

질의분해 적합성 피드백을 이용한 정보검색에 관한 연구

A Study on Information Retrieval Using Query Splitting Relevance Feedback

김영천** · 박병권* · 이성주**

Young-cheon Kim**, Byung-gweun Park*, and Sung-joo Lee**

*서강정보대학 정보통신과

**조선대학교 전자계산학과

요 약

순수한 부울 검색 시스템은 문서와 질의 사이의 유사도를 나타내는 문서값을 계산할 수 없기 때문에, 검색된 문서들을 질의를 만족하는 정보에 따라 정렬할 수 없다. 부울 검색 시스템의 이러한 단점을 보완하는 방법으로 MMM 모델, Paice 모델, P-norm 모델이 개발되었다. 이러한 방법들은 부울 연산자를 유연하게 연산하는 공통된 특성을 지니고 있다. 본 논문에서는 높은 검색 효과를 제공하는 질의분해 적합성 피드백(QSRF)을 이용한 정보 검색 모델을 제안한다. 질의 분해 적합성 피드백 모델의 연산 특성이 MMM, Paice, P-norm 모델보다 우수함을 설명하고, 또한 성능 비교를 통하여 이를 입증한다.

Abstract

In conventional boolean retrieval systems, document ranking is not supported and similarity coefficients cannot be computed between queries and documents. The MMM, Paice and P-norm models have been proposed in the past to support the ranking facility for boolean retrieval systems. They have common properties of interpreting boolean operators softly. In this paper we propose a new soft evaluation method for Information retrieval using query splitting relevance feedback model. We also show through performance comparison that query splitting relevance feedback(QSRF) is more efficient and effective than MMM, Paice and P-norm.

Key Words : 정보검색, 적합성 피드백, 질의 분해, 부울 검색, 유사도

1. 서 론

순수한 부울 검색 모델은 문서와 질의 사이의 유사도를 나타내는 문서값을 계산할 수 없기 때문에, 검색된 문서들을 질의를 만족하는 정도에 따라 정렬할 수 없다는 단점을 지니고 있다.

순수한 부울 검색 시스템의 단점을 보완하기 위하여 퍼지 집합 모델(Fuzzy Set Model)이 개발되었다. 퍼지 집합 모델은 색인어가 문서 내에서 갖는 중요성을 반영하는 색인어가 중치를 이용하여 문서값을 계산함으로써 부울 검색 시스템의 문제점을 극복하였다. 그러나 퍼지 집합 모델은 많은 경우에 부정확한 문서값을 생성하기 때문에 정보 검색 모델로서 부적합하다고 비판되어 왔다. 이것은 AND와 OR 연산을 위하여 사용하는 MIN과 MAX 연산자가 단일 피연산자 의존 문제(Single Operand Dependency Problem)를 발생시키기 때문이다.

퍼지 집합 모델의 단일 피연산자 의존 문제를 극복하기 위하여 MMM 모델, Paice 모델, P-norm 모델이 개발되었

다. 이들 모델들은 AND와 OR 연산을 위하여 MIN과 MAX 대신에 부울 연산자를 유연하게 연산하는 새로운 연산자를 사용함으로써 단일 피연산자 의존 문제를 극복하였을 지라도 다음과 같은 단점을 지니고 있다. 첫째, MMM 모델은 빠른 검색 시간을 제공할 지라도, 부정확한 문서값을 생성할 수 있는 요인을 지니고 있다. 둘째, MMM 모델의 문제점이 Paice, P-norm 모델에서는 발생하지 않을 지라도, 이들 모델은 검색 시간이 느리다는 단점을 가지고 있다.

정보검색에서 가장 중요하면서도 어려운 문제 중의 하나는 사용자가 원하는 정보를 찾기 위한 효율적인 질의를 작성하는 일이다. 하지만 전체 문서집합의 구성에 대해 미리 알고 있지 않는 한 이상적인 최적의 질의는 작성할 수 없다. 대신 최초에는 시험적 질의(tentative query)로 검색을 수행한 후, 이전의 검색 결과에 대한 평가에 기반하여 다음 번 검색의 질의를 개선시키는 방법이 적합성 피드백(relevance feedback)이다.

적합성 피드백은 특정 질의와 적합한 문서들은 유사한 벡터로 표현된다고 가정한다. 따라서, 어떤 문서가 주어진 질의에 적합하다고 판단되면 질의를 적합한 문서와의 유사도가 증가하도록 변환하여 질의를 개선시킨다. 이렇게 개선된 질의는 최초 적합하다고 판단된 문서와 유사한 문서들을 추가적으로 검색하여 더 많은 양의 적합 문서를 검색해 낼 수 있다.

접수일자 : 2001년 4월 21일

완료일자 : 2001년 5월 30일

실제 이 적합성 피드백을 이용할 때에는 전체 문서집합에 대해 적합문서와 부적합문서를 미리 알 수 없으므로, 이미 적합성이 알려져 있는 문서들의 정보에 기반 하여 질의확장을 수행한다. 이때의 적합성을 사용자가 알려주는 방법을 사용자 적합성 피드백(user relevance feedback)이라 하고, 사용자의 개입 없이 초기질의로 검색된 결과 문서 중 상위 문서를 적합한 문서로 간주하여 적합성 피드백을 적용하는 방법을 의사 적합성 피드백(pseudo relevance feedback)이라 한다[2],[7].

본 논문에서는 질의 분해 적합성 피드백 모델을 제안한다. 본 논문의 구성은 다음과 같다. 2장에서는 부울 연산자를 유연하게 연산하는 기존의 방법들 MMM, Paice, P-norm 모델에 대하여 기술한다. 3장에서는 높은 검색 효과를 제공하는 질의 분해 적합성 피드백 모델을 제안한다. 4장에서는 질의 분해 적합성 피드백 모델과 MMM, Paice, P-norm 모델의 성능을 비교한다. 마지막으로 5장에서 결론 및 앞으로의 연구 방향을 제시한다.

2. 부울 연산자를 유연하게 연산하는 기존의 방법들

퍼지 집합 모델의 단일 피연산자 의존 문제를 극복하기 위해 MMM 모델, Paice 모델, P-norm 모델이 개발되었다. 이들 모델들은 AND와 OR 연산을 위하여 MIN과 MAX 대신에 부울 연산자를 유연하게 연산하는 새로운 연산자를 사용함으로써 퍼지 집합 모델, MMM 모델, Paice 모델, P-norm 모델을 기반으로 하는 정보 검색 시스템은 <T, Q, D, F>로 정의되는 확장된 부울 검색 체계(Extended Boolean Retrieval Framework) 내에서 설명될 수 있다.

- ① T는 질의와 문서를 표현하기 위해 사용되는 색인어 들의 집합이다.
- ② Q는 시스템이 인식할 수 있는 질의들의 집합이다.Q에 속하는 각각의 질의 q는 색인어들과 부울 연산 자 AND, OR, NOT으로 구성된 부울 수식이다.
- ③ D는 문서들의 집합이다. D에 속하는 각각의 문서 d 는 w_i 가 색인어 t_i 의 가중치일 때, $\{(t_1, w_1), \dots, (t_n, w_n)\}$ 와 같이 표현된다. 색인어 가중치 w_1 는 0부터 1 사이의 값을 갖는다.
- ④ F는 문서값을 계산하는 순위 결정 함수(Ranking Function)로서 다음과 같이 정의된다.

$$F: D \times Q \rightarrow [0, 1]$$

검색함수 F는 각 쌍의 (d, q)에 0부터 1 사이의 값을 지정한다. 이 값은 문서 d와 질의 q 사이의 유사도를 의미하며, 질의 q에 대한 문서 d의 문서값이다.

검색 함수 F(d, q)는 다음과 같은 2단계 과정을 거쳐서 계산된다.

- (i) 질의에 나타난 각각의 색인어 t_i 에 대하여, F(d, t_i)는 문서 d에서 색인어 t_i 의 가중치 w_i 로 정의 된다.
- (ii) 부울 연산자 AND와 OR는 (a), (b), (c), (d)에서주어진 식들을 이용하여 계산되고, NOT은 $F(d, NOT t_i)=1-w_i$ 로 계산된다.

두 개 이상의 부울 연산자를 포함하는 부울 질의는 가장 안쪽에 위치하는 절부터 순환적으로 계산된다.

퍼지 집합 모델의 부울 연산자 계산식 (a)는 두 개의 피연산자를 갖는 이항연산이고, MMM, Paice, P-norm 모델의 연산자 계산식은 2개 이상의 피연산자를 갖는 다항연산이다. 이것은 퍼지 집합 모델의 MIN과 MAX 연산자가 결합법칙을 만족하는데 비하여 MMM, Paice, P-norm 모델의 연산자는 결합법칙을 만족하지 못하기 때문이다. 결합법칙을 만족하지 못할 경우, 임의의 문서에 대하여 두 개의 동일한 질의($(t_1 AND t_2) AND t_3$ 와 $t_1 AND(t_2 AND t_3)$)의 문서값이 서로 다르다. MMM, Paice, P-norm 모델은 이러한 문제점을 다항연산을 가능하게 함으로써 극복하였다.

$$F(d, t_1 AND t_2) = MIN(w_1, w_2) \quad (1)$$

$$F(d, t_1 OR t_2) = MAX(w_1, w_2) \quad (2)$$

(a) 퍼지 집합 모델

$$F(d, t_1 AND \dots AND t_n) = r \cdot MAX(w_1 \dots w_n) + (1-r) \cdot MIN(w_1 \dots w_n) \quad (3)$$

$0 \leq r \leq 0.5$

$$F(d, t_1 OR \dots OR t_n) = r \cdot MIN(w_1 \dots w_n) + (1-r) \cdot MAX(w_1 \dots w_n) \quad (4)$$

$0.5 \leq r \leq 1$

(b) MMM 모델

$$F(d, t_1 AND \dots AND t_n) = \frac{\sum_{i=1}^n (r^{i-1} \cdot w_i)}{\sum_{i=1}^n (r^{i-1})} \quad (5)$$

($0 \leq r \leq 1$, w_i '는 오름차순정렬)

$$F(d, t_1 OR \dots OR t_n) = \frac{\sum_{i=1}^n (r^{i-1} \cdot w_i)}{\sum_{i=1}^n (r^{i-1})} \quad (6)$$

($0 \leq r \leq 1$, w_i '는 내림차순정렬)

(c) Paice 모델

$$F(d, t_1 AND \dots AND t_n) = 1 - \left[\frac{(1-w_1)^p + \dots + (1-w_n)^p}{n} \right]^{\frac{1}{p}} \quad (7)$$

($1 \leq p \leq \infty$)

$$F(d, t_1 OR \dots OR t_n) = \left[\frac{w_1^p + \dots + w_n^p}{n} \right]^{\frac{1}{p}} \quad (8)$$

($1 \leq p \leq \infty$)

(d) P-norm 모델

1. 정규화

확장된 부울 검색 체계를 기반으로 하는 검색 모델은 문서값을 계산하기 위하여 색인어 가중치를 사용한다. 색인어 가중치는 역문헌빈도(Inverse Document Frequency)와 색인어 출현빈도(Term Frequency)로부터 유도될 수 있다. N이

$$Q^{new} = \alpha Q^{old} + \frac{\beta}{|R|} \sum_{D_i \in R} D_i - \frac{\gamma}{|N|} \sum_{D_i \in N} D_i \quad (13)$$

Ide Regular는 식(14)와 같다.

$$Q^{new} = \alpha Q^{old} + \beta \sum_{D_i \in R} D_i - \gamma \sum_{D_i \in N} D_i \quad (14)$$

Ide Dec Hi는 식(15)와 같다.

$$Q^{new} = \alpha Q^{old} + \beta \sum_{D_i \in R} D_i - \gamma \max_{n-r}(D_i) \quad (15)$$

여기서 Q^{new} 는 새로 확장된 피드백 질의 벡터를 Q^{old} 는 확장되기 전 단계의 질의 벡터를 의미한다. R과 N은 각각 초기 검색된 문서집합 중에서 적합하다고 판단된 문서집합과 부적합하다고 판단된 문서집합을 |R|과 |N|은 각각 해당 문서집합의 문서 개수를 뜻한다. $\max_{n-r}(D_i)$ 는 부적합문서 중 최상위 문서를 나타낸다. α, β, γ 는 이전 단계의 질의, 적합 문서집합, 부적합문서집합 간의 중요도를 조율하는 상수이다.

식 (13)은 적합문서나 부적합문서의 정보를 각 문서집합의 크기로 정규화하여 질의 확장에 적용하는 방법이고, 식 (14)는 정규화 과정 없이 조율 상수만으로 중요도를 조정하는 방법이며, 식 (15)는 질의를 확장하는데 적합문서는 모두 사용하지만 부적합문서에 대해서는 최상위 문서 하나만 사용하는 방법이다. 여러 가지 실험 결과 위 3가지 식 모두 비슷한 결과를 보인다고 알려져 있다.

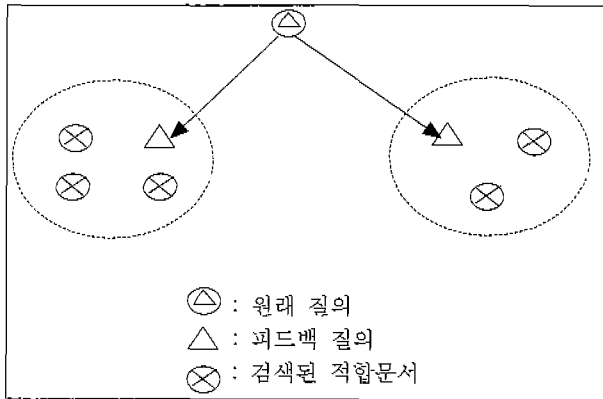


그림 3. 질의 분해를 이용한 적합성 피드백
Fig 3. Relevance Feedback Using Query Splitting

특히 $\gamma=0$ 인 경우, 즉 부적합문서집합의 정보는 사용하지 않고 적합문서 집합의 정보만을 사용하여 질의확장을 하는 경우를 양성 피드백(positive feedback)이라 한다. 양성 피드백은 의사 적합성 피드백에 자주 이용된다.

3. 적합 문서 추출

초기질의와 각 문서 사이의 유사도 계산은 정보검색에서 많이 사용하는 코사인 유사도(cosine similarity) 식 (16)을 이용한다.

$$\text{sim}(S_j, Q^0) = \frac{S_j \cdot Q^0}{|S_j| \times |Q^0|} = \frac{\sum_{i=1}^t w_{ij} \times w_{iq}}{\sqrt{\sum_{i=1}^t w_{ij}^2} \times \sqrt{\sum_{i=1}^t w_{iq}^2}} \quad (16)$$

여기서, S_j 는 각 문서 벡터, Q^0 는 초기질의의 벡터를 의미하고, w_{ij} 와 w_{iq} 는 단어 i 가 각각 문서와 초기질의에서 갖는 가중치이다. t 는 각 문서 벡터와 초기질의의 벡터를 생성하는데 사용된 단어의 총 개수이다. 식 (16)에 의한 유사도값에 따라 문서를 내림차순 정렬한 후 유사도 값이 큰 상위 k 개의 문서를 적합 문서로 간주한다.

4. 실험 및 결과

본 논문에서는 성능 평가 자료로서 CHODIC를 사용하였다. CHODIC은 500개의 문서와 21개의 질의로 구성되어 있다. 문서와 질의 사이의 연관성 평가는 문서 제목을 기준으로 설정하였다.

정보 검색 시스템의 검색 효과는 그림 4에서 정의되는 재현능력도(Recall)와 검색정밀도(Precision)를 이용하여 평가된다.

재현능력도는 문서집합에서 사용자가 원하는 문서를 어느 정도 검색하였는가를 나타내고, 검색정밀도는 검색된 문서들 중에서 사용자가 원하는 문서가 얼마나 포함되어 있는가를 나타낸다.

예를 들어 200개의 문서로 구성된 문서 집합과 관련된 문서의 수가 5개인 질의를 가정하자. 사용자가 검색 시스템을 사용하여 6개의 문서를 검색하였을 때 검색된 문서 중에서 질의에 관련된 문서가 4개라 하면 재현능력도와 검색정밀도는 각각 0.8과 0.67이 된다.

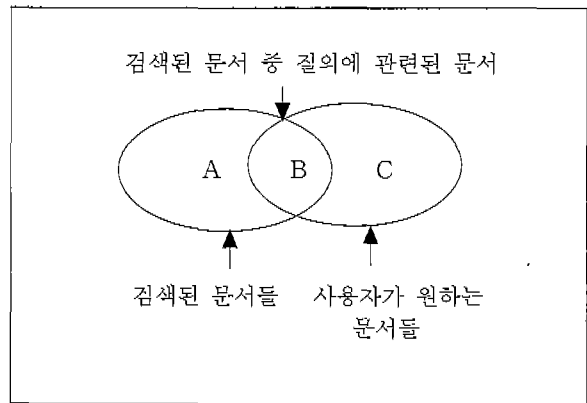


그림 4. 검색 효과 측정 방법
Fig 4. Measure Method of Retrieval Effect

검색정밀도(P)와 재현능력도(R)의 식(17), (18)과 같다.

$$P = \frac{B}{A+B} \quad (17)$$

$$R = \frac{B}{B+C} \quad (18)$$

표 1은 MMM, Paice, P-norm, RF 모델의 검색효과를 보여준다. 본 논문에서는 검색 효과를 평가하기 위하여 질의들에 대한 평균 검색정밀도를 계산한다. 각각의 질의에 대한 검색정밀도는 재현능력도를 0.25, 0.5, 0.75에 고정시켜 계산된 검색정밀도들의 평균값이다. 또한 표에 나타난 검색 효과는 가장 높은 검색 효과를 나타내는 매개변수에 대한 것이다. RF(Relevance Feedback), P-norm 모델이 MMM,

Paice 모델보다 높은 검색 효과를 제공한다.

MMM 모델은 Paice 모델보다 높은 검색 효과를 나타내고 있다. 이는 색인어 가중치가 역문헌빈도와 출현빈도로부터 유도하기 때문에 검색 효과를 저하시키는 요인이 발생하기 않았기 때문이다.

표 1. 검색 효과 비교(단위 : 검색정밀도)

Table 1. Retrieval Effect Comparison(Unit: Precision)

	Average
MMM	0.327
Paice	0.318
P-norm	0.362
RF	0.602

표 2에서는 질의분해 적합성 피드백(Query Splitting Relevance Feedback)의 검색결과가 적합성 피드백(Relevance Feedback) 결과에 비해 검색정밀도가 3.49% 향상을 보여주고 있다.

표 2. 질의 분해 피드백(단위 : 검색정밀도)

Table 2. Query Splitting Feedback(Unit:Precision)

	Average
Relevance Feedback	0.602
QSRF	0.623

표 3과 표 4에서는 검색문서수를 다르게 제한했을 때 검색효율이 어떻게 달라지는가를 분석한 것이다. P-norm, 적합성 피드백(RF), 질의분해 적합성 피드백(QSRF)를 이용하여 검색문서수를 10건과 20건으로 제한한 경우 재현능력과 검색정밀도를 보여 주고 있다.

검색문서수를 10건으로 제한하였을 때 적합성 피드백 검색 결과는 P-norm 결과에 비해 재현능력과 검색정밀도가 각각 56%, 50% 향상되었고, 질의분해 적합성 피드백 결과는 적합성 피드백 결과에 비해 재현능력과 검색정밀도가 각각 2.6%, 3.2%가 향상되었다.

검색문서수를 20건으로 제한하였을 때 적합성 피드백 검색 결과는 P-norm 결과에 비해 재현능력과 검색정밀도가 각각 46.81%, 48.72% 향상되었고, 질의분해 적합성 피드백 결과는 적합성 피드백 결과에 비해 재현능력과 검색정밀도가 각각 1.45%, 5.17%가 향상되었다.

표 3. P-norm, RF, QSRF의 재현능력도 비교

Table 3. Recall Comparison of P-norm, RF, QSRF

검색문서수 \ 척도	재현능력도		
	P-norm	RF	QSRF
문서수≤10	0.25	0.39	0.40
문서수≤20	0.47	0.69	0.70

표 4. P-norm, RF, QSRF의 검색정밀도 비교

Table 4. Precision Comparison of P-norm, RF, QSRF

검색문서수 \ 척도	검색정밀도		
	P-norm	RF	QSRF
문서수≤10	0.42	0.63	0.65
문서수≤20	0.39	0.58	0.61

그림 5는 검색 문서수 20건으로 제한하였을 때 P-norm 초기검색과 적합성 피드백, 질의분해 적합성 피드백 검색결과를 재현능력도와 검색정밀도로 표현한 성능곡선으로 비교한 것이다. 재현능력도의 증가에 따른 검색정밀도의 하강 현상이 두드러지지 않고 있음을 볼 수 있다.

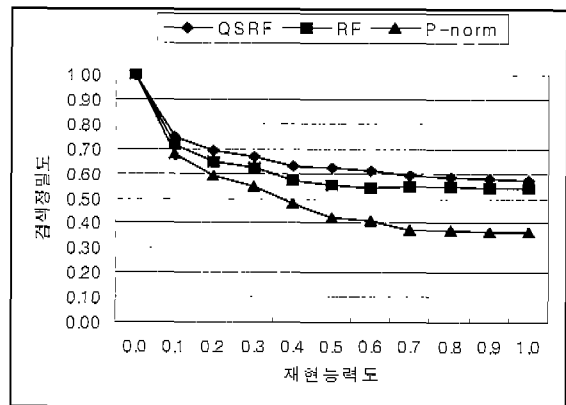


그림 5. P-norm 초기와 RF, QSRF의 실험 결과
Fig 5. Experimentation Result of P-norm Initial and RF, QSRF

5. 결론

정보 검색 시스템의 중요한 목적중의 하나는 단순히 사용자 질의를 만족하는 문서들의 집합을 검색하는 것이 아니라, 질의를 만족하는 정도에 따라 검색된 문서들에 순위를 부여함으로써 사용자가 필요한 정보를 얻는데 소모되는 시간을 최소화시키는 것이다.

이러한 부울 검색 시스템의 단점을 보완하기 위해 퍼지집합 모델이 제안되었다. 그러나 퍼지 집합 모델은 단일 피연산자 의존 문제로 인하여 많은 경우에 부정확한 문서값을 생성하는 것으로 알려져 왔다.

퍼지 집합 모델의 문제점을 개선하는 방법으로서 MMM 모델, Paice 모델, P-norm 모델이 개발되었다.

본 논문에서는 정보 검색 분야에서 사용되는 적합성 피드백에 기초하여 높은 검색 효과를 제공하는 질의 분해 적합성 피드백(QS-RF) 모델을 제안하였다. 실험 결과 제안하는 질의 분해 적합성 피드백 방법으로 정보검색하는 경우가 질의 분해 적합성 피드백을 이용하지 않는 경우 보다 더 좋은 정밀도를 보였다.

QS-RF 모델은 기존의 방법들보다 높은 검색 효과를 제공함을 성능 비교를 통하여 입증되었다.

참 고 문 헌

[1] Ricardo Baeza-Yates and Berthier Ribeiro-Neto, Modern Information Retrieval, Addison-Wesley Publishing Company, 1999.

[2] Daniel Marcu, Discourse trees are good indicators of importance in text, In Inderjeet Mani and Mark Maybury, eds, *Advances in Automatic Text Summarization*, pp. 123-136, The MIT Press, 1999.

[3] Mark Sanderson, Accurate User Directed Summarization from Existing Tools, In *Proceedings of the 7th International Conference on Information and Knowledge Management*, pp. 45-51, 1998.

[4] Regina Barzilay and Michael Elhadad, Using Lexical Chains for Text Summarization, In Inderjeet Mani and Mark Maybury, eds, *Advances in Automatic Text Summarization*, pp. 111-121, The MIT Press, 1999.

[5] Anastasios Tombros and Mark Sanderson, Advantages of Query Biased Summaries in Information Retrieval, In *Proceeding of ACM-SIGIR'98*, pp. 2-10, 1998.

[6] J.H. Lee, M.H. Kim and Y.J. Lee, Information Retrieval Based on Conceptual Distance in Is a Hierarchics, *Journal of Documentation*, vol. 49, no. 2, pp. 188-207, 1993.

[7] M.H. Kim and J.H. Lee and Y.J. Lee, Analysis of Fuzzy Operators for High Quality Information Retrieval, *Information Processing Letters*, vol. 46, no. 5, pp. 251-256, 1993.

[8] G.Salton, *Automatic Text Processing : The Transformation, Analysis, and Retrieval of Information by Computer*, Addison Wesley, 1989.

[9] J.H.Lee, W.Y. Kim, M.H. Kim and Y.J. Lee, Enhancing the Fuzzy Set Model with Positively Compensatory Operators, *Proceedings of the 3rd International Symposium on Database Systems on Advanced Applications*, Taejon, Korea, pp. 368-375, 1993.

저 자 소 개



이 성 주(Sung-Joo Lee)

1970년 : 한남대학교 물리학과 (이학사)
 1992년 : 광운대학교 전자계산학과(이학석사)
 1998년 : 대구가톨릭대학교 전자계산학과 (이학박사)
 1988년~1990년 : 조선대학교 전자계산소 소장
 1995년~1997년 : 조선대학교 정보과학대학장
 1981년~현재 : 조선대학교 컴퓨터공학부 교수

관심분야 : 소프트웨어 공학, 프로그래밍 언어, 객체지향 시스템, 리프 집합



박 병 권(Byung-Gweun Park)

1988년 : 조선대학교 전자계산학과 (학사)
 1990년 : 조선대학교 전자계산학과 (석사)
 2000년 : 조선대학교 전자계산학과 (박사)
 1991년~1994년 : 광주은행 전산업무부
 1995년~현재 : 서강정보대학 정보통신과 조교수

관심분야 : 소프트웨어공학, 객체지향시스템, 퍼지집합, 리프집합



김 영 천(Young-Chon Kim)

1992년 : 광주대 전자계산학과 졸업
 1996년 : 조선대 컴퓨터공학과 졸업 (공학석사)
 1998년~현재 : 조선대 전자계산학과 박사 과정

관심분야 : 객체지향시스템, 소프트웨어 공학, 유전자알고리즘, 정보검색