# NATURAL LANGUAGE PROCESSING CS491

# PROJECT REPORT 2

## Group Chat Text Segmentation using Topic Modelling

**GROUP NUMBER- G5**

**TEAM MEMBERS DETAILS**

**1) KUMARI RENUKA (U101115FCS111)**

**2) RISHABH KUMAR KANDOI (U101115FCS283)**

**3) SOUMI PAL (U101115FCS158)**

# TABLE OF CONTENTS

# PROJECT DETAILS

**GROUP MEMBERS-**

| S.NO | NAME | EMAIL | E .NUMBER |
|------|------|-------|-----------|
| 1 | Kumari Renuka | kumari.renuka@st.niituniversity.in | U101115FCS111 |
| 2 | Rishabh Kandoi | rishabhk.kandoi@st.niituniversity.in | U101115FCS283 |
| 3 | Soumi Pal | Soumi.pal@st.niituniversity.in | U101115FCS158 |

## ASSIGNED PAPER

*Topic Segmentation using bottleneck method*

## CURRENT PAPER

*. Topic segmentation with a structured topic model*

# PAPER SUMMARY

## DESCRIPTION OF THE PAPER ASSIGNED AND THE PROBLEM FACED

**Information Bottleneck (IB)** method has been successfully applied to clustering in the NLP domain. Specifically, IB attempts to balance the trade-off between accuracy and compression (or complexity) while clustering the target variable, given a joint probability distribution between the target variable and an observed relevant variable. Similar to clustering, this paper interprets the task of text segmentation as a compression task with a constraint that allows only contiguous text snippets to be in a group. In the proposed technique the information bottleneck method augmented with **sequential continuity constraints**. Furthermore, they have utilized critical **non-textual clues** such as time between two consecutive posts and people mentions within the posts.

However, **information bottleneck** is recently proved to a theoretical foundation for deep learning and very complex in its own to be integrated with the natural language processing. Also, we are not be able to find sufficient resources and research papers related to chat text segmentation using information bottleneck.

*Reference-* Vishal, S., Yadav, M., Vig, L. and Shroff, G., 2017. Information Bottleneck Inspired Method For Chat Text Segmentation. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (Vol. 1, pp. 194-203).

## SUMMARY: GROUP CHAT TEXT SEGMENTATION USING TOPIC MODELLING

The conversation platforms have now become prevalent for both personal and professional usage. For instance, in a large enterprise scenario, project managers can utilize these platforms for various tasks such as decision auditing and dynamic responsibility allocation. Logs of such conversations offer potentially valuable information for various other applications such as automatic assessment of possible collaborative work among people etc. It is thus vital to invent effective segmentation methods that can separate discussions into small granules of independent conversational snippets. Independent means that a segment should be self-contained and discussing the same topic. Topic modelling is one such method which can be utilized to separate discussions into small granules and provide a topic name to each segment.

In the paper named "*Topic Segmentation with a Structured Topic Model*" presents a new hierarchical Bayesian unsupervised topic segmentation model, integrating a point-wise boundary sampling algorithm with a structured topic model. This new model takes advantage of the high modelling accuracy of structured topic models to produce a topic segmentation based on the distribution of latent topics. This model provides high quality segmentation

performance on Choi's dataset, as well as two sets of meeting transcripts and written texts. After reading the paper we have implemented the topic segmentation in python.

**Reference-** Du, L., Buntine, W. and Johnson, M., 2013. Topic segmentation with a structured topic model. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies* (pp. 190-200).
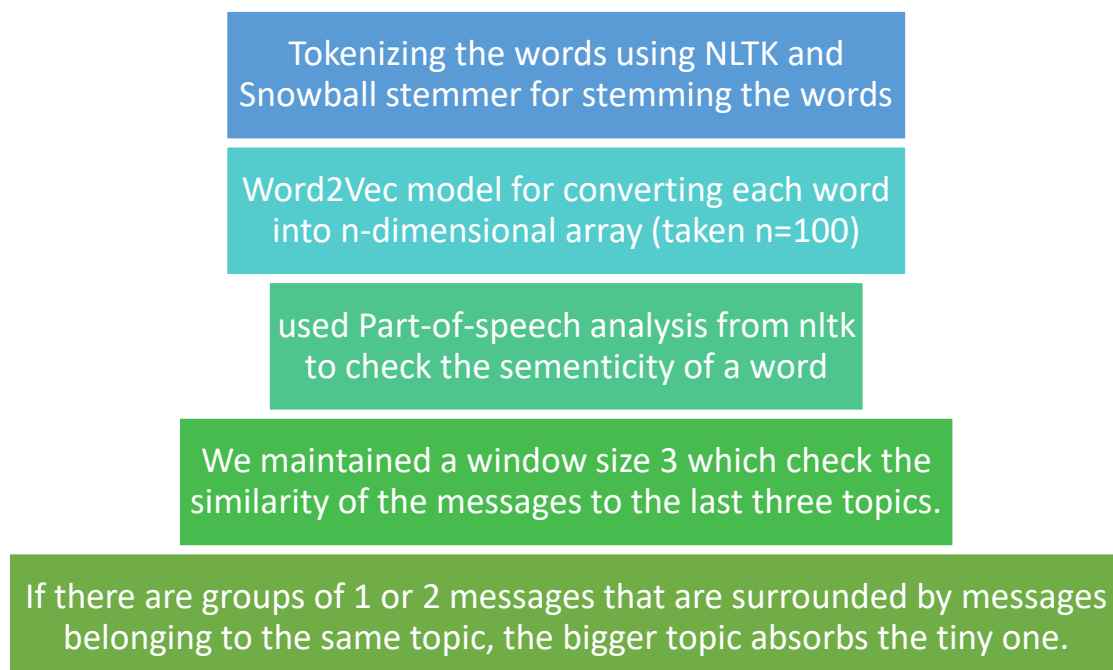
## TOPIC SEGMENTER

This tool segment group chat into several conversations and each conversation represents a different subject. We have used the knowledge of Natural Language Processing and Neural Networks.

## IMPORTANT POINTS TO CONSIDERED

1. Messages from a same group and related to a particular topic many not be sequential.
2. We have segmented a single large conversation related to the same topics into smaller subtopics.
3. The reply texts such as 'OK', 'nice ', 'I agree' are treated as the category of the same topic.

### FIG1: WORKING OF A TOPIC SEGMENTER

Tokenizing the words using NLTK and Snowball stemmer for stemming the words

Word2Vec model for converting each word into n-dimensional array (taken n=100)

used Part-of-speech analysis from nltk to check the sementicity of a word

We maintained a window size 3 which check the similarity of the messages to the last three topics.

If there are groups of 1 or 2 messages that are surrounded by messages belonging to the same topic, the bigger topic absorbs the tiny one.

# WORK DONE SO FAR

After going through the reference paper provided and analysing the feasibility of the implementation of the paper, we found that the method involved in doing chat text segmentation is beyond our scope to implement it in the same way as described in the paper. The unavailability of in-depth information about implementing Information Bottleneck method for Text Segmentation made us search for some other possible methods to do so.

We found there are other numerous ways that can be adopted to get the same results with rather less complexity like Naïve Bayes, SciPy, etc. One such method is using nltk library which is built for NLP in Python. We chose this library because it can work simply over unsupervised data as well. This means we don't have to spend more time in training the dataset rigorously to get the valid output and thus more time can be spent in analysing and refining the output as per the requirements.

## *Till now, our method is designed to work as follows:*

1. Built the dataset as described below.
2. Use "nltk" to tokenize each message removing stop words (e.g. 'to', 'and', 'a') and punctuation signs, except in cases like "It's", where the apostrophe is necessary.
3. Then we use a Snowball stemmer to reduce each word to its root form (e.g. 'agrees' and 'agreed' to 'agree').
4. Now we have a simplified vector text for each message.
5. Then we use a Word2Vec model to convert each word into an n-dimensional vector (by default n is 100). The advantage of this model is that algebraic operations on these vectors carry semantic meaning (e.g. vector('king') - vector('man') ~= vector('queen') - vector('woman')). We use the whole set of messages as the training data of the model.
6. As a post-step, if there are groups of 1 or 2 messages that are surrounded by messages belonging to the same topic, the bigger topic absorbs the tiny one. This happens when one or two messages are short and so grammatically incoherent that it's not possible to related them to any existing topic.

## *Contribution of each team member:*

1. Rishabh: Dataset conversion and tokenization.
2. Renuka: Stemming and vectorizing.
3. Soumi: Integrating all to convert raw data to topics. Researched and came out with the idea adopted and explained it briefly to the group members.

# SCREENSHOTS OF OUTPUT

**FIG2: RAW DATASET**



**FIG3: JSON DATA 1**

**FIG4: JSON DATA 2**



**FIG5: JSON DATA 3**

**FIG6: JSON DATA 4**



**FIG5: JSON DATA 3**

**DATASET DESCRIPTION:**

We have used the same dataset that is being provided to us along with the reference paper. The format of the dataset is as follows:

1. Each text file contains around 25 chat conversations.
2. Each conversation consists of the text being written by the user along with the User ID and timestamp.
3. Each conversation comes in a single line.

This dataset is then converted to JSON format for further processing. The JSON format thus formed is as follows:

1. There are 4 objects named as "anon_text", "user", "ts" and "type".
2. Extract all chat text from each conversation and store it to "anon_text".
3. Extract all User IDs from each conversation and store it to "user".
4. Extract all timestamp information from each conversation and store it to "ts".
5. Since currently each conversation if by default of type "message", so for each conversation, store the string "message" to the object "type".

# PLAN FOR THE REST OF THE SEMESTER

The work remaining are to refine the output in the best way possible. To do so, the following are planned to be implemented till the next/final submission:

1) To identify reply objects. There a considerable number of messages that are simply reply messages (e.g. 'Ok. let's do it', 'I agree', 'Mine is better', etc.) which don't contribute to the topic but are part of the same topic. Semantically, these messages by themselves carry no meaning unless they belong to a bigger topic. We call these message reply objects, and we treat them as a particular case. According to the paper "Topic Detection in Group Chat Based on Implicit Reply" more than 90% of the messages are reply objects in a group conversation.

2) Complete Incorporate a message containing <@user_id> also as a reply message.

3) Introduce a dynamic window size feature to make the machine remember the last topic for better topic segmentation.

4) Understanding the dataset used and the different ways to process it.

5) Introduce a concept of Similarity Distance. Since each conversation is already converted to vector, so we can apply POS tagging to identify semantic meaning for each conversation and provide a score to be compared to other conversation. This helps to segment the texts more efficiently.

**These works are planned to be distributed in following manner:**

1. **Rishabh:** Identify Reply Objects.
2. **Renuka:** Find Similarity distance.
3. **Soumi:** Incorporate dynamic window size functionality