

令和3年2月4日

景気ウォッチャー調査を用いたテキスト分析の方法について（第1回）

茨城県政策企画部統計課 企画分析グループ

1 はじめに

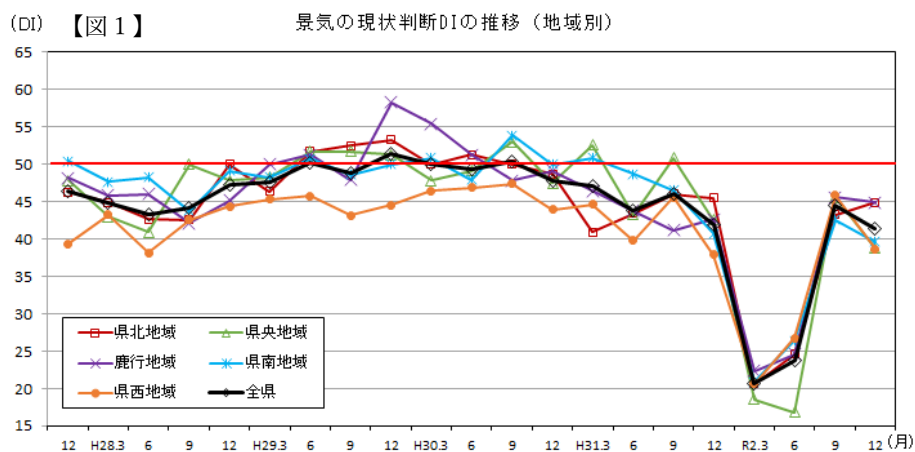
人々の景況感を掴むことを目的に、茨城県統計課（以下「当課」という。）が実施している「茨城県景気ウォッチャー調査」（以下「本調査」という。）は、約300名の調査客体の景気判断理由をテキストベースで得られるという特徴がある。今回は、このような数値化できないデータを視覚的に分かりやすくして俯瞰できるようにするテキスト分析の方法として、単語間の共起（ある2つの語が同じ文章中に出現すること。）関係を可視化する「共起ネットワーク図」を中心に紹介し、調査客体の景気判断理由を考察する。

なお、分析にあたってはプログラムを書かなくても共起ネットワーク図を作成することができる「KH Coder」²を用いる。また、今回の分析に用いるデータは当課公式サイト³で公表している。

2 調査結果からの推測

景気の状態判断DI（調査時点の景況感を示す指標。以下「DI」という。）の推移を見ると、令和2年以降の動きに特徴があることがわかる（図1）。令和2年3月調査では、DI値は大きく低下しており、リーマン・ショック期のDI値⁴と同程度にまで落ち込んだ。その半年後の9月調査では、未だDI値が50を下回っているものの、大きく低下する前の令和元年12月調査と同程度のDI値まで上昇した。

令和2年3月頃は、国内で新型コロナウイルスの感染が拡大し始めた頃である。令和2年3月調査の調査客体の景気判断理由のコメントでも、新型コロナウイルスに関するものが多く見られたことから、新型コロナウイルスがDI値の落ち込みに繋がったと推測できる。



¹ 調査概要：<https://www.pref.ibaraki.jp/kikaku/tokei/fukyu/tokei/betsu/bukka/watch/chogai.html>

² テキストデータを統計的に分析するためのフリーソフトウェア。（<http://khcoder.net> よりダウンロード可能）

³ <https://www.pref.ibaraki.jp/kikaku/tokei/fukyu/tokei/betsu/bukka/watch/bunseki/index.html>

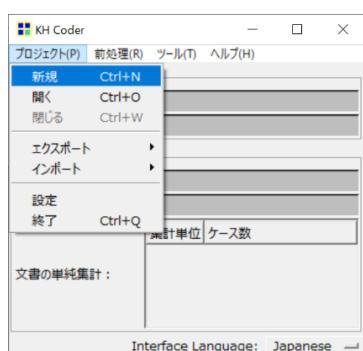
⁴ リーマン・ショック期はDI値が大きく低下し、平成21年3月には本調査開始以来過去最低の18.5となった。

3 分析データの読み込みと前処理

「KH Coder」を使って、景気判断理由に新型コロナウイルスに関するコメントがどのくらい多いか可視化してみよう。

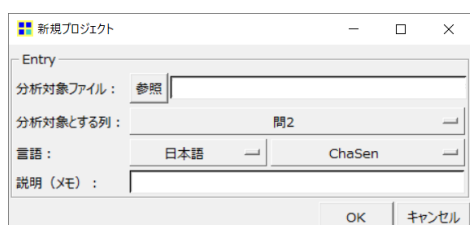
(1) 新規プロジェクトの作成

KH Coder を起動し、[プロジェクト(P) - 新規]から新規プロジェクトを作成する。



(2) 分析対象ファイルの読み込み

「参照」ボタンから、分析対象のデータファイルとして令和2年3月調査の景気判断理由のテキストデータをまとめた Excel ファイル「2020.03 茨城県景気ウォッチャーコメント.xlsx」を選択して読み込む。



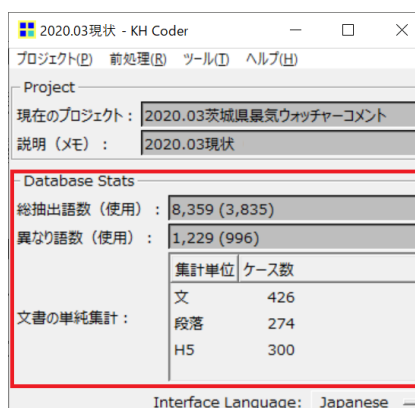
この際、分析対象とする列を指定する。ここでは、現状の景気判断理由が書かれている「問2」列を選択する。

(3) 前処理の設定・実行

分析に先立っては、「形態素解析」と呼ばれる前処理が必要である。この前処理によって、分析対象のテキストデータから語が切り出され、頻出語はどのようなものがあるかなどの確認ができるようになる。前処

理には、形態素解析のための辞書が必要であり、KH Coder では、「MeCab」か「ChaSen」の2種類が用意されているが、ここでは「MeCab」を選択し、「OK」ボタンを押す。

次に、[前処理(R) - 前処理の実行]から前処理を実行する。前処理が完了すると、「Database Stats」内に前処理の結果が表示される。



(4) 抽出語リストの確認

[ツール(T) - 抽出語 - 抽出語リスト]から、頻出上位の語を確認する。すると、「ウイルス」、「コロナ」、「新型」などの語が頻出していることが分かる。

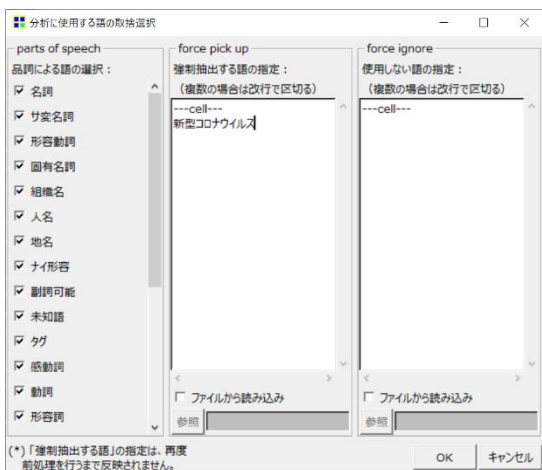


ここで、「新型コロナウイルス」といった最近になって使われ始めたような語は、先述した辞書には登録されていないため、「新型コロナウイルス」1語では切り出されない。

このような場合には、「新型コロナウイルス」を1語で強制的に抽出させる設定を加える。

(5) 強制抽出する語の指定

[前処理(R) - 語の取捨選択]から、強制的に抽出させたい語を左枠「force pick up」内に入力する。



強制抽出する語の指定が完了したら、[前処理(R) - 前処理の実行]から再度前処理を実行する。前処理の完了後、抽出語リストに「新型コロナウイルス」の語が正しく強制抽出されていることを確認する。

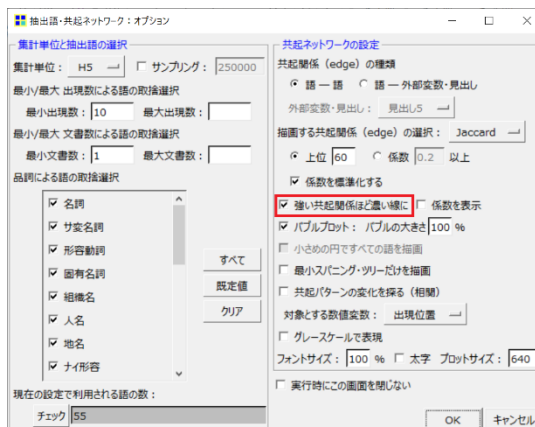
#	抽出語	品詞/活用	頻度
1	新型コロナウイルス	タグ	188
2	影響	サ変名詞	150
3	減少	サ変名詞	54
4	キャンセル	サ変名詞	41
5	状況	名詞	33
6	客	名詞C	31
7	客数	名詞	29
8	悪い	形容詞	27
9	売上げ	名詞	27
10	減る	動詞	26

4 語と語の共起関係

前処理が完了すると、共起ネットワーク図の作成などの分析ができるようになる。

[ツール(T) - 抽出語 - 共起ネットワーク]から、共起ネットワーク図の設定を行う。

ここでは、「品詞による語の取捨選択」などの項目は既定値のままで行う。ただし、共起度の強弱を可視化するため、「強い共起関係ほど濃い線に」にはチェックを入れることを推奨する⁵。



以上の設定で共起ネットワーク図を作成した結果は【図2】⁶のとおり。この図を見ると、「新型コロナウイルス」という語に、「キャンセル」や「減少」といった語が共起している。「減少」には「客数」も共起していることから、新型コロナウイルスの影響により、キャンセルや客数の減少が起こっていることが読み取れる。実際のテキストデータを覗いても、そのようなコメントが複数見られた。

⁵ 共起の強弱をはかる指標として Jaccard 係数などが用いられる。ここでは、既定の設定である Jaccard 係数により共起の強弱を計算させている。

⁶ この図では比較的強い共起関係にある8つのグループが8色で色分けされているが、色そのものには意味はない。

<調査客体のコメント> (一部抜粋)

・ゴルフ場 (県央地域)

例年3月は来場客数がピークとなるが、新型コロナウイルスの影響から、キャンセル数も多く、前年比で20%以上の減少が予想される。

・旅行代理店 (鹿行地域)

新型コロナウイルスの影響により、予約のキャンセルが相次ぎ、客数が減った。

5 語と景気評価の共起関係

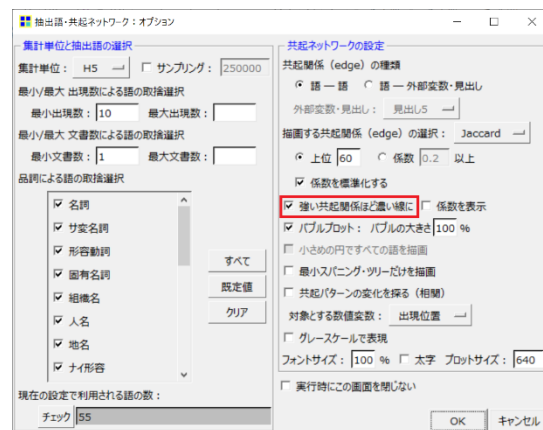
「4 単語間の共起関係」では、語一語での共起関係を図にしたが、今度は語一外部変数・見出しでの共起関係を見ていく。

「外部変数・見出し」で、令和2年3月調査の「問1」(調査月における景況感についての評価を問うもの)を選択すると、問1の各回答(「1:良くなっている」「2:やや良くなっている」「3:変わらない」「4:やや悪くなっている」「5:悪くなっている」と共起している語の様子を見ることができる。

例えば、ある語「word」が「4:やや悪くなっている」や「5:悪くなっている」と共起していれば、調査月における景況感で、4や5のマイナスな評価をしている人のコメントに「word」という語が多く見受けられると言うことができ、「word」がDI値の低下要因になっている可能性が高い。

では、[ツール(T) - 抽出語 - 共起ネットワーク]から、共起ネットワーク図の設定を再度行ってみる。

「共起関係(edge)の種類」を「語一外部変数・見出し」に変更し、「外部変数・見出し」には「問1」を選択する。



作図した結果は、【図3】⁷のとおり。赤い四角は「外部変数・見出し」に指定した問1の各回答番号で、四角内の数字が大きいほど調査月の景況感に対する評価が低い。この図から、次の2点に注目したい。

- ① 調査月の景況感の評価が「3:変わらない」以下の評価において「新型コロナウイルス」と共起しており、特に、4や5の評価で共起度が高い。
- ② 「4:やや悪くなっている」や「5:悪くなっている」といった低い評価において「キャンセル」と共起している。

これらの結果から、新型コロナウイルスについては、「2 調査結果からの推測」で述べた「新型コロナウイルスがDI値の落ち込みに繋がった」とする推測のとおりであることが見て取れる。また、「キャンセル」という語は出現頻度が比較的高かった(41件)ものの、4や5といった低い評価にのみ共起関係が見られたという点から、「キャンセル」が調査客体の景況感をより下向きにさせ、令和2年3月調査のDI値の下げ幅を

⁷ この図では語が3色で色分けされているが、その色分け規則は語から出ている線の本数によって異なる。

大きくさせたことが考えられる。

6 まとめ

「4 語と語の共起関係」と「5 語と景気評価の共起関係」の2つの側面から共起ネットワーク図を作成した結果、新型コロナウイルスが景況感に悪影響を及ぼし、特に、キャンセルやそれに伴う客数の減少がDI値の大幅低下を後押ししたという示唆が得られた。

以上のように、テキストデータを視覚的に分かりやすくして俯瞰するには、今回紹介した方法が有効である。とりわけ、本調査はテキストの量が多く、全て読むには時間を要するため、内容や特徴を大まかに把握したい場合には特に有効となる。

