

MATHEMATICAL MODEL OF INFORMATION RETRIEVAL SYSTEM WITH FEEDBACK PROCESS

Sakai, Hiroshi

Department of Information Systems, Interdisciplinary Graduate School of Engineering Sciences,
Kyushu University

Haraguchi, Makoto

Research Institute of Fundamental Information Science, Kyushu University | Department of
Information Systems, Interdisciplinary Graduate School of Engineering Sciences, Kyushu
University

Takeya, Shunichi

Research Institute of Fundamental Information Science, Kyushu University | Department of
Information Systems, Interdisciplinary Graduate School of Engineering Sciences, Kyushu
University

Kano, Seigo

Research Institute of Fundamental Information Science, Kyushu University | Department of
Information Systems, Interdisciplinary Graduate School of Engineering Sciences, Kyushu
University

<https://doi.org/10.5109/13357>

出版情報 : Bulletin of informatics and cybernetics. 21 (1/2), pp.65-76, 1984-03. 統計科学研究会
バージョン :
権利関係 :



MATHEMATICAL MODEL OF INFORMATION RETRIEVAL SYSTEM WITH FEEDBACK PROCESS

By

Hiroshi SAKAI*, **Makoto HARAGUCHI****
Shun-ichi TAKEYA** and **Seigo KANŌ****

Abstract

A mathematical model of information retrieval system with feedback process is described. The model is characterized by learning process in which the system gradually guesses the user's interest of retrieval in terms of keywords. Precise formulation of the total system based on probability theory and several results are presented.

1. Introduction

There are generally two types of document retrieval system when classified by the retrieval method. The one is Boolean type which is common in most commercial systems, and the other is document ranking type using inner product value as in SMART system [1]. In both types, a user usually retrieves documents through trials and errors, and finally gets an optimum result. From user's point of view, the fundamental problems of retrieval are how to find the appropriate keywords for his interest of retrieval, and, if it is possible, how to make system understand his interest.

In this paper, we deal with the ranking type method, and formulate a feedback system including users. We try to express the retrieval process uniformly on the basis of probability theory. In practice, the content of this paper is a system's learning process of keywords which appropriately reflect the user's interest of retrieval.

In the field of information retrieval, such query reforming problem is generally considered as query adjustment problem [1, 2, 3, 4]. There are several methods for it, for example, relevance feedback [1], but its formula includes indefinite parameters which may be decided by experience. Chow's feedback query [3] is another example, but it is only one time feedback and has some assumptions about parametric distribution on terms. We describe the total system's flow including users without using indefinite parameters or special parametric distributions.

* Department of Information Systems, Interdisciplinary Graduate School of Engineering Sciences, Kyushu University 39, Kasuga, Fukuoka, 816 Japan.

** Research Institute of Fundamental Information Science, Kyushu University 33, Fukuoka, 812 Japan.

2. Necessary Factors for Feedback Process

First, we define several factors for our feedback system, where D and T stand for the document set and the index term (keyword) set, respectively.

DEFINITION 2.1 (*User participation I*). Assume that a user can mark $R_\omega(d)$ -value for documents, where “ ω ” is one of the user’s interest of retrieval. If d_i is more relevant to the user’s interest than d_j , then $R_\omega(d_i)$ is greater than $R_\omega(d_j)$.

$$(2.1) \quad R_\omega : D \longrightarrow [0, 1].$$

The relation between documents and terms used in the retrieval system is given beforehand as Boolean type matrix Γ . But each term for a document doesn’t seem to describe it uniformly. For this reason, we have the next definition.

DEFINITION 2.2 (*Contribution value*). Each document is assumed to have initially a probability distribution of terms. That is, for $\forall d \in D$,

$$(2.2) \quad M_d : \text{supp}(d) \longrightarrow [0, 1],$$

$$\sum_{t \in \text{supp}(d)} M_d(t) = 1, \quad \text{supp}(d) = \{t \in T \mid t \text{ can retrieve } d\}.$$

We call $M_d(t)$ a contribution value of term t for document d . By (2.2), a contribution matrix Γ' is made, and it is not known well to user.

DEFINITION 2.3 (*Entropy of term over d*). $M_d(t)$ has probability over a document d . The entropy $H(d)$ of term is,

$$(2.3) \quad H(d) = \sum_{t \in \text{supp}(d)} -M_d(t) \times \log M_d(t).$$

Since an entropy expresses disorder, the less $H(d)$ becomes, the better the document d will be recognized.

Our feedback system Ψ is roughly sketched in Fig. 2.1. In each step, a document set presented by the system will be evaluated by the user. We call the set the system request.

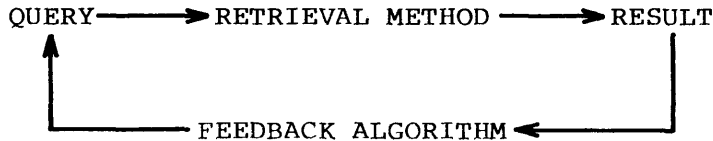


Fig. 2.1. Feedback process

Except the definition of the convergence of feedback process, our total system is denoted as follows ;

$$(2.4) \quad \Psi = (T, D, \Gamma', G, F, A, R_\omega(\cdot)),$$

G : feedback algorithm,

F : inner product retrieval method,

A : system request series,

$R_\omega(\cdot)$: user participation.

We regard Ψ as a retrieval process.

3. Definitions of the Feedback System

DEFINITION 3.1 (*Query*). The query of n -step is,

$$(3.1) \quad Q_n = (q_1^n, \dots, q_{\#T}^n),$$

q_i^n : weight of term t_i , $\#T$: number of terms.

DEFINITION 3.2 (*Occupancy of term*). Occupancy of term t_i in n -step query is,

$$(3.2) \quad P_n(t_i) = q_i^n / \sum_{j=1}^{\#T} q_j^n.$$

DEFINITION 3.3 (*Efficiency of term*). After a document is evaluated by user, the efficiency of term t_i is given as follows;

$$(3.3) \quad M(t_i|d) = R_\omega(d) \times M_d(t_i),$$

$$(3.4) \quad M(t_i|X_n) = (1/\#X_{ni}) \times \sum_{d \in X_{ni}} M(t_i|d),$$

X_n : system request, $X_{ni} = \{d \in X_n | M_d(t_i) > 0\}$ ($\neq \phi$).

$M(t_i|d)$ doesn't depend on query and is stationary with respect to steps.

DEFINITION 3.4 (*Value of retrieval*). N -step query Q_n gives documents those values of retrieval;

$$(3.5) \quad P_{Q_n}(d) = \sum_{t \in \text{supp}(d)} M_d(t) \times P_n(t),$$

$$(3.6) \quad \text{supp}(d) = \{t \in T | M_d(t) > 0\} \quad (\neq \phi),$$

if $\text{supp}(d) = \phi$ then $P_{Q_n}(d) = 0$.

Now, we will use (3.5) as F of the system Ψ .

DEFINITION 3.5 (*Term probability*). When document d_j is evaluated by query Q_n , the post probability of term t_i over d_j is,

$$(3.7) \quad P_{Q_n}(t_i|d_j) = M_{d_j}(t_i) \times P_n(t_i) / P_{Q_n}(d_j), \quad (P_{Q_n}(d_j) \neq 0).$$

These definitions are all given as probabilities and expected values.

PROPOSITION 3.1.

$$(3.8) \quad 1) \quad 0 \leq P_{Q_n}(d) \leq \max_{t \in \text{supp}(d)} M_d(t).$$

$$(3.9) \quad 2) \quad R_\omega(d_i) > R_\omega(d_j) \Leftrightarrow \sum_{t \in \text{supp}(d_i)} M(t|d_i) > \sum_{t \in \text{supp}(d_j)} M(t|d_j).$$

PROOF.

$$\begin{aligned} 1) \quad P_{Q_n}(d) &= \sum_{t \in \text{supp}(d)} M_d(t) \times P_n(t) \\ &\leq \sum_{t \in \text{supp}(d)} \{ \max_{t \in \text{supp}(d)} M_d(t) \} \times P_n(t) \\ &= \max_{t \in \text{supp}(d)} M_d(t) \times \sum_{t \in \text{supp}(d)} P_n(t) \end{aligned}$$

$$\begin{aligned}
&\leq \max_{t \in \text{supp}(d)} M_d(t). \\
2) \quad \sum_{t \in \text{supp}(d)} M(t|d) &= \sum_{t \in \text{supp}(d)} R_\omega(d) \times M_d(t) \\
&= R_\omega(d) \times \sum_{t \in \text{supp}(d)} M_d(t) \\
&= R_\omega(d).
\end{aligned}$$

PROPOSITION 3.2.

$$(3.10) \quad 1) \quad \sum_{t \in \text{supp}(d)} P_{Q_n}(t|d) = 1.$$

$$2) \quad \text{If } M_d(t_i) = M_d(t_j) \text{ for } \forall t_i, t_j \in \text{supp}(d) \text{ and } \sum_{t \in \text{supp}(d)} P_n(t) = 1 \text{ then,}$$

$$(3.11) \quad P_{Q_n}(t_i|d) = P_n(t_i).$$

PROOF.

$$1) \quad \sum_{t \in \text{supp}(d)} P_{Q_n}(t|d) = \sum_{t \in \text{supp}(d)} \{M_d(t) \times P_n(t) / P_{Q_n}(d)\} = 1.$$

$$\begin{aligned}
2) \quad P_{Q_n}(t_i|d) &= M_d(t_i) \times P_n(t_i) / \sum_{t \in \text{supp}(d)} M_d(t) \times P_n(t), \\
\sum_{t \in \text{supp}(d)} M_d(t) \times P_n(t) &= M_d(t_i) \times \sum_{t \in \text{supp}(d)} P_n(t) \\
&= M_d(t_i),
\end{aligned}$$

$$\begin{aligned}
\text{i. e. } P_{Q_n}(t_i|d) &= M_d(t_i) \times P_n(t_i) / M_d(t_i) \\
&= P_n(t_i).
\end{aligned}$$

$P_{Q_n}(t|d)$ has a probability over a document d , and $P_n(t)$ has a probability over a term set. After the document is evaluated, probability $P_{Q_n}(t|d)$ is equal to the prior term probability, if Condition 2) is satisfied.

PROPOSITION 3.3 (*Shannon's inequality*). If $P_1 = (p_{11}, \dots, p_{1n})$ and $P_2 = (p_{21}, \dots, p_{2n})$ are two discrete probability distributions and $0 \times \log 0 = 0$ then

$$(3.12) \quad \sum_{i=1}^n -p_{1i} \times \log p_{1i} \leq \sum_{i=1}^n -p_{1i} \times \log p_{2i}.$$

Equality is given when two distributions are the same.

PROOF. In general, $\log x \leq x - 1$ ($x > 0$), so

$$\begin{aligned}
\sum_{i=1}^n p_{1i} \times \log (p_{2i} / p_{1i}) &\leq \sum_{i=1}^n p_{1i} \times \{(p_{2i} / p_{1i}) - 1\} \\
&= \sum_{i=1}^n (p_{2i} - p_{1i}) \\
&= 0. \\
\sum_{i=1}^n p_{1i} \times \log (p_{2i} / p_{1i}) &= \sum_{i=1}^n -p_{1i} \times \log p_{1i} + \sum_{i=1}^n p_{1i} \times \log p_{2i} \\
&\leq 0.
\end{aligned}$$

PROPOSITION 3.4.

$$(3.13) \quad 0 \leq H(d) \leq \log(\#supp(d)).$$

PROOF. By Proposition 3.3,

$$\begin{aligned} H(d) &= \sum_{t \in supp(d)} -M_d(t) \times \log M_d(t) \\ &\leq \sum_{t \in supp(d)} -M_d(t) \times \log(1/\#supp(d)) \\ &= \log(\#supp(d)). \end{aligned}$$

Next, $-M_d(t) \times \log M_d(t) > 0$ in $(0, 1)$, so if $\#supp(d) \geq 2$ then $H(d) > 0$. If $\#supp(d) = 1$ then $H(d) = 0$.

THEOREM 3.1. $P_{Q_n}(t|d_i)$ has a probability over d_i by (3.10), so we define an entropy of terms, and write it $H_{Q_n}(d_i)$. Then,

$$(3.14) \quad H_{Q_n}(d_i) \times P_{Q_n}(d_i) < H(d_i), \quad (H(d_i) > 0 \text{ only}).$$

PROOF.

$$\begin{aligned} H_{Q_n}(d_i) &= \sum_{t \in supp(d_i)} -P_{Q_n}(t|d_i) \times \log P_{Q_n}(t|d_i) \\ &\leq \sum_{t \in supp(d_i)} -P_{Q_n}(t|d_i) \times \log M_{d_i}(t) \text{ by (3.12)} \\ &= (1/P_{Q_n}(d_i)) \times \left\{ \sum_{t \in supp(d_i)} -M_{d_i}(t) \times P_n(t) \times \log M_{d_i}(t) \right\} \end{aligned}$$

$0 \leq P_n(t) \leq 1$ and $H(d) > 0$, so $\#supp(d) \geq 2$ (by Proposition 3.4)

$$\begin{aligned} &< (1/P_{Q_n}(d_i)) \times \left\{ \sum_{t \in supp(d_i)} -M_{d_i}(t) \times \log M_{d_i}(t) \right\} \\ &= H(d_i)/P_{Q_n}(d_i). \end{aligned}$$

As $H(d_i)$ is fixed, according to (3.14) an increment of retrieval value causes the concentration of term (precisely speaking, term probability). The relation among the factors defined above is shown in Fig. 3.1.

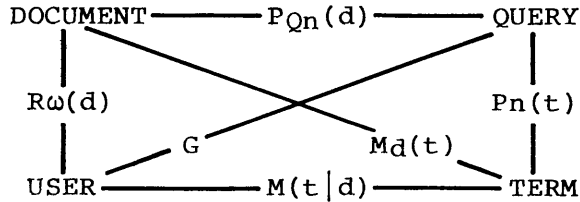


Fig. 3.1. Factors of the system Ψ

Next we define G of the system Ψ .

DEFINITION 3.6. After the user participation for n -step system request X_n , the weight of term t_i in n -step is defined as follows;

$$(3.15) \quad \begin{aligned} G : q_i^n(X_n) &= (1/\#X_{ni}) \left\{ \sum_{d \in X_{ni}} R_\omega(d) \times P_{Q_{n-1}}(t_i|d) \right\} \quad (X_{ni} \neq \phi), \\ G : q_i^n(X_n) &= 0 \quad (X_{ni} = \phi). \end{aligned}$$

For the formulation of the transition of the system Ψ' , (3.15) will be the most important formula. It shows that term t_i will possess the relevancy of a document d with the probability $P_{Q_{n-1}}(t_i|d)$ that is the probability of term t_i over d after evaluated by Q_{n-1} . Thus we use the term weights as to express the relevancy of documents. Therefore an evaluation of a document will eventually gives the same effect to documents that have similar terms. This property comes from the special condition of document-term matrix I' .

4. Definition of Classes by Document-Term Matrix

Before retrieval, we (the system) could know the relation between documents and terms according to I' . We define several classes on documents and terms using I' .

DEFINITION 4.1.

$$(4.1) \quad \text{supp}(d_i) = \{t \in T \mid M_{d_i}(t) > 0\},$$

$$(4.2) \quad \text{supp}(D_i) = \bigcup_{d \in D_i} \text{supp}(d).$$

DEFINITION 4.2.

$$(4.3) \quad C^1(d_i) = \bigcup_{t \in \text{supp}(d_i)} \{d \in D \mid M_d(t) > 0\},$$

$$(4.4) \quad C^n(d_i) = \bigcup_{t \in \text{supp}(C^{n-1}(d_i))} \{d \in D \mid M_d(t) > 0\}.$$

$C^n(d_i)$ is not monotone decreasing and $C^n(d_i) \subset D$, so $C^n(d_i)$ converges to $C(d_i)$. These classes are called as follows;

$C(d_i)$: document class of d_i -analogy,

$\text{supp}(C(d_i))$: term class of d_i -analogy.

PROPOSITION 4.1. *Any one document belonging to one class analogizes with documents in the same class and specifically if $C(d_i) = D$ then any one document analogizes with the total document set.*

PROOF. Let d be an arbitrary document in $C(d_i)$. Then d_i analogizes with d , so inversely d analogizes with d_i . Therefore $C(d) \supset C(d_i)$ and $d_i \in C(d)$. Equally, $C(d) \subset C(d_i)$ will be derived. So $C(d) = C(d_i)$.

PROPOSITION 4.2. *If $d \in C(d_i)$ then $\text{supp}(d) \cap \text{supp}(C(d_i)) = \phi$.*

PROOF. Suppose $\text{supp}(d) \cap \text{supp}(C(d_i)) \neq \phi$, then there is at least one term $t (\in T)$ and it satisfies the following;

$$M_d(t) > 0 \quad \text{and} \quad M_{d'}(t) > 0 \quad (d' \in C(d_i)).$$

The condition shows that d is analogized by d' through term t . It is contradictory to $d \in C(d_i)$.

PROPOSITION 4.3. *$C(d_i) \cap C(d_j) = \phi$ or $C(d_i) = C(d_j)$.*

PROOF. If $d_i \in C(d_j)$ then the proof is equal to the one of Proposition 4.1. Therefore, suppose as follows;

$$d_i \notin C(d_j) \quad \text{and} \quad C(d_i) \cap C(d_j) \neq \phi,$$

then $\exists d \in C(d_i)$ and $\exists d' \in C(d_j)$. That is, d_i is analogized by d_j through d , which contradicts to $d_i \notin C(d_j)$.

PROPOSITION 4.4. $\text{supp}(C(d_i)) \cap \text{supp}(C(d_j)) = \phi$ or $\text{supp}(C(d_i)) = \text{supp}(C(d_j))$.

PROOF. If $C(d_i) = C(d_j)$ then $\text{supp}(C(d_i)) = \text{supp}(C(d_j))$. This is trivial. Next, according to the result of Proposition 4.2 $\text{supp}(d) \cap \text{supp}(C(d_j)) = \phi$ for $\forall d \in C(d_i)$, so

$$\bigcup_{d \in C(d_i)} \text{supp}(d) \cap \text{supp}(C(d_j)) = \phi.$$

It shows that $\text{supp}(C(d_i)) \cap \text{supp}(C(d_j)) = \phi$.

These definitions are similar to those of transitive closure in graph theory. According to these propositions and definitions, we get the decomposition classes of D and T as follows;

$$(4.5) \quad D = \bigcup_i C(d_i), \quad T = \bigcup_i \text{supp}(C(d_i)).$$

THEOREM 4.1. Formula (3.15) is well-defined on the only one class of D and T . Namely, the classes are independent with respect to retrieval.

PROOF. Let $C(d_i)$ and $C(d_j)$ are two classes. By (4.5), $\text{supp}(C(d_i)) \cap \text{supp}(C(d_j)) = \phi$, so the system request $X_n \subset C(d_i)$ gives no information about $\text{supp}(C(d_j))$. Therefore $P_n(t) = 0$ for $\forall t \in \text{supp}(C(d_j))$ and $P_{Q_n}(d) = 0$ for $\forall d \in C(d_j)$.

Theorem 4.1 clarifies the sets which Ψ can deal with. By Theorem 4.1 and (3.15) the system Ψ has a property of inference in a class. If we (the system) have any information about one class, we (the system) can reach to any document in the class by $R_\omega(d)$ -value, but can never reach to documents of other classes. In practice, the classes are decided by matrix I , so the indexing of document is especially important in Ψ .

In general, when we construct a feedback system, we must prepare both ways of concentration and expansion of the retrieved document set. But in Ψ , the expansion of term set is decided in the initial step according to Theorem 4.1, and its concentration is done in sequence. We call here the process 'feedback'.

5. Some Results and Total System Algorithm

THEOREM 5.1. The occupancy of term t_i in n -step ($n \geq 2$) is,

$$(5.1) \quad \begin{aligned} P_n(t_i) &= Z_{in} / \sum_{j=1}^{\#T} Z_{jn}, \\ Z_{in} &= P_1(t_i) \times \prod_{m=1}^{n-1} E_{Q_m}[t_i | X_{m+1}], \\ E_{Q_m}[t_i | X_{m+1}] &= (1/\#X_{m+1,i}) \times \left\{ \sum_{d \in X_{m+1,i}} M(t_i | d) / P_{Q_m}(d) \right\}. \end{aligned}$$

PROOF. By induction,

$$\begin{aligned} q_i^2(X_2) &= (1/\#X_{2i}) \times \left\{ \sum_{d \in X_{2i}} R_\omega(d) \times P_{Q_1}(t_i | d) \right\} \\ &= (1/\#X_{2i}) \times \left\{ \sum_{d \in X_{2i}} R_\omega(d) \times (M_d(t_i) \times P_1(t_i) / P_{Q_1}(d)) \right\} \\ &= P_1(t_i) \times (1/\#X_{2i}) \times \left\{ \sum_{d \in X_{2i}} R_\omega(d) \times M_d(t_i) / P_{Q_1}(d) \right\} \\ &= P_1(t_i) \times E_{Q_1}[t_i | X_2]. \end{aligned}$$

$$\begin{aligned}
P_2(t_i) &= q_i^2 / \sum_{j=1}^{\#T} q_j^2 \\
&= \{P_1(t_i) \times E_{Q_1}[t_i | X_2]\} / \sum_{j=1}^{\#T} \{P_1(t_j) \times E_{Q_1}[t_j | X_2]\} \\
&= Z_{i2} / \sum_{j=1}^{\#T} Z_{j2}.
\end{aligned}$$

Next,

$$\begin{aligned}
q_i^n(X_n) &= (1/\#X_{ni}) \times \left\{ \sum_{d \in X_{ni}} R_\omega(d) \times P_{Q_{n-1}}(t_i | d) \right\} \\
&= (1/\#X_{ni}) \times \left\{ \sum_{d \in X_{ni}} R_\omega(d) \times (M_d(t_i) \times P_{n-1}(t_i) / P_{Q_{n-1}}(d)) \right\} \\
&= (1/\#X_{ni}) \times \left\{ \sum_{d \in X_{ni}} R_\omega(d) \times M_d(t_i) / P_{Q_{n-1}}(d) \right\} \\
&\quad \times \left\{ Z_{i, n-1} / \sum_{j=1}^{\#T} Z_{j, n-1} \right\} \\
&= \left(1 / \sum_{j=1}^{\#T} Z_{j, n-1} \right) \times (1/\#X_{ni}) \\
&\quad \times \left\{ \sum_{d \in X_{ni}} R_\omega(d) \times M_d(t_i) \times Z_{i, n-1} / P_{Q_{n-1}}(d) \right\} \\
&= \left(1 / \sum_{j=1}^{\#T} Z_{j, n-1} \right) \times E_{Q_{n-1}}[t_i | X_n] \times Z_{i, n-1} \\
&= \left(1 / \sum_{j=1}^{\#T} Z_{j, n-1} \right) \times Z_{in}.
\end{aligned}$$

$$\begin{aligned}
P_n(t_i) &= q_i^n / \sum_{j=1}^{\#T} q_j^n \\
&= Z_{in} / \sum_{j=1}^{\#T} Z_{jn}.
\end{aligned}$$

In system Ψ , the occupancy of term t_i for query Q_n is sequentially reformed by (5.1) using system requests and $R_\omega(d)$ -values. We want to define the end of retrieval with the convergence of term occupancy, but it is not trivial.

PROPOSITION 5.1. *If $P_n(t_i)=0$ then $P_m(t_i)=0$ ($m \geq n$).*

PROOF. By (3.7), $P_{Q_n}(t_i | d) = M_d(t_i) \times P_n(t_i) / P_{Q_n}(d)$. $P_n(t_i)=0$ is a condition, so $P_{Q_n}(t_i | d)=0$ for $\forall d \in D$. By (3.15), $q_i^{n+1}=0$, so $P_{n+1}(t_i)=0$. We can prove it, sequentially.

According to Proposition 5.1, once the system regards a term as useless, the term doesn't cause any effect, so it can be eliminated from the term set.

PROPOSITION 5.2. *The process such as $P_n(t_i) > 0$ and $P_{n+1}(t_i) = 0$ occurs in either one of the following two cases.*

$$(5.2) \quad 1) \quad R_\omega(d) = 0 \quad \text{for } \forall d \in X_{n+1, i} \quad (\neq \phi),$$

$$(5.3) \quad 2) \quad X_{n+1, i} = \phi.$$

PROOF. By (3.2), $P_{n+1}(t_i) = q_i^{n+1} / \sum_{j=1}^{\#T} q_j^{n+1}$, and in order to $P_{n+1}(t_i) = 0$, q_i^{n+1} must be 0.

According to (3.15), q_i^{n+1} is 0, when 1) or 2) is satisfied.

The case 1) is reasonable, and term's concentration will depend on the degree of elimination of terms by 1). The case 2) expresses a unique condition that system requests must satisfy. In each step, a system request must be decided such that the term set of requested documents covers the term set of 1-step before, and this condition excludes unreasonable elimination of terms. The problem of the system request seems to be solved by using a thesaurus [5, 6] and it is very important and difficult. We will consider the problem next time.

In fact, the term set, the document set and $supp(d)$ depend on the steps. So we denote the substances of them as T_n , D_n and $supp_n(d)$, respectively.

PROPOSITION 5.3. T_n , D_n and $supp_n(d)$ are not monotone increasing in inclusion relation. A term is eliminated when one of the conditions of Proposition 5.2 is satisfied, and a document is eliminated when $supp_n(d)(=supp(d) \cap T_n) = \phi$ is satisfied.

PROOF. If $q_i^n = 0$ then $T_n = T_{n-1} - \{t_i\}$, so $T_n \subset T_{n-1}$ and never become $T_n \supset T_{n-1}$. It shows that T_n is not monotone increasing. By (3.6), if $supp_n(d) = \phi$ then $P_{Q_n}(d) = 0$, so $D_n \subset D_{n-1}$ and never become $D_n \supset D_{n-1}$. It shows D_n is not monotone increasing, too.

Substantial sets are decided by the next algorithm.

ALGORITHM 5.1.

```

begin
  calculate  $q_1^n, \dots, q_{\#T_{n-1}}^n$ ;
   $T_n := T_{n-1}$ ;
   $D_n := D_{n-1}$ ;
  for  $i := 1$  until  $\#T_{n-1}$  do
    if  $q_i^n = 0$  then  $T_n := T_n - \{t_i\}$ ;
  for  $j := 1$  until  $\#D_{n-1}$  do
    begin
       $supp_n(d_j) := supp(d_j) \cap T_n$ ;
      if  $supp_n(d_j) = \phi$  then  $D_n := D_n - \{d_j\}$ 
    end;
  end

```

PROPOSITION 5.4. When $M_{d_i}(t_j) = \delta_{ij}$ (complete information) then one document makes one class, and can't analogize with other classes. So we eventually evaluate the total document set. In this case, system Ψ is useless.

PROOF. Trivial.

We defined classes before retrieval, but in practice, we deal with the subset $D^*(\subset D)$. The total algorithm is shown next.

TOTAL ALGORITHM.

$$D_1 (= \bigcup_i D(i) : D(i) \text{ is a subset of } C(d), d \in D(i)), \quad (*)$$

$$T_1 (= supp(D_1)), P_1(t_i) \ (i=1, \dots, \#T_1). \quad (**)$$

(*) and (**) are given initially.

begin

query-type: $Q := \bigcup_i Q_i$ (Q_i : query for $supp(D(i))$);

$n := 2$;

```

user: for  $i := 1$  until  $\#X_n$  do decide  $R_\omega(d)$ -value ;
       $sum := 0$  ;
term-set:  $T_n := T_{n-1}$  ;
query: for  $j := 1$  until  $\#T_{n-1}$  do
      begin
        calculate  $q_j^n(X_n)$  ;
        if  $q_j^n(X_n) = 0$  then  $T_n := T_n - \{t_j\}$  ;
         $sum := sum + q_j^n(X_n)$ 
      end ;
occupancy: for  $j := 1$  until  $\#T_n$  do  $P_n(t_j) := q_j^n(X_n) / sum$  ;
value: for  $i := 1$  until  $\#D_1$  do calculate  $P_{Q_n}(d_i)$  ;
      comment  $D_n$  and  $supp_n(d)$  are decided by Algorithm 5.1 ;
      if continue then  $n := n + 1$  else stop ;
      comment System's convergence check is not given ;
      goto user
end

```

6. Convergence of Retrieval

It is natural to think that if term occupancies converge under the influence of each step's information, then the final occupancies will be the optimum term occupancies. By (5.1), each step's information is $E_{Q_n}[t | X_{n+1}]$. It is concerned with Q_n , $M_d(t)$, $R_\omega(d)$ and the system requests. The expansion of term set is decided in the initial step, and the system request is given in such a way that the support of documents covers substantial term set, so the process doesn't seem to depend on the system requests too much.

Next, we describe how to mark $R_\omega(d)$ -value in order to converge the process.

DEFINITION 6.1. (*User participation II*). Suppose that a user doesn't mark precise $R_\omega(d)$ -value for system requests but he only classifies the system request roughly in the order of relevancy.

PROPOSITION 6.1. *If a user only divides a system request into two, namely, he classifies a system request as " $X_n = X_n(1) \cup X_n(2)$, $X_n(1) \cap X_n(2) = \emptyset$ and $X_n(1)$ is more relevant than $X_n(2)$ ", then $R_\omega(d)$ -values can be decided as to satisfy the following ;*

$$(6.1) \quad P_{Q_n}(d_1) > P_{Q_n}(d_2) \quad \text{for } \forall d_1 \in X_n(1), \forall d_2 \in X_n(2).$$

PROOF.

$$\begin{aligned} P_{Q_n}(d) &= \sum_{t_m \in \text{supp}(d)} M_d(t_m) \times P_n(t_m) \\ &= \sum_{t_m \in \text{supp}(d)} M_d(t_m) \times \left\{ q_m^n / \sum_{j=1}^{\#T_n} q_j^n \right\}. \end{aligned}$$

So, (6.1) becomes

$$\sum_{t_m \in \text{supp}(d_1)} M_{d_1}(t_m) \times q_m^n > \sum_{t_k \in \text{supp}(d_2)} M_{d_2}(t_k) \times q_k^n.$$

By (3.15),

$$(6.2) \quad \sum_{t_m \in \text{supp}(d_1)} (M_{d_1}(t_m) / \#X_{nm}) \times \left\{ \sum_{d \in X_{nm}} R_\omega(d) \times P_{Q_{n-1}}(t_m | d) \right\}$$

$$> \sum_{t_k \in \text{supp}(d_2)} (M_{d_2}(t_k)/\#X_{nk}) \times \left\{ \sum_{d \in K_{nk}} R_\omega(d) \times P_{Q_{n-1}}(t_k|d) \right\}.$$

Here, the same $R_\omega(d)$ -value is given for documents in the same class, and it must satisfy $R_\omega(d_1) > R_\omega(d_2)$, so there are two different values of $R_\omega(d)$. (6.2) is changed to the next.

$$\begin{aligned} & R_\omega(d_1) \times \left\{ \sum_{t_m \in \text{supp}(d_1)} (M_{d_1}(t_m)/\#X_{nm}) \times \sum_{d \in X_{nm}} P_{Q_{n-1}}(t_m|d) \right\} \\ & > R_\omega(d_2) \times \left\{ \sum_{t_k \in \text{supp}(d_2)} (M_{d_2}(t_k)/\#X_{nk}) \times \sum_{d \in X_{nk}} P_{Q_{n-1}}(t_k|d) \right\}. \end{aligned}$$

Let $h_{ni} = \sum_{t_j \in \text{supp}(d_i)} (M_{d_i}(t_j)/\#X_{nj}) \times \sum_{d \in X_{nj}} P_{Q_{n-1}}(t_j|d)$. Then,

$$(6.3) \quad R_\omega(d_1)/R_\omega(d_2) > h_{n2}/h_{n1}.$$

In (6.3), if X_n is given then the right side is fixed for d_1 and d_2 . Next we substitute the maximum of numerator and minimum of denominator for each of the right side. We write it $K_n(X_n(1), X_n(2))$, and finally get the next formula.

$$(6.4) \quad R_\omega(d_1)/R_\omega(d_2) > \max\{1, K_n(X_n(1), X_n(2))\}.$$

If the system marks $R_\omega(d)$ -values according to (6.4), then (6.1) is concluded.

THEOREM 6.1. *If a user classifies X_n as $\bigcup_i X_n(i)$ ($X_n(i)$ more relevant than $X_n(i+1)$), then $R_\omega(d)$ -values can be decided to satisfy the following;*

$$(6.5) \quad P_{Q_n}(d_i) > P_{Q_n}(d_j) \quad \text{for } \forall d_i \in X_n(i), \quad \forall d_j \in X_n(j) \quad (i < j).$$

PROOF. Sequentially, system can decide conditions of $X_n(i)$ and $X_n(i+1)$ by Proposition 6.1. The following conditions are decided.

$$\begin{aligned} & R_\omega(d_1)/R_\omega(d_2) > \max\{1, K_n(X_n(1), X_n(2))\}, \\ & R_\omega(d_2)/R_\omega(d_3) > \max\{1, K_n(X_n(2), X_n(3))\}, \\ & \quad \vdots \\ & R_\omega(d_{m-1})/R_\omega(d_m) > \max\{1, K_n(X_n(m-1), X_n(m))\}. \end{aligned}$$

If $R_\omega(d_1)$ is fixed then the other $R_\omega(d)$ -values are decided by these conditions, and they cause (6.5).

Theorem 6.1 needs enormous calculation. But the ranking of X_n influences the ranking of total document set. We finally find that a user need not give precise $R_\omega(d)$ -values for a system request but just classify it. The method has merits of Boolean type and ranking type. Moreover if we sequentially adopt documents that are ranked as upper middle, the process converges steadily.

7. Conclusion

We defined several factors necessary for information retrieval system with feedback and formulated the retrieval process in contrast with Salton's relevance feedback formula. For the sake of the classes defined by matrix I' , we cleared substantial sets that are dealt with by the system Ψ in the present step. (3.15) is the most important formula, and expresses the analogy of document in one class.

The user participation I and II are defined. Both of them will give the direction of retrieval, but the effect of II is much more steady, clear and realistic.

The authors hope that these discussions will clear the total retrieval process and give vivid influence to formulate it more in detail. The practical implementation of our system is not considered at present because of enormous calculation.

References

- [1] SALTON, G. and MCGILL, M.J.: Introduction to modern information retrieval, McGraw-Hill, New York (1983).
- [2] VERNIMB, C.: *Automatic query adjustment in document retrieval*, IPM, 13 (1977), 339-353.
- [3] CHOW, D. and YU, C.T.: *On the construction of feedback query*, JACM, 29 (1982), 127-151.
- [4] DRIYANSKII, V.M.: *Retrieval models in on-line documentary information systems: An analytic review*, Cybernetics, 17 (1981), 269-287.
- [5] RAS, Z.W.: *An algebraic approach to information retrieval systems*, IJCIS, 11 (1982), 275-293.
- [6] MAZUR, Z.: *Inverted organization in the information retrieval system based on thesaurus with weights*, IPM, 15 (1979), 227-234.

Received September 20, 1983