

索引의 理論에 대해서

B. C. Landry & J. E. Rush 著

南 台 祐 譯

(전남대학교 강사)

本 論文은 Proceedings of the ASIS, 5, (1968)에
게재된 "Toward a Theory of Indexing"을 完역한
것이다.

이 論文의 目的은 情報의 蓄積과 檢索의 理論을
위한 基本的인 要素를 제시하는데 있다. 이러한 理
論은 索引作成의 一般的인 理論을 證立함이 최선의
方法이라 믿어진다. 이론에 대한 基本的인 前提를
언급하고, 필수 불가결한 概念을 정리한 後에 索
引의 理論과 情報蓄積과 檢索과의 關係성을 고려해
본 것이다.

索引作成과정과 一般的인 커뮤니케이션과정과의
유사성이 논의되어졌으며, 索引作成은 정보를 증가
(즉, 엔트로피의 감소)시키는 작용의 절차로서 간
주하였다.

索引의 理論的인 概念은 現實世界에서의 索引作成
시스템과 關係되어 발전해 왔으며, 역조현상을 나
타내기도 하였다. 질문의 공식, 檢索과 檢索과의 關
連성이 논의되어졌으며, 이들의 概念은 人間組織의
效率性과의 關係성을 갖고 있다.

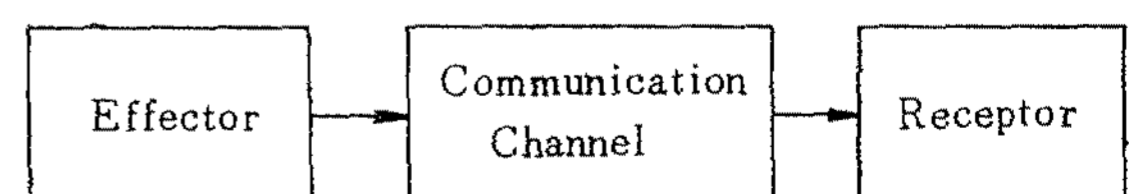
이 論文에서 提示한 觀點은 索引作成 과정속에서
더욱 상세한 調査 관찰을 위한 유용한 골격(frame-
work)을 제시하는데 주目的이 있다. <訳者註>

I. 索引의 理論에 대해서

情報를 蓄積시키고, 檢索하는 것에 대한
研究目的은 擴大되는 人間의 知識을 效果的
으로 体系化시키는 問題에 대한 解決策을
摸索하는데 있다. 이러한 問題를 解決하기

위한 多樣的인 試圖가 이루어 졌지만, 論理가
整然한 學問에도 基礎가 되는 問題解決에
대한 基本的인 模型이 없는 까닭에 原則的인
인 面에서 問題가 발생되고 있다. 現在 매
우 集中的으로 努力이 展開되고 있는 그 自
체도 情報의 蓄積과 檢索에 대한 理論的인
基礎의 認識이 否足하여 方向設定이 잘못되
고 있는 것이다. 이러한 緣由로써 본 論文
은 情報의 蓄積과 檢索에 대한 原則的인 理
論의 基礎를 마련하는데 그 目的을 두고 있
다.

情報의 流通 또는 커뮤니케이션은 送信者
(effector)와 受容者와의 커뮤니케이션링크
(communication link)의 形式으로 構成된
다고 할 수 있다(그림 1). 이러한 模型에서
情報는 送信者로 부터 커뮤니케이션의 채널
(channel)을 통해서 受容者에게로 傳達된다.
이러한 傳達速度는 時間當 傳達되는 情報의
要素로서 측정된다. 情報의 要素를 잘못 表
現하든지 혹은 流通을 制限하여 情報要素에
손실을 가져오는 어떠한 理論이나 實際도



<그림 1> 일반적인 情報傳達의 模型

부적당한 것으로 간주되어야 한다. 이러한 것을 前提로 할 때 본 論文은 앞으로의 研究調査에 基礎와 展望을 提供해 주기 위해서 情報의 蓄積과 檢索의 意味를 再評價해 보고자 한다.

II. 커뮤니케이션과 索引作成

情報의 蓄積과 檢索은 커뮤니케이션 過程과 관련성이 있다. 發送者와 受容者와의 커뮤니케이션은 그들이 서로 관련성을 갖게 되어 어떤 結果가 발생된 것이다. 文档테이션의 경우에도(여기에서 文档테이션은 情報의 創案에서부터 情報를 利用하기까지 커뮤니케이션의 모든 過程을 말한다) 이것은 情報를 획득하려는 사람이 情報組織에 대해 予備知識과 利用者가 願하는 情報를 檢索할 수 있는 記号(表現物)를 갖고 있어야 한다는 것을 意味한다.

이러한 最近의 分析은 우리가 手作業을 考慮하거나 컴퓨터시스템을 고려하든 간에 우리는 이것들과 같은 문제에 부딪치게 된다. 主된 要点은 情報의 流通과 그리고 궁극적으로 정보검색이 索引作成의 과정에 달려있는 것이다. 索引作成은 情報流通過程에서 發生하는 概念의 순서를 나타내며 그들 情報相互間(즉, 情報流通過程에서 發生되는 情報要素의 認知模型)의 關係를 決定해야 한

다. 그러나 우리는 現在의 索引理論이 이러한 基準에 미치지 못한 것으로 믿고 있다. 즉, 文档에서 概念間의 關係나 순서를 명확히 하여 적절한 情報에 接近하도록 하는데 失敗하고 있다는 것이다(여기에서 文档은 정돈이 잘된 情報要素의 集書를 의미한다).

최근의 索引作成의 過程에서 文档은 소수의 “중요한” 概念의 集書로 간주된다. 이 “중요한”이란 낱말은 매우 주관적인 概念이다. 이러한 중요한이라는 概念을 부여 받으면, 그러한 情報는 그것과 비슷한 개념들과 함께 文档에서 따로 분리되어 蓄積된다. 이러한 處理過程은 文档의 개념을 파괴할 뿐만 아니라 索引을 위해 選択된 概念들간의 關係도 파괴한다. 따라서 索引을 무시하는 그러한 原則은 그것이 文档에서 情報要素間의 유일한 組織을 模糊하게 하는 단점이 있다.

이 分野를 좀 더 理解하기 위해서는 현재의 情報의 蓄積과 檢索에 대한 理論은 現代數學의 거의 모든 분야를 導入해야 한다.¹⁻⁵⁾

그렇지만 情報의 蓄積과 檢索을 위해 기초가 되는 理論과 概念的인 運作的 發展을 統一化시키고 다가오는 미래를 연구하는 노력을 통해서라기 보다 情報의 蓄積과 檢索에 대한 理論을 定立시키려는 시도가 여러가지 문제점을 개별적으로 설명하려는 이론을 통해서 나타난 것이다. 따라서 特定문

1) A. G. Dale and N. Dale, "Some Clumping Experiments for Associative Document Retrieval", Amer. Doc., 16, 1, 5-9 (1965).
 2) F. J. Damerau, "Experiment in Automatic Indexing", Amer. Doc., 16, 4, 283-289 (1965).
 3) D. J. Hillman, "Study of Theories and Models of Information Storage and Retrieval", Center for the Information Sciences-

Lehigh University, Bethlehem, Pennsylvania (1963-1967).

4) M. E. Maron, "Automatic Indexing: An Experimental Inquiry", J. Assoc. Computing Machinery, 8, 4, 404-417 (1961).

5) H. E. Stiles, "The Association Factor in Information Retrieval", J. Assoc. Computing Machinery, 8, 2, 271-279 (1961).

제를 다루기전에 一般的인 概念들이 定立되어야 한다. 우리는 이러한 개념의 정립이 索引의 개념을 다시 계통적으로 설명함으로써 가능하다고 생각된다.

III. 情報伝達에 있어서의 索引의 役割

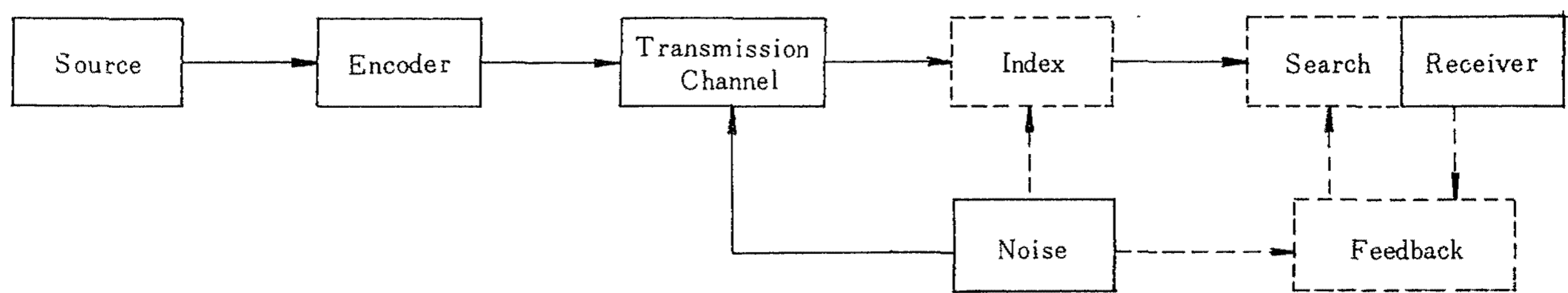
實際상의 모델로서, 다음과 같은 Shannon의 一般화된 커뮤니케이션 圖式⁶⁾(그림 2)에서 考慮해 보면, 이 모델에서 우리들은 情報의 伝達過程에서의 索引의 役割과 그 位置를 알 수 있다.

만일 우리가 情報伝達過程에서 어떤 시점에서 횡단면적으로 본다면, 文獻중의 各개의 情報은 이러한 각각의 시점에서 한 개의 點으로 볼 수 있다. 索引은 情報要素들간의 類以點과 差異點을 明確히 해야 하며, 각 시점간의 안에서 그들사이의 연결關係를 인지해야만 한다. 다시말해서 索引은 다양한 위치에서 情報의 伝達을 可能케 해야 하며, 상관적인 方向과 情報要素의 유통량도 이해할 수 있어야 한다. 차후에 좀더 명확하게 되겠지만, 索引作成은 어떤 시점에서의

情報의 流通상태를 完全하게 표시해 주는 相關성의 集合體로 볼 수 있다. 또한 원칙적으로 이들의 關係성은 역으로도 되어야 하는데 즉, 그들은 文獻도큐멘트의 유통량을 재조직할 수 있어야 한다는 것이다. 索引은 또한 情報의 伝達과 受信사이의 여과장치로도 說明된다. 즉 라디오의 주파수와 같은 비슷한 역할을 수행한 것이다. 여과장치로서 索引의 獨특한 特성은 索引이 복잡한 정보들로 부터 전달할 부분을 선택해야 하고, 수신과정은 고려해야 한다는 것이다.

IV. 索引의 節次上 特性

Yovits와 Ernst는 情報를 意思決定에 있어서 가치가 있는 데이터(data)로써 定義를 내리고 있다⁷⁾ 分明하게 말해서 情報要素(data)는 意思決定過程에서 最적의 時間에 最적의 형태를 갖추어서 이용가능한 것이어야 한다는 것이다. 이것이 바로 索引의 機能인 것이다. 우리가 실제 현실생활의 경험에서 얻는 데이터는 보통 잘 정리된 報告書나 文獻이다. 이러한 잘 정리된 데이터(즉, 文獻)의 集合은 소위 Document space



<그림 2> Shannon의 커뮤니케이션의 모델에서는 索引(---)을 포함시켜서 적용하고 있다. 雜音(noise)은 정보를 부정확하게 分析하게 하는 것이다.

6) C. E. Shannon and W. Weaver, The Mathematical Theory of Communication, University of Illinois Press, 1964.

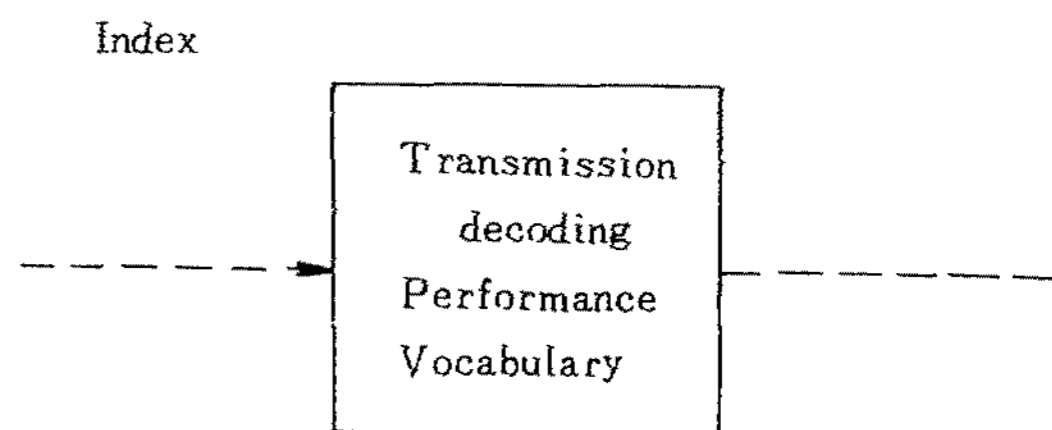
7) M. C. Yovits and R. L. Ernst, "Generalized Information Systems", in

Electronic Handling of Information: Testing and Evaluation, A. Kent, ed., Thompson Book Co., Washigton, D. C., 1967. 279-290.

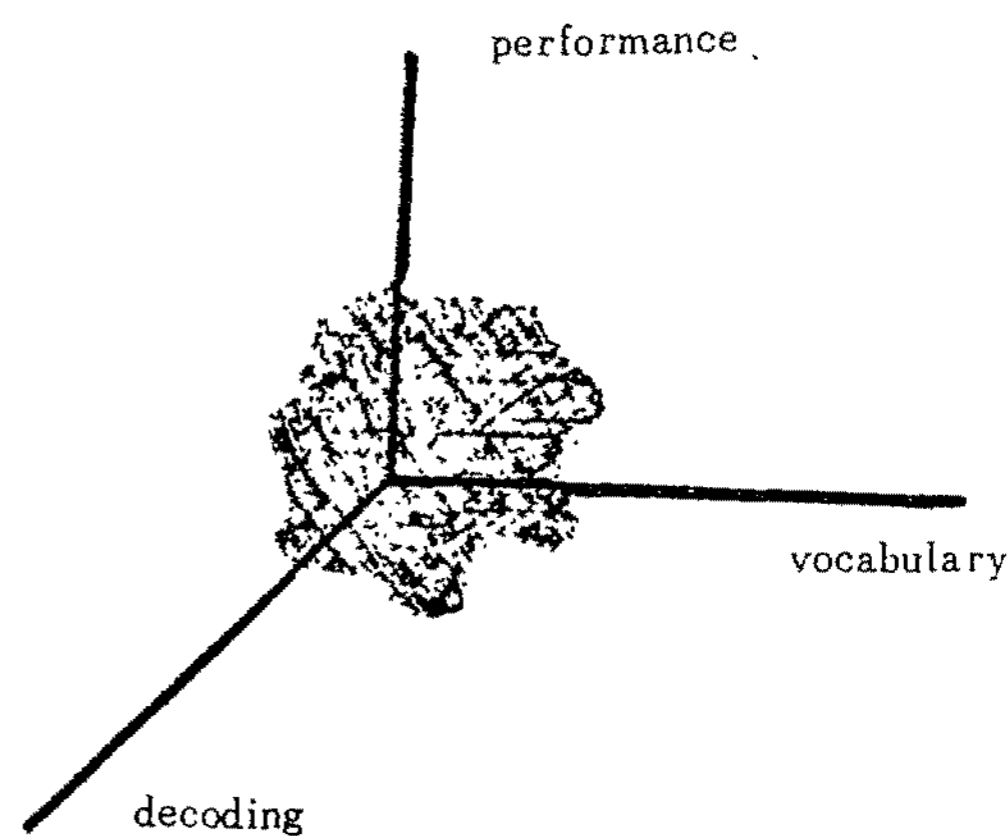
를 형성하게 된다. Document space의 要素는 資料들간의 뒤섞인 배열순서를 보여줄 뿐이다. 巨視眼的으로 볼 때 도큐먼트사이 에 존재하는 순서는 다만 색인작성시스템에서 도착되는 각각의 시간을 의미한다. 따라서 索引作成의 시스템은 도착되는 도큐먼트의 모든 정보요소들의 순서와 상호관련성을 결정해야 한다. 그러한 조직과정을 受信者 모두도 알아야 한다.

Document space의 概念과 대등한 것으로 우리들은 Index space⁸⁾를 Performance, Vocabulary 그리고 Transmission decoding의 3가지 集合 Cartesian積으로 定義한다. Performance 集合의 要素는 이미 索引되어 있는 것에 의존하는 1) 再現率(recall)과 개개 探索者에 의존하는 2) 適合率(relevance), 그리고 그 댓가에 대한 費用으로서 또는 意思決定過程에서의 가치로서 表現되는 3) 利益(benefits)에 기초를 두고 索引作成⁹⁾의 척도가 된다. Vocabulary 요소의 集合은 索引語가 되며, 자주 差異性, 同等性, 一般性的의 集合과 함께 語彙의 순서표의 형식으로서 그 자체가 명백하게 나타난다. 語彙는 情報檢索에서 또한 중요한 구성요소가 된다. 왜냐하면 모든 檢索은 이들과 똑같은 단어와의 相關성으로 이루어져 있기 때문이다. Transmission decoding 은 실제 索引作成의 段階 또는 情報要素間的 關係를 나타내는데 필요로 하는 技法의 집합이다.

Index space의 개념으로 우리들은 다음과 같은 結論을 導出해 낼 수 있다. 檢索過程은 Index space로 볼 수 있고 수신자가 요구하는 특수한 정보는 Index space의 요소



〈그림 3〉 索引의 要素



〈그림 4〉 Index-space

로서 구성되어야만 한다는 것이다. 이러한 概念으로 推理해 볼 때, 檢索의 過程은 요구되는 정보요소를 Index space의 요소안에 同形異義語로 배치하는 것으로 볼 수 있다. 理想的인 경우로서는 적어도 개개의 探索者만큼 많은 수의 順次的인 배치가 있어야 한다. 이러한 것으로 볼 때 索引의 本質은 情報의 傳達과 探索者와의 사이에 존재하는 公有영역이다(그림 3참조). 제 3자에게는 索引의 구성요소들이 불확실한 범위의 영역에서 나타나는 것처럼 보이는 반면에(그림 4를 보라) 개별적인 受容者와 相關되는

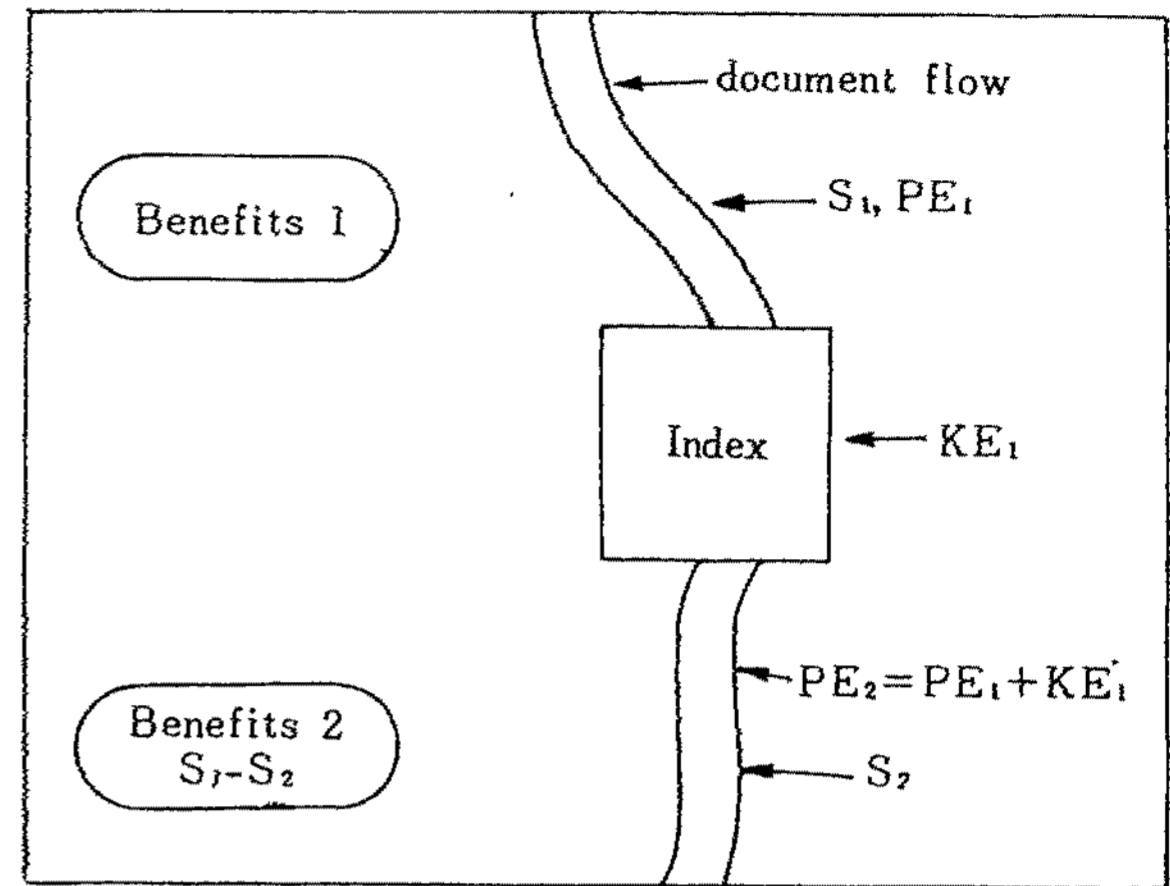
8) J. Rothstein, "An Overview from an Informational-Organizational-Operational Viewpoint," in Communication: Concepts and Perspectives, L. Thayer, ed., Spartan

Books, New York, N. Y., 1967, 397-423.

9) J. A. Swets, "Effectiveness of Information Retrieval Methods", AD-656-340, (1967)

Index space의 요소의 위치는 명확하다.

이러한 영역은 개개의 受容者를 나타내는 점들의 전체적인 平均이며, 그리고 空間에서의 각 점은 Performance와 Vocabulary, 그리고 decoding의 관계의 한 기준이 되는 것이다. 平面에서의 한 점은 제 3자가 受容者에 의해 決定되는 위 3가지의 기준에 따르는 중요성을 이해만 한다면 매우 명확해질 것이다. 따라서 제 3자가 檢索을 不分明한 受容영역의 연속과정이라고 보는 이유를 理解하게 될 것이다.



〈그림 5〉 도큐멘테이션 시스템

V. 索引과 엔트로피(Indexing and Entropy)

情報(도큐멘트)는 索引作成過程에 매우 무질서하게 접수된다(그림 2).

한 질의 도큐멘트에서 情報要素 사이에서는 分명한 관계가 존재하지 않는다. 이와같이 索引作成시스템에서는 다양한 도큐멘트들간의 순서를 결정시키거나 또는 정보의 엔트로피를 감소시켜준다. 그렇지만 열역학적으로 볼 때 이러한 엔트로피 감소는 다른 도큐멘테이션 시스템에서의 엔트로피의 증가를 가져온다. 그러한 엔트로피의 增加는 索引作成過程에서 얻어지는 혜택으로 볼 수 있다¹⁰⁾.

그림 5에서는 도큐멘테이션 시스템에서 도큐멘트의 유통과정을 나타내 주고 있다.

索引作成前에는 도큐멘트의 集合은 잠재적인 에너지 PE_1 과 무질서의 단위(엔트로피) S_1 을 갖는다. 索引作成過程은 운동에너지

의 KE_1 의 入力を 필요로 하는 업무로서 설명된다. 이와 같이 索引作成後(도큐멘트의 정보요소의 순서)에는 도큐멘트수집의 잠재적인 에너지가

$$PE_2 = PE_1 + KE_1 \dots \dots \dots (1)$$

수준으로 증가한다. 이러한 잠재적인 에너지의 증가는 $S_2 < S_1$ 등식과 같이 S_1 에서 S_2 로의 엔트로피에서 중요한 감소로 영향을 미친다.

시스템에서의 엔트로피는 변함이 없이 남아 있어야 하는데 이르기 위해서는 잠재적인 혜택의 엔트로피가 $S_1 - S_2$ 만큼 증가해야 한다. 즉 광범위한 혜택이 이용될 수 있어야 하며, 정보요소의 전체가 意思決定과정에서 보다 유용해져야 한다. 그리고 엔트로피가 보다 큰 가치를 가져야 한다 [Yovits와 Ernst⁷⁾의 관념].

이러한 혜택을 측정하는 한 척도는 索引作成過程에서 생성되는 경제적인 힘이다. 이 방법에서 혜택을 B, 정보생산비용을 Cd,

10) J. W. Murdock and D. M. Liston, "A General Model of Information Transfer : Theme Paper 1968 Annual Convention", Amer. Doc., 18, 4, 197-208(1967).

11) D. J. de Solla Price, "Research on Rese-

arch", in Journeys in Science : Small Steps-Great Strides, D. L. Arm, ed., University of New Mexico Press, Albuquerque, New Mexico, 1967, 12-13

索引作成費用을 C_i 이라고 가정하면

$$B = Cd + C_i \dots\dots\dots(2)$$

가 된다. 索引作成으로 부터 얻는 惠澤은 상당한 時間差를 습하면 된다. 즉 전체의 혜택 B_T 는 각 시간에서 혜택의 습 t_i 가 된다. 이것은

$$B_T = \sum_{i=1}^N B_{t_i} \dots\dots\dots(3)$$

가 된다.

N 이란 $N < i$ 일 경우 $B_{t_i} = 0$ 으로서 최대량의 시간이 된다. 좀더 간결하게 해 보면, B_T 가 어떤 극한값에 수렴한다는 것이다. 그리고 어느 시간이고 주어진 시간 t 에서 惠澤은 상수이므로 정보생산에서의 비용을 감소시킨다면, 索引으로부터의 效率性은 提高되는 것이다.

VI. 再現性 語義論과 構文論

完全한 索引作成시스템은 完全한 再現과 適合性에 있다. 즉, 索引作成시스템은 探索者로 하여금 모든 適用可能한 文獻를 檢索할 수 있도록 하게 한 것이다. 文獻를 探索한 결과에 대한 절대적인 適合性은 보장될 수가 없다. 適合性은 개별 探索者의 相對的인 만족도의 側面에서만 測定될 수 있다. 따라서 探索者의 特別한 要求條件이 明確하게 公式化되어 표시되어 있다면 再現性의 意味가 明白해진다. 現實 世界에서 檢

索시스템이 이러한 요구조건에 부응치 못하게 될 理由는 情報의 索引作成이 不完全하기 때문이다.

最近에서야 再現性(recall)과 適合性(relevance), 그리고 再現성과 適合성에 대한 制約性에 상당한 注意力을 기울이기 시작했다.¹²⁻¹⁶⁾ 그러나 더욱 重要한 概念인 構文論과 語義論^{17,18)}은 不幸히도 索引作成에 종사하는 사람들이 새로운 관심을 거의 기울이지 않고 있다. 이러한 概念들은 索引作成의 過程이 어떻게 커뮤니케이션에 있어서 정보요소들의 體系化를 불명료하게 하는가에 대한 좋은 例가 된다. 情報는 自然語의 형태에서 索引으로 도달되는데, 이것은 構文論과 語義論의 관계가 강한 것을 보여준다. 그렇지만 그러한 관계는 檢索結果가 빈곤하다는 것에서 알 수 있는 것처럼 索引作成이후에 이들의 관계는 상당히 弱化된다. 링크(link)나 룰(role)과 같은 分類시스템은 단지 情報要素를 위해서 構文論 또는 語義論의 작성자의 제한된 인원으로 제공해 주는 이유로 다만, 다양한 수준에서 문제점의 所在를 제공해주고 있을 뿐이다.^{15,19)} 그러므로 정보의 단위가 링크와 룰과 같은 그러한 종류에 정확하게 속하지 않고, 커뮤니케이션에 있어서 構文論이나 語義論의 잠재적인 오류가 檢索을 실패로 이끌지도 모른다는 점은 명확하다. 이러한 난관은

12) W. Goffman, "On Relevance as a Measure", Inform. Storage Retrieval, 2, 3, 217-220 (1964)
13) C. J. Maloney, "Semantic Information", AD-636-871, (1962).
14) M. Rubinoff, "Semantic Tools in Information Retrieval", AD-660-087, (1967).
15) J. D. Sinnett, "An Evaluation of Links and Roles Used in Information Retrieval", AD-432-198, (1963).
16) M. Taube, "A Note on the Pseudo-Mathematics

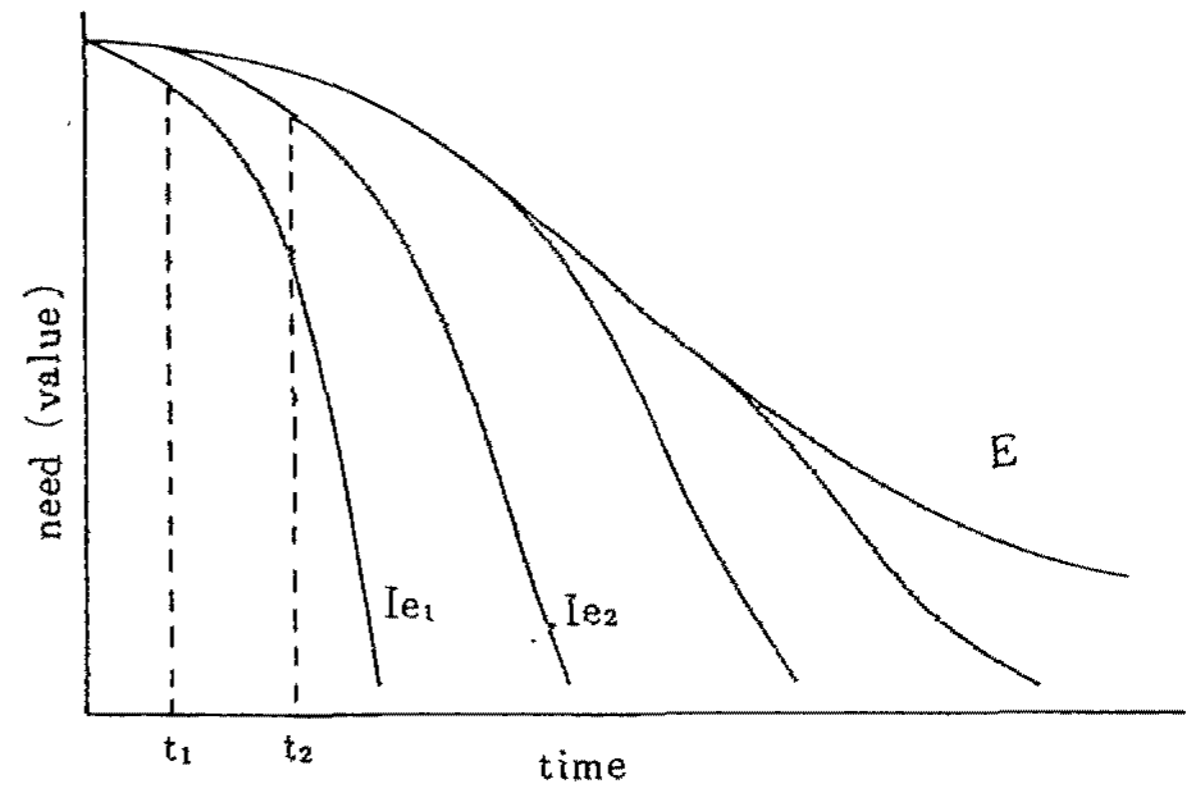
of Relevance", Amer. Doc., 16, 2, 69-72 (1965)
17) L. B. Doyle, "The Microsyntax of Text", System Development Corporation, SP-1083, (1963).
18) L. B. Doyle, "Semantic Road Maps for Literature Searchers", J. Assoc. Computing Machinery 8, 4, 553-578 (1961)
19) M. Taube, "Notes on the Use of Roles and Links in Coordinate Indexing", Amer. Doc., 12, 2, 98-100 (1961)

正本의 文獻에 있어서 句讀符號와 語의 關係의 索引으로써 제거될 수 있다.

VII. 人間의 限界性(The Human Limitation)

情報蓄積과 檢索은 全적으로 人間이 組織하는 領域의 중심이 된다. 따라서 情報의 蓄積과 檢索은 人間행위 그 자체와 關連되어 있는 것이다. 더욱기, 索引作成과 情報檢索의 궁극적인 目的은 人間組織 시스템에 그 基礎를 두고 있고 따라서, 어떠한 索引作成시스템에서도 人間組織의 結晶을 반영하게 된다. 아직까지 行動科學이 人間行爲에 對한 精確한 理論을 定立하지 못했을 뿐만 아니라 人間心理를 完全히 파악하지 못했기 때문에, 우리는 더욱 情報의 精確한 檢索과 蓄積의 理論을 기대할 수가 없다. 특히 索引作成에서 그 精確도를 기대하기에는 더욱 어려운 事實이다. L. A. Zadeh²⁰⁾ fuzz집합이론의 精神에서 우리는 fuzz이론에 對해서 현실에 對한 조사치에 만족해야만 한다.

情報과 人間行爲를 考慮하는 定義에 비추어 볼 때 단순한 檢索作業이 一般적으로 意思決定을 하는데 요구하는 情報요인의 모든 것들을 나타내는데 충분하지는 않다. 一般적으로, 첫번째 檢索되는 情報요인은 단지 요구하는 것에 對한 접근에 불과하다. 다음에는 索引과의 相互作用에 의해 더욱 잘 定義된 接近치를 제공해 줄 수 있다. 意思決定은 본질적으로 質問과 對答의 結果이므로, 意思決定過程에서 運用할 수 있는 檢索된 情報요인의 精進적인 습을 意味한다. 아래 D.



〈그림 6〉 情報가치 쇠퇴

C. Stone²¹⁾의 그림에서, 우리는 意思決定에 있어서 그들의 役割에 對해서 行動의 Poisson型을 나타내는 것과 같은 情報요인으로 간주할 수 있다. 即, 意思決定에 있어서 어떤 情報要因의 가치는 시간이 경과됨에 따라 감소된다. 그렇지만 檢索이 더욱 特수화됨에 따라서, 가치(요구정보)의 증가는 반감된다.

그림 6은 각각의 새로이 檢索된 要素의 가정된 행위, 즉 Ie_i ($i=1, 2, 3, \dots$)를 나타낸 것이고, 또한 각 情報單位는 다음 질문의 필요성에 對한 情報가치의 반감비율을 나타내 주고 있는 것이다.

포괄선(envelope curve) E는 要因의 特수성에 對한 逆임을 나타내 주고 있다. 바꿔 말하면 情報要因이 더욱 精確해 짐에 따라 索引에 對한 情報적인 놀라운 가치의 量이 제로(Zero)에 가까운 것이다.

VIII. 結 言

두 가지 主要한 論點은 이제 前章의 論議에서 明白해졌다.

20) L. A. Zadeh, "Fuzzy Sets", Inform. Control, 8, 338-353 (1965).

21) D. C. Stone, "Word Statistics in the

Generation of Semantic Tools for Information Systems", AD-664-915, (1967)

- 1) 索引作成시스템의 여과를 거친 情報는 증가된 組織과 더욱 밀접한 상호관계를 나타내야만 한다.
- 2) 索引作成시스템에서 얻어진 어떤 혜택을 실현키 위해서는 情報를 組織化하는 비용과 정보를 생산하는 비용은 같아야 한다. 색인작성시스템에서 얻어진 혜택이 색인에 의해 확장된 능력을 인정한다면, 아마도 이러한 상충되는 개념들은 정말로 상호 공존할 수 있다.

정보를 生産하는 비용이 통제없이 증가되지 않는다는 것은 인용주 (3)에서와 마찬가지로 명백하다. 索引은 통제되지 않는 비용 조건 아래에서는 경제적으로 불가능한데 이는 모든 혜택이 잠정적인 혜택으로 남아 있기 때문이다. 情報處理過程에 대한 現在의 狀況은 정보생산의 호황으로 유추해 볼 수 있다. 그래서 더욱 양질의 정보를 檢索하는 일은 어렵게 되었다. 따라서 우리는 情報檢索에 대한 더욱 효율적인 方法을 찾는다. 우리는 索引作成과 檢索의 효율성을

높이기 위하여 전송된 정보를 집중시키는 더 좋은 方法을 摸索하여야 한다.

최대한의 조직화와 혜택의 目的을 달성하기 위해서는 본 논문에서 논의된 개념에 따라 정보시스템이 발달되어야만 하는 것은 명백한 사실로 입증된 것이다. 그렇지만, 索引作成의 요소(효용성, 어휘, 記号化, 적합성, 재현성, 혜택과 필요성)등은 당연히 비슷한 실체이므로 우리는 索引의 fuzzy 이론 밖에 가질 수 없도록 제한성이 있다고 강조되고 있다. 관련된 인간조직 영역이 학습과 의사결정과정에 기반을 두고 있기 때문에 색인작성시스템이 적용적이고 상호의존적²²⁾ 이어야만 하는 사실은 이제 명백하다. 索引作業이 Ohio 주립대학교에서 행해지고 있는 것은 이러한 관점에서이다. 우리의 노력은 기본적인 索引理論, 自動化索引절차, 데이터축소와 索引作業에 대한 컴퓨터 시뮬레이션의 업무에 직접 관련되어 있는 것이다.

22) H. Boroko, "Interactive Displays for

Document Retrieval", AD-661-657,