

仮名の Ngram を用いた源氏物語写本の系統分類の試み

齊藤 鉄也 (淑徳大学 経営学部)

本論文では、源氏物語の短い本文を持つ巻の、写本本文の仮名の Ngram を用いた系統分類の調査結果を述べる。表記が異なる本文データの Ngram を用いた巻ごとの分類結果からは、3gram から 5gram の調査データで、青表紙本系統と河内本系統の分類ができること、系統分類が明瞭な巻と不明瞭な巻があることを明らかにした。加えて、文学や文献学の系統分類の研究成果と本調査結果を比較し、その分類結果が一致している点があることから、本調査方法に一定の有効性がある可能性を示した。

An Attempt at Stemmatic Classification of Manuscripts of the *Tale of Genji* Using *Kana* Ngrams

Tetsuya Saito (College of Business Administration, Shukutoku University)

In this paper, I report the results of an experiment at stemmatic classification, using *kana* Ngrams to analyze textual and orthographic variations between multiple manuscripts of individual chapters within the *Tale of Genji*. Classification results were obtained for several selected short-length chapters using Ngrams of 3~5 grams. The results for manuscripts of these chapters showed that (1) the method was able to identify which stemma (Aobyoshi-bon or Kawachi-bon) a given chapter belonged to, and that (2) stemmatic classification was easier for some chapters of the work than for others. Furthermore, a comparison of these results with stemmatic classification results obtained by literary and philological methods of research showed that the conclusions reached were identical, suggesting that the method used in this experiment may have a certain degree of validity.

1. はじめに

本論文では、仮名の Ngram を用いた源氏物語写本の系統分類の調査結果を述べる。これまで文学や文献学分野では、本文の内容や、語や文節ごとの意味の相違に着目して本文の分類を行ってきた。これに対し、本調査では、本文の表記の相違に基づいて分類を試みたことが異なっている。

調査対象とした写本は、出版されている源氏物語のうち、短い本文を持つ巻を中心に「空蟬」「花宴」「花散里」「閑屋」「篝火」「鈴虫」の全文と「柏木」の一部の、計 145 写本である。その分類方法は、通行仮名に翻字した写本の本行本文データの Ngram に対し、クラスター分析、主成分分析、系統学的手法、ヒートマップを用いた。

その結果、3gram から 5gram の調査データで、青表紙本と河内本の写本の系統分類ができること、巻ごとに分類結果に差があり、分類結果、明瞭な巻と不明瞭な巻があることを明らかにした。加えて、文学や文献学の研究成果と本調査結果を比較し、その分類結果が一致している点が多いことから、本調査方法に一定の有効性を示している可能性を示した。

2. 本研究の目的と関連研究

本研究の目的は、表記の相違を持つ本文データを用いて、写本の系統分類を計量的な処理のみで行い、その結果を可視化することと、その結果を

文学や文献学の研究成果と比較し、本方法の有効性と限界を明らかにすることである。その最初の段階として、本調査では、源氏物語の短編を中心に選択し、本行本文の仮名データの Ngram を用いて系統分類を試みた。

これまで文学や文献学での源氏物語本文の系統分類は、語や文節といった単位に対し、主として意味の相違に基づいて行われてきた。そのため、漢字と仮名、仮名遣いや音便の相違といった意味には影響を与えない表記の相違は本文異同として扱われてこなかった[1]。これに対し、本調査で用いた写本の本文のデータでは、これらの表記の異同も相違として扱っている。また、本文の修正や挿入といった傍記は対象とせず、本文の傷も相違として扱っている。これまでも、特定の漢字と仮名の表記の相違に着目した本文の類似性の指摘は存在する[2]。本調査は、特定の漢字と仮名の相違だけではなく、仮名遣いや音便、本文の傷も含めた本文の表記の相違を対象に系統分類を行う初めての試みである。

源氏物語写本の本文は、これまで青表紙本系統、河内本系統という二系統と、それ以外の別本群に分類されてきた[1]。この分類枠組みを検討し[3]、河内本群と別本群の二分類とする議論[4]や、既存の分類の有効性を指摘する議論[5]もある。本文の異同に基づき、その系統は分類されているが、本文異同の多様性が原因となり、研究者により分類

結果が異なることもある。この問題に対して、源氏物語の系統分類へ計量的な手法の導入を行った研究がある[5]。この研究では、本文異同を語単位で数値化したデータを用い、系統分類を試みている。この本文異同の数値化には文学や文献学の知見が必要であり、その知見のない情報学の研究者が系統分類を行う際の参入障壁となっている。

そこで、本論文では、この参入障壁を越えて系統分類を行うために、写本本文の仮名を対象に Ngram を用いた系統分類を試みた。写本本文の変体仮名を通行仮名に翻字したデータから Ngram を生成して、教師なし分類手法を用いることで、系統分類の簡易化と結果の可視化を行う。この本文データを用いて系統分類することが可能であれば、研究の効率化とその進展が期待できる。また、研究成果が蓄積され、複数の手法に基づく調査結果の比較が可能であれば、その結果の信頼性の向上も期待できる。

加えて、近年、写本の出版やインターネット上で公開が進みつつあることから、本研究は、公開されている写本画像から文字認識が実用化された際の応用事例として位置付けることができる。この場合、本手法の結果と文献学の結果との比較と検証を行わなければ、その分類結果の妥当性が不明であり、結果として本手法の有効性も不明である。この点で、本研究は、その検証を行い、この応用事例の基礎的な知見を得る研究と言える。写本の文字認識から分類結果の可視化まで自動化が可能であれば、調査や研究が効率化され、文学や文献学への貢献が期待できる。

3. 調査対象と生成した本文データ

調査対象には、源氏物語のうち、短い本文を持つ巻を選択した。その理由は、源氏物語の写本が出版またはインターネット上で公開され入手し易いこと、調査対象である写本が大量に存在するため、調査の効率化が必要であること、本研究を今後の長い本文を持つ巻に調査対象を拡大する際の初期段階として位置付けていることがある。

加えて、源氏物語の青表紙本系統と河内本系統間の本文異同は、54 帖全体で一定ではなく、巻ごとに異なり、前半の巻に多いことが指摘されている[6]。そこで、調査対象とした巻は、本文異同が大きい巻から少ない巻までを対象とした。

結果として、調査対象とした巻と写本の数は「空蟬」が 19 本、「花宴」が 20 本、「花散里」が 20 本、「閑屋」が 21 本、「篝火」が 20 本、「柏木」24 本、「鈴虫」が 21 本である。「柏木」以外の巻に関しては全文を調査している。「柏木」は長文であるため、尊経閣文庫本「柏木」の藤原定家筆部の本文の文字数 3247 字を基準として本文範

囲を選択し、各写本の該当する範囲を調査対象とした。

調査対象としたデータは、本行本文の変体仮名を通行仮名に変換した本文である。本文の漢字は対象としていない。挿入されている本文や、誤字や脱字を修正した文字といった傍記は対象としていない。繰り返し符号である踊り字も対象としていない。この理由は、通行仮名のみを用いた本文を分類した場合、通行仮名と常用漢字のみを用いた本文を分類した場合、通行仮名と常用漢字と、繰り返し符号を対応する文字列に変換した本文を分類した場合の分類結果を比較したが、それらの分類結果は、どれも似た結果となったからである。そのため、最も単純な通行仮名のみを用いている。

4. 調査結果の概要と考察

最初に分類の概要を把握するために、調査対象とした各巻に対して、2gram から 7gram の文字列を生成し、階層的クラスタ分析を行った。階層的クラスタ分析では、写本間の距離には平方ユークリッド距離の 1/2 を用い、クラスタ間の距離にはワード法を用いた。調査結果を判断する基準として、各写本の解題の(推定)書写年代と(推定)系統を採用した。これに基づき本調査の分類結果と比較している。階層的クラスタ分析の結果、青表紙本系統と河内本系統の異同が大きい巻は、「空蟬」「花宴」「花散里」「閑屋」であった。これらの巻では、3gram から 5gram において、青表紙本系統と河内本系統で異なるクラスタを構成することが明らかになった。「柏木」は、二つの系統に分類できたが、その距離は近かった。「篝火」「鈴虫」は、これらの系統分類の結果は明確ではなく、同一のクラスタ内に二つの系統の写本が混在する結果となった。

次に主成分分析を用いて階層的クラスタ分析の裏付けを試みた。5gram のデータに対して、主成分分析を行い、さらに k 平均法を用いた非階層的クラスタ分析を行った。5gram を選択した理由は、本文異同を比較するためには、本文中の文字列の位置が特定できる、より長い文字列が望ましいからである。以下の各図は、その分類結果である。どの結果も第一主成分と第二主成分を軸として用いているが、それらの軸の意味は明らかではない。図では各写本を表す文字列に、写本名称と書写年代、系統を用いている。書写年代は解題に基づく時代と時代区分を用いている。時代区分の記載がない写本も存在する。系統には、青表紙本系統は「青」、河内本系統は「河」、別本群は「別」として表している。写本名称とその出典は注にまとめた。

4.1 「空蟬」の分類結果

図1は「空蟬」の分類結果である。結果は、左上の楕円の位置にこれまでの研究により河内本系統とされる写本が集まる。この楕円には、別本である陽明文庫本と玉里文庫本も所属している。ほぼ中央左に存在する飯島本は河内本に分類されているが、この楕円の中には存在せず、穂久邇文庫本と同じ群に分類されている。

この他の写本は、解題によれば全て青表紙本に分類されている。これらの写本は、さらにいくつかの群に分類されている。このうち一つの群は中央下の楕円で表されているが、鎌倉時代と室町時代の写本が混在している。

中院文庫本は日大三条西家本を書写し、両写本は本文中の漢字の位置が同じであることが知られている。この図においても、これらの写本は右上の近い位置に存在している。

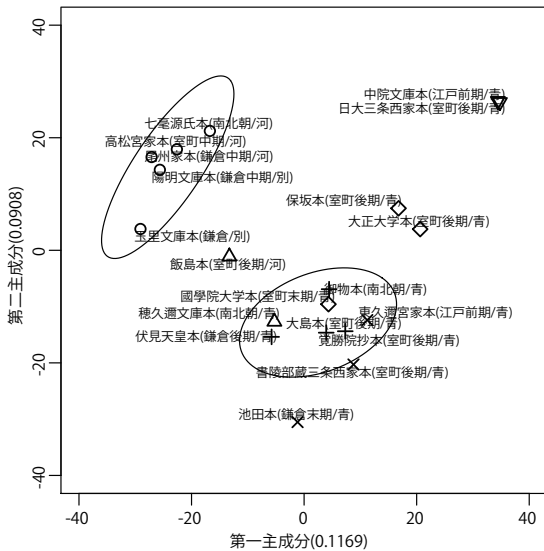


図1 「空蟬」5gramの主成分分析結果

4.2 「花宴」の分類結果

図2は「花宴」の分類結果である。図では、左側の楕円の位置に河内本系統とされる写本が集まる。中央に位置する別本である御物本は青表紙本系統の陽明文庫本と同じ群に分類されている。

中央右にある楕円には、青表紙本系統の写本である、鎌倉時代と室町時代の写本が混在している。このうち、池田本と明融臨模本は、本文のある一行が同じ位置で改行されていることから、その関係が論じられている写本である。この図においても、これらの写本は近い位置に存在することから、本文、よく似た写本であることを表している。また、共に室町時代書写の大正大学本と保坂本もよく似た写本であることを表している。

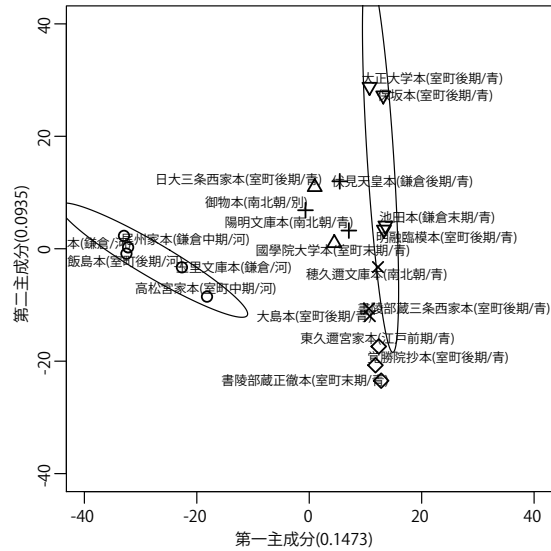


図2 「花宴」5gramの主成分分析結果

4.3 「花散里」の分類結果

図3は「花散里」の分類結果である。左側の楕円に河内本系統とされる写本が集まる。「花宴」と同様に、中央に位置する別本は御物本と陽明文庫本が同じ群にある。この他に、青表紙本系統の写本が集まる、二つの楕円が存在し、中央の楕円には鎌倉時代から室町時代に書写された写本が混在し、右の楕円には室町時代の写本、存在する。

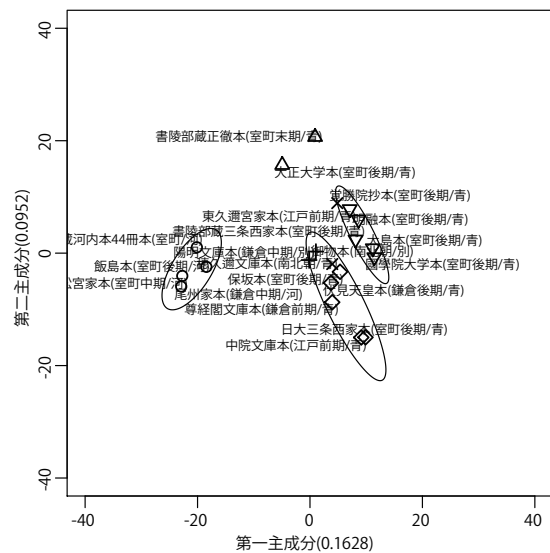


図3 「花散里」5gramの主成分分析結果

4.4 「関屋」の分類結果

図4は「関屋」の分類結果である。左側の楕円には河内本系統とされる写本が集まる。別本である陽明文庫本と玉里文庫本は、異なる群に存在する。それ以外の青表紙本系統の写本は、中央の楕円に鎌倉時代の写本が多く集まり、その右に室町

時代の写本が集まっている。

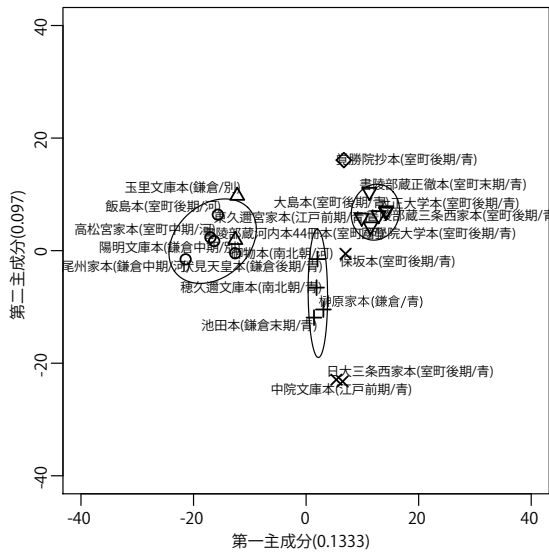


図4「関屋」5gramの主成分分析結果

4.5 「篝火」の分類結果

図5は「篝火」の分類結果である。その分類結果はこれまでと異なり、中央に楕円が集まる。これまでの巻と比較して、河内本系統の写本が存在する中央上の楕円に青表紙本系統である池田本が存在する。その他の写本の中では、左側の楕円には主として鎌倉時代の写本が集まり、主として右側の楕円には室町時代の写本が集まる。

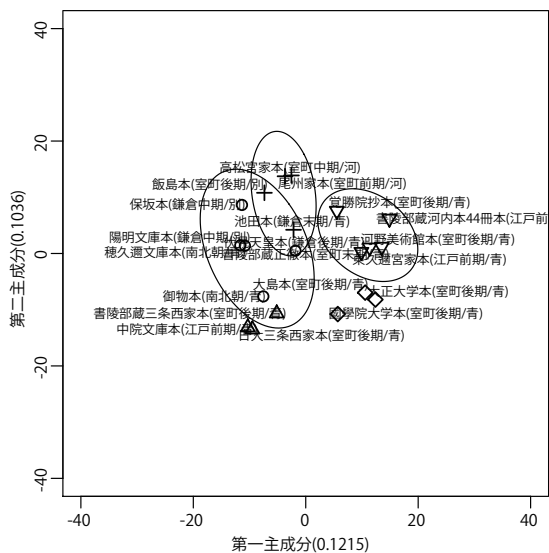


図5「篝火」5gramの主成分分析結果

4.6 「柏木」の分類結果

図6は「柏木」の分類結果である。「柏木」は長編であるため、その本文のうちの3000字程度を対象とした分類結果である。左側の楕円に主として河内本系統の写本が所属している。この群に

は別本である御物本が属している。下には主として鎌倉時代の別本が位置する。一部が重なる二つの楕円には主として青表紙本系統の写本が存在するが、これらの共通の性質は明らかではない。

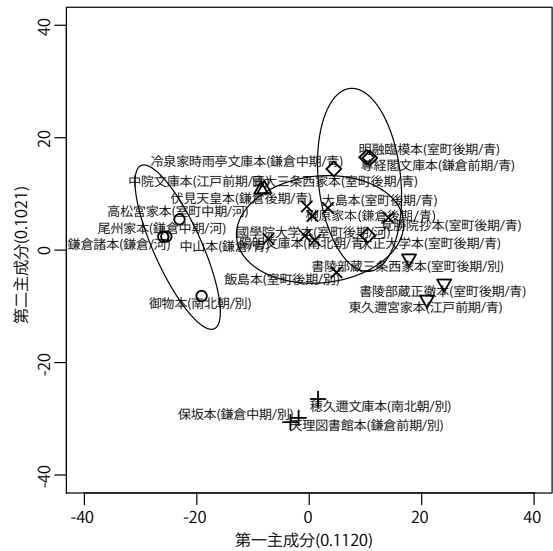


図6「柏木」5gramの主成分分析結果

4.7 「鈴虫」の分類結果

図7の「鈴虫」の分類結果からは、主に3つの群に分類できることを示している。左上の楕円は青表紙本系統と河内本系統の写本が混在している。下の楕円には別本が属している。右の楕円には室町時代に書写された青表紙本系統の写本が属している。

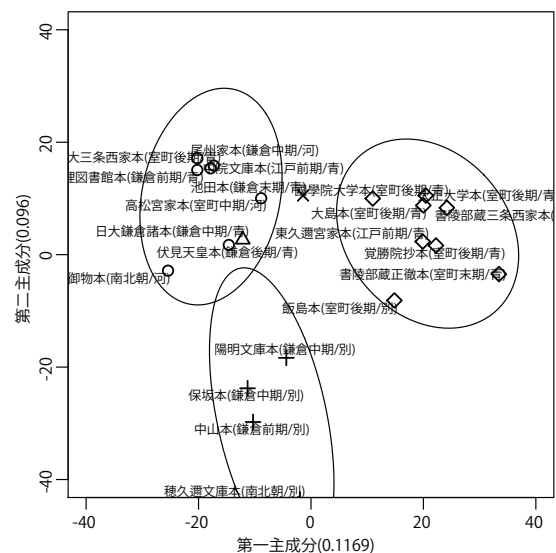


図7「鈴虫」5gramの主成分分析結果

5. 文献学の系統分類結果との比較

上記で取り上げた巻のうち、既存研究において

系統分類結果が発表されている「空蟬」「花散里」「柏木」「鈴虫」について、本調査結果の比較し、その妥当性を検討する。ここでは写本間の関係を明らかにするために、系統学的手法を用いた。

5.1 「空蟬」の系統分類結果の比較

「空蟬」は、新美[5]において、系統学的手法を用いて、系統分類が行われている。そこでは、写本から語を単位として本文異同を数値化し、距離の計算方法にハミング距離を用いて、複数の系統学的手法を用いて分類している。

この結果と比較するために、本調査の 5gram の本文データに対し、距離の計算方法としてコサイン距離を用いた Neighbor-net[7]ネットワークを図 8 に示す。

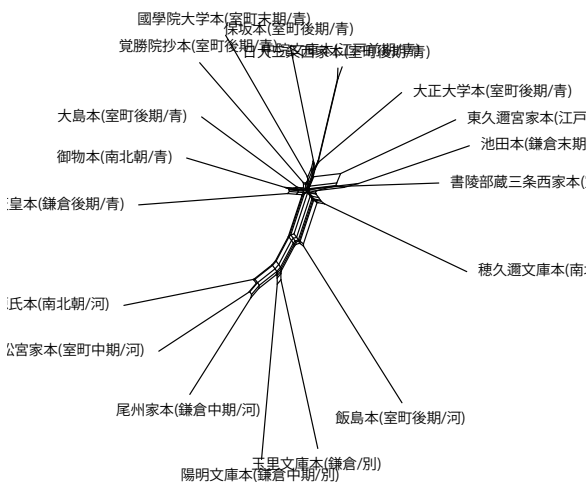


図 8 「空蟬」の Neighbor-net ネットワーク

図 8 の左下に河内本系統とそれらと本文が似ていると考えられる別本が集まっている。それ以外は全て青表紙本系統である。

飯島本を除く河内本系統の写本と二つの別本の関係は、新美[5]の結果と同様である。その調査では、飯島本は調査対象でないため、その分類結果は不明である。青表紙本系統のうち、中央左に位置する大島本、御物本、伏見天皇本の群と、上に位置する、中院文庫本と重複している日大三条西家本と保坂本の群が、それぞれ近い位置に存在すること、右に位置する穂久邇文庫本と書陵部蔵三条西家本が離れて位置することも同様の結果となっている。

新美[5]の結果と本調査とは、調査対象とした写本も異なり、その結果の単純な比較はできないが、同一の調査対象である一部の写本については、同様の結果となっている。この点では、文献学の知見を用いた語の本文異同の分類結果と、表記の

相違を持つ 5gram の仮名文字列の分類結果の一部が合致している。このことは、文献学の知見を用いずに Ngram により系統分類ができる可能性を示している、と言える。

5.2 「花散里」の系統分類結果の比較

「花散里」は伊藤[8]と大内[9]において、写本分類が検討されている。伊藤[8]においては、源氏物語の写本を河内本群と別本群として分類し、さらに、その中をいくつかの小群に分類した分類案を提案している。

「空蟬」と同様に「花散里」においても、本調査の 5gram の本文データにコサイン距離を用いて、Neighbor-net ネットワークを作成した。この結果を図 9 に示す。

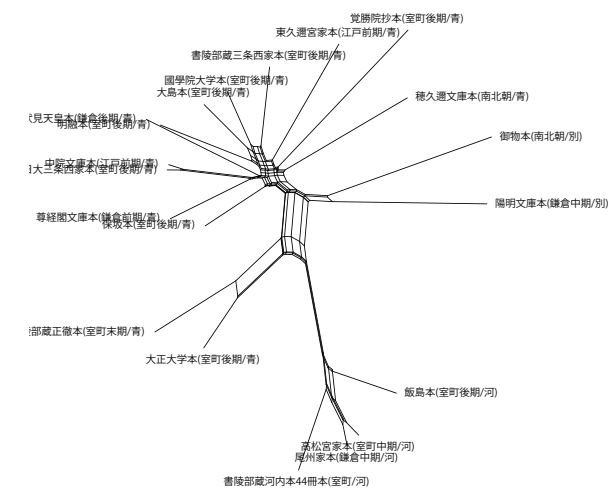


図 9 「花散里」の Neighbor-net ネットワーク

図 9 では、下に河内本系統の写本が集まっている。中央右に位置する別本である御物本と陽明文庫本は互いに近い関係あることを示している。また、残る青表紙本のうち、左下に位置する書陵部蔵正徹本と大正大学本は互いに近い関係にあることを示している。それ以外の青表紙本は左上に集まっている。

伊藤[8]の結果は、25 写本を物語の内容に沿って分類し、河内本群と別本群に分けている。本調査とは、そのうちの 13 写本が一致している。それらの写本のうち、下に位置する尾州家本と高松宮家本、右に位置する御物本と陽明文庫本が小群を構成することに関しては一致している。図 9 では、日大三条西家本と中院文庫本、尊経閣文庫本、伏見天皇本が互いに近いことを示しているが、この結果は、伊藤[8]の結果とは異なっている。この相違は分類の視点の相違によることが考えられ、その原因分析は今後の課題である。

5.3 「柏木」の系統分類結果の比較

「柏木」は阿部[3]と加藤[10]において青表紙本系統の写本の中の分類を論じている。阿部[3]では、尊経閣文庫本、陽明文庫本、日大三条西家本を例に挙げて、青表紙本系統の写本は、さらに小群に分類できることを指摘している。加藤[10]では、別本とされる書陵部蔵三条西家本と、大正大学本と書陵部蔵正徹本といった室町時代の青表紙本系統の写本との類似性の指摘がある。

本調査では、これまでと同様のデータと距離の計算方法を用いて系統樹ネットワークを作成し、これらの指摘の確認を行った。

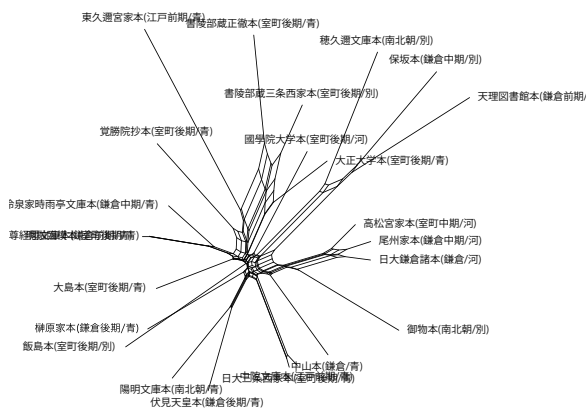


図 10 「柏木」の Neighbor-net ネットワーク

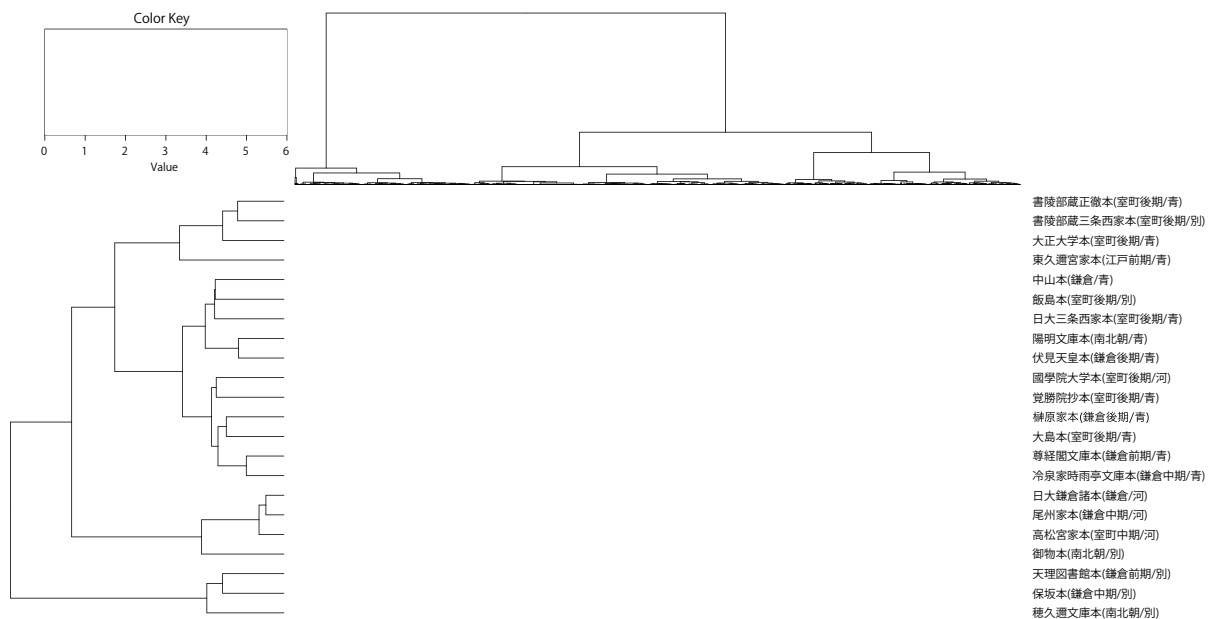


図 11 「柏木」の階層的クラスター分析を伴うヒートマップ

図 10 では、中央右に河内本系統の写本が集まっている。別本である御物本は、これらの写本に近い位置に存在する。加藤[10]では、御物本は「河内本のくずれたもの」とされ、河内本との類似性が指摘されている。この他の別本である徳久邇文庫本と保坂本、天理図書館本は右上に群を構成し、そのうち保坂本と天理図書館本は互いにより近い関係にあることを示している。

残る青表紙本のうち、下に位置する文字列が重複する写本は日大三条西家本と中院文庫本である。左に位置する文字列が重複する写本は尊経閣文庫本と明融臨模本である。これらは互いに本文が類似していることを示している。図からは阿部[3]が指摘する尊経閣文庫本と明融臨模本、大島本は近いこと、陽明文庫本と伏見天皇本、榊原家本は近いこと、日大三条西家本は尊経閣文庫本よりも陽明文庫本に近いことが指摘できる。また、中央上に位置する加藤[10]が指摘する室町時代の三つの写本は互いに近いことも指摘できる。この様に、図 10 からは、これまでの文学や文献学分野で指摘されてきた写本間の関係が、本調査で用いた手法においても可視化できている。これも、本手法を用いて系統分類ができる可能性を示していると言える。

さらに、本文の異同箇所の出現傾向と系統分類の関係を調査するために、図 11 に、階層的クラスター分析を伴うヒートマップ[11]を示した。階層的クラスター分析の距離の計算方法は平方ユークリッド距離の 1/2、クラスター構成方法はウォード法を用いた。

6. まとめ

本論文では、複数の手法を用いて、源氏物語の短編写本を中心に、仮名の Ngram を用いた系統分類の可能性の調査報告をした。調査対象とした本文は写本の本行本文の仮名である。そのため、本文としては誤脱があり、漢字と仮名や音便といった意味に影響を与えない表記の相違を持つ本文データを調査対象としている。この本文データに対し、計量的な分類手法という統一的な手法を用いて系統分類を行い、その結果を既存の文学や文献学の結果と比較し検討した。調査結果からは、これまでの文献学の知見と同等の結果も存在することから、表記に基づいて系統分類するという本手法の有効性を示している可能性。ある。また、本調査結果は異なる本文データと計量的な手法を用いた結果であることから、一致する点に関しては、既存の系統分類結果の蓋然性を高めているということができる。

今後の課題としては、既存の研究成果との相違の原因を分析するために、調査対象の写本数を増やすこと、本文内容を考慮したデータや表記が異なるいくつかの本文データを作成し、異なる本文データを用いた場合の系統分類結果を比較することがある。

注

本調査で対象とした古典籍は、公刊された複製本または影印本、デジタルデータに依拠している。

- 1)陽明文庫本, 陽明叢書国書篇源氏物語, 思文閣出版.
- 2)保坂本, 保坂本源氏物語, おうふう.
- 3)尾州家本, 尾州家河内本源氏物語, 八木書店.
- 4)伏見天皇本, 源氏物語伏見天皇本, 古典文庫.
- 5)池田本, 新天理図書館善本叢書池田本, 八木書店.
- 6)御物本, 御物各筆源氏, 貴重本刊行会.
- 7)穂久邇文庫本, 日本古典文学影印叢刊, 貴重本刊行会.
- 8)高松宮家本, 高松宮御蔵河内本源氏物語, 臨川書店.
- 9)書陵部蔵三条西家本, 宮内庁書陵部蔵青表紙源氏物語, 新典社.
- 10)大正大学本, 大正大学附属図書館源氏物語写本.
- 11)日大三条西家本, 日大鎌倉諸本, 日本大学蔵源氏物語, 八木書店.
- 12)飯島本, 書芸文化院春敬記念書道文庫蔵飯島本源氏物語, 笠間書院.
- 13)大島本, 大島本源氏物語, 角川書店.
- 14)覚勝院抄, 源氏物語聞書覚勝院抄, 汲古書院.
- 15)國學院大學本, 國學院大學伝為家本, 國學院大学図書館デジタルライブラリー.
- 16)書陵部蔵正徹本, 宮内庁書陵部蔵書寮文庫 554-14, 「源氏物語(正徹本・初音首欠)」.
- 17)中院文庫本, 京都大学附属図書館所蔵貴重書源氏物語.
- 18)東久邇宮家旧蔵本, 国立国会図書館デジタルコレクション東久邇宮家旧蔵本源氏物語.
- 19)玉里文庫古筆源氏, 鹿児島大学附属図書館, 天

213-1361.

- 20)桃園文庫蔵明融臨模本, 東海大学出版会.
- 21)尊経閣文庫本, 雄松堂書店.
- 22)書陵部蔵河内本 44 冊本, 宮内庁書陵部図書寮文庫 150-744, 「源氏物語(有欠・河内本)」.
- 23)榊原家本, 国文学研究資料館.
- 25)天理図書館, 天理図書館善本叢書和書之部源氏物語諸本集, 八木書店.
- 26)冷泉家時雨亭文庫本, 冷泉家時雨亭叢書, 朝日新聞社.
- 27)中山本, 国立歴史民俗博物館蔵貴重典籍叢書文学篇物語, 臨川書店.

参考文献

- 1) 池田亀鑑(編著): 源氏物語大成研究編, 中央公論社 (1985).
- 2) 中古文学会関西西部会(編): 大島本源氏物語の再検討, 加藤洋介: 大島本源氏物語の本文成立事情 - 若菜下巻の場合 -, pp.167-208, 和泉書院 (2009).
- 3) 阿部秋生: 源氏物語の本文, 岩波書店 (1986).
- 4) 伊藤鉄也: 源氏物語本文の研究, おうふう (2002).
- 5) 新美哲彦: 源氏物語の受容と生成, 武蔵野書院 (2008).
- 6) 増田繁夫, 鈴木日出男, 伊井春樹(編): 源氏物語研究集成第十三巻源氏物語の本文, 加藤洋介: 河内本の成立とその本文 - 源親行の源氏物語本文校訂 -, pp.113-143, 笠間書院 (2000).
- 7) Bryant, D. and Moulton, V.: NeighborNet: An agglomerative method for the construction of planar phylogenetic networks, *Molecular Biology and Evolution*, Vol.21, No.2, pp.255-265 (2012).
- 8) 鈴木一雄(監), 秋山虔, 室伏信助(編): 源氏物語の鑑賞と基礎知識 No.29 花散里, 伊藤鉄也: 『源氏物語』の諸本別本について, pp.213-222, 至文堂 (2003).
- 9) 大内英範: 源氏物語本文資料整理の方法, 情報知識学会誌, Vol.13, No.2, pp.32-40 (2003).
- 10) 加藤洋介: 定家本源氏物語の復元とその限界, *国語と国文学*, Vol.82, No.5, pp.126-141 (2005).
- 11) Kobayashi, Y.: HeatMap with Hierarchical Clustering: Multivariate Visualization Method for Corpus-based Language Studies, *NINJAL Research Papers*, Vol.11, No.1, pp.25-36 (2016).
- 12) 池田利夫: 源氏物語の文献学的研究序説, 笠間書院 (1988).
- 13) 伊藤鉄也: 『源氏物語』の異本を読む - 「鈴虫」の場合, 臨川書店 (2015).
- 14) 中尾央, 三中信宏(編著): 文化系統学への招待, 矢野環: 『老葉』に対する系統学的アプローチ - 宗祇による連歌の系譜, pp.35-63, 勁草書房 (2012).