

# Opinion Summarization using Entity Features and Probabilistic Sentence Coherence Optimization: UIUC at TAC 2008 Opinion Summarization Pilot

Hyun Duk Kim, Dae Hoon Park, V.G.Vinod Vydiswaran, ChengXiang Zhai  
{hkim277, dpark34, vgvinodv, czhai}@illinois.edu  
Department of Computer Science  
University of Illinois at Urbana-Champaign  
Urbana, IL 61801

## Abstract

This paper talks about participation of the University of Illinois at Urbana-Champaign (UIUC) in TAC 2008 Opinion Summarization pilot. We mainly explored two ideas: (1) use of entity recognition and parsing to enhance a standard retrieval method for opinion retrieval, and (2) use of a coherence language model to optimize the ordering of sentences in a summary. Our result showed that use of entity recognition during retrieval led to mixed results and re-ordering with coherence language model was not as good as heuristic polarity-based ordering using guiding phrases. Our additional experiments showed that the performance of coherence language model can be different depending on probability function and word selection.

## 1 Introduction

As more people post their opinions on the Web, the Internet is fast becoming a popular and dynamic source for opinions. Due to the large volume and wide range of data, there is a growing need to summarize opinionated documents. The opinion summarization pilot task of the 2008 Text Analysis Conference (TAC) offered a good opportunity to study and evaluate methods for summarizing opinionated text documents. We participated in this task and studied two problems in opinion summarization: (1) Can we improve sentence retrieval by assigning more weights to noun phrases or entities in the query? (2) How can we optimize the ordering of sentences to increase the coherence of a summary?

The pilot Opinion Summarization Task in TAC 2008 asked the participating teams to summarize the opin-

ions expressed in the blog documents. The data used was BLOG06 corpus that had been used earlier in 2006 and 2007 in many TREC tracks. There was a set of target topics on which the participants were evaluated. For each target topic, there was a set of questions, usually dealing with the opinion expressed about the target topic. Also, a set of relevant document ids was released for every target topic. The final goal was to generate a summary of the opinions expressed in the given document set about the particular aspects of the target topic as referred to in the given set of questions. Additionally, there were some sample answer snippets released for every target topic that participating teams could optionally use.

Each summary had a restriction on the length that was based on the number of questions in the target topic. The submitted summaries were evaluated on how many relevant opinions were covered, the length of the summary generated, and the coherence of the summary.

Our overall approach for solving this problem was to break the task into four steps, viz. data preprocessing, sentence retrieval, sentence filtering, and summary generation. Fig. 1 shows an overview of the system. In the data preprocessing step, we eliminated useless tags from the target documents and identified meaningful keywords in the query questions. In the sentence retrieval step, we extracted relevant sentences from target documents with the queries from the previous step. In the sentence filtering step, we filtered the retrieved sentences based on the question and the target sentence polarity. We also removed redundant sentences from the list. In the summary generation step, we ordered retrieved sentences based on polarity and a statistical coherence language model.

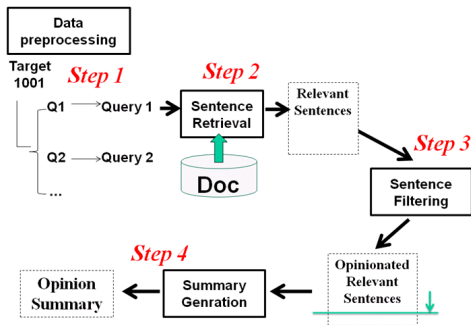


Figure 1: System overview

Specifically, to retrieve relevant sentences, we used standard language modeling approach and inference network retrieval frameworks. We explored an idea of assigning different weights on words in the query question based on linguistic analysis. Since a large portion of questions asked about people’s opinion on some target topics, we hypothesized that named entities and noun phrases should probably have higher weights than other words in the query. So, we used natural language processing techniques to identify named entities and noun phrases, and studied whether assigning more weights on them would improve retrieval accuracy.

After retrieving a set of relevant sentences about the target, we used opinion orientation features to further filter out retrieved sentences that do not match the desired opinion polarity. The most important difference between our opinion summarization task and a normal summarization task is that the query question and target documents are opinionated. If the question asked for positive aspect of one object, the answer also should follow the same orientation. To facilitate this, we used a heuristic dictionary-based method for predicting opinion polarity, and filtered out sentences that do not match the polarity orientation of the question.

The next step in our method was to select a set of representative sentences by removing redundancy. We used some simple methods to assess sentence similarity and used the maximal marginal relevance (MMR) ranking method [1] to select up to a certain number of sentences that were not redundant with respect to each other.

The last step was to generate a readable coherent summary based on the selected non-redundant sentences. The main question we studied in this step was how to optimize the coherence of a summary. We proposed a statistical coherence language model to order sentences. Specifically, we estimated the probability of two words,  $u$  and  $v$ , occurring in a pair of adjacent sentences, using

some natural text as training data. We could then use this model to compute a coherence score of a candidate-ordering of two sentences included in a summary. An optimal order of multiple sentences can be generated using a greedy algorithm by maximizing the pair-wise coherence scores. In addition to this statistical method, we also explored using polarity to separate opinions and heuristically adding additional words to improve readability.

We submitted three runs: a polarity based run (UIUC1), a coherence-based run (UIUC2), and a baseline run (UIUC3). Specific information of each run will be explained in later sections and is summarized in Table 1. None of these runs used the answer snippets provided along with the data because we believe that in a real application of opinion summarization, we would not have such information. All these runs followed the general method described above with key variations in some of the steps. The results of these runs and some follow-up experiments show that increasing the weights of entities and noun phrases in the query improved performance for some queries but could not make significant improvement for some queries. Also, simple statistical coherence optimization did not perform as well as heuristic polarity ordering method with guiding phrases.

In the rest of the paper, we will first describe each step of our method in more detail and then discuss results with our additional experiments.

## 2 Data preprocessing

### 2.1 Document processing

**Remove tags and special characters:** We used the In-dri toolkit<sup>1</sup> for our experiments. We, first, eliminated UNICODE special characters from the documents.

**Extract main contents:** In a blog document, we are mainly interested in the textual contents of the blog entry. So, in addition to eliminating the HTML tags, we also delete the text between two tags if it is too short and does not belong to the main part of the blog. In addition, we also filter out the sections generated by a javascript.

**Sentence splitting:** We use a sentence as the basic unit of retrieval. We believe that a sentence is a natural unit for summarization because it is a smallest unit that can deliver a coherent message, and a summary can be composed of a sequence of sentences. We applied a Sentence Segmentation tool<sup>2</sup> to each preprocessed doc-

<sup>1</sup>From <http://www.lemurproject.org/>

<sup>2</sup>From <http://L2R.cs.uiuc.edu/~cogcomp/atool.php?tkey=SS>

ument. The Sentence Segmentation tool handles honorifics and initials within names, so we were able to obtain a clean sentence boundary using the tool.

## 2.2 Question processing

**Common word elimination:** Before we retrieve sentences using questions as queries, we removed common words from the questions. As most query questions have similar formats, there are words that are not related to the topic, such as “people”, “opinion”, and function words. We applied a simple heuristic on the frequency of words: we counted all the words in the given questions and we filtered out words having a frequency count more than a threshold = 3.

**Find more important words:** We considered named entities and noun phrases as more important words because they convey distinct information compared to common words, and most of questions were asking for opinions about a specific entity. We parsed the question to find all the named entities. We built an NE recognizer using Learning Based Java (LBJ)<sup>3</sup> to identify PERSON, ORGANIZATION, and LOCATION entities. LBJ is a modeling language for the rapid development of learning-based systems, designed for use with the Java programming language. Further, we parsed each question description with a POS tagger<sup>3</sup> to identify noun phrases. For example, in the sentence “The new mobile phone was released today at the tech expo.”, “mobile phone” and “tech expo” are identified as noun phrases.

## 3 Sentence Retrieval

For sentence retrieval, we used the basic retrieval model in Indri, which is a combination of the language modeling and inference network retrieval frameworks, with term weighting [7]. We considered each individual sentence as a document, and built an Indri index over all sentence-documents. We then retrieved the sentences using the filtered list of words from the questions as queries.

With IndriQuery, we can assign different weights for individual words / groups of words. We built an Indri query using a weighted combination of the named entities, noun phrases, and other words from the query to improve the recall. The retrieved sentences formed the initial ranked candidate-list. For the run using different weighting (UIUC1), we used 10:2:1 as weights for named entity : noun phrase : other words. Our approach

gave named entities a much higher weight compared to other terms and also gave noun phrases a slightly higher weight than other words. The weighting scheme was decided based on preliminary experiments of relevance judgment results from TREC Blog retrieval task. We checked the MAP score for different weighting schemes, such as 2:1:1, 5:2:1, 10:2:1, etc., and chose 10:2:1 because it had the highest score. Although both experiments used the same BLOG06 data, we cannot be sure that this is the optimal weighting for the current task at hand because the evaluation was done on a different query set. It is possible that another combination of weights may improve the performance for this task. In fact, our UIUC1 run did not perform as well as the other run (UIUC2) which uses uniform weighting.

## 4 Sentence Filtering

### 4.1 Polarity filtering

We implemented simple dictionary-based polarity decision module. We first obtained lists of words expressing positive and negative sentiments from the following sources: 1) sentiment words<sup>4</sup> and 2) good and bad (positive and negative) sentiment adjectives<sup>5</sup>. We also added some positive and negative words that commonly occur in questions to generate the final positive/negative word list. The following is the complete list of words we added to the existing lists: Positive words: support, like, great, suggest, approve, want; Negative words: negative, complain, complaint, object, oppose. These word lists were then used by the polarity decision module to decide one of four polarities: viz. positively-opinionated, negatively-opinionated, mixed-opinionated, and non opinionated. For an input text segment, the polarity decision module counted the number of words that overlapped with the positive/negative word lists. Depending on the number of positively and negatively inclined words in the target text segment, the polarity module decided whether the text segment is positively, negatively, or mixed opinionated. If the number of positive and negative words were quite skewed, we classified the document as either positive or negative. If we find a negation in a sentence, such as “not” or “never”, we reversed the polarity. If the numbers of positive and negative words were quite even, then we classified the document as a mixed-opinionated document. If there were no common words between the text segment and the posi-

<sup>4</sup>From <http://eqi.org/index.htm>

<sup>5</sup>From <http://www.keepandshare.com/doc/view.php?u=12894>

<sup>3</sup>From <http://L2R.cs.uiuc.edu/~cogcomp/software.php>

tive/negative lists, then we classified the text segment as non opinionated.

We applied the polarity decision module in two steps. First, based on the assumption that the task wanted us to summarize “opinions”, we simply filtered out all non-opinionated sentences from our initial candidate list. Second, we applied opinion decision module to the question as well, and then based on the question polarity, we selected sentences which have the same polarity as the question. For example, if the question asked for “people’s positive opinion about” a particular target, only positively opinionated sentences about the target were retained.

Polarity decision module was also used for paragraph division (See Sec. 5).

## 4.2 Remove redundancy

After the initial retrieval and filtering based on polarity, there can still be many similar sentences in the candidate list. In such a case, we need to keep only one representative and eliminate others from the candidate list. We used the Text Similarity tool<sup>6</sup> to find pairwise sentence similarity. The tool provides lesk value [6] based on counting common words and phrases between two strings. Phrasal matching gets higher scores than single word match. Since longer sentences tend to have more matches, the final scores are normalized by the length of strings.

Based on this similarity score, we applied MMR technique [1] to eliminate redundant sentences. For each sentence  $s$ , we calculated similarity score with each sentence ranked higher than  $s$ , and if the sentence similarity score of  $s$  with any of them is above a set threshold, we removed  $s$  from the candidate list. This way all “duplicate” low-ranked sentences are filtered out. For our submissions, we set the similarity threshold to 0.8.

## 4.3 Cut sentence list up to the limit

Finally, we retained all the sentences in the rank order and pruned the list when the total length of sentences reached the summary size limit (7000 non-white spaced characters per question, each target has one to three questions). Among summaries for 22 targets, 14 summaries reached the limit, and the average length of our final summary results was 11,080 non-white spaced characters (accumulated over all sub-questions in one target).

<sup>6</sup>From <http://www.d.umn.edu/~tpederse/text-similarity.html>

# 5 Summary Generation

## 5.1 Result post-processing

Based on the preliminary experiment results, we noticed that a few noisy phrases such as date and time of blog entry occurred frequently in the blog article collection. Since these did not contribute to a blog opinion, we removed them by matching against the following simple regular expression:

```
Posted_by: ([A-Za-z0-9()]+) *? |\n\\w\\w\\w\\d+, \\d\\d\\d\\d\\d\\d?: \\d\\d: \\d\\d\n[AP]M
```

## 5.2 Paragraph division

Even though we had a good list of relevant sentences as summary, we needed additional structure to make the summary more coherent. We divided the summary into paragraphs based on the questions. The answer to each question was grouped together in one paragraph. Additionally, we needed to distinguish the positive opinions from the negative opinions on a topic. Hence, we further organized the paragraphs by polarity. We kept all sentences with a positive polarity together at the start of the paragraph, followed by all sentences with negative and mixed polarity. This way, we avoided changing the topic of interest and the sentiment orientation too often in the final summary and hence increased coherence.

To make paragraph division more clear and the change of topic more lucid, we also inserted guiding text snippets. At the beginning of each paragraph, we added guiding sentences like “The first question is ...” or “Following are positive opinions ...”.

## 5.3 Statistical coherence optimization

We also experimented with ordering the sentences so as to optimize the statistical coherence. We considered the following simple statistical approach: Suppose  $S = \{s_1, \dots, s_n\}$  is the set of candidate sentences in the summary. Let  $c(s_i, s_j)$  be an asymmetric coherence measure between two sentences  $s_i$  and  $s_j$ . We can then generate a summary by enumerating all the possible orders of these  $n$  sentences and choosing the one that maximizes the overall “pairwise coherence”. That is, we seek for an ordering of sentences  $s'_1, \dots, s'_n$  that maximizes the overall sum  $c(s'_1, s'_2) + c(s'_2, s'_3) + \dots + c(s'_{n-1}, s'_n)$ .

Because enumerating all the permutations is infeasible, we use a greedy algorithm to approximate it. We try to build a sequence of sentences, one pair at a time.

Among all the pairs, we selected the one that has the highest coherence score. Let it be  $\langle s_i, s_j \rangle$ . This forms the first unit in the sentence sequence. Then, we find the highest coherence pair among the ones that can be used to extend the sentence sequence in either direction. In other words, if the current sentence sequence is  $\langle s_i, \dots, s_j \rangle$ , then among all the pairs of the form  $\langle \cdot, s_i \rangle$  or  $\langle s_j, \cdot \rangle$ , we selected the one with the highest coherence score and added it to the sentence sequence.

The coherence measure  $c(s_i, s_j)$  can be defined in many ways. In our experiments, we defined it as the average pointwise mutual information of a word in the first sentence,  $s_i$ , and a word in the second sentence,  $s_j (= s_{i+1})$ . For every word pair  $\langle u, v \rangle$ , we first find  $p_{pmi}(u, v)$ , the pointwise mutual information of word  $u$  in a sentence  $s_i$  and word  $v$  in the following sentence  $s_j$ . Since this is sensitive to the order of sentences,  $p_{pmi}(u, v)$  is asymmetric. We used simple additive smoothing.

$$p_{pmi}(u, v) = \frac{cnt(u, v) + 0.001}{freq(u)freq(v) + 1.0} \quad (1)$$

where  $cnt(u, v)$  is the count of  $u$  and  $v$  in two adjacent sentences in order ( $u$  in  $s_i$  and  $v$  in  $s_j$ ), and  $freq(u)$  is the frequency of  $u$  in the corpus.

Then, we define the asymmetric coherence measure  $c(s_i, s_j)$  over two sentences  $s_i, s_j$ , as

$$c(s_i, s_j) = \sum_{u \in s_i, v \in s_j} \frac{p_{pmi}(u, v)}{|s_i||s_j|}$$

where  $p_{pmi}(u, v)$  is as defined above and estimated using all adjacent sentences in the training data. Computationally, we get all the words in adjacent sentences  $s_i, s_{i+1}$  and accumulate the number of times we observe each combination of words  $(u, v)$  with  $u$  coming from sentence  $s_i$  and  $v$  from sentence  $s_{i+1}$  over the training data. We then normalize the counts to get the pointwise mutual information  $p_{pmi}(u, v)$ .

This approach is similar to the Mirella Lapata’s method [5] with some differences in the definition of the scoring function and the technique used to search for the best order.

## 6 Evaluation

### 6.1 Submission

We submitted three runs, labeled UIUC1, UIUC2, and UIUC3. To evaluate the effect of different techniques,

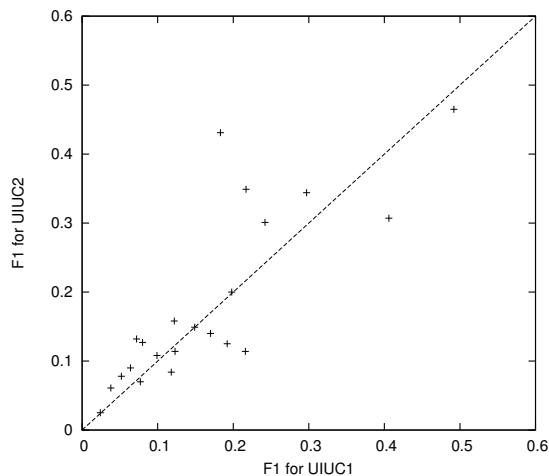


Figure 2: Scatter-plot comparing F1 scores of UIUC1 with UIUC2

we applied different combination of techniques detailed above, for each run. They have been summarized in Table 1.

Techniques	UIUC1 Polarity	UIUC2 Coherence	UIUC3 Baseline
Filter non-opinionated sentences	Off	<b>On</b>	<b>On</b>
Filter sentences having same polarity as questions	<b>On</b>	Off	Off
Paragraphs divided on question	<b>On</b>	Off	Off
Paragraphs divided on polarity (pos/neg/mixed)	<b>On</b>	Off	<b>On</b>
Statistical coherence optimization	Off	<b>On</b>	Off

Table 1: Techniques enabled in the submitted runs. On indicates active and Off indicates inactive features.

## 6.2 Evaluation results and Discussion

### 6.2.1 Official evaluation results

Among the submissions, two runs, UIUC1 and UIUC2, were evaluated in the official TAC 2008 evaluation phase. The official results are summarized in Table 2.

	Recall	Precision	F, with $\beta = 1$
UIUC1	0.309181818	0.125090909	0.165045455
UIUC2	0.390136364	0.124818182	0.180545455

Table 2: Recall, Precision, and F with  $\beta = 1$

	F-Score	Grammaticality	Non-redundancy	Structure/Coherence	Fluency/Readability	Responsiveness
UIUC1	18	31	30	12	17	21
UIUC2	15	20	19	30	33	14

Table 3: Rank among all runs (Total: 36 runs)

	F-Score	Grammaticality	Non-redundancy	Structure/Coherence	Fluency/Readability	Responsiveness
UIUC1	6	15	15	5	6	8
UIUC2	3	9	10	15	16	4

Table 4: Rank among runs that did not use answer-snippet (Total: 19 runs)

Table 2 shows the average scores over 22 summaries based on the pyramid nugget method. Among the two evaluated runs, UIUC2 showed better recall, and both runs had similar precision. The F-score showed UIUC2 had better performance. Fig. 2 shows a scatter-plot of the F1 score for all queries over the two runs.

The pyramid score based evaluation is mainly performed over the contents. On the issue of content selection alone, more sophisticated techniques were applied to UIUC1, such as different weighting for named entities and noun phrases and sentence filtering using polarity matching. However, the performance evaluation result showed that methods of UIUC1 were not as effective as those of UIUC2.

We hypothesize that two factors could contribute to degraded retrieval performance in UIUC1: (1) using high weights for named entity/noun phrases might add noise and harm the retrieval performance, or (2) if the polarity decision module made a wrong decision, it could degrade the overall performance. Since there can be many implicitly opinionated expressions in a blog document, we feel that the polarity decision module did not work as well as it worked with other better-formatted text corpora.

Table 3 and Table 4 show the average score rankings for our runs in six evaluation subcategories. The ranks mentioned in Table 3 are among all the evaluated runs and the ranks in Table 4 are among the submitted runs that did not use the answer-snippets. Since we did not use the provided answer-snippets, we analyze our performance based on the ranks given in Table 4.

Our F-score ranking was about average among all the submitted runs. However, among the runs without using answer-snippet, our runs showed pretty good performance especially on content retrieval. As mentioned above, the F-score measure from pyramid nuggets computes the performance based on content selection. Since the answer-snippet information provides significant clues for relevant information, it is more appropriate to compare content selection performance among runs that did

not use answer snippet. The run, UIUC2, was ranked 3<sup>rd</sup> in F-score among runs that do not use answer-snippets. Therefore, we can say that the content retrieval module in our system has pretty good performance.

In both runs, the non-redundancy scores are low. This is because we applied redundancy removal procedure to sentence list based on each question individually. We extracted sentences for each question separately and did not merge them all for one target until the final summary creation stage. This led to presence of duplicate sentences in one target topic summary which was generated by merging relevant sentences of all the questions in the target. To remove duplicates which came from different question’s relevant sentence list, we need to apply redundancy removal again after merging sentences from each question.

UIUC1 obtained high score in the structure/coherence criterion. Among runs without using answer-snippet, UIUC1 was ranked 5<sup>th</sup>. UIUC1 adopted both question and polarity paragraph division techniques. Moreover, guiding text snippets for each paragraph were also added in UIUC1. Compared to UIUC1, UIUC2 had low structure/coherence score. This suggests that our statistical coherence optimization technique did not work as well as the heuristic polarity ordering with guiding phrases. Using a greedy algorithm that looks at only the next sentence, limits the algorithm from obtaining clues from non-neighboring sentences, and this could harm the coherence score considerably.

Finally, UIUC2 had a better overall responsiveness score than UIUC1. In terms of ranking, UIUC1 showed average performance and UIUC2 showed fairly good performance. Among the runs that did not use answer-snippets, UIUC2 was ranked 4<sup>th</sup>.

## 6.2.2 Additional experiments on entity feature weighting

To examine the relationship between content retrieval performance and weighting on entity features such as

named entities and noun phrases, we performed some additional experiments. There has been previous work on trying differential weighting for named entities. Hassel’s experiment on Swedish text [2] showed the importance of careful usage of named entity feature in automatic summarization. He pointed out that high weighting on named entities can degrade information retrieval performance because it can cause reference errors and prioritize elaborate sentences over sentences having general information.

In our experiments, we evaluated the effect of weighting named entities by computing how well the system retrieved relevant nuggets. We calculated precision, recall, and F-scores based on the pyramid nuggets on the intermediate sentence extraction results before applying the polarity module so as to exclude the effect of the polarity module. Instead of applying our redundancy check, we eliminated obviously redundant sentences by removing adjacent duplicates. For each sentence in the ranked sentence list, we checked for presence of any of the nuggets. Since the nugget text given as part of pyramid nuggets’ data is a description rather than the exact nugget, it does not exactly match with the string in the extracted sentences. We approximated the presence of nuggets by calculating similarity between the sentence and provided nugget text. After stemming and removing stopwords from the target sentence and nugget text, we used the Text Similarity tool<sup>6</sup> to count the overlap of words. If the similarity score between the extracted sentence and nuggets in the same target was above the certain threshold, we believed that the sentence matched the corresponding nugget.

The similarity score threshold was tuned with the evaluation results of our TAC submission that listed which nuggets were found in our submission result. Based on the preliminary experiment with the first two target results, we tuned the similarity threshold in increments of 0.05 and decided the threshold to be 0.35.

We experimented with 2 sample settings of weights: in Fig. 3 and Fig. 4, “1-1-1” indicates equal weights to all components and “10-2-1” indicates that the named entities were weighted 5 times higher than noun phrases, which were weighted twice as much as other words in the query. Fig. 3 shows the scatter plot comparing the F1 scores of retrieval runs with 1-1-1 and 10-2-1 weighting schemes. We see that for many queries, the F-scores are higher when the named entities and noun phrases are weighted higher. However, the improvement is not uniform across queries and does not result in any significant improvement. Further, for very few queries, the performance degrades when named entities are weighted

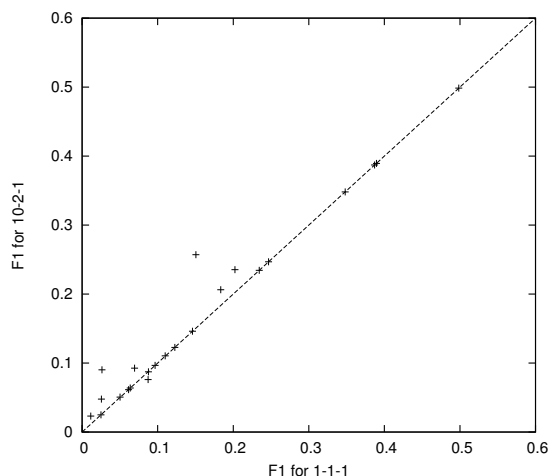


Figure 3: Scatter-plot comparing F1 scores of retrieval runs with 1-1-1 and 10-2-1 weighting schemes

higher. Fig. 4 shows the F-scores of nuggets found in two sample queries, 1022 and 1027. For these queries, the non-uniform weighting scheme performed better than the uniform scheme. The plots show that more sentences with relevant nuggets are retrieved when the named entities are given higher weights. Typical summary lengths are 4000 characters. The summaries were truncated based on the maximum allowable length in TAC (i.e. 7000× number of questions in a target). Based on these experiments, it seems that weighting named entities higher than other words improves performance to some extent, but the overall improvement is not very significant.

### 6.2.3 Additional experiments on polarity module

To check the performance of our polarity module, we performed additional experiments on the tagged data set used in Minqing Hu and Bing Liu’s previous work [3, 4]. They used 14 product review from Amazon<sup>7</sup>, and the sentences in the data set have manually generated tags indicating features and sentiments (positive and negative). We collected all the positive and negative sentences, without segregating product and feature-specific reviews, and checked if our polarity module can classify those sentences correctly. Because Amazon reviews are free-format text posted by web users, we can assume that they have similar characteristics to blog posts. Table 5 shows the classification results. The true classification rate is 38% and the exact opposite misclassification rate is 16%. These results show that the performance of po-

<sup>7</sup><http://www.amazon.com>

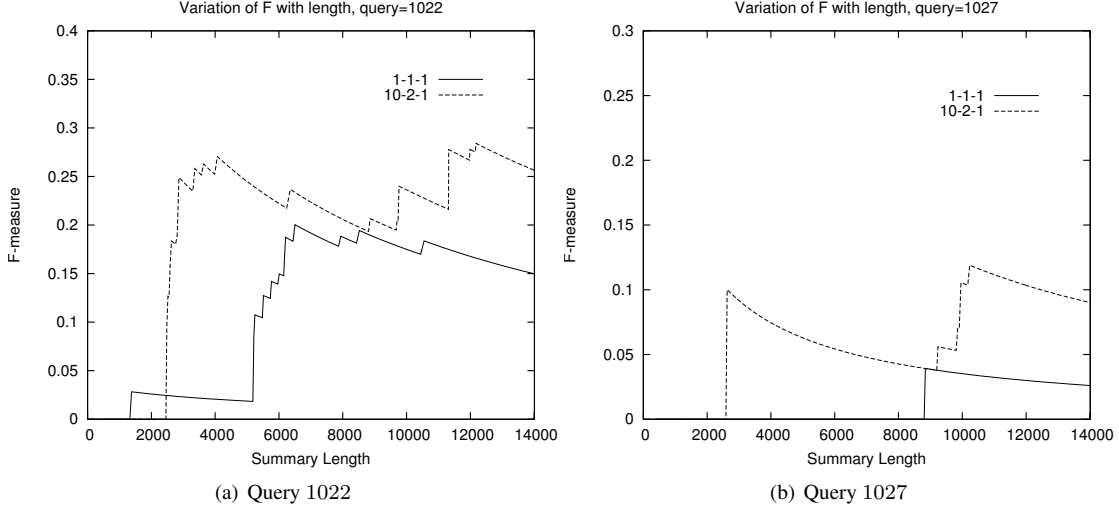


Figure 4: Variation of F-scores with length

larity module was low and may have caused some useful sentences to be filtered out in our final summary.

Classification result	Positive	Negative
NonOpinionated	1063	598
Positive	1363	371
Negative	383	412
Mixed	296	210
Total	3105	1591
Exact Match	1363/3105=0.44	412/1591=0.26
	(1363+412)/(3105+1591)=0.38	
Exact Opposite	383/3105=0.12	371/1591=0.23
	(383+371)/(3105+1591)=0.16	

Table 5: Polarity module performance (# sentences)

## 6.2.4 Additional experiments on coherence optimization

We conducted additional experiments on coherence optimization module to check its performance. We measured the performance of coherence module by comparing the order of sentences given by the coherence module with “ideally-ordered” sentences. Since there is no “ideal” order in a summary, we assumed that the sentence ordering in the original blog document is ideal. Table 6 shows the split of the sentences from relevant blog documents into training set (used to train the coherence module) and the test set.

In addition to the basic optimization that we applied to the submitted run, we tried different coherence functions. Instead of using pointwise mutual information,

	Number of sentences
Given data set	52383
Training set	38561 (73.6%)
Test set	13822 (26.4%)

Table 6: Coherence optimization experiment data set

$p_{pmi}(u, v)$  (Eq. 1), we first tried to use strict joint probability,  $p_{sjp}$ , defined as

$$p_{sjp}(u, v) = \frac{cnt(u, v) + 1}{\sum_{u, v} cnt(u, v) + |\text{dictionary}|^2} \quad (2)$$

Secondly, we suspected that the bias of pointwise mutual information caused the low performance of the original algorithm, we tried to use general mutual information,  $p_{gmi}$ , defined as

$$p_{gmi}(u, v) = p(u, v) \log \frac{p(u, v)}{p(u)p(v)} + p(\neg u, v) \log \frac{p(\neg u, v)}{p(\neg u)p(v)} + p(u, \neg v) \log \frac{p(u, \neg v)}{p(u)p(\neg v)} + p(\neg u, \neg v) \log \frac{p(\neg u, \neg v)}{p(\neg u)p(\neg v)} \quad (3)$$

where  $(\neg u, v)$  means “ $u$  does not exist in sentence  $i$  and  $v$  exists in sentence  $i + 1$ ” (similarly  $(u, \neg v)$  and  $(\neg u, \neg v)$ ), and

$$p(u, v) = \frac{cnt(u, v) + 0.25}{N + 1}$$

$$p(\neg u, v) = \frac{cnt(\neg u, v) + 0.25}{N + 1}$$

$$p(u, \neg v) = \frac{cnt(u, \neg v) + 0.25}{N + 1}$$



Selection of training words	Coherence score		
	Strict Joint Probability	Mutual Information	Pointwise Mutual Information
No omission	0.022259	0.041651	<b>0.056063</b>
Omitted 5% most frequent words (counts > 33)	0.031314	0.049498	0.051460
Omitted 10% most frequent words (counts > 11)	0.028899	0.046103	0.045725
Omitted 20% most frequent words (counts > 6)	0.020373	0.034785	0.032219
Omitted 40% most frequent words (counts > 2)	0.019845	0.021882	0.022259
Omitted 50% most frequent words (counts > 1)	0.018486	0.020599	0.019694
Omitted 95% least frequent words (counts < 33)	0.022108	0.032898	0.044065
Omitted 90% least frequent words (counts < 14)	0.022108	0.032823	0.049045
Omitted 80% least frequent words (counts < 6)	0.020599	0.037048	0.054931
Omitted 60% least frequent words (counts < 2)	0.022410	0.041802	0.056591
Omitted stopwords <sup>a</sup>	0.036746	0.054554	<b>0.057119</b>
Omitted non-stopwords	0.017958	0.022863	0.020750
Omitted connecting words <sup>b</sup>	0.022108	0.043235	0.057421
Omitted non-connecting words	0.017958	0.019241	0.017958
Omitted transitive words <sup>c</sup>	0.020599	0.042405	0.056968
Omitted non-transitive words	0.018939	0.017279	0.019694

<sup>a</sup><http://www.lextek.com/manuals/onix/stopwords1.html>

<sup>b</sup><http://www.rsc.ccc.tn.us/owl/writingcenter/OWL/Connect.html>

<sup>c</sup><http://www.virtualsalt.com/transits.htm>

Table 7: Coherence optimization result. The coherence score of the baseline random ordering was 0.01586.

$$p(\neg u, \neg v) = \frac{cnt(\neg u, \neg v) + 0.25}{N + 1}$$

$$N = cnt(u, v) + cnt(\neg u, v) + cnt(u, \neg v) + cnt(\neg u, \neg v)$$

For unseen pairs,

$$p_{gmi}(u, v) = 0.5 \times \min(\text{pairs seen in training})$$

We also ran the optimization experiments after first excluding frequent words, then excluding rare words. Words having high frequency tend to be stopwords and are not related to sentence ordering, so they can be removed. Some rare signal words, such as “first” or “second”, may be significant clues for deciding sentence coherence. Therefore, ignoring such rare clues can affect the performance. Based on this hypothesis, we excluded frequent and rare words when we train our coherence model.

To evaluate the model, we used a very-conservative strict pair matching formulation. We set the evaluation score to be the ratio of number of correct adjacent sentence pair to the number of total adjacent sentence pair. The range of the used measurement is 0 to 1, and 1 means exactly same ordering as the ideal document.

Table 7 shows the experiment results. For the basic setup, we obtained the coherence score of 0.056063, which is 3.5 times larger than that of random ordering, 0.01586. However, it is still a low score as the range of measure is [0, 1]. This result shows that the performance of current probabilistic coherence optimization module

is poor. This corresponds with the low coherence score results of the UIUC2 run.

Among different coherence functions, pointwise mutual information showed the best performance. Excluding some words in sentence ordering also showed some positive results. By omitting frequent words, we could obtain a higher score in strict joint probability and mutual information functions. After omitting stopwords, we could obtain performance increase in all three functions. Similar to stopwords, eliminating connecting words and transitive words also helped to improve performance generally. However, because connecting words and transitive words contain key words for sentence ordering, such as “first”, “second” or “next”, performance improvement was lower than that when eliminating stopwords. The performance even decreased in strict joint probability when eliminating transitive words.

We can infer that the superiority of pointwise mutual information came from its effective handling on frequent words. When  $u$  and  $v$  are frequent, the denominator in Eq. 1 increases, and it penalizes the pointwise mutual information,  $p_{pmi}(u, v)$ . However, the denominator in Eq. 2 for strict joint probability,  $p_{sjp}(u, v)$ , is fixed regardless of frequency of  $u$  and  $v$ . Similarly, mutual information,  $p_{gmi}(u, v)$ , is also unaffected by the frequency of  $u$  and  $v$ , since  $N$  (in Eq. 3) sums up all cases of  $u$  and  $v$  being present. Therefore, when eliminating frequent words, we could see performance improvement in strict joint probability and mutual information, while we could not see it in pointwise mutual infor-

$u$	$v$	$p_{sjp}(u, v)$
the	the	4.40E-006
to	the	3.07E-006
the	to	3.06E-006
the	of	2.87E-006
the	and	2.87E-006
of	the	2.85E-006
and	the	2.83E-006
a	the	2.68E-006
the	a	2.67E-006
...	...	...

Table 8: Top ranked pairs using strict joint probability

mation which already penalized frequent words. When eliminating stopwords, performance improved in all functions including pointwise mutual information. This is because stopwords list contains not-meaningful words which were not frequent in the training set. The performance improvements of strict joint probability and mutual information when eliminating stopwords are still larger than that of pointwise mutual information. Table 8 confirms our intuition that the top ranked  $p_{sjp}(u, v)$  pairs are indeed stopword pairs.

## 7 Conclusion

Our submission to the opinion summarization pilot task at TAC 2008 gave us a great opportunity to try out various techniques and ideas to tackle the problem. Contrary to our initial expectation, the run having higher weighting on named entities and noun phrases, UIUC1, showed lower retrieval performance. Our experiments showed that differential weighting of named entities and noun phrases helped retrieve more relevant sentences and improve the nugget score in some cases. By evaluating our polarity module, we showed the possibility that using imperfect polarity decision model is the reason that UIUC1 had a degraded performance in content selection. The run that adopted statistics coherence optimization module showed low coherence evaluation. The experiment results on the performance of statistical optimization also showed that sentence ordering based on only pointwise mutual information between words is not effective. However, by comparing coherence ordering results using strict joint probability and mutual information, we could learn that using pointwise mutual information was the better choice. Moreover, additional experiments with selected words based on frequency of words showed the possibility of better performance by tuning characteristics of training set.

As part of future analysis, we can consider using less

strict coherence measures and different ordering methods. We also need to examine more factors affecting retrieval and sentence ordering performance. It will also be instructive to analyze what weighting scheme would perform well based on the query.

Our result analysis showed that we can improve further by using summary structuring methods and content selection methods from UIUC2 with additional redundancy removal after merging sentences from each question. We also need to improve the precision of the polarity decision module by preparing better-selected positive/negative word lists or applying different decision methods.

## References

- [1] J. Carbonell and J. Goldstein. The use of MMR, diversity-based reranking for reordering documents and producing summaries. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 335–336, 1998.
- [2] M. Hassel. Exploitation of named entities in automatic text summarization for swedish. In *Proceedings of NODALIDA '03- 14th Nordic Conference on Computational Linguistics*, 2003.
- [3] M. Hu and B. Liu. Mining and summarizing customer reviews. In *Proceedings of International Conference on Knowledge Discovery and Data Mining*, 2004.
- [4] M. Hu and B. Liu. Mining opinion features in customer reviews. In *Proceedings of National Conference on Artificial Intelligence*, 2004.
- [5] M. Lapata. Probabilistic text structuring: Experiments with sentence ordering. In *Proceedings of the 41st Meeting of the Association of Computational Linguistics*, 2003.
- [6] M. Lesk. Automatic Sense Disambiguation using Machine Readable Dictionaries: How to tell a pine cone from a ice cream cone. In *Proceedings of the 5th Annual International Conference on Systems Documentation*, 1986.
- [7] H. T. Trevor Strohman, Donald Metzler and W. B. Croft. Indri: A language-model based search engine for complex queries. In *Proceedings of the International Conference on Intelligent Analysis*, 2005.