

Question Answering Using “Common Sense” and Utility Maximization Principle

Tomoyosi Akiba

Department of Information and Computer Sciences, Toyohashi University of Technology
1-1-1 Hibarigaoka, Tenpaku-cho, Toyohashi-shi, 441-8580, JAPAN
akiba@cl.ics.tut.ac.jp

Katunobu Itou

Graduate School of Information Science, Nagoya University
1 Furo-cho, Nagoya, 464-8603, JAPAN

Atsushi Fujii

Graduate School of Library, Information and Media Studies, University of Tsukuba
1-2 Kasuga, Tsukuba, 305-8550, JAPAN

Abstract

In this paper, we propose two new methods targeting NTCIR-4 QAC2. First, we use knowledge resembling “common sense” for question answering purposes. For example, the length of a runway in an airport must be a few kilometers, but a few centimeters. In practice, we use specific types of information latent in document collections to verify the correctness of each answer candidate. Second, we use the utility maximization principle to determine the appropriate number of answers for a list question. We estimate the expected value of the evaluation score, on the basis of the probability scores for multiple answer candidates. We show the effectiveness of our methods by means of experiments.

1 Introduction

This paper describes our question answering systems participated in all of the subtasks, i.e., subtasks 1, 2, and 3, of Question Answering Challenge (QAC) 2 carried out at NTCIR-4. In order to participate QAC2, our systems have been developed from the scratch with several new methods. Among them, two outstanding methods are proposed.

Human commonly uses a kind of knowledge called “common sense” to solve problems. For example, the length of an airport’s runway should be some kilometers and should not be some centimeters. That kind of knowledge can be used to help selecting the appropriate answers for question answering. Common sense is based on human experience. Because large-scale document collections, or corpora, include many cases

about the world, it can be used as the knowledge resource including common sense. One of our methods utilizes corpora as such a knowledge resource without any preprocessing of knowledge extraction.

Another novel method is about selecting a set of answers for list questions, which are dealt in QAC2 subtask 2 and 3. The method applies the decision theory to select the optimal set of answers that maximize the resulting utility function.

Section 2 describes our definition of question answering as a search problem. Section 3 describes the method of utilizing common sense in corpora. Section 4 describes our method to deal with the context of answer candidates. Section 5 describes the method of selecting the set of answers for list questions. Section 6 describes the method to deal with the series of related questions. Section 7 describes some experimental results of our proposed methods.

2 Question Answering as a Search Problem

The question answering process is often seen as the sequence of the question analysis, the relevant document (or passage) retrieval, answer extraction and answer selection processes. In this paper, we recast these processes as a search problem.

Question Answering (1) Given query q and document set D , from all substrings in D , $S = \{(d, p_s, p_f) | d \in D, p_s < p_f; p_s$ and p_f are positions in $d\}$, by using an evaluation function $L(a|q)$ defined on $a \in S$, select \hat{a} that maximizes $L(a|q)$.

Question Answering (2) Given query q and document set D , from all substrings in D , $S = \{(d, p_s, p_f) | d \in D, p_s < p_f; p_s \text{ and } p_f \text{ are positions in } d\}$, by using an evaluation function $L(A|q)$ defined on $A \in 2^S$, select \hat{A} that maximizes $L(A|q)$.

Question Answering (1) is the problem of finding a single best answer, which corresponds to the factoid question in TREC [15] and the subtask 1 in NTCIR Question Answering Challenge (QAC). **Question Answering (2)** is the problem of finding one or more answers exhaustively and exactly, which corresponds to the list question in TREC and the subtask 2 in NTCIR QAC. Because the search space of question answering is vast, an approximation technique is needed to limit the search space. We search only the document fragments that are relevant to the question.

In existing question answering systems, the evaluation function L is constructed by one or more properties as below:

- a. the property for each answer candidate, and
- b. the property for context (surrounding text) of each answer candidate.

For the property (a), many systems today examine the agreement between semantic category from the question and that from the answer candidate. These two categories are typically extracted by using question analysis and the named entity extraction respectively.

For the property (b), the similarity between the question and the context of the answer candidates is examined. This process can be seen as selecting an appropriate passage to extract the answer candidate. This passage retrieval is one of the common research topics for question answering [14].

The next two subsections will explain our approach for constructing the two properties, respectively.

3 The Property concerning Answer Candidates

3.1 Previous Work

For the property (a) in Section 2, a number of existing QA systems evaluate the agreement between the semantic category determined by a question and that for each answer candidate. Named entity (NE) extraction is commonly used to obtain the category of the answer candidate. The categorization adopted by the named entity extraction is an important function for question answering. In general, the more detailed categorizations a system uses, the better performance it achieves. For example, the system proposed by Lee and Lee [9], which used 62 categories, performed best among the participants in QAC1 subtask 1. However, the NE-based method is associated with the following problems.

- The development of knowledge base used for the named entity extraction is expensive.
- The accuracy of the named entity extraction affects on the performance of the question answering. An excessively detailed categorization potentially can reduce the performance of question answering.

In view of these problems, we propose an NE-free method¹ to evaluate the agreement between the semantic categories of the question and an answer candidate. We elaborate on our method in the rest of this section.

3.2 Evaluating Semantic Relations using a Corpus

In TREC and NTCIR, a question often consists of the word or phrase that directly express the semantic constraint for the possible answer. For example, the query “2000 nen no NHK taiga dorama wa nan desu ka?” (What was the NHK roman-fleuve TV drama broadcasted in 2000?) implies that the answer should be the instance of “NHK taiga dorama” (the NHK roman-fleuve TV drama). In addition the phrase “kioku youryou” (memory capacity) in the question “ZIP no kioku youryou wa ikutsu desu ka?” (What is the capacity of ZIP?) implies that the answer should be a numerical expression followed by a unit expression, such as “MB”(mega byte) and “GB”(giga byte). We shall call these central words and phrases that directly express the semantic constraint for the possible answers “Question Focus (QF).”

Our method examines the presence of the semantic relation between the QF extracted from a question and each answer candidate. The result is reflected to the evaluation function L of the QA search problem (see Section 2). The examination is performed with the specific language patterns. Figure 1 illustrates the process of the method.

Our method is advantageous, because any preprocessing to extract knowledge (e.g. one for NE extraction). In addition, because we do not predefine the semantic categories, the granularity of semantic constraints is determined dynamically driven by each question.

3.3 Implementation

We first retrieve the documents including both the QF and an answer candidate and identify the specific linguistic patterns including them in the retrieved documents. For this we use hand-crafted regular expressions based on the surface expression and lexico-syntactic information obtained by Japanese morphological analysis.

¹We used the method instead of using the named entity extraction with detailed categories in our system participated in QAC2. However, the method can be used together with the conventional named entity extraction with detailed categories and it may improve the performance of question answering.

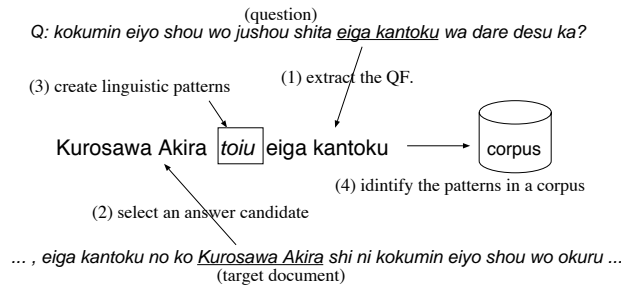


Figure 1. The process of the proposed method.

Various language units can be considered as QFs, which are nouns compound nouns, phrases, and complex phrases including relatives. For example, the question “kokumin eiyo shou wo jushou shita eiga kantoku wa dare desu ka?” (Who is the film director who received the national honorary prize?) includes multiple candidates of the QF, which are “kantoku” (director), “eiga kantoku” (film director), “jushou shita eiga kantoku” (the film director who have received a prize), and “kokumin eiyo shou wo jushou shita eiga kantoku” (the film director who have received the national honorary prize). In general, using a larger unit of word sequence as a QF leads richer information about the category of the answer and, therefore, improves the precision of the answer extraction, although the coverage is decreased. If we use complex phrases including relatives as QF’s, the total performance is dependent of deep NLP, such as parsing.

In view of the above problem, we experimentally use the following unit as a QF.

- a largest (compound) noun, or a longest sequence of nouns.
- exclude vague and non-informative nouns that have the same specificity with WH-words , e.g., “mono” (thing), “namae” (name), “jinbutsu” (person), and “basho” (place).

For example, we select “eiga kantoku”(film director) as a QF from the question above.

3.4 Evaluating Name Expressions

A factoid question expects either the name or numerical expression as the answer. For the answers that are name expressions, the hypernym-hyponym relation between the QF and each answer candidate is examined. We use the lexico-syntactic patterns for this test, e.g., “AC *toiu* QF” (QF such as AC), “AC *igaino* QF” (QF other than AC), “QF · AC”, in which QF and AC are the surface expressions of the question focus and the answer candidate, respectively. Hearst [7] and other related works used similar patterns for extracting semantic relations from corpora. However, we

use such patterns not for extracting the relations itself but for directly examining the relation using corpora without extraction.

3.5 Evaluating Numerical Expression

A Japanese numerical expression consists of the sequence of numbers followed by the unit expression, such as “250 yen”. For each part our method examines their validity as the answer.

3.5.1 Evaluating Unit Expression

The lexico-syntactic patterns can also be used for examining the semantic relation between the QF and the unit expression used in the answer candidate. We use the regular expression “QF AUX* num UNIT”, in which AUX is an auxiliary word, num is a number (a sequence of digits) and UNIT is a unit expression in an answer candidate. Murata et al. [11] used similar patterns for extracting semantic relations between the QF and the unit expression from corpora, though, again, we used such patterns not for extracting but for examining the relations.

3.6 Additional String Based Method to Evaluate Relations

3.6.1 Evaluating Numbers

The set of numbers appeared with a topic (QF) in corpora can be considered as the common cases of values about the topic. Thus, by examining the proximity between the number part of the answer candidate and that in corpora appeared in the similar context, we can check whether the number of the answer candidate is appropriate for the topic. Using the patterns of evaluating the unit expression mentioned above, we extracted the set of numbers. Because this process can be seen as random sampling, the set of the numbers follows the Gaussian distribution. We made and tested the hypothesis that the number part of the answer candidate is a sample from the distribution, in order to select the appropriate answers. To put it more concretely, we calculated the minimum critical rate from the Gaussian distribution that the hypothesis would not be rejected and reflected it to the evaluation function L . Figure 2 illustrates the process of the method.

In addition to the corpus based method described above, the following string based method is also utilized.

For evaluating name expression as an answer, the appearance of the QF in each answer candidate is investigated. A Japanese compound noun often has its head at the end of the noun sequence. For example, When the string “yama” (mountain) is extracted as a QF, an answer candidate “Fuji yama” (Mt. Fuji) can be judged as an instance of the QF, because it has the substring equal to the QF in its end.

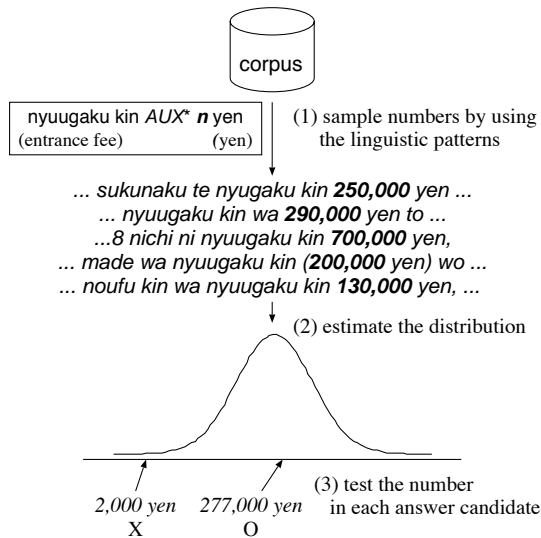


Figure 2. The process of evaluating numbers.

For evaluating numerical expression as an answer, the agreement between the unit expression appeared in the question and that in each answer candidate is investigated. A Japanese question asked for a numerical expression often includes the unit expression following a WH-word. For example, since the question “*tyunijia no jinkou wa nan nin desu ka?*” (How is the population in Tunisia?) has the unit expression “nin”(the suffix of numerical expression counting humans) following the WH-word “*nan*” (what), the answer should be a numerical expression whose unit expression is “nin”.

3.7 Related Works

Acquisition and Verification

A lot of works initiated by Hearst [7] has been focused on extraction semantic relations from unstructured text. In particular, Fleischman et al.[6] utilized the extracted relations for question answering. These previous works was the method of “acquisition” of knowledge from corpora, while our method was “verification” of the specific relations using corpora.

Generally speaking, “acquisition” is the process that seeks for all the pairs of objects that fulfill the given constraints, while “verification” is the process that seeks if one specific pair of objects fulfill the constraints. One of the problems of “acquisition” in practical use is that it needs a great deal of computational and spatial costs. Some limitation on the extent of the acquisition is indispensable in practical use. For example, the unit of just one word rather than any word sequence is often used as the target of acquisition. On the other hand, “verification” is much less expensive when the specific pair of objects is already known. Because the QFs and the answer candidates are known,

the “verification” can be applied effectively in the process of question answering.

Question Focus

The notion of the question focus was first introduced by Moldovan et al.[10]. They utilized the QFs for answering the query that has “what” as the query term and is ambiguous in extracting the answer type. Ittycheriah et al.[8] emphasized the answers who had hypernym or hyponym relationship in WordNet with the QF. Prager et al.[12] focused on answering “What is X?” question. The WordNet was consulted from the extracted QF and the hypernyms were considered as the answer candidates of the what-is question.

Using World Wide Web

Several works made use of a large-scale text collection, namely the World Wide Web, for question answering[4, 5]. These works took advantage of the vast amount of text as the target of extracting answers. On the other hand, our method utilized corpora as general knowledge resources. Therefore the method using WWW can be applied with our method to improve the performance of question answering.

4 The Propertie concerning the Context of Answer Candidates

4.1 Selecting Optimal Context

Selecting the length of the context, or selecting passage in other words, is one of the common research topics for question answering[14]. The context is used to calculate the similarity against the query. Some systems use a sentence as the context, while other systems use a paragraph. The longer the context is selected, the more candidates can be picked up and be considered as the answer. It raises the recall of the answer, while it reduces the precision because the more wrong candidates are also picked up.

Another difficulty arises if we look into headlines of newspaper articles to extract answer candidates in addition to contents of the articles. Because a headline of an article is apart from the content, it does not have the neighbor sentences. Whole the content can be considered as the context of the headline, though using such a long context (whole the article) reduces the precision of the answers.

Considering the examination above, we adopted dynamic passage selection used for selecting the optimal context. Suppose we are going to select the context of an answer candidate a , who belongs to a sentence s_i of a document (a content of an article) $d = s_1 s_2 \cdots s_i \cdots s_n$. Let $s'_i = s_i - \{a\}^2$, h be the headline of d , and t be the string “今年今月今日” (this

²We approximated $s'_i \approx s_i$ in order to reduce the cost of calculation in our participated system.

year, this month, today). Given a number $k > 0$, let $S_i = \{h, t, s_{i-k}, \dots, s_{i-1}, s'_i, s_{i+1}, \dots, s_{i+k}\}$. The optimal context \hat{C}_i is selected from $C_i \in 2^{S_i}$ by maximizing the following evaluation measure $F(C_i)$.

$$R(C_i) = \frac{\text{Score}(q \wedge C_i)}{\text{Score}(q)}$$

$$P(C_i) = \frac{\text{Score}(q \wedge C_i)}{\text{Score}(C_i)}$$

$$F(C_i) = \frac{1 + \beta^2}{\frac{\beta^2}{R} + \frac{1}{P}}$$

where $\text{Score}(A)$ is a sum of the IDF's of uni-gram and bi-gram in the word sequence A and $\text{Score}(A \wedge B)$ is a sum of the IDF's of uni-gram and bi-gram appeared commonly in A and B , which will be defined in the next subsection. The context of headline is selected from C_i for $i = 1 \dots n$ that maximize $F(C_i)$.

We used $k = 1$ for our system participated in QAC2. The evaluation measure F corresponds to the (weighted) F-measure often used in IR research. We chose $\beta > 1$ to emphasize the recall for the selection.

4.2 Similarity Calculation using Content Word bi-gram

In order to select the appropriate passage, the measure of similarity between the query and the passage must be constructed. The most basic measure used in many QA systems is word-based, which counts the number, or sums up the weighted values like TF-IDFs, of common words that appear in both the query and the passage, like a document retrieval manner. However, it fails to capture the similarity of the higher order relations of the word sequences. On the other hand, some systems[13] adopt the measure of similarity between the syntactic structures of the query and the passage. The disadvantages of such an approach include that the measure needs expensive syntactic parsing and that the accuracy of the parsing becomes critical for the result.

We extended the simple word-base similarity measure to utilize the content word bi-gram. In addition to the sum of the IDF's of the common words (uni-gram) both in the query and the passage, the extra IDF's of neighboring content words (bi-gram), allowing some sort of functional word like “ \emptyset ” or symbols like “ \cdot ” between them, are given if these word sequence is commonly appeared both in the query and the passage. An example of the calculation is shown in figure 3. The advantages of this measure include that it can capture some higher order relations of the word sequences including word orders, and that it does not need expensive NLP components like parsing.

Q: 2004年の大河ドラマは何ですか？

case1: “今年大河ドラマ「新選組」...”

IDF(“今年”(=“2004年”)) + IDF(“大河”) + IDF(“ドラマ”) + IDF(“今年大河”) + IDF(“大河ドラマ”)

case2: “ドラマ「大河の一滴」は2004年...”

IDF(“ドラマ”) + IDF(“大河”) + IDF(“2004年”)

Figure 3. Similarity using content word bi-gram

5 Extracting Set of Answers for subtask 2 and 3

5.1 Removing Duplication from the Answers

In QAC2 subtask 2 and 3, it is required to extract a set of answers that has no duplication. If the set include the n duplicated answers, $n - 1$ answers are considered to be incorrect answers.

Our system adopted two methods to remove the duplications. The one is the character-based method to find the answer candidates that are the abbreviation of another candidate. In special case, it also removes the candidates that have same expression of another candidate.

The other is the pattern matching based method to find the pair of candidates that indicate same object. The patterns like “ $AC1(AC2)$ ” are used to find the pair from the target articles that the pair has been extracted.

In both case, the top scored candidate was survived if there found the set of duplicated answer candidates.

5.2 Selecting Set of Answers by using Expected Utility

In order to select the set of answers, we calculated expected utility of the evaluation measure used in the subtasks, i.e., F-measure, and select the best strategy that maximize the expected utility.

Suppose the extracted answer candidates from the query q are $C = \{c_1, c_2, \dots, c_n\}$, each of which has the plausibility score $L(c_i|q)$ calculated by the evaluation function mentioned in section 2. Suppose also that the sequence $c_1 \dots c_n$ is sorted in descending order by the score $L(c_i|q)$. Let A be the set of correct answers. We make an assumption that all the answers are included in C , i.e., $A \subset C$. This assumption is approximately fulfilled when sufficiently large n is selected.

Suppose the number of correct answers $|A|$ is known to be i . Let a set of answers $C_s \subset C$ be selected for evaluation. Using the number of correct answers $|A|$, the number of selected answers $|C_s|$, and the number of selected correct answers $|A \cap C_s|$, the

F-measure $F(|A|, |C_s|, |A \cap C_s|)$ is calculated as follows.

$$F(|A|, |C_s|, |A \cap C_s|) = \frac{2 \cdot \frac{|A \cap C_s|}{|A|} \cdot \frac{|A \cap C_s|}{|C_s|}}{\frac{|A \cap C_s|}{|A|} + \frac{|A \cap C_s|}{|C_s|}}$$

Therefore the expected value of the F-measure $E(C_s | |A| = i)$ when selecting the answer set C_s given $|A| = i$ can be calculated as follows.

$$E(C_s | |A| = i) = \sum_{k=1}^i P(C_s, k | |A| = i) F(i, |C_s|, k)$$

where $P(C_s, k | |A| = i)$ is the conditional probability that the just k correct answers are included in the set C_s given that the number of correct answers $|A|$ is i .

The conditional probability $P(C_s, k | |A| = i)$ can be approximately calculated as an extension of the Hypergeometric Distribution by following formula.

$$P(C_s, k | |A| = i) = \frac{\sum_{E \in \text{sel}(C_s, k)} \sum_{F \in \text{sel}(C - C_s, i - k)} p(E \cup F)}{\sum_{D \in \text{sel}(C, i)} p(D)}$$

where $\text{sel}(D, i)$ is the set of the combination of selecting i elements from the set D , and $p(D)$ is calculated as follows.

$$p(D) = \sum_{x \in D} f(L(x|q))$$

where f is a non-decrement function defined in $x \geq 0$ that is introduced to revise the value of evaluation function. The values of evaluation function L are meaningful in their ordering but not in their quantities, thus the revision of the values is indispensable.

Until here, we suppose the number of correct answer $|A|$ is known. Using the prior probability $P(|A| = i)$, the expected value $E(C_s)$ can be calculated as follows.

$$E(C_s) = \begin{cases} P(|A| = 0) \cdot 1 & \text{if } C_s = \{\} \\ \sum_{i \geq 1} P(|A| = i) E(C_s | |A| = i) & \\ \text{otherwise} & \end{cases}$$

The best answer set \hat{C}_s can be selected by using $E(C_s)$ as follows.

$$\hat{C}_s = \text{argmax}_{C_s \subset C} E(C_s)$$

Note that because the probability $P(C_s, k | |A| = i)$ that approximately calculated above is independent against the combination of the elements of C_s , the possible best strategy can be obtained among the j -best candidates in their scores, i.e., either $C_s = \{\}$ or $C_s = \{c_1 \cdots c_j\}$ for $j \geq 1$. The different selection of the probability model, including the dependent model against the combination, would result in the different selection of C_s .

In the calculation above, the revision function f and the prior probability $P(|A| = i)$ must be specified. Additionally we gave the upper limit of the number of the selected answer J where $C_s = \{c_1 \cdots c_j\} (1 \leq j \leq J)$. Our two systems participated in QAC2 subtask 2 differed with these parameters. The first system used $f(x) = x^2$ and $J = 5$. The query analysis module was used to expect the number of correct answers e . If it found a number in the QF or more than one wh-words in a given question, e was expected to be the number. Otherwise, e was expected to be 1. The prior probability was obtained by using e as follows.

$$P(|A| = i) = \begin{cases} 1 & \text{if } i = e \\ 0 & \text{otherwise} \end{cases}$$

The second system used $f(x) = x^4$, which was chosen from $f(x) = x^n (n > 0)$ that performed best using QAC1 formalrun test collection, and $J = 10$. The prior probability was defined as follows.

$$P(|A| = i) = \begin{cases} \alpha & \text{if } i = 0 \\ 1 - \alpha & \text{if } i = 1 \\ 0 & \text{otherwise} \end{cases}$$

where the constant α is the prior probability that the query has no answers. We set $\alpha = 4/200$ that was selected from the practical value of QAC1, which has 4 no answer questions out of total 200 questions. The two systems selected average 3.573 and 2.784 answers from the QAC2 subtask 2 formalrun test collection, respectively. Both systems performed almost the top among the QAC2 subtask 2 participant systems. This result indicated the effectiveness of our approach of using expected values.

We would like to note that our systems participated in QAC2 subtask 2 and 3 adopted F-measure as the evaluation measure used for calculating the expected values, because we knew the evaluation of the subtasks would be made by it. The system can also use any other evaluation measures dependent on the purpose. For example, the weighted F-measure can be used as the evaluation measure in order to obtain the answers emphasized on either recall or precision.

6 Answering a Series of Questions for subtask 3

In QAC2 subtask 3, the system is required to answer a series of related questions. We constructed three systems for this subtask.

The first and second system were the simple extension of the two systems participated in the subtask 2 described in last section. In these system, the questions in the series are simply combined and treated as a single input to the systems, except that the query type and the question focus are extracted from the question currently being handled. For example, suppose a series of questions is $q_1 q_2 \cdots q_i$, in which the question

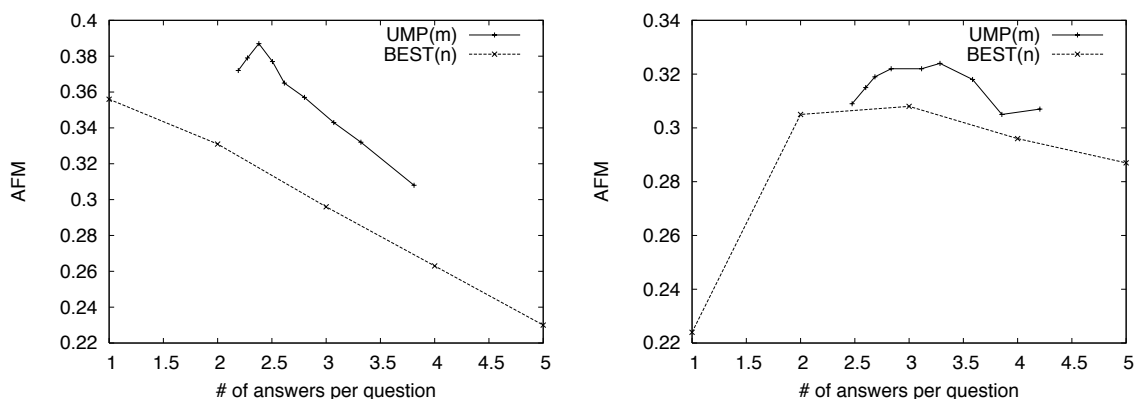


Figure 4. The relation between the average number of answer per question and AFM with respect to QAC1 (left) and QAC2 (right) subtask2 test collections.

currently asked is q_i , the query type and the QF are extracted only from q_i , while the other clues, including the content words using for the passage selection process, are extracted from all the questions $q_1, q_2 \dots q_i$.

The third system was constructed by extending the second system above so as adding the system's answers for previous questions to the input. Suppose a series of questions is $q_1 q_2 \dots q_i$ and the system have returned a series of answer sets $(a_{1,1} a_{1,2} \dots a_{1,j_1}), (a_{2,1} a_{2,2} \dots a_{2,j_2}), \dots, (a_{i-1,1} a_{i-1,2} \dots a_{i-1,j_{i-1}})$ that are extracted from $q_1, q_2, \dots q_{i-1}$, respectively. The union of the queries and the answers $q_1 \dots q_i, a_{1,1}, \dots, a_{i-1,j_{i-1}}$ are used as a single input to the system, except the query type and QF extraction that are extracted only from q_i .

7 Experiments

7.1 Evaluating Semantic Relations using Corpora

The evaluation of the method described in section 3 took place by using QAC1 and QAC2 test collections. The detailed experimental results using QAC1 test collection are found in [3]. In this paper, we examined only the total performance of question answering.

The accuracy of extracting the question focuses by our query analysis module was shown in Table 1. The question focuses from about 70 % out of all the queries in QAC1 subtask 1 were correctly extracted, while the QFs from about only 40 % out of all in QAC2 subtask 2. One of the reason why the accuracy in QAC2 was considerably reduced compared with that in QAC1 was that the queries in QAC2 has more variety of expression than that in QAC1.

Table 2 shows the total performance of our question answering systems. The 'BASE' indicates the result by the system that does not use the proposed method mentioned in section 3 but only use the string based

Table 1. The accuracy of QF extraction with respect to subtask1 of QAC1 and 2.

collection	QAC1	QAC2
total # of queries	196	195
QF exists (upper bound)	154 (77.0%)	131 (65.5%)
successfully extracted	139 (69.5%)	79 (39.5%)

method mentioned in section 3.6, while the '+pattern' indicates the result by the system participated in QAC2 that uses the proposed method.

With respect to QAC1, by using the proposed method, there observed considerable improvement, where the MRR for subtask 1 and AFM for subtask 2 increased +0.058 and +0.062, respectively. On the other hand, there observed smaller improvement with respect to QAC2, where the MRR and AFM increased +0.015 and +0.035, respectively. It was because the low accuracy on the QFs extraction in QAC2.

7.2 Selecting Set of Answers by using Expected Utility

Two strategies of selecting set of answers were compared.

BEST(n) select n best scored answers for each question ($n = 1 \dots 5$).

UMP(m) select by using Utility Maximization Principle proposed in Section 5.

The first system participated in QAC2 (mentioned in Section 5) was basically used for the experiment for

Table 2. The performance of question Answering with respect to QAC test collection

collection subtask	QAC1		QAC2	
	1(MRR)	2(AFM)	1(MRR)	2(AFM)
BASE	0.458	0.322	0.480	0.283
+pattern	0.516	0.384	0.495	0.318

Table 3. Effects of the proposed methods with respect to QAC2 test collection.

method used	TKBQ-2	-DP	-BG	+QF	+CD	+QF,+CD
subtask1 (MRR)	0.498	0.474	0.502	0.523	0.642	0.666

TKBQ-2 = the baseline system that corresponds to our QAC2 participated system for subtask 1, -DP = not used the dynamic passage selection (see Section 4.1) but static neighboring 3 sentences as a passage, -BG = not used the context word bi-grams for calculating similarity of passages (see Section 4.2), +QF = correct Question Focuses were given (corresponding to no QF extraction failure), +CD = correct articles that include a correct answer were given (corresponding to no IR failure)

UMP(m), excepted that the parameter m was introduced to revise the value of evaluation function as $f(x) = x^m$.

Figure 4 shows the relations between the average number of answers per question and AFM with respect to each strategy and test collection. It showed that our method using Utility Maximization Principle considerably outperformed the basic n-best selection strategy.

7.3 Effects of the Proposed Methods

Table 3 shows the contributions of the proposed methods toward the total performance of question answering with respect to the QAC2 test collection³.

8 Conclusion

Novel methods, each of which was used as a component of question answering, was proposed, including the method utilizing semantic relations in corpora, the method of dynamically selecting the optimal context of the answer candidates, the method of measuring the similarity between the query and the context by using the content word bi-gram, the method of selecting the set of answers for list questions, and so on.

We would also like to note that we have great interest in developing and evaluating the speech-driven question answering system. The detailed report of our system can be found in [2]. We are also interested in making the system accept spontaneously spoken queries [1].

References

[1] T. Akiba, A. Fujii, and K. Itou. Collecting spontaneously spoken queries for information retrieval. In *Proceedings of 4th*

³The parameter balancing the elements consists of the evaluation function $L(a|q)$ have been modified a little after participating in QAC2. This modification resulted in a small difference between the final MRR values (from 0.495 to 0.498).

International Conference on Language Resources and Evaluation, 2004.

- [2] T. Akiba, K. Itou, and A. Fujii. Adapting language models for frequent fixed phrases by emphasizing n-gram subsets. In *Proceedings of European Conference on Speech Communication and Technology*, pages 1469–1472, 2003.
- [3] T. Akiba, K. Itou, and A. Fujii. Selecting answers of common sense by using text collections for question answering. In *Proceedings of the Tenth Annual Meeting of The Association for Natural Language Processing*, pages 297–300, 2004. (in Japanese).
- [4] E. Brill, J. Lin, M. Banko, S. Dumais, and A. Ng. Data-intensive question answering. In *Proceedings of TREC-10*, 2001.
- [5] C.L.A. Clarke, G.V. Cormack, T.R. Lynam, G.M. Li, and G.L. McLearn. Web reinforced question answering. In *Proceedings of TREC-10*, 2001.
- [6] M. Fleischman, E. Hovy, and A. Echihiabi. Offline strategies for online question answering: Answering questions before they are asked. In *Proceedings of Annual Meeting of the Association for Computational Linguistics*, pages 1–7, 2003.
- [7] M. A. Hearst. Automatic acquisition of hyponyms from large text corpora. In *Proceedings of International Conference on Computational Linguistics*, pages 539–545, 1992.
- [8] A. Ittycheriah and S. Roukos. IBM's statistical question answering system – TREC-10. In *Proceedings of TREC-11*, 2001.
- [9] S. Lee and G. G. Lee. SiteQ/J: A question answering system for Japanese. In *Proceedings of The third NTCIR Workshop*, 2003.
- [10] D. Moldovan, S. Harabagiu, M. Pasca, R. Mihalcea, R. Goodrum, R. Girju, and V. Rus. LASSO: A tool for surfing the answer net. In *Proceedings of TREC-8*, pages 65–73, 1999.
- [11] M. Murata, M. Utiyama, and H. Isahara. A question-answering system using unit estimation and probabilistic near-terms IR. In *Proceedings of The third NTCIR Workshop*, 2003.
- [12] J. Prager and J. Chu-Carroll. Answering what-is question by virtual annotation. In *Proceedings of Human Language Technology Conference*, pages 26–30, 2001.
- [13] T. Takahashi, K. Hawata, S. Kouda, and K. Inui. Seeking answers by structural matching and paraphrasing. pages 87–94, 2003.
- [14] S. Tellex, B. Katz, J. Lin, A. Fernandes, and G. Marton. Quantitative evaluation of passage retrieval algorithms for question answering. In *Proceedings of ACM SIGIR*, pages 41–47, 2003.
- [15] E. Voorhees and D. Tice. The TREC-8 question answering track evaluation. In *Proceedings of the 8th Text Retrieval Conference*, pages 83–106, Gaithersburg, Maryland, 1999.