

Implicit Communication Detection Using Topics Model on Asynchronous Communication Data

Charles Panaccione

Pearson Knowledge Technologies
4940 Pearl East Circle, Suite 200
Boulder, CO 80301
charles.panaccione@pearsonkt.com

Peter Foltz

Pearson Knowledge Technologies
4940 Pearl East Circle, Suite 200
Boulder, CO 80301
peter.foltz@pearsonkt.com

Abstract

Detecting implicit communication is a difficult problem since it involves uncovering a person's intentions rather than accept the explicit communication of individuals. We used various forms of topic models on a large dataset of asynchronous communication to detect aspects of implicit communication. We show that by analyzing rates of change of words in topics and looking at the divergence of words in topics, we can detect anomalous behavior of words which shows what words and topics are active. The detection of such activity is one way to uncover implicit communication in messages.

1 Introduction

Language conveys a great deal of information about people and their audience, and also provides a view into their situation. The primary goal in communicating with another person is to share knowledge or information relevant to the particular context. Most of the information exchanged in communication is explicitly stated. Indeed, pragmatically, there is a strong assumption among parties communicating that what is being stated is true [3]. However, along with the explicit content of a message, individuals' choices of words and how the information is conveyed can sometimes hint at a much richer insight into their motives, a background signal that may elaborate or even contradict the explicit message. Thus, detecting and understanding the implicit information in communications is critical to comprehending the complete nature of an information interchange.

2 Dataset

We used data from a study conducted at the Army War College. The Army War College has a large distance learning program with online discussions actively monitored by course instructors. The data was collected from an asynchronous distance learning planning activity offered to several hundred senior officers in January 2004 over a ten day period. There were 20 separate discussion groups with 12 to 15 participants per group. Participants were U.S. government personnel from all over the world—from Kuwait to Kenya to Kansas. The activity was titled “Interagency Process Simulation” and dealt with U.S. foreign and security policy and the future of NATO [4]. Students addressed the following issues during the simulation:

- Continued U.S. engagement in Europe through NATO
- Russian membership in NATO
- European Union security and defense development

Participants were given unique roles, such as Deputy Secretary of State, and assigned to

departments and committees, e.g., State Department, Policy Coordinating Committee, Deputies Committee.

The discussion generated 1829 comments with a few hundred comments per group. Comments were quite long and well-thought out, averaging around 150 words, indicating that students were heavily engaged in the simulation. Because the data included instructor ratings of performance, measures of on-topic and off-topic discussion, assigned roles and tasking of participants, and time stamped indications of events, we were able to compare model predictions to measures of ground truth.

3 Approach

We initially used LDA [2] to derive an overall picture of the topics and then we developed and tested other topic model methods including author-topic modeling [5] and discrete dynamic topic modeling [2] on the War College data. As part of the work, we extended some functionality in the discrete dynamic topic model. First, we examined how fast certain words move in and out of favor over time. By looking at the rate of change of words, we can detect if certain words or topics are being promoted or suppressed. This technique can also be combined with the author-topic model to see which individuals are responsible for pushing certain words or topics over time in order to determine how they may be influencing the discussion.

Besides using the currently available topic model methods, we extended the topics model further by measuring the divergence of topics and words over time relative to average strength of a topic or word. In our data when the relative strength of a word exceeds a threshold, the model would signal that a word or topic is exhibiting anomalous behavior.

To train the topic models, we used a corpus containing NATO and DARPA documents we collected from previous studies.

4 Results

4.1 LDA

The first application was to use LDA to perform general clustering of words and topics. In doing this analysis, we noticed three interesting clusters. The first cluster referred to a tasking that evaluated continued viability such as whether difficulties over Iraq indicate lasting divergence. The second cluster consisted of five topics that describe proposals for future military use such as coalitions of the willing, or new military capabilities. The third cluster consisted of three topics and described NATO and EU enlargement to Central and Eastern Europe.

4.2 Author-Topics Model

The next step in our topic modeling was to incorporate authors with topics. The War College data contain 128 authors over 1829 posts. The Author-Topic model correctly identified the author 18% of the time for authors who posted 10 or more times. The random rate for this analysis is 0.8% so the Author-Topic model did over 20 times better than random chance. For comparison, we used LSA along with Random Forest to do the same analysis and the success rate for that method was under 10%.

4.3 Discrete Dynamic Topics Model

The last method we examined was dynamic topic models. The War College data consists of asynchronous online communication between individuals where the content or topic of conversation can change quickly. The initial work with topic models successfully showed the most important topics over the entire War College data but a finer grain approach was needed to extract a deeper understanding of the data. The messages were segmented into 10 time slices with each slice representing a day in the class. A discrete dynamic topic model was applied to the War College Data.

4.4 Extensions to Discrete Dynamic Topics Model

The dynamic topics model is extended to show the strength of a particular word across all topics. For a given word, the joint probability of the word strength and topic strength is calculated to get an overall strength of word for a time slice. The strength of a word is the sum of the joint probabilities of a word across all topics. For example, if the word “NATO” appears in 3 topics, the strength of the word NATO is the probability of the word “NATO” in a particular topic times the probability of that topic occurring. This is performed for all three topics and the probabilities are summed to get an aggregate probability or strength. In our model, we are assuming independence between topics. The strength of a word across all 10 time slices is smoothed. This smoothing process dampens some of the extreme high and low values and has been used in dynamic topic models before [1]. When the strength of the word is plotted over time, we can see how the influence of a word changes. In Figure 1, we see four words: war, allies, consensus and economic. Some words show decreasing strength over time such as “consensus”; ironically, this result could be due to individuals and teams failing to reach consensus. Words can exhibit wild fluctuations such as “economic” while other words are more stable such as “war”.

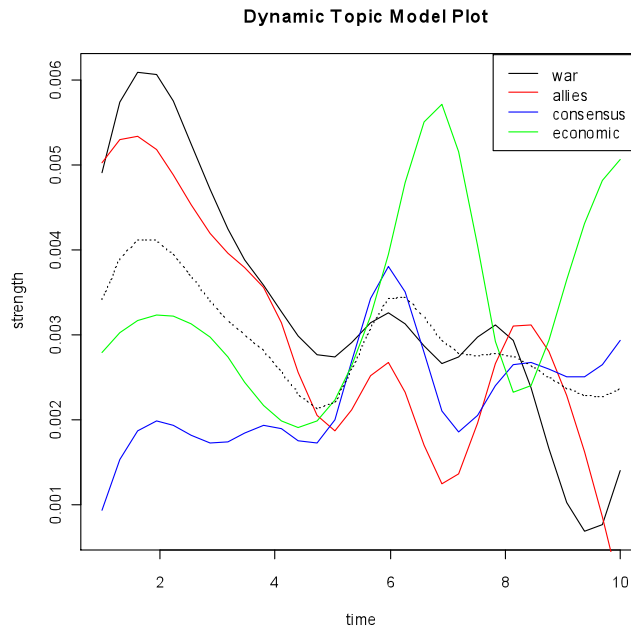


Figure 1: Dynamic Topic Model Plot. The dashed line is the mean strength.

While the dynamic topic model graph is interesting on its own, we can derive new metrics to detect implicit communication. The first of these metrics is the rate of change (slope) of the relative strength of a word across all topics. This metric shows how fast a word is moving with respect to other words. A fast moving word shows that a topic is being strongly pushed or repressed while a slow moving word shows that a topic is slow to change.

In Figure 1, “war” and “allies” look to have the same slope while “consensus” and “economic” look to have different slopes. However, Figure 2 tells a different story. In that plot, “war”, “allies” and “consensus” all have similar slopes while “economic” has a different slope. That tells us that “consensus” is moving at the same rate even though its relative strength is different from “war” and “allies”. We can also deduce that these three words are topically related due to their similar rates of change. “Economic”, on the other hand, moves in the opposite direction compared to the other three words. This effect could be the result of downplaying “economic” when discussing “war”, “allies” or consensus and vice versa.

The last metric is divergence from word or topic norms as shown in Figure 2. This plot shows how far away a particular word's curve is from the average curve; that is, it plots the difference between a word curve and the mean word curve (the dashed line in Figure 1). This metric is useful as it can be used as a signal strength detector. We can create confidence bounds for this plot and whenever a curve falls outside of the confidence bounds, the relative strength of that word is beyond what is normally expected.

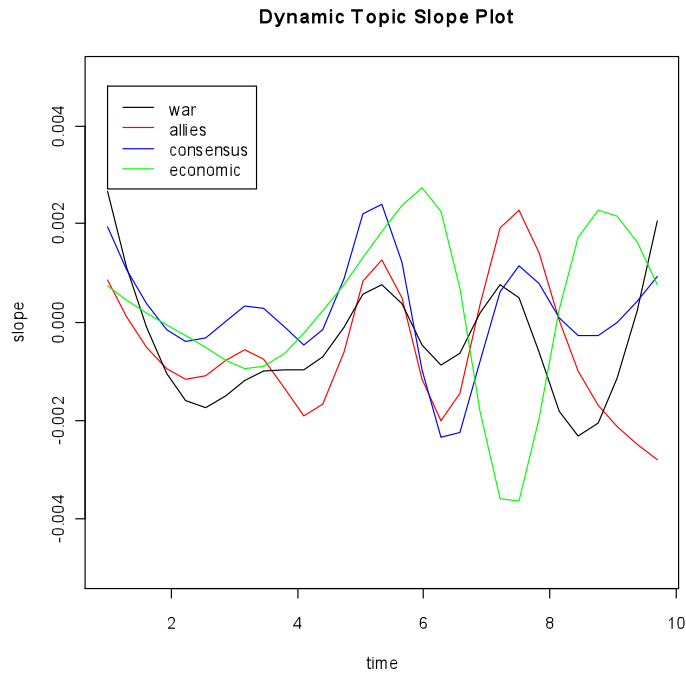


Figure 2: Dynamic Topic Model Slope Plot showing the changes in relative strength for words across all topics.

4 Conclusion

We have introduced a method for detecting implicit communications using topic models on asynchronous communication data. Our approach uses a variation of dynamic topic model to show the relative strength of a word in a topic. By comparing relative strength over time, we can see which words and topics are being communicated more or less strongly.

References

- [1] Blei, D. M. and Lafferty, J. D. 2006. Dynamic topic models. In Proceedings of the 23rd international Conference on Machine Learning (Pittsburgh, Pennsylvania, June 25 - 29, 2006). ICML '06, vol. 148. ACM, New York, NY, 113-120.
- [2] Blei, D.M., NG, A.Y. and Jordan, M.I. (2003). Latent Dirichlet Allocation. Journal of Machine Learning Research, Vol. 3, pp. 993-1022.
- [3] Grice HP (1975) Logic and conversation. Syntax & Semantics, eds Cole P, Morgan JL (Academic, New York), Vol 3.
- [4] Krupnick, C. (2004) U.S. Army War College's interagency process simulation. Interactive Technologies Conference, Arlington, VA.
- [5] Rosen-Zvi, M., Griffiths, T., Steyvers, M., and Smyth, P. (2004). The author-topic model for authors and documents. In Proceedings of the 20th conference on Uncertainty in Artificial Intelligence, pages 487-494, Arlington, VA, USA. AUAI Press.