

クラウドソーシングによる語義調査

人間文化研究機構 国立国語研究所 加藤 祥

1. はじめに

本発表では、多義語の意味調査例を紹介し、意味の調査における一般的な解釈としてのクラウドソーシング実験の活用可能性を考えたい。文脈における多義語の意味判定は、読み手によって、あるいは一人の読み手であっても揺れが生じるものである。そのため、用例を用いて多義語の意味判定のグラデーションを調査し、多義のネットワークと派生関係の解明を試みている。被験者実験では、調査対象とする多義語について、実験協力者に指標文と判定文を提示し、用例中の意味の類似度を判定してもらう。この結果、語義の関係がどのように用例に現れているのか整理することが可能となる。これらの調査手法と「一般的な」意味解釈を用いた調査事例を示す。

2. 語義調査を行う背景

国立国語研究所では、『現代日本語書き言葉均衡コーパス』（以降BCCWJ）に『分類語彙表¹』番号の付与を行っている。分類語彙表番号は表1のような構成であり、BCCWJの国語研短単位（概ね語に相当する）に対し、人手によって文脈上該当すると判断された5桁の情報が付与される。すなわち、分類語彙表番号付与済みBCCWJ（以降BCCWJ-WLSP）は、均衡性のある新聞・雑誌・書籍（BCCWJ中の約35万語）に含まれる自立語²のすべてに意味的な情報付与を行ったデータである（加藤ら2019aなど）。

表1 分類番号の構造（例：この（分類番号：3.1010））

類	部門	中項目	分類項目
相 (3)	関係 (.1)	真偽 (.10)	こそあど (.1010)

BCCWJ-WLSPを用い、多義語を語義別に頻度調査することが可能となった。そこで、代表義（山崎・柏野，2017）、基本義などの情報と組み合わせた頻度調査や派生関係の調査などが試みられている（加藤ら，2019c）。しかし、BCCWJ-WLSPの整備にあたっては、文脈による多義語の意味判定が行われており、読み手の違いによって、あるいは一人の読み手の中でも場合によって、意味判定の結果に揺れの生じることが考えられる。

よって、たとえば(1)のように、読み手（BCCWJ-WLSPの情報付与作業員）がどの語義と判断するか迷ったとする例などは、付与された意味情報を一語義の頻度と数えて良いのかという疑問が生じる。また、類似の用例においては、付与済みの判断にも、作業員によって、あるいは場合によって揺れの生じている可能性がある。

(1) 両手を胸の前でぎゅっと握りしめて、加護は遠い目をした。

（サンプルID：PB39_00372，下線は著者による。以下同様）

BCCWJ-WLSPの付与情報の精度を高めるためには、『分類語彙表』で複数の番号（語義）を有する語（多義語）について、揺れの生じるような文脈においてどの語義が最も適しているのかを客観的に判定したい。そのため、分類語彙表番号体系を把握している作業員数名の判断に限らず、広く一般的にどの語義と判定されるのかという観点を取り入れることにした。

3. 調査の設計

読み手が語の意味を判断する際に揺れや迷いの生じる用例と、語義的に類似度の高い用例を収集し、用例間の語義の関係性を調査することで、用例の整理を進めることができる。そのため、実験協力者に指標文と判定文を提示し、指標文と対照して各判定文の類似性を判断してもらうという被験者実験を考えた。

¹ 国立国語研究所（編）（2004）. 国立国語研究所資料集14『分類語彙表一増補改訂版一』. 大日本図書.

² 同対象範囲には、助動詞の用法も別途付与されている（加藤ら，2019b）。

3.1 調査対象

まず、BCCWJ-WLSPの相の類（形容詞・形容動詞・副詞などの修飾語句を主に含む）の語³を語義別に集計したデータ（加藤ら，2019c）を参照し、多義語として複数語義が付与され、かつ各語義の頻度が5以上取得できる語を選んだ。次に、調査対象とした語の各語義において典型的と考えられる用例と、同コロケーションや類似文脈ながら付与済みの語義に判定揺れが見られる用例（前掲の(1)参照）を収集した。このほか、先行研究などに語義の記述があるもののBCCWJ-WLSPから取得できなかった種類の用例や、語義の付与判定が困難と考えられる用例も、BCCWJから追加収集した。

3.2 調査方法

先に収集した用例群のすべてが指標文となり得、かつすべての用例と類似度判定を行う組み合わせを作成し、全用例の組み合わせについて実験協力者に意味的な類似度の判定を依頼する。

実験協力者は、提示された文（用例）中の対象語との意味的な類似度を判定し、0（まったく違う）から5（まったく同じ）のいずれかを選択する。画面の関係上、1用例（指標文）に対し、判定する用例（判定文）をランダムに5例表示し、選択肢としてラジオボタンを表示した（画面例は図1）。提示する用例は、基本的に調査対象語の前後20語程度を含む一文以上としたが、画面表示の関係上、一文全てが提示できない場合には、読点を基準として区切ることにより一部を省略したことがある（よって、表示範囲では語義の決定が文脈上困難と考えられる場合、調査対象用例にできないという制限がある）。

実験協力者は、Yahoo! クラウドソーシング（<https://crowdsourcing.yahoo.co.jp/>）によって募集した Yahoo! 日本語 ID を有する20歳以上の男女である。用例1（例文）×1（判定文）の各組合せにつき20名以上を募集した（募集画面例は図2）。

図1 作業画面例

図2 実験協力者募集画面例（募集終了時のもの）

4. 調査例

本調査の手順によって、高い類似度と判定された用例群を、語義的に類似したまとまりとして可視化することができる。また、用例と語義との対応を確認し、『分類語彙表』や先行研究の記述を検証可能である。いずれの用例とも類似度が低く判定された用例は、一般的な読み手の判断に揺れや迷いの生じる用例といえるが、比較的類似度の高い用例を取得することによって語義の関係を整理できる。あるいは、用例を語義的に分類することで分析の手がかりとすることも可能であろう。

以下では、具体的な調査例を示す。4.1 節は、多義語の語義を調査した例である。対義語の対照についてもあわせて紹介する。4.2 節は、詳細な語義分類の可能性のある表記に着目し、分析のために語義的な用例分類を用い

³ 他の類（体：主に名詞，用：主に動詞など）についても順次調査を進める予定である。

た例である。

4.1 多義語の語義的な用例分布調査例

調査対象語について、用例の類似度分布から、用例を語義的に分類することが可能である。ここでは、位置を示す「遠い」(33 用例, 56,100 組, 実験協力者 560 名(異なり))と「近い」(42 用例, 90,300 組, 実験協力者 1,453 名(異なり))の調査結果を例として、高類似度用例群の可視化を試みる⁴。

調査手順は3節の通りである。結果、語義的に高類似度と判定された用例の分布を図3に示す。

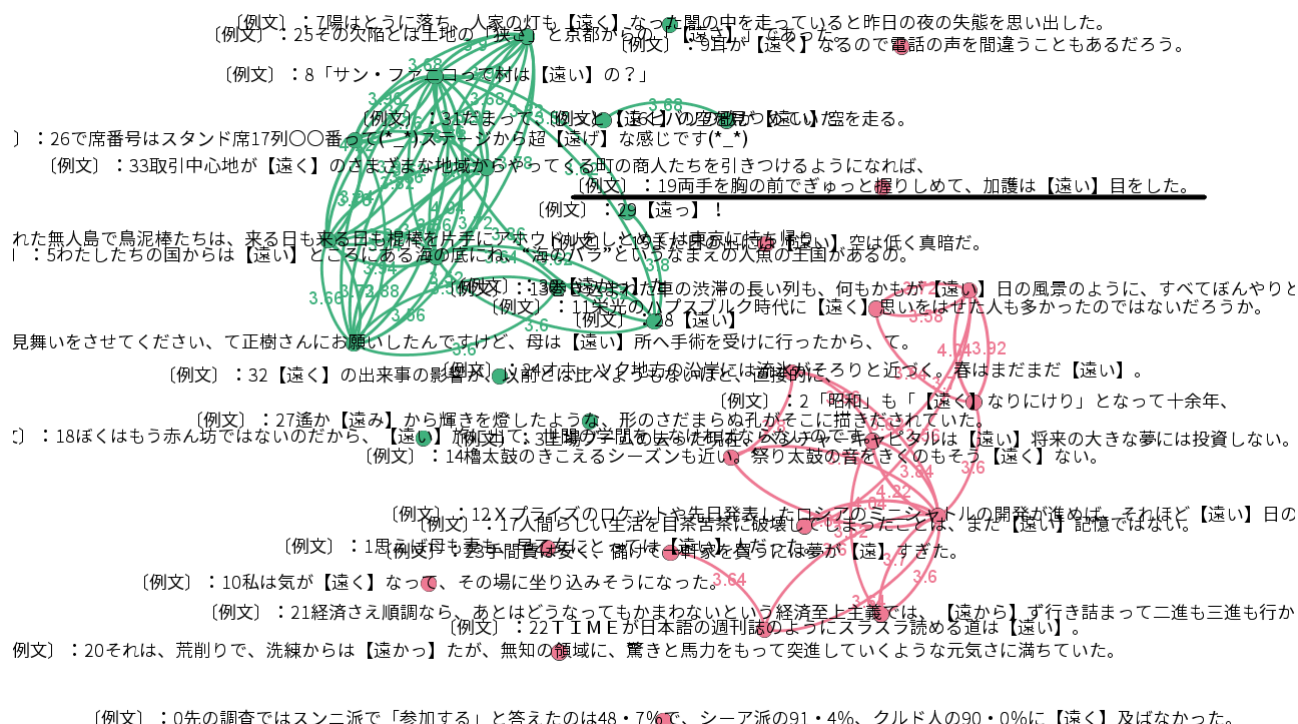


図3 「遠い」の語義的高類似度用例群(類似度 3.5 以上を線と数値で示す)

「遠い」は、類似度の高い(判定平均 3.5 以上)用例を見ると、空間(分類番号 3.1911(遠近):以下同様)、時間(3.1600)の2種類の語義が大きく分類可能であった。図3では、線でつながれた用例が大きく2つのクラスタに分かれており、左のクラスタが空間を、右のクラスタが時間を表す。また、図3で高類似度の線を表示していない(図3で点在する)用例についても、類似度を見ることで、関係(3.1110)、生理(3.5710(「耳が遠い」「気が遠くなる」)の分類が、概ね類似度が高い用例群となっていることがわかった。

(1)のような例(図3中下線部)は、高類似度の用例群のいずれにも属さないが、(1)の用例と類似度の高い判定結果となる用例群を調べると、空間の解釈、時間の解釈、「気が遠くなる」のいずれもが同程度に類似した意味と判断されていることがわかった。また、(1)と類似度の高い(判定平均 3.42)(2)をはじめ、時間の解釈の場合は過去のことに限定されており、「遠い目」「気が遠くなる」などが、過去と関わる解釈がなされている可能性などが明らかとなった。

(2) 栄光のハプスブルク時一代に遠く思いをはせた人も多かったのではないだろうか。

(サンプル ID : PB23_00052)

「遠い」の対義語と考えられる「近い」では、「遠い」と同様の意味的な分布が期待される。そこで、図4は、語義的に高類似度と判定された「近い」の用例分布を示す。図4の用例分布を見ると、「近い」は、空間(3.1911)

⁴ 分析や整理した語義のグラデーション作成の試みなどについては、加藤ら(2019d)を参照されたい。

表記でも同様に出現している。個人の使い分け意識があるとしても、一般的な使い分けは不明であった。

そこで、一般的な読み手⁵による語義的な類似度判定結果によって、各表記の用例の語義的な分布を整理することができれば、読み手がどのように読み分けているのかがわかると考えた。すなわち、用例を語義的に分類することで、表記の使い分けの効果を解明できる可能性があると期待した。

この調査では、71 例について、異なり 2,387 人による延べ 42,650 件 (8,530 画面) の回答を得た。図 5 に、4.1 節同様、語義的に高類似度と判定された〈あたたかい〉の用例分布を示す。

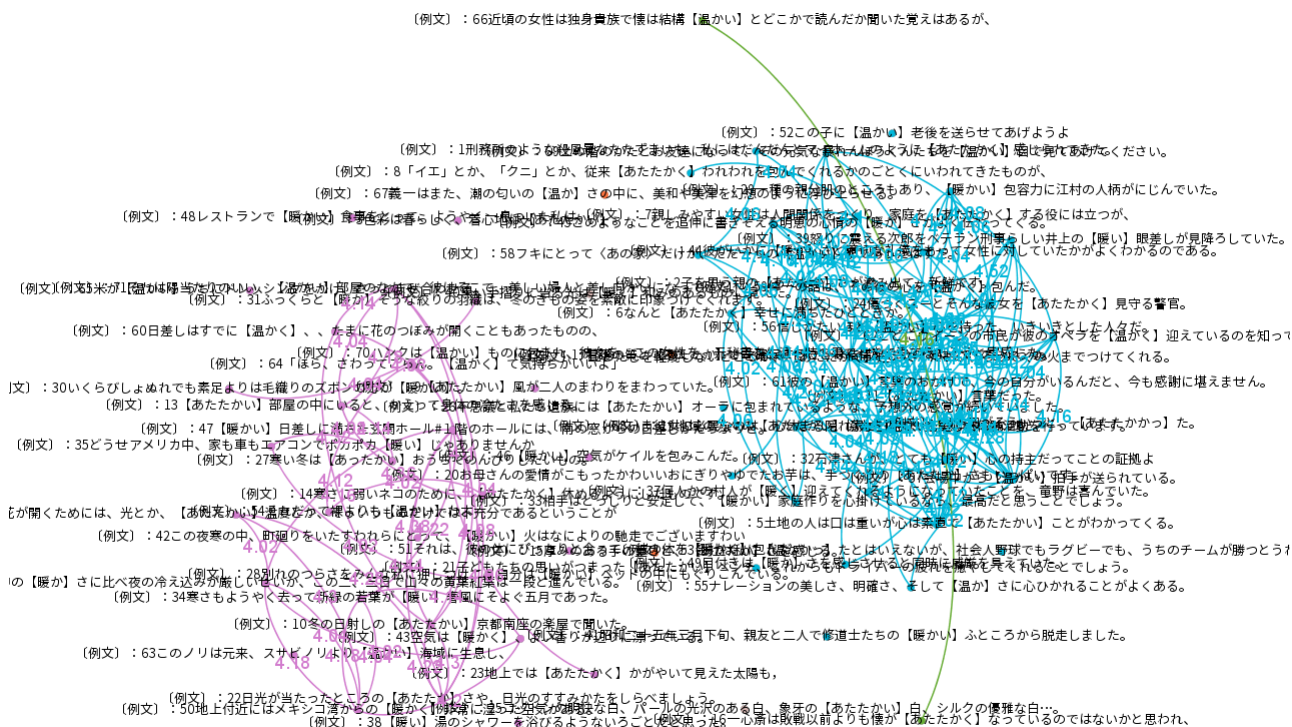


図 5 〈あたたかい〉の語義的高類似度用例群 (類似度 4.0 以上を線と数値で示す)

〈あたたかい〉について、高類似度 (4.0 以上) と判定された用例は、大きく 2 つに分類された。右のクラスターは比喩的な用例群であり、左のクラスターは比喩でない用例群である。

比喩でない用例群に表記差は確認できなかったものの、図 5 では高類似度の線を表示できていないが「衣類」「場所」「温度」などが類似性の高い用例群と分類された。また、比喩的な用例群では、「あたたかい」が 13 件と最も多く出現し、「温かい」が 9 件、「暖かい」が 8 件という、表記的に若干ひらがなの多い傾向が見られた。比喩的な意味であるために、「温かい」「暖かい」の書き分けでない「あたたかい」において類似度が高く判定された可能性が考えられた。このほかの類似性のまとまりとしては、「温かい目」「暖かい目」「あたたかく見守る」、「あたたかいマイホーム」「温かい家」「暖かい家庭」「あたたかい家庭」などの類似文脈の用例がつながっているほか、「あたたかい言葉」「温かい拍手」「暖かい心情」などが高い類似度でつながっている傾向が見られる。いずれもコロケーションや文脈において語義的に類似したカテゴリでまとまる傾向であった。

以上のような読み手が一般に着目する言語形式を手掛かりに、身体そのものあるいは身体と直接的に接触して取得可能な温度について、「温かい」を用いられやすい傾向が見られ、「暖かい」は身体に直接的に接しないが体感可能な場合に用いられやすい傾向となっていることがわかった。

また、分類の区別を見ると、〈あたたかい〉とを感じる身体がいずれであるのかという視点が関わっていた。た

⁵ 3.2 の同画面同手順により、別途〈あたたかい〉の語のみ表示した調査も試みた。文脈と関わらない場合の読み手の判断を調査しておきたかったためである。この結果では、「あたたかい」「アタタカイ」はいずれの表記とも高い類似度 (3.6~4.1) とされ、ひらがなやカタカナの影響は見られなかった。「暖かい」と「温かい」、「温かい」と「暖かい」を対照すると、類似度の判定は下がった (2.5~2.8)。語のみの表示では書き手としての意識が影響するだろうが、読み手の意識としても「温かい」と「暖かい」には意味的な差のある可能性が考えられた。

例えば、(3)の「暖かい眼差し」は、着点（差される対象）の視点と読み取ることができるが、(4)の「温かい目」は、起点（目を持つ対象）の視点と読み取れるという違いがあった。このほかの用例でも、どちらの視点で文脈を読み取ったのかという判断が、読み手の意味解釈に関わっていると考えられた。読み手によって分類された用例分析の結果として、「温かい」「暖かい」の表記は、身体との距離感によって使い分けられる傾向が確認されたのである。

(3) 怒りに震える次郎をベテラン刑事らしい井上の暖い眼差しが見降ろしていた。

(サンプル ID : LBe9_00011)

(4) 上の階のかたとお友達になって、その元気な暴れんぼうくんたちを温かい目で見てあげてください。

(サンプル ID : PM51_00423)

類似度で用例が分類されたことにより、類似性の高い用例がまとまり、詳細な用例分析がしやすくなる。このように、クラウドソーシング実験によって、多数の一般的な判断結果を収集し、用例の語義的な分類を行うことができる。一般的な見地の分類が、詳細な用例分析のきっかけとなり得る。

5. まとめと展望

本発表では、クラウドソーシングを用いた実験を考えるに至った背景と実験設定を紹介した。また、実際にクラウドソーシングを用いた調査の試みとして、多義語の語義分布調査例と、語義的な用例分類を分析に使用した例（表記を観点とした語義分布調査例）を示した。

クラウドソーシングを用いることで、一般的な判断として多数の集合結果が短期間に収集可能であるため、我々が目的としているような意味的情報付与データの整備に多数決的な結果を反映させることをはじめ、理論の検証を含めた実態調査や、分析の困難なデータの初期分類などにも活用可能性が考えられる。

今後も、BCCWJ-WLSP の整備を進めるとともに、語義研究への活用を図りたい。

謝辞

本研究は、国立国語研究所コーパス開発センター共同研究プロジェクト、JSPS 科研費 18K00530, 18K00634, 18K18519, 19K00591, 19K00655, 国立国語研究所所長裁量経費 2018 によるものです。

参考文献

- 加藤祥, 浅原正幸, 山崎誠(2019b). 分類語彙表番号を付与した『現代日本語書き言葉均衡コーパス』の書籍・新聞・雑誌データ. 日本語の研究, 15 (2), p. To Appear.
- 加藤祥, 浅原正幸, 山崎誠(2019b). 『現代日本語書き言葉均衡コーパス』新聞・書籍・雑誌データの助動詞に対する用法情報付与. 日本語学会 2019 年度春季大会予稿集, p. To Appear.
- 加藤祥, 田邊絢, 浅原正幸, 古宮嘉那子, 新納浩幸(2019c). 多義語の語義分布と語義間の派生関係調査の試みー相の類を中心に. 言語処理学会第 25 回年次大会
- 加藤祥, 西内沙恵, 浅原正幸(2019d). 多義語用例の類似度による語義の分類; 「遠い」と「近い」を例に. 日本認知言語学会第 20 回全国大会予稿集, p. To Appear.
- 文化庁(2015). 言葉に関する問答集 総集編. (「言葉に関する問答集」10, 1984 年)
- 八亀裕美(2015). 〈関係〉を表す形容詞の意味と用法: 「近い」と「遠い」. 甲南大學紀要, 文学編, pp.11-22.
- 山崎誠, 柏野和佳子(2017). 『分類語彙表』の多義語に対する代表義情報のアノテーション. 言語処理学会第 23 回年次大会発表論文集, pp.302-305.
- 吉村弓子(1981). 同音語の用法ー『温かい』と『暖かい』ー. 日本語と日本文学, 1, pp.47-56.