left(\left(\sum_{k=1}^{2} (\beta_k + m_{g_{st}hk}^{(b)}) - x \right) \right)$$

and,

$$P(t_{b}|\mathbf{V}, \mathbf{g}, \mathbf{w}, \mathbf{t}_{-b}, \alpha, \beta, \eta)$$

$$= \frac{P(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t}|\alpha, \beta, \eta)}{P(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t}_{-b}|\alpha, \beta, \eta)}$$

$$\propto \frac{P(\mathbf{g}, \mathbf{V}, \mathbf{w}, \mathbf{t}|\alpha, \beta, \eta)}{P(\mathbf{g}, \mathbf{V}_{-b}, \mathbf{w}_{-b}, \mathbf{t}_{-b}|\alpha, \beta, \eta)}$$

$$\propto \frac{\prod_{v=1}^{V} \prod_{x=1}^{e_{v}^{(b)}} (\eta_{v} + c_{t_{b}v} - x)}{\prod_{x=1}^{\sum_{v=1}^{V} e_{v}^{(b)}} \left(\sum_{v=1}^{V} (\eta_{v} + c_{t_{b}v}) - x\right)}$$

$$\times \prod_{g=1}^{G} \prod_{h=g}^{G} \frac{\prod_{k=1}^{2} \Gamma(\beta_{k} + m_{ghk}^{(b)})}{\Gamma(\sum_{k=1}^{2} (\beta_{k} + m_{ghk}^{(b)}))}$$