

고문헌 저작자의 진위판별을 지원하는 시스템 설계

이근무*, 이근우**

* 경주대학교 컴퓨터공학과

** 부경대학교 사학과

The Study on Author's Determination Supporting

System Design to Ancient Literature

Rhee KunMoo, Rhee KunWoo

*Dept of Computer Science, Kyongju University

**Dept of History, Bukyung University

요약

이 논문에서는 현재 그 진위가 한국 고대사 연구의 초미의 관심이 되고 있는 화랑세기 등 고대사 저술의 진위를 판단하는 과학적 방법론을 제안하는 데 있다. 이런 방법론의 전통은 성서의 바울서한의 저자의 진위 논쟁, 셰익스피어 작품의 진위논쟁 등 세계적 관심이 되는 것에서부터 정치적 사건 및 개인의 송사에까지 다양한 스펙트럼에서 논구될 수 있으며 이런 결과들은 현재 우리의 인문학계 특히 고대사학과 민족 정체성에 대한 거대 담론들에 대한 여러 형태의 유용한 실증적 증거를 마련해주게 될 것이다. 또한 다학문적, 학제적 연구의 새로운 모멘텀이 될 수 있을 것이다.

1. 서론

이 논문에서는 화랑세기 등 상고 및 고대사의 여러 역사 저술의 진위를 결정하는 과학적 방법론을 제안하는 데 있다. 이런 방법론의 전통은 성서의 바울서한의 저자의 진위 논쟁, 셰익스피어 작품의 진위논쟁 등 세계적 관심이 되는 것에서부터 정치적 사건 및 개인의 송사에까지 다양한 스펙트럼에서 논구될 수 있다. 문헌의 작자에 대한 진위 판정은 지금까지 첫째 기록 내용에 대한 검토와 역사적 사실에 대한 고증, 둘째 원본 자료가 있을 경우 그 필적

의 감정, 셋째 사용된 종이의 지질, 잉크 혹은 먹물 등의 화학적 분석에 의한 연대추정에 의존하는 것이 보통이었다. 그러나 원저자의 원고가 남아있지 않을 경우에는 필적감정이나 화학적 분석을 할 수 없거나 역사적 사실이나 고증을 할 수 없는 경우가 존재한다. 특히 고대 문헌의 경우에는 필사본으로서 전해오는 경우에는 필적감정이나 화학적 분석이 불가능하다. 또한 근년에 와서는 컴퓨터를 이용한 편집작업이 진행되면서 이러한 전통적인 방법으로는 문헌의 진위판단에 한계가 있다. 따라서 이러한 문제해결을 위한 새로운 과학적 방법이 필요하게 되었다. 여기에서 새롭게 대두된 방법론들이 한 문장의

길이, 단어의 문자수, 품사의 출현율, 특수한 말의 출현율 등의 문장의 수량적 특성을 조사하여 분석함으로써 저자의 문장스타일과 습관 등을 계량적으로 파악하여 문헌의 진위 여부를 판정하는 방법론이 제안되고 있다. 문헌의 스타일과 습관을 확인한다는 것은 개인에게 지문이 개인을 구별하는 증거가 되는 것처럼 개인의 문장의 특성 역시 진위 판단의 근거가 될 수 있다는 가정에 근거하고 있다.

이러한 계량적 방법에 의한 작자 진위에 대한 연구는 유럽의 경우 100여 년 이상의 전통을 가지고 있으나 우리나라에서는 연구자들의 식견으로는 아직 일천한 것이 사실이다. 현재의 고대사와 관련되어 그 진위가 의심받는 화랑세기와 같은 문체 저작들에 대한 진위판정에 효과적으로 이용할 수 있는 방법론과 시스템을 설계하고자 한다. 이 논문은 제1절, 제2 절에서는 제안된 진위 판정의 분석 기법들을 제시하고 제3절에서는 관련 연구를 살펴보고, 제4 절에서는 저작자 판별지원 시스템을 제안하고 제5절에서 결론을 맺는다.

2. 저작 스타일과 습관의 진위 검증을 위한 방법

1) 비율

특정 단어의 출현율, 품사의 출현율, 총단어 수에 대한 각기 다른 단어의 비율 등

$p =$ 특정단어/총단어의 수

예를 들어 명사의 출현비율=문헌중에 나타난 명사의 수/ 문헌중의 총 단어의 수

2) 평균치(mean), 최빈치(mode), 중앙값(median) 범위(range), 표준편차(standard deviation), 분산(variance)

작가의 한 문장의 길이가 길 수도 짧을 수도 있다. 이럴 경우 예를 들어 한 작가의 한 문장의 평균단어 수 혹은 문자의 수의 평균값을 이용해 비교해 볼 수 있다. 이외에도 최빈값, 중앙값 등을 이용해 비교해 볼 수 있다. 또한 범위, 분산, 표준편차 등도 문장의 습관을 구별하는 좋은 도구가 될 수 있다.

3) 상관계수 (correlation coefficient)

$$\rho_{XY} = \frac{\sigma_{XY}}{\sqrt{\sigma_X^2 \sigma_Y^2}}, \quad -1 \leq \rho \leq 1$$

$$\sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)]$$

저자의 특성 중에는 연령에 따라 시기에 따라 변화하는 특징들이 있을 수 있다. 이런 경우 상관계수를 이용하여 그 변화의 특징을 확인할 수 있다.

4) 다변량 분석 방법(판별 분석, 주성분분석, 군집분석)

판별분석은 저자가 불분명한 문헌에 대해 저자의 가능성이 있는 인물을 판별하는 기법이다.

$$(\bar{X}_1 - \bar{X}_2)' S_p^{-1} x_0 - \frac{1}{2} (\bar{X}_1 - \bar{X}_2)' S_p^{-1} (\bar{X}_1 + \bar{X}_2) \geq \ln \left[\frac{\alpha(1/2)}{\alpha(2/1)} \cdot \frac{p_2}{p_1} \right]$$

이때 x_0 를 G_1 으로 분류하고
아니면 x_0 를 G_2 로 분류한다.

여기서, 공동된 표본분산은

$$S_p = \frac{(n_1 - 1)S_1 + (n_2 - 1)S_2}{n_1 + n_2 - 2}$$

이다.

주성분분석은 다음의 조건을 만족시키는 선형결합 함수 Y_1, Y_2, \dots, Y_p 를 구하고자 하는 것이다.

① $Var(Y_1) \geq Var(Y_2) \geq \dots \geq Var(Y_p)$

② $\sum_{i=1}^k Var(Y_i) = \sum_{i=1}^k Var(X_i)$

③ $Cov(Y_i, Y_j) = 0, \quad i \neq j = 1, 2, \dots, p$

이렇게 구한 선형결합함수 중 $k(k < p)$ 개를 변수집단 X 를 축약시킨 주성분으로 결정하게 된다. 이는 많은 변수를 가진 정보를 주성분이라 부르는 몇 개의 합성성분으로 줄여서 분류하는 방법이다.

군집 분석은 data matrix에서 Euclid's distance Mahalanobis's distance의 유사성 척도로 집단을 분류하는 기법이다.

집필자가 불분명한 문헌에 대해 집필자의 가능성이 있는 인물이 둘 이상 있을 때, 이들에 대한 여러 관찰된 변수 값이 주어지면 판별분석, 요인분석 클러스터링 등의 방법 등을 통해 작가를 판별하고, 분류해 줄 수 있다.

3. 관련연구

1851년 런던대학 수학교수이고 논리대수의 창안자인 드 모르간(Augustes De Morgan)은 캠브리지의 목사 헤랄드에게 편지를 보내서 신약성경 중에서

바울서간 집필자에 대한 추정문제와 관련하여 각 서한에 대한 단어당 평균 문자수를 이용하여[1] 작자의 진위를 찾아낼 수가 있음을 밝혔다. 이 서한은 1882년 사후 그 부인에 의해 출판되었다.[2] 이를 본 오하이오대학의 물리학자 멘덴홀(T.C. Mendenhall)은 집필자를 추정하는데 모르간이 제안한 단어당 평균 문자수보다는 단어 길이의 범위(range)와 분포를 이용하는 것이 더욱 효과적이라는 제안을 하였다. 이 방법을 이용하여 멘덴홀은 1887년 SCIENCE 지에 동일한 집필자가 작성한 문장에 나타나는 단어길이의 도수분포를 WATT SPECTRUM이라 이름하고 이를 이용하여 집필자를 추정할 수 있는가를 검토하였다.

1901년 멘델홀은 Shakespeare작품의 진위논쟁에 뛰어 들어 그때까지 Shakespeare의 생애에 대해 불분명한 것이 많아 그의 작품이 동시대의 유명한 철학자이며 정치가인 Francis Bacon을 위시해서 여러 동시대 인물들이 쓴 것이라는 설이 퍼져 있었다. 멘델홀은 Shakespeare와 Bacon의 작품에서 각각 40만 단어와 20 만 단어를 표집하여 단어의 길이의 분포를 측정하였다. 그 결과 Bacon은 1단어 당 3문자에 최빈값(mode)을, Shakespeare는 4문자라는 것을 보여 주어 Bacon이 Shakespeare작품의 저자라는 설을 부정하였다.

그 이외에도 수학자 코시(A. L. Cauchy)는 평균, 중앙값, 사분위범위를, “그리스도에 대해서”(The Imitation Christ:De Imitatione Christi)의 저자를 찾는데 이용하였으며, 위에서 살펴 본 문장의 길이에 관한 척도에 대한 연구에서 나아가서 1944년 울(G. U. Yule)은 저자의 어휘량에 주목하여 그는 K 특성치를 제안하였다.[3]

작품 중에 단어의 나타나는 단어 F_i 가 X_i 번 나타난다고 할 때

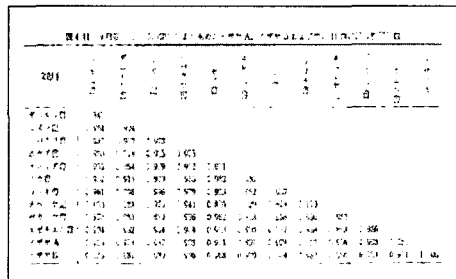
$$S1 = \sum XiFi, S2 = \sum X^2Fi,$$

$$K = 10^4(S2 - S1/S1^2)$$

1986년 藤泥偉作은 울의 데이터를 정준 상관을 이용하여 울의 연구결과를 지지하였다.[4]

이외에도 구약의 이사야서에 대한 진위 논의 역시 많은 시사점을 얻을 수 있다. 이사야서는 여러 사실적 의문으로 1970년대까지만 해도 39장을 전후로 전 반부와 후반부가 다른 작자에 의해 쓰여진 것으로

이해되었다. 그러나 1973년 Adams와 Rencher는 헤브라이어의 접두어의 사용비율을 이사야서에서 조사한 결과를 분석하여 이사야서의 전후반부 간에는 접두어의 사용율에서 높은 상관관계를 보이고 있어 이사야서는 1인에 의한 저작일 가능성이 높다는 연구를 발표하였다.[5] 그러나 1994년 村上征勝은 아래 <그림1>과 같이 이들의 방법론을 다른 구약성경에 대한 무선표집을 통해 상관분석을 한 결과 접두어를 이용한 상관분석의 결과는 표집한 모든 구약의 부분들에서 높은 상관을 보여 이를 통해 구약전체가 한 사람의 집필자에 의한 저작으로 해석할 여지가 있음을 보여 Adams와 Rencher의 연구에 회의를 표시하고 이사야서 집필자에 대한 추정에 새로운 문제를 제기하였다.[6]



<그림 1 > 구약성경 여러저자들간의 correlation

중국의 경우에도 중국의 대표적인 고대소설인 홍루몽(紅樓夢)이라는 소설의 경우 작품의 일부가 다른 사람에 의해 쓰여졌을 것이라는 논쟁이 계속 되어 오고 있었다. 1987년 夏旦大學 李賢平은 <그림 2>와 같이 120장의 각 장에 출현하는 47개의 虛字를 이용하여 주성분분석, 정준상관분석, 군집분석(clustering analysis)을 실시하였다. 그 결과 다음과 같이 결론을 맺었다. “홍루몽의 80회까지는 조설근이 내려오는 石頭記”를 수정하여 쓴 것이고 그 이후는 120 회까지는 조설근의 자신의 초기소설 “風月寶玉”을 새롭게 추가하여 쓴 것을, 그 동료들이 수정 가필한 것이라는 결론이었다.[7]

4. 저자 판별 지원 시스템 제안

저자 판별을 위한 지원 시스템은 다음과 요건을 갖추어야 한다.

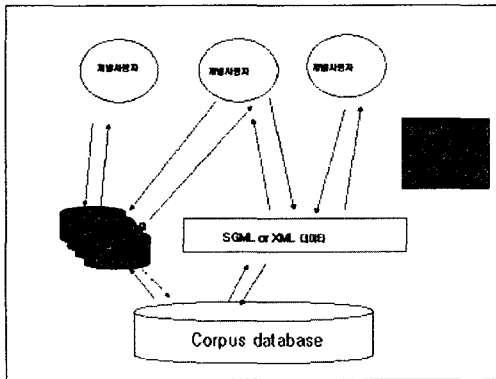
- 1) 고문헌 자체를 자유로 입력하고 검색할 수 있는 편집기가 구현되어야한다.
- 2) 고문헌은 주로 한자체계로 되어 있고 현재 우리

가 쓰고 있는 상용 혹은 확장한자에는 없는 자들이 많아 이를 처리하기 위한 처리기가 구현되어야 한다.

3.) 한 문장 혹은 단어의 길이 등 기본통계와 필요한 자료를 출력할 수 있어야 한다.

4) 단일 문헌 혹은 사료에 대한 입력자료 및 중간처리 단계의 자료들이 데이터베이스화되어 향후 시대별, 작자별, 기록의 종류별 등 필요한 검색 조건별로 검색이 가능하여야 한다.

위와 같은 특성을 고려한 시스템의 기본 구성도는 다음과 같다.



<그림 3> 판별지원시스템 기본구성도

5. 결론

본 연구는 최근 역사학계의 중심논쟁으로 되고 있는 화랑세기 등의 저작에 대한 진위를 밝히기 위한 과학적인 방법론을 제안하고자 하는 것이다. 이를 통하여 우리 역사의 지평과 전망을 다시 살펴보는 계기를 마련할 수 있을 것이다. 저작의 진위를 확인하는 연구에는 문학, 언어학, 통계학, 문헌정보학, 컴퓨터공학 등 다양한 학문분야에 의한 다학문적 접근을 필요로 한다.

이러한 연구의 성공을 위해서는 관련분야 연구자가 서로 지식이 부족한 부분에 대한 보완을 하면서 조직적 다면적 접근을 통한 광범위한 시야를 가지고 연구를 진행하여야 한다. 이를 위해서는 우선 각분야 연구자가 서로의 의견을 나눌 수 있는 기회를 만들어 다른 분야의 연구에 관한 지식을 이해하는 동시에 공통의 문제 의식을 가지고 공동 연구를 진행할 수 있는 체제를 확립하는 노력이 기울여야 할 것

이다.

저자의 진위를 판정하는 문제는 단순히 문헌에만 한정되는 것은 아니다. 국내의 실정은 알 수 없으나 동경미술관 감정위원회에서 감정한 591 점의 감정 미술 작품 가운데에서 38%에 해당하는 228 점이 위작인 것으로 발표되었다. 이와 같이 글, 그림, 도자기 등의 작품에도 문헌이상의 위작이 있어 진위논쟁은 앞으로 계속될 것이다. [8] 이러한 학문간 교류에 의한 연구는 이외에도 다양한 여러 영역에 그 연구의 위력을 더해 갈 수 있을 것이다.

참고문헌

- [1] R. D. Lord, Studies in the history of probability and statistics VIII. De Morgan and the statistical study of literature style. *Biometrika*, 45, 1958
- [2] S. E De Morgan, Memoir of Augustus de Morgan by his wife Sophie Elizabeth De Morgan with selection from his letters. Longman, Green and Co., 1882.
- [3] G. U. Yule, The statistical study of Literacy Vocabulary. Cambridge University press, 1944
- [4] 藤退偉作, 未知の著作権の推定について. BASIC 數學 2 月号, 1986
- [5] L. L. Adams and A. C. Rencher, The popular critical view of the Isaiah problem in light of statistical style analysis. *Computer Studies in the Humanities and Verbal Behavior*. 7(3-4), 1973.
- [6] 村上征勝, 眞贋の科學, 朝倉書店, 1994.
- [7] 李賢平, 紅樓夢成書新說, 夏旦大學, 1987年 第5期
- [8] 種村秀弘, 贋作者列傳, 青土社, 1992.