# Query Based Summarized Email Extraction

## GROUP 07

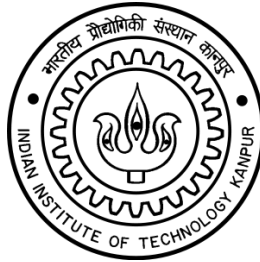[1]Keshaw Singh      [2]Roshan Kumar      [3]Sandipan Mandal

[1]13347, skeshaw@iitk.ac.in
[2]13590, roshan@iitk.ac.in
[3]13616, mandals@iitk.ac.in

## CS671A: Final Report

# Contents

# 1 Introduction

## 1.1 Problem Statement

Our Problem statement can be summarised in one line as follows

**"Given a query, extract 'topic-wise' summarized email content from various conversations"**

To explain further, our system takes a query as input and returns summarised content of relevant mails in the database. If our query can correspond to more than one meaning, it returns seperate summary corresponding to each.

## 1.2 Dataset

The dataset we used is **Enron Email Dataset**
available at `https://www.kaggle.com/wcukierski/enron-email-dataset`.

Brief description of the dataset:-

- Consists of a single 'emails.csv' file.

- Data from 150 users, mostly senior management of enron.

- Has approximately 0.5M emails over course of 2 years.

A small preview of the **emails.csv** file.

| file | message |
|---|---|
| allen-p/_sent_mail/1. | Message-ID: <18782981.1075855378110.JavaMail.evans@thyme> Date: Mon, 14 May 2001 16:39:00 -0700 (PDT) From: phillip.allen@enron.com To: tim.belden@enron.com Subject: Mime-Version: 1.0 Content-Type: text/plain; charset=us-ascii Content-Transfer-Encoding: 7bit X-From: Phillip K Allen X-To: Tim Belden <Tim Belden/Enron@EnronXGate> X-cc: X-bcc: X-Folder: \Phillip_Allen_Jan2002_1\Allen, Phillip K.\'Sent Mail X-Origin: Allen-P X-FileName: pallen (Non-Privileged).pst Here is our forecast |
| allen-p/_sent_mail/10. | Message-ID: <15464986.1075855378456.JavaMail.evans@thyme> Date: Fri, 4 May 2001 13:51:00 -0700 (PDT) From: phillip.allen@enron.com To: john.lavorato@enron.com Subject: Re: Mime-Version: 1.0 Content-Type: text/plain; charset=us-ascii Content-Transfer-Encoding: 7bit X-From: Phillip K Allen X-To: John J Lavorato <John J Lavorato/ENRON@enronXgate@ENRON> X-cc: X-bcc: X-Folder: \Phillip_Allen_Jan2002_1\Allen, Phillip K.\'Sent Mail X-Origin: Allen-P X-FileName: pallen (Non-Privileged).pst Traveling to have a business meeting takes the fun out of the trip. Especially if you have to prepare a presentation. I would suggest holding the business plan meetings here then take a trip without any formal business meetings. I would even try and get some honest opinions on whether a trip is even desired or necessary. As far as the business meetings, I think it would be more productive to try and stimulate discussions across the different groups about what is working and what is not. Too often the presenter speaks and the others are quiet just waiting for their turn. The meetings might be better if held in a round table discussion format. My suggestion for where to go is Austin. Play golf and rent a ski boat and jet ski's. Flying somewhere takes too much time. |
| allen-p/_sent_mail/100. | Message-ID: <24216240.1075855687451.JavaMail.evans@thyme> Date: Wed, 18 Oct 2000 03:00:00 -0700 (PDT) From: phillip.allen@enron.com To: leah.arsdall@enron.com Subject: Re: test Mime-Version: 1.0 Content-Type: text/plain; charset=us-ascii Content-Transfer-Encoding: 7bit X-From: Phillip K Allen X-To: Leah Van Arsdall X-cc: X-bcc: X-Folder: \Phillip_Allen_Dec2000\Notes Folders\'sent mail X-Origin: Allen-P X-FileName: pallen.nsf test successful. way to go!!! |
| allen-p/_sent_mail/1000. | Message-ID: <13505866.1075863688222.JavaMail.evans@thyme> Date: Mon, 23 Oct 2000 06:13:00 -0700 (PDT) From: phillip.allen@enron.com To: randall.gay@enron.com Subject: Mime-Version: 1.0 Content-Type: text/plain; charset=us-ascii Content-Transfer-Encoding: 7bit X-From: Phillip K Allen X-To: Randall L Gay X-cc: X-bcc: X-Folder: \Phillip_Allen_Dec2000\Notes Folders\'sent mail X-Origin: Allen-P X-FileName: pallen.nsf Randy, Can you send me a schedule of the salary and level of everyone in the scheduling group. Plus your thoughts on any changes that need to be made. (Patti S for example) Phillip |

# 2 Pre-processing

## 2.1 Reasons for Pre-processing

- Difficult to use in the format provided.

- Flooded with irrelevant information.

- Inclusion of forwarded mails result in repetitive dataset.

- Unexpected tokens, numeric need to be appropriately handled.

## 2.2 Method of Pre-processing

- Used **pandas** library to properly organize the mail for later use. After a few processing steps any field of the mail can be accessed easily using **pandas**. For example suppose in python, we want to access the content, subject, to and from field of $i^{th}$ mail

```
content = pandas_frame['content'][i]
subject = pandas_frame['Subject'][i]
To = pandas_frame['To'][i]
From = pandas_frame['From'][i]
```

- Wrote Regular Expressions to remove appended 'Forwarded' message.

# 3 Initial Steps

## 3.1 Using LDA

Latent Dirichlet Allocation(LDA) [2] represents documents as mixtures of topics that spit out words with certain probabilities.
Suppose we are given a set of documents. We decide upon the number of topics, say, K. Then learning happens as follows:

- Go through each document, and randomly assign each word in the document to one of the K topics.

- This random assignment already gives both topics representations of all the documents and the word distributions of all topics, though not very good ones.

- To improve upon these distributions, for each document d:

  – Go through each word w in d and for each topic t, compute two things:
    1. p(topic t | document d) = the proportion of words in document d that are currently assigned to topic t
    2. p(word w | topic t) = the proportion of assignments to topic t over all documents that come from this word w.

  Reassign w a new topic, where we choose topic t with probability p(topic t | document d) * p(word w | topic t) (according to our generative model, this is essentially the probability that topic t generated word w, so it makes sense that we resample the current word's topic with this probability). In other words, in this step, we're assuming that all topic assignments except for the current word in question are correct, and then updating the assignment of the current word using our model of how documents are generated.

– After repeating the previous step a large number of times, we'll eventually reach a roughly steady state where our assignments are pretty good. So use these assignments to estimate the topic mixtures of each document (by counting the proportion of words assigned to each topic within that document) and the words associated to each topic (by counting the proportion of words assigned to each topic overall).

## 3.2   Simple Mail Grouping

- Assume forwarded or replied mails to have same subject appended with 'FW:' or 'RE:'

- Find relevant emails this way and arrange them in order of time sent to get the whole string

- Follow this by topic analysis and summarization

- People often modify subject and/or don't continue the conversation by replying back

- As many details are lost this way summarization and query search won't give results.

## 3.3   Trying Different Mail Representation

Various representations for email were tried, which include:

- Bag of Words

- Weighted Bag of Words

- Average of Bag of Words

- Bag of Vectors

- Weighted Bag of Vectors

- Doc2Vec [4]

In the end, we settled for Doc2Vec representation for our model. Its working has been explained below.

### 3.3.1   Distributed Representations of Documents(Doc2Vec)

We use Doc2vec which is an unsupervised model that learns fixed-length feature representations from variable pieces of text, such as sentences, paragraphs or documents. This has the potential to overcome the weakness of the bag-of-words model. This has been known to perform excellently on text as well as sentiment classification tasks. The two models of training tried are:

**3.3.1.1   Paragraph Vector: A distributed memory model**   The key points to note how the model is trained are:

1. Every paragraph is mapped to a unique vector, represented by a column in matrix D and every word is also mapped to a unique vector, represented by a column in matrix W.

2. The paragraph vector and word vectors are concatenated to predict the next word in a context.

3. The model is the trained in a similar manner to word2vec with an additional document token along with word tokens to predict the next word in a sequence

This technique is shown in Figure 1

Figure 1: In this model, the concatenation or average of this vector with a context of three words is used to predict the fourth word. The paragraph vector represents the missing information from the current context and can act as a memory of the topic of the paragraph.



Figure 2: Distributed Bag of Words version of paragraph vectors. In this version, the paragraph vector is trained to predict the words in a small window.

**3.3.1.2 Paragraph Vector without word ordering: Distributed bag of words** The key points to note how the model is trained are:

1. We ignore the context words in the input, and force the model to predict words randomly sampled from the paragraph in the output.

2. We sample a text window, then sample a random word from the text window and form a classification task given the Paragraph Vector.

This technique is shown in Figure 2

## 3.4 Using K-Means

- Attempted clustering by K-Means using representations described before

- K is brought down from no. of groups obtained previously

- So obtained clusters are assigned "key-phrases" which serve to address our query

    The nature of emails causes clustering to perform poorly.

5

Some representations did provide decent results but not what we wanted

# 4 Algorithm Followed

## 4.1 Algorithm Details

Pseudocode of the algorithm is mentioned in **Algorithm 1**

---
**Algorithm 1** Get Summarized Mails
---
1: **procedure** SUMMARIZEDMAILS(Query, MailCorpus, StopWords, Idf)     ▷ Mail corpus is pre-processed
2:     PreprocessedQuery ← REMOVESTOPWORDSANDSTEM(Query)
3:     ImportantTerms ← **k** Terms in **PreprocessedQuery** having highest **idf**.
4:     RelevantMails ← All Mails in MailCorpus having any word matching any terms from **ImportantTerms**
5:     Clusters ← BIRCHCLUSTERING(RelevantMails)
6:     **for each** cluster ∈ Clusters **do**
7:         summary ← LSASUMMARIZATION(cluster)
8:         topic ← EXTRACTTOPIC(summary)
9:         **print** topic
10:         **print** summary
---

## 4.2 Details of Steps Involved

### 4.2.1 Stemming and Stop words Removal

The first step in our pipeline is to remove stopwords and stem the contents and queries. It is required because of the following reasons:

- Users may query for some stopword, which don't carry much information and occurs in many mails.

- Users may query for different form of the same word. Without stemming, we might miss out on some form.

### 4.2.2 Email Extraction

- Use idf value of words to get the 3 most important words from the query. It is an important step because some words appear in many mails and hence don't carry any specific information. High idf implies appearance in lesser number of mails and hence carrying more information.

- Now use these 3(max) words to extract the mails where they occur.

### 4.2.3 Clustering

**BIRCH** (balanced iterative reducing and clustering using hierarchies) [1] is an unsupervised data mining algorithm used to perform hierarchical clustering over particularly large data-sets. An advantage of BIRCH is its ability to incrementally and dynamically cluster incoming, multi-dimensional metric data points in an attempt to produce the best quality clustering for a given set of resources.

- Used BIRCH, with email vectors trained with Doc2Vec. Clustering was better at this stage because number of clusters was limited. No need to specify the number of clusters.

- Another reason for clustering was that there may be different instances of the same event. We need to differentiate between those instances.

Consider the search term **survey**. As we can see from the results shown below, the clustering algorithm is able to properly differentiate between all the senses.

### 4.2.4 Summarization

We then executed summarizer on the cluster. We cleaned the emails for forwarded mails and other text irrelevant at this stage to come up with proper summary.

We used four summarizers for this purpose. The first three come as a part of the **pyTLDR** [8] package in python. All three of these implementations are extractive - that is, they simply extract and display the most relevant sentences from the input text. PyTLDR comes with a built-in sentence tokenizer that is used for summarization. The tokenizer performs stemming in several languages as well as stop-word removal. The last one comes as a part of the **sumy** [7] package.

- TextRankSummarizer [pyTLDR]

  This approach uses PageRank [5] algorithm, where "votes" or "in-links" are represented by words shared between sentences. The length of the summary can be specified either as a number of sentences, or a percentage of the original text.



- LsaSummarizer [pyTLDR]

  This method reduces the dimensionality of the article into several "topic" clusters using singular value decomposition, and selects the sentences that are most relevant to these topics. This is a rather more abstract summarization algorithm.



- RelevanceSummarizer [pyTLDR]

  This method computes and ranks the cosine similarity between each sentence vector and the overall document, removing the most relevant sentence at each iteration.

- LsaSummarizer [sumy]

  Algorithm is same as second but is taken from different library.



### 4.2.5   Topic Extraction

Next we extracted the most probable topic/keyphrase from the summaries corresponding to each cluster. We used 'RAKE' [6] library for this purpose. This was done to identify the clusters. We used the highest scoring keyphrase as topic for the cluster.

From the results, we observed that the topic extractor was not performing well. In many cases, the topic was unable to describe the content. It was happening because often mail contents spans over multiple small topics. Hence it becomes difficult to extract a single topic to describe the content.

| Search Term | Number of relevant mail | Number of clusters | Purity |
|---|---|---|---|
| Survey | 27 | 4 | 0.94 |
| Deal | 38 | 14 | 0.90 |
| Project | 21 | 7 | 0.90 |
| Birthday | 12 | 5 | 0.97 |
| Meeting | 51 | 17 | 0.92 |

Table 1: Performance of clustering algorithm for various search terms

| Search Term | Algorithm | Important Points | Redundancy | Coherence |
|---|---|---|---|---|
| Survey | TextRank[pyTLDR] | 1 | 1 | 1 |
| Survey | LSA[pyTLDR] | 1 | 1 | 1 |
| Survey | RelevanceScore[pyTLDR] | 1 | 1 | 1 |
| Survey | LSA[sumy] | 1 | 1 | 1 |
| Deal | TextRank[pyTLDR] | 0 | 1 | 1 |
| Deal | LSA[pyTLDR] | 0 | 1 | 0 |
| Deal | RelevanceScore[pyTLDR] | 0 | 1 | 0 |
| Deal | LSA[sumy] | 1 | 1 | 1 |
| Project | TextRank[pyTLDR] | 1 | 1 | 0 |
| Project | LSA[pyTLDR] | 1 | 1 | 0 |
| Project | RelevanceScore[pyTLDR] | 0 | 1 | 0 |
| Project | LSA[sumy] | 1 | 0 | 0 |

Table 2: Performance of Summarization algorithms for various search terms

# 5  Evaluation of Algorithm

## 5.1  Evaluation of Clustering Algorithm

We evaluated the performance of Clustering algorithm using **Purity**[3] of clusters.
Some results are tabulated in Table 1.

As we can see from Table 1, even though the number of clusters is quite small the **Purity** is very close to **1** in every case.
From this we can conclude that the clustering algorithm is performing pretty well in clustering the mails by different instances of the search term.

## 5.2  Evaluation of Summarization

We evaluated the summarization algorithms w.r.t to **inclusion of important points, non-redundancy and coherence**.
But our principal concern in this project is **inclusion of important points.**
We had 5 human annotators evaluate our summary on a scale of 0-1 (whether it is unacceptable or acceptable respectively) under each point. The majority opinion is tabulated in Table 2.

From the Table 2, if we consider the column corresponding to **Important Points**, we observe that LSA from sumy library has a score of 1 (Acceptable) for all the three search terms. Since our main concern was inclusion of important points, we can conclude that LSA algorithm from sumy is performing better than others in the sense that it is always capturing the main points from the cluster.

## 5.3   Evaluation of Topic Extraction

We did not explicitly evaluate the performance of Topic Extraction because we could see from the screenshots (of topic extraction), the performance is very poor.

# 6   Problems faced

- No general format of writing mails, depends on user to user.

- No systematic way to keep track of threads.

- Modifying subjects while continuing on the same thread - becomes difficult to keep track of the threads by the Subject.

- Lots of stuffs from various domains - becomes difficult to extract a unique topic from the mails.

# 7   Suggested Improvements

- Semantic analysis may be used to extract the main information asked for by the user.

- Some sort of data structure can be used to properly organize the data.

- All the pre-processing can be done before hand to efficiently reply to the queries.

# References

[1]   *BIRCH*. URL: https://en.wikipedia.org/wiki/BIRCH.

[2]   Edwin Chen. *Introduction to Latent Dirichlet Allocation*. URL: http://blog.echen.me/2011/08/22/introduction-to-latent-dirichlet-allocation/.

[3]   *Evaluation of clustering*. URL: http://nlp.stanford.edu/IR-book/html/htmledition/evaluation-of-clustering-1.html.

[4]   *Gensim Doc2Vec*. URL: https://radimrehurek.com/gensim/models/doc2vec.html.

[5]   Larry Page. *Page Rank Algorithm*. URL: https://en.wikipedia.org/wiki/PageRank.

[6]   *Python RAKE*. URL: https://pypi.python.org/pypi/python-rake/1.0.5.

[7]   *Python sumy*. URL: https://pypi.python.org/pypi/sumy.

[8]   *Python TLDR*. URL: https://pypi.python.org/pypi/PyTLDR.