

ModDBS- X^M : A Diversity Based Summarizer for DUC

Tadashi Nomoto

National Institute of Japanese Literature

1-16-10 Yutaka Shinagawa

Tokyo 142-8585 Japan

nomoto@nijl.ac.jp

1 Description of the system

ModDBS- X^M is a clustering based single document summarizer. It is a generic extractive summarizer which requires of the input nothing more than the availability of basic IR statistics such as term and document frequency. Therefore it could be adapted for any language and domain without much effort.

The system goes through three major states to generate a summary: data preparation, summarization, and post-summarization. An input text is first examined for its conformity to the XML syntax; some portions of it are extracted for use in summarization, which are passed on to the sentence selection step, which in turn builds diverse topical clusters over the input and chooses representative sentences thereof. The selected sentences are then put through a post-summarization process, where parenthetical expressions are identified and removed.

1.1 Data preparation

The system takes as input an XML document and runs a sequence of operations to prepare data for summarization. These include:

1. validating the document with an XML parser (Derksen, 2000)
2. isolating particular XML elements for use in summarization, i.e. subtrees rooted at <TEXT>, <HEADLINE> (or <HL>) and <LEADPARA> (or <LP>).
3. identifying sentences in selected elements using MXTERMINATOR (Ratnaparkhi, 1997)
4. assigning POS tags to word tokens found in the sentences, which is done with LimaTK (Yamashita, 1999).

Texts not compliant with the XML 1.0 standards are rejected. (There were a few of them in the DUC test data.) Also rejected are texts with more than one <TEXT> element since we took <TEXT> to mean the main body of the text, as was the case with the majority of texts in the test data. (In the FBIS articles, we took [Text] as <TEXT>.)

The purpose of the initial sequence of operations is to identify sentences in particular portions of a

well-formed XML document and to produce POS information on each word in the sentences. Let us call the set of sentences generated by the preparation process, a *source text*.

The initial sequence is followed by the creation of a table holding the tfidf weight for each word type in the source text, which is defined as: for a given term x ,

$$T(x) = (1 + \log(tf(x))) \cdot idf(x)$$

$tf(x)$ is the frequency of term x in a document, $idf(x)$ is the inverse document frequency of x . The document frequency of a term was determined with a reference to the test data set alone.

For index terms, we used everything except for punctuation marks, non-linguistic symbols, particles such as case marker. We did not use a stoplist except for those elements already excluded from the set of index terms.

1.2 Summarization

The summarization step here uses a slightly modified version of DBS/ X^M (Nomoto and Matsumoto, 2001). (Hence the name ModDBS- X^M .) It takes as input a source text with supplementary information such as POS assignments and the weight table, and performs the following operations:

1. **Find-Diversity**: Find diverse topic clusters in the source text.
2. **Reduce-Redundancy**: For each topic cluster, locate most important sentences and take them as representative of that cluster.

The goal of the summarization here is to extract sentences in such a way that they collectively create a general picture of what the source text is about. The following look at some details of each operation.

1.2.1 Find-Diversity

Find-Diversity is a clustering algorithm built upon the K -means extended with Minimum Description Length Principle (MDL) (Rissanen, 1997; Li, 1998). In fact the algorithm here is an MDL-version of X -means (Pelleg and Moore, 2000), an extension of the K -means clustering algorithm with the functionality

of estimating K , the number of clusters which otherwise needs to be supplied by the user. We call our adaptation of X -means ‘ X^M means.’

K -means is a hard clustering algorithm that produces a clustering of input data points into K disjoint subsets. It dynamically redefines a clustering by relocating each centroid to the center of mass of points associated with it and re-associating the centroid with points closest to it.

In an effort to improve the time-efficiency and scalability of clustering, we have modified **Find-Diversity** as formulated as part of DBS/ X^M to incorporate a strategy for rapidly refining initial points for clustering (Bradley and Fayyad, 1998). The modified **Find-Diversity** determines initial points by clustering, using K -means, *centroids* of a set of subsamples drawn randomly from the source text (each subsample containing 10% of the source), and selecting, among them, those which give a least distorted clustering for the centroids.

1.2.2 Reduce-Redundancy

Reduce-Redundancy is a simple sentence ranking algorithm based on Zechner (1996), where one takes the weight of a given sentence as the sum of tfidf scores of index terms occurring in that sentence. The weight W of sentence s is given by:

$$W(s) = \sum_{x \in s} T(x)$$

where x denotes an index term. **Reduce-Redundancy** applies itself to each topic cluster from **Find-Diversity**, ranks sentences there according to the weighting model above, and selects those with highest scores. As is apparent from the above formula, the sentence weight is *not* normalized for length. The idea of choosing as representative best scoring sentences is intended to minimize the loss of the resulting summary’s relevance to a potential query. (The exact number of sentences to select from each cluster depends on the desired length of a summary.)

1.3 Post-Summarization Process

The post-summarization process operates on each of the sentences from the summarization step. What it does is a simple removal of parenthetical expressions from a sentence. Other than that, no paraphrasing or transformation of a sentence is performed; conjunctives such as *but* or *and* occurring sentence initially are also left intact. (Incidentally, we have found that about a half of the DUC per-doc summaries which have a low rating (“2”) for grammaticality had some parentheses incorrectly removed.) Finally, the sentences are put in the order in which they appear in the source text.

2 Some Remarks

The average number of peer units (PUs) for (per-doc) summaries generated by ModDBS was 3.41 while that of model units (MUs) was 6.39. Since the precise definition of model or peer units was not available at the time of writing, we simply assumed them to be some clause-like elements. The DUC results indicate that model summaries are about twice as long as those generated by the system; this would mean that in order to make system summaries comparable in unit length to model summaries, one has to make them twice as long, i.e. 200-word long. The number of PUs for human created (per-doc) summaries was found to be 5.53 on average, which is far closer to the average number of MUs.

The reason why the system summaries tend to contain less units apparently has to do with the system’s inability to perform within sentence (or unit) reductions; since a per-doc summary is limited to 100 words in length, to increase the number of units included in a summary demands reduction of words contained in a unit, which does not happen with the present system except for a simple removal of parentheses.

References

- P. S. Bradley and Usama M. Fayyad. 1998. Refining initial points for k-means clustering. In *Proceedings of the Fifteenth International Conference on Machine Learning (ICML98)*, pages 91–99. Morgan Kaufmann.
- Enno Derksen. 2000. Perl module: libxml-enno. <http://www.cpan.org>.
- Hang Li. 1998. *A Probabilistic Approach to Lexical Semantic Knowledge Acquisition and Structural Disambiguation*. Ph.D. thesis, University of Tokyo, Tokyo.
- Tadashi Nomoto and Yuji Matsumoto. 2001. A new approach to unsupervised text summarization. In *Proceedings of the 24th International ACM/SIGIR Conference on Research and Development in Information Retrieval*, New Orleans, September. ACM.
- Dan Pelleg and Andrew Moore. 2000. X -means: Extending K -means with Efficient Estimation of the Number of Clusters. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML2000)*, pages 727–734, Stanford University, CA, June-July. Morgan Kaufmann.
- Adwait Ratnaparkhi. 1997. MXTERMINATOR. <http://www.cis.upenn.edu/~adwait/statnlp.html>.
- Jorma Rissanen. 1997. Stochastic complexity in learning. *Journal of Computer and System Sciences*, 55:89–95.
- Tatuo Yamashita. 1999. LimaTK. <http://cl.aist-nara.ac.jp/~tatuo-y/ma>.
- Klaus Zechner. 1996. Fast Generation of Abstracts from General Domain Text Corpora by Extracting Relevant Sentences. In *Proceedings of the 16th International Conference on Computational Linguistics*, pages 986–989. Copenhagen, Denmark.