

看護研究のためのテキストマイニング (質的/量的研究現場での利用事例)

2015.11.20

東京医科歯科大学

大河原知嘉子

テキストマイニングとは

- 質的データである文字テキストを量的分析手法である統計や多変量解析などによって分析する手法

方法	データ	分析方法
量的研究	数値(量的データ)	量的分析(統計)
質的研究	文字(質的データ)	質的分析
データマイニング	数値(量的データ)	量的分析(統計)
テキストマイニング	文字(質的データ)	量的分析(統計)

- ✓ テキストマイニングは質的分析と量的分析との両方の特徴を持つ。
- ✓ Mix Method(混合研究法)としても有用。

分析対象

1. 情動的・情報的・文字資料

- Ciniiなどの文献データベース、webサイト
- ガイドブック、ハンドブック
- 情報中心の記事
- 新聞記事

2. 物語(ナラティブ)的・文字資料

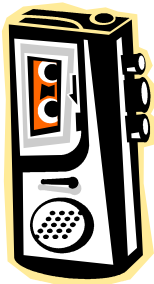
- 闘病記、感想文、レポート、小説、書籍
- ブログ、ツイッター
- 体験談、インタビューの逐語録

テキストマイニング

自由記載のアンケート

ご自由にお書きください

面接での逐語録



今までは



人の手でテキストを集計
人がテキストを読み、
内容を理解し、情報を得る

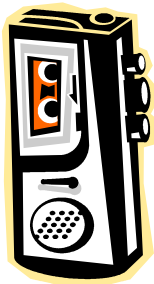
- 利点：人が理解した意味内容を分析に反映できる
- 欠点：大量のデータの分析困難
分析が主観的
情報処理に人手と時間がかかる

テキストマイニング

自由記載のアンケート

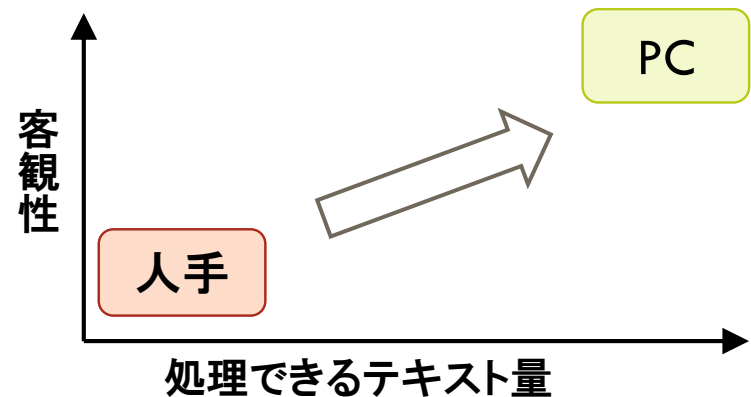
ご自由にお書きください

面接での逐語録



テキスト
マイニング

パソコンを利用して、
テキストを数値データと
同じように処理できる



利点：分析が客観的、再現可能
大量のデータの分析が可能
時間短縮が出来る

欠点：意味内容にまで踏み込んだ
分析が出来ない
ことばの落とし穴

テキストマイニングの効果

- 大量のテキストデータを分析できる
- テキストデータの特徴や傾向を量的に示せる
→ 文章基本情報、単語頻度分析など
- テキストデータの特徴を視覚的に示せる
→ 言葉ネットワーク分析、対応バブル分析など
- 文章と属性との関係を定量的に説明できる
- 再現性があり誰でも同じ基準で分析できる
- 分析時間を大幅に削減できる

人が読むだけでは得られない情報を
獲得できる

テキストマイニングの限界

- 文章の意味内容や行間を加味した分析が困難
- カテゴリーを作成してもことばの繋がりや使われ方、頻度などにより文章が分類される
 - ✓ ことば自体の意味で文章を分類するため、文章全体の意味内容での分類は困難
- 日本語の落とし穴
 - ① 英語と異なり、ことばの切れ目が不明確
 - ② 同じことばでも意味がさまざま→同音異義語の存在
 - ③ 会話分析ではことばの置換がおきるため関係性の分析が出来ない・・・意味的分類、会話の落とし穴

辞書作りが分析の精度を左右する

①不明瞭なことばの切れ目

- 英語と異なりことばの切れ目が不明確

E) It is my special day in my life.



E) It / is / my / special / day / in / my / life.

J) 優秀な日本人の看護師は美しい



J) 優秀な / 日本人の看護師は / 美しい

J) 優秀な日本人の / 看護師は / 美しい

優秀なのは誰？
日本人の看護師？
日本人？

文脈の中での理解が必要

どこで区切りますか？

- わかちがき=文章を意味のあるまとまりに区切ること
ex) 私は毎朝テレビで天気予報を確認する。
 - 私は毎朝テレビで天気予報を確認する。・・・形態素解析
 - 私は毎朝テレビで天気予報を確認する。・・・分かち書き

ユーザ辞書登録

辞書

単語

ユーザ辞書 類義語辞書 分割辞書

新規ユーザ辞書--*--

見出し語	品詞	読み
毎朝テレビ	名詞 固有名詞組織	
▶ 看護師長	名詞 一般	
*	名詞 一般	

▼ 新規ユーザ辞書--*--

↑

↓

名前変更

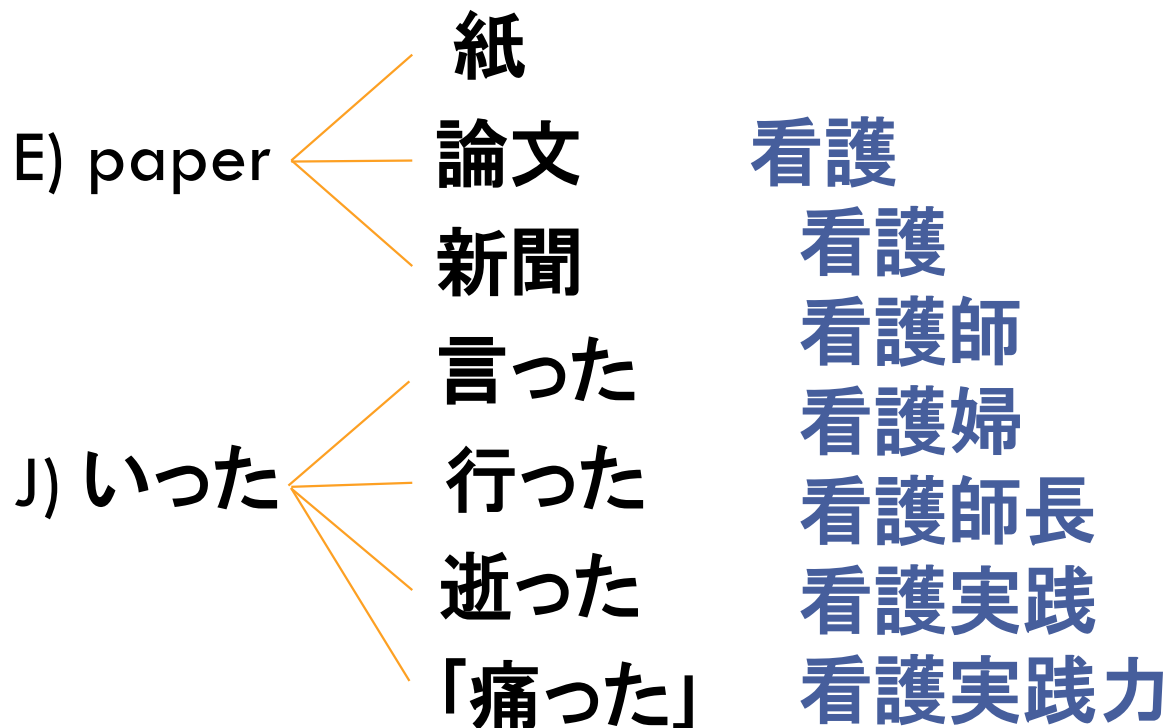
削除

辞書環境

新規登録 読み込み 書き出し 一括出力

②同音異義語

- 同じことばでも意味や構成がさまざま



辞書登録

□ 類義語辞書

辞書

ユーザ辞書 類義語辞書 分割辞書

新規類義語辞書--*--

代表語	品詞
看護師	名詞 一般

類義語

看護師
看護婦
ナース
看護師
看護婦さん
看護師さん
▶ 看護師さん
*

代表語	品詞
*	名詞 一般

新規類義語辞書
 形容詞辞書

↑
↓
名前変更
削除

辞書環

新規登録 読み込み 書き

□ ユーザ辞書

辞書

ユーザ辞書 類義語辞書 分割辞書

新規ユーザ辞書--*--

見出し語	品詞	読み
毎朝テレ	名詞 固有名詞組織	
▶ 看護師長	名詞 一般	
*	名詞 一般	

新規ユ

↑
↓
名前変更
削除

辞書環

新規登録 読み込

③-1 意味的分類

- 「看護師の対応がよい」
- 「対応した人の態度が良かった」
- 「気持ちの良い対応だった」

看護師について聞いている調査であればどれも
看護師の対応の仕方が良かったと言っている



ある程度分類は自動的にできるものの簡単には
いかず、分析過程で文章をどう扱うかがポイントとなる

③-2 会話の落とし穴

- 会話の録音データから逐語録を作成し、それをデータとして分析する場合におこる落とし穴
 - 会話には句読点がないため、“、”や“。”をつける位置は研究者
 - 聞いた会話は、漢字になっていないため、漢字変換をするのは研究者
 - 「あー」、「えーつと」等、意味をなさないつなぎ言葉の扱い
 - 「ハハハ」「あはは」「ウフフ」等笑い声の扱い

録音データ

□ 録音データ

- えーまあでもわたしに**その**かんがえなさいっていうんきかいとじかんをちゃんとあたえてくれてこたえがでたときもちゃんとそれをきいてくれるひとだったみたい**いなあはは**

□ テープ起こし・・・「素起こし」「**ケバ取り**」「**整文**」

- でも私に考えなさいっていう**機会**と時間をちゃんと与えてくれて、答えが出た時**も**ちゃんとそれを聞いてくれる、その聞いてくれる人だった**みたい**な。

③-2 会話の落とし穴:ことばの置換

- 会話分析ではことばの置換がおきるため
関係性の分析が出来ない

例:もう嫌になっちゃって...私



私はもう嫌になった

私はもう嫌になっちゃって〇〇

分析

文章の基本情報

項目	値
総行数	300
平均行長(文字数)	14.2
総文数	80
平均文長(文字数)	53.3
延べ単語数	1,634
単語種別数	814

- 分析対象となっている文章の特徴を量的にとらえるために最初に行う分析
- その文章の行数や行長、延べ単語数、単語種別数などから、文章の長さや、使用していることばの豊富さなどを明らかにする

品詞の出現頻度

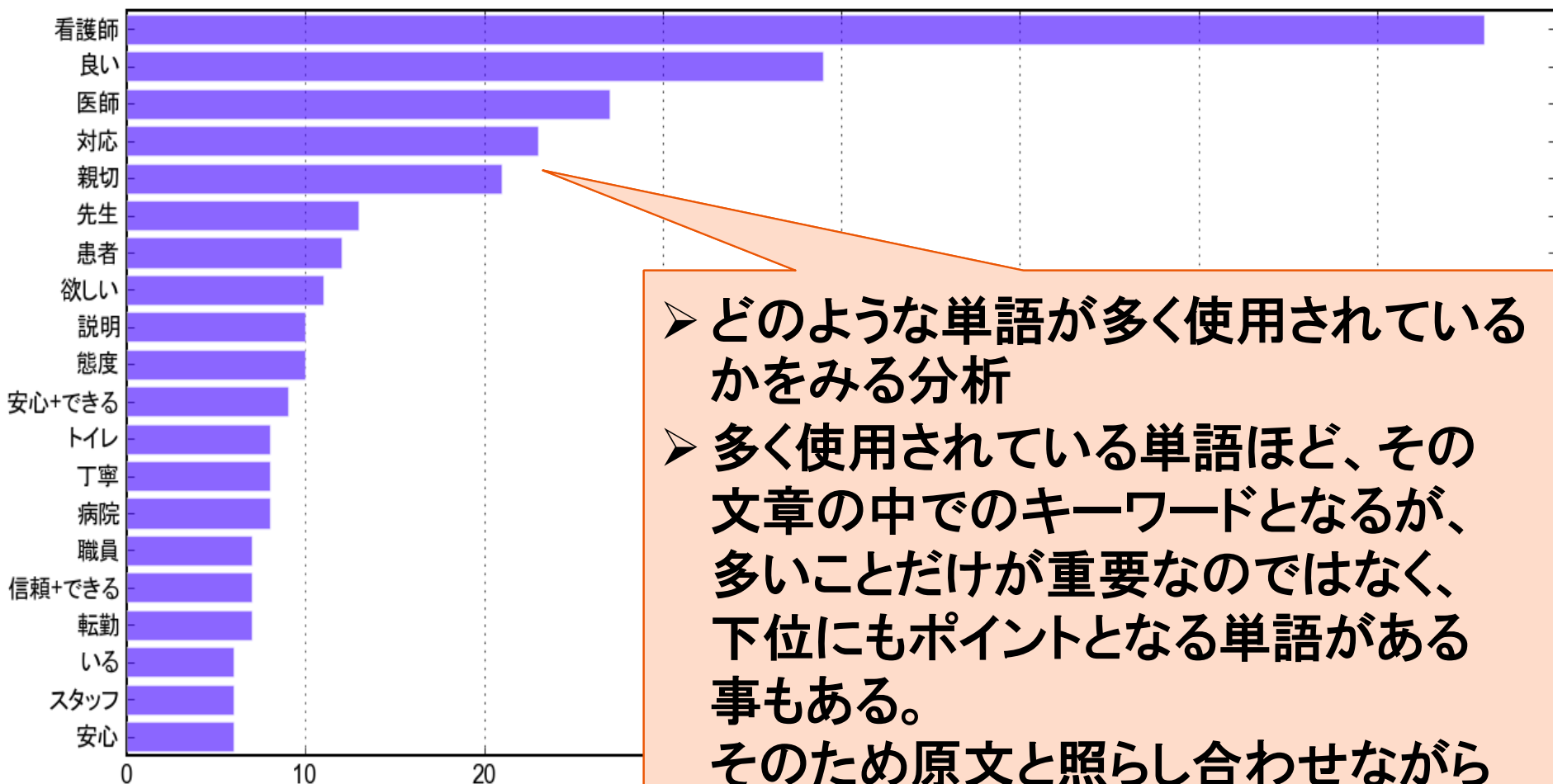
品詞	出現回数
名詞	4,406
動詞	3,318
形容詞	790
副詞	298
連体詞	245
接続詞	45

- その文章内で、どのような品詞の単語が多く使用されていたかを表す
- この文章だと名詞、動詞が多く使用されているが、それ以外の品詞はあまり使用されていない



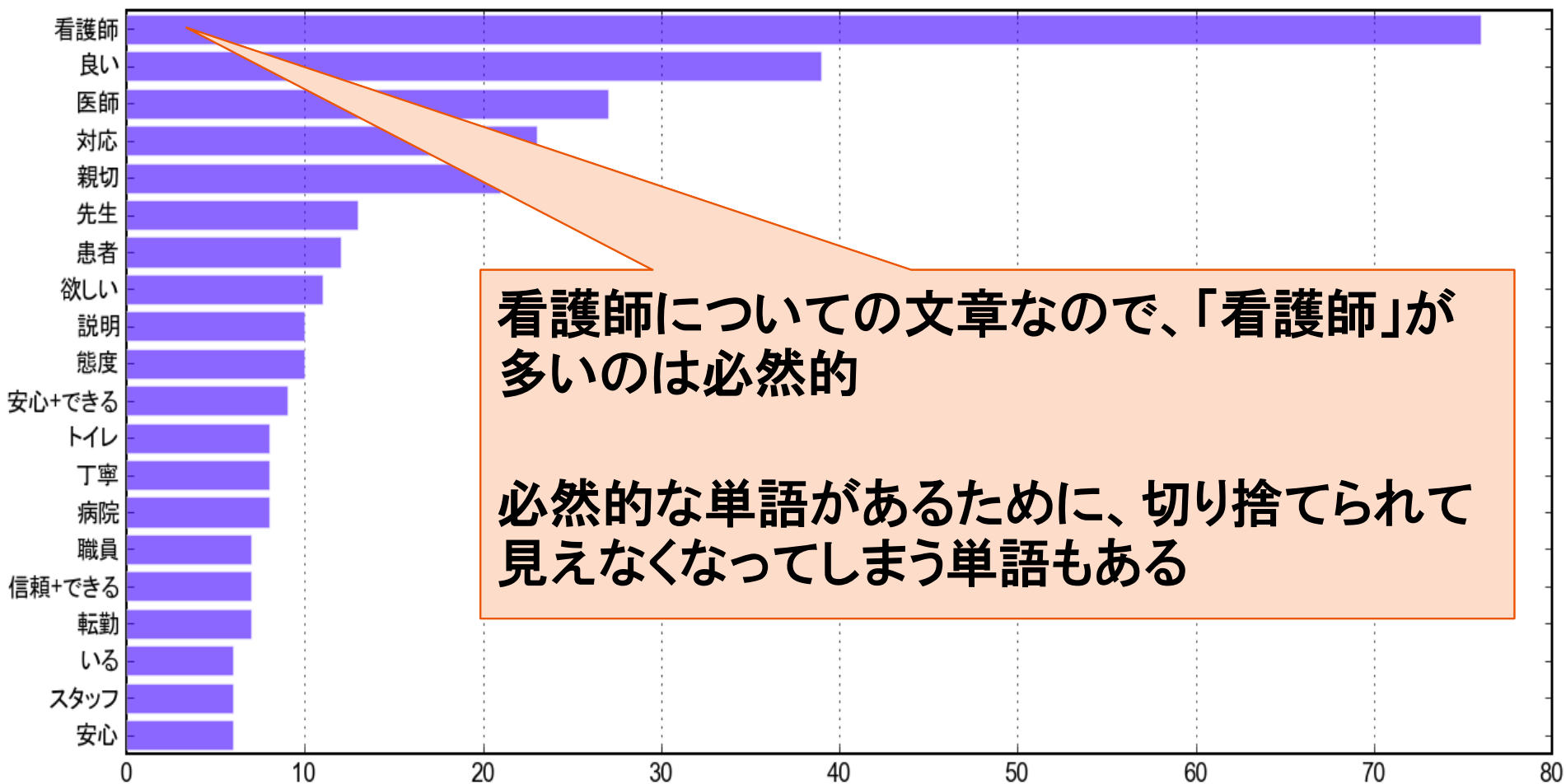
分析対象を決定するのに重要となる分析

単語頻度解析

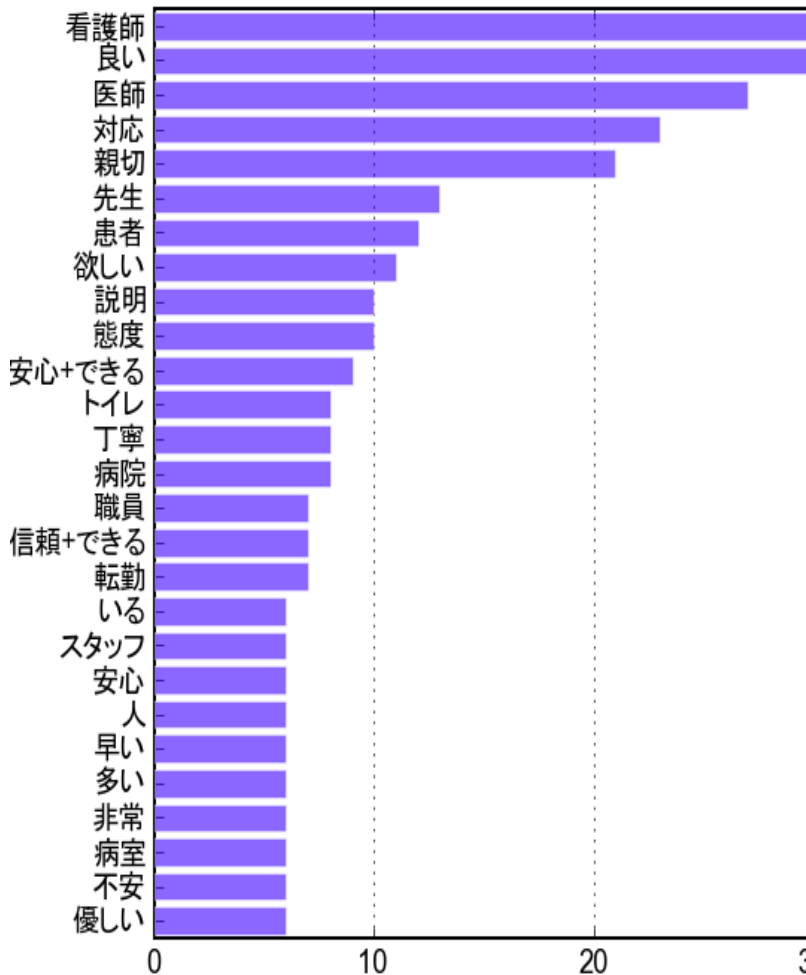


- どのような単語が多く使用されているかをみる分析
- 多く使用されている単語ほど、その文章の中でのキーワードとなるが、多いことだけが重要なのではなく、下位にもポイントとなる単語がある事もある。
そのため原文と照らし合わせながら注目する

単語頻度解析



単語頻度解析



単語頻度解析

フィルタ条件 属性別集計の設定

品詞

名詞・形容詞・動詞

オリジナル設定

単語

属性 述語属性を区別する

頻度 回 以上 回 以下

文字数 文字 以上 文字 以下

行中に現れる重複単語のカウントを1とする

上記の条件を満たすもののうち、上位 件を抽出する

順位が同じものは上位件数を超えても出力する

「その他」をカウントする 「合計」をカウントする

「合計」からの割合も抽出する

係り受け頻度解析～看護師～

The image displays three overlapping screenshots of a software interface for dependency frequency analysis, specifically for the word '看護師' (Nurse).

Leftmost Screenshot (Filter Settings):

- 品詞フィルタ: イメージ (名詞 - 形容詞・形), 行動 (名詞 - 動詞・サ変), 話題一般 (名詞 - 形容詞・形, 動詞・サ変), オリジナル
- 係り元品詞: [] 係り先品詞: []
- 頻度: 2 回 以上
- 行中に現れる重複表現のカウントを1とする
- これらの条件を満たすもののうち、上位順位が同じものは上位件数をこえても
- 「その他」をカウントする 「合計」
- 「合計」からの割合も抽出する
- Buttons: 初期値として使用, OK

Middle Screenshot (Filter Settings):

- 単語フィルタ: 係り元単語 [] 係り先単語 []
- 述語属性: 係り元述語属性 [] 係り先述語属性 []
- 文字数フィルタ: 係り元 1 文字 以上 [] 文字, 係り先 1 文字 以上 [] 文字
- 独立した単語もカウントする
- 述語属性を区別する
- 係り元/係り先を入れ替えたフィルタ条件でフィルタリングした結果も表示する
- Buttons: 初期値として使用, OK

Rightmost Screenshot (Single Word Filter):

単語抽出条件

単語	条件
* 看護師	を含む
*	を含む

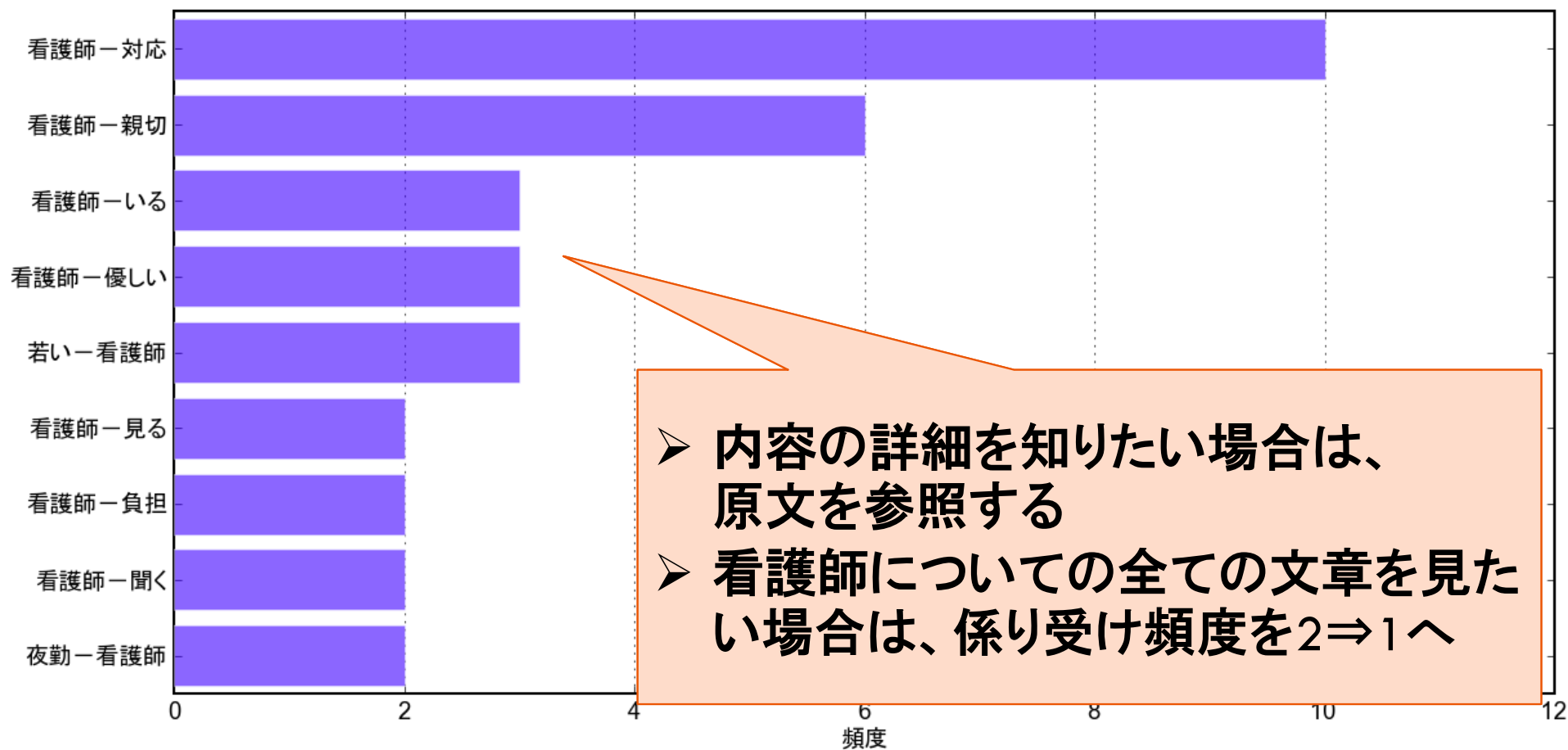
選択された行の条件をまとめて変更します

[を含む] [比一致する] [含まない] [比一致しない]

全ての条件を満たす (AND)
 いずれかの条件を満たす (OR)

Buttons: OK, キャンセル

係り受け頻度解析～看護師～



共起関係と係り受け関係

共起関係

- 文章中にある言葉が同時に出現するか否かという確率をもとにした言語現象



ことばの方向性はない

係り受け関係

- 主語述語の関係、修飾と被修飾の関係、補助の関係、並立の関係など文章の中で単語と単語がどのようにつながっているかを示す関係



ことばの方向性がある

テキストマイニングの基本は共起関係

身近な共起関係

Google

検索

すべて

画像

地図

動画

ニュース

ショッピング

もっと見る

東京都文京区
場所を変更

東京医科歯科大学

東京医科歯科大学

東京医科歯科大学病院

東京医科歯科大学 偏差値

東京医科歯科大学 図書館

東京医科歯科大学 歯学部

東京医科歯科大学大学院

東京医科歯科大学歯学部附属病院

東京医科歯科大学 教養部

他のキーワード: [東京医科歯科大学偏差値](#) [東京医科歯科大学図書館](#)

[東京医科歯科大学教養部](#) [東京医科歯科大学学生協](#)

[東京医科歯科大学難治疾患研究所](#)

[国立大学法人 東京医科歯科大学](#)

www.tmd.ac.jp/ - キャッシュ

東京医科歯科大学 TMDU の公式ホームページです。大学案内、入学案内、学部・大学院・附属病院等の紹介、研究活動、産学連携、国際交流など、東京医科歯科大学に関する情報をご覧ください。

[Google+ ページ](#)

東京

高文

文
外

東洋

日大
学

©201

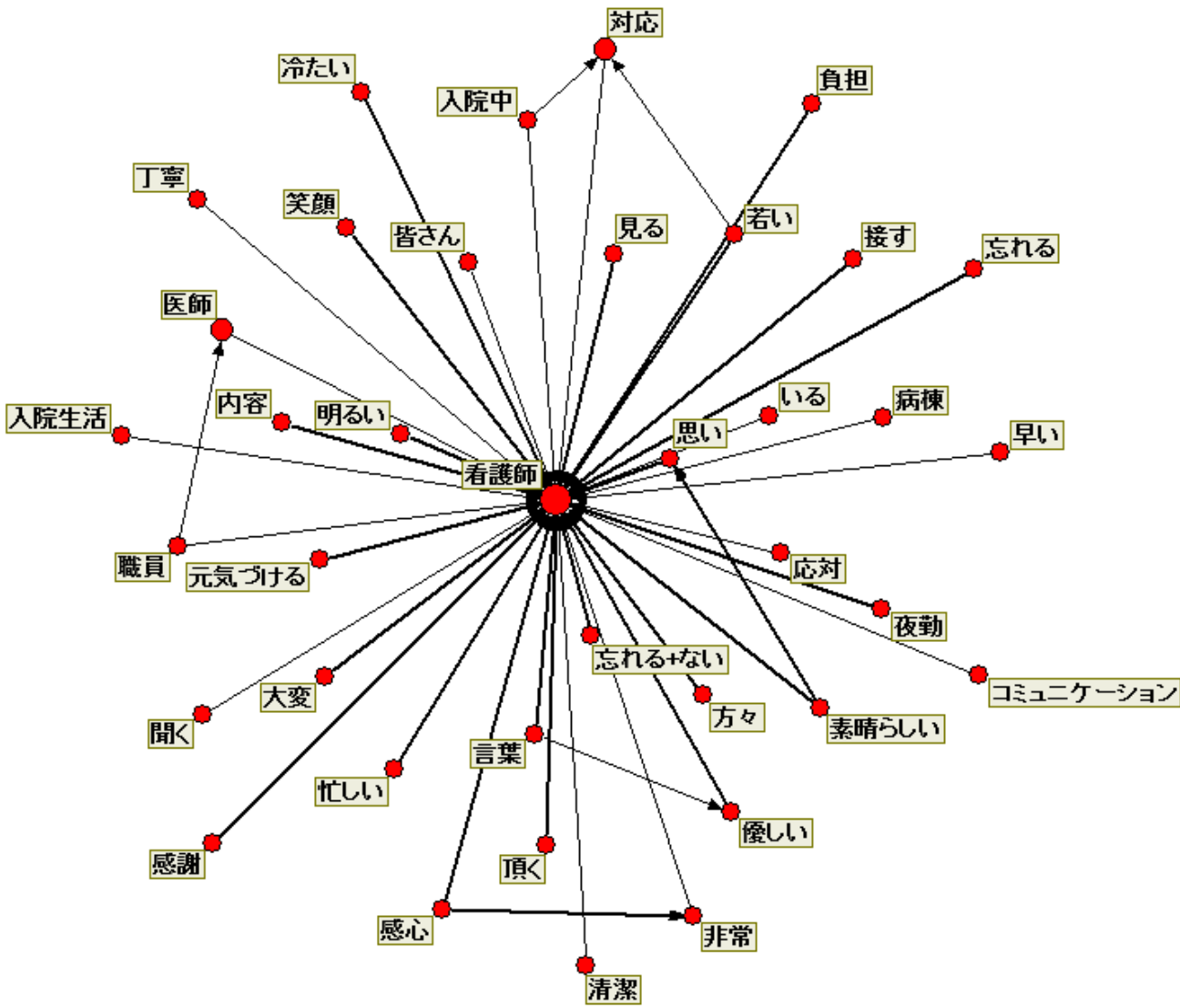


特徴

注目語情報

- ある特定の単語に注目し、出現頻度や使われ方等の情報から、その単語がテキストの中でどのように使われているかを調べることが出来る
 - 他のどのような単語・属性と同時に出現しているか
→ 共起関係⇒ **ネットワーク図**で表示
 - 注目した単語がどのような表現の中で使われているか
→ 係り受け関係⇒ **注目語表現情報**で表示

注目語情報～ネットワーク図～

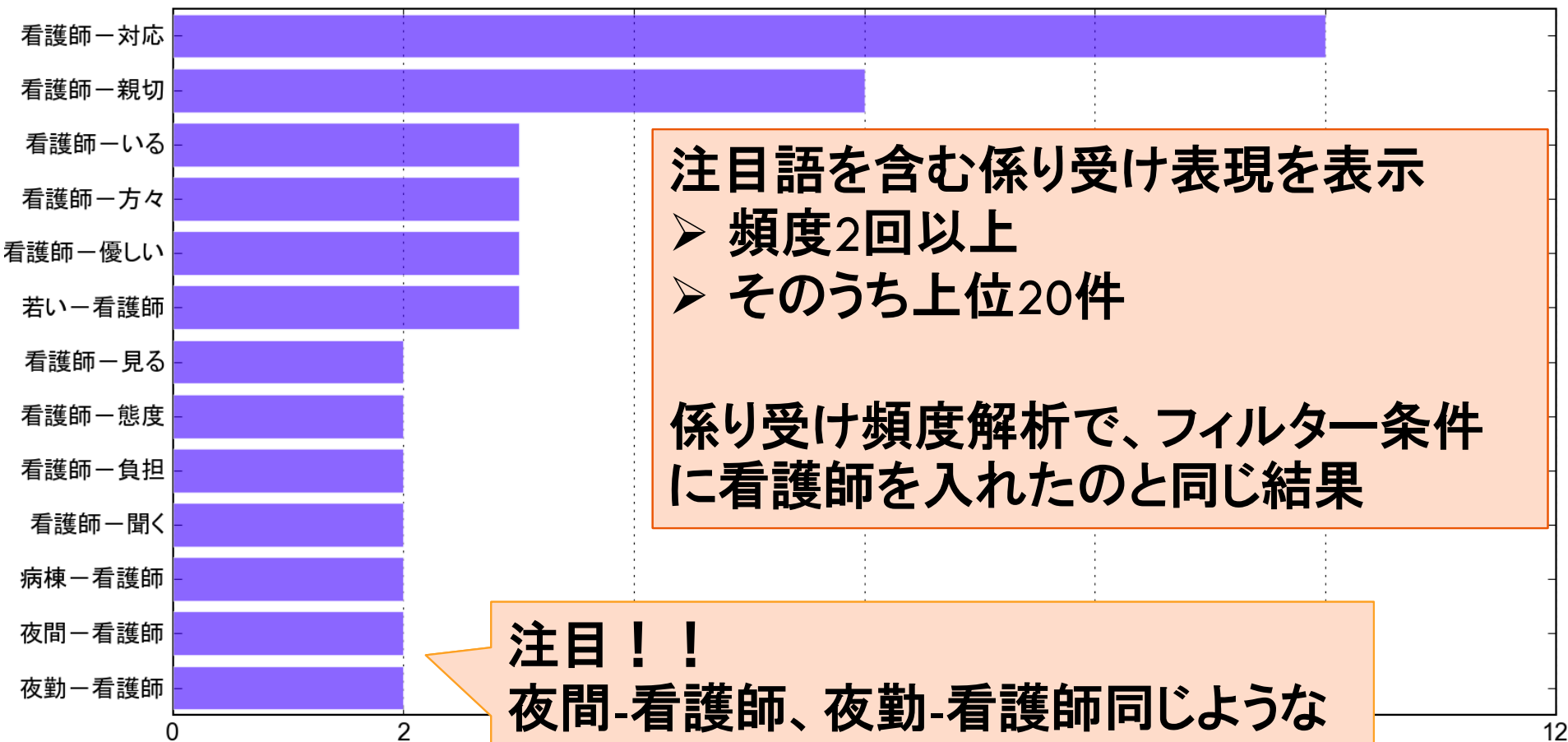


下記の共起ルールに
当てはまるものを
最大100ルール図示

- 注目語:看護師
- 出現頻度2回以上
- 行単位での共起

看護師と共起性が
強い単語は太い線で、
弱い単語は細い線で
結ばれる

注目語情報～注目語表現情報～



注目語を含む係り受け表現を表示

- 頻度2回以上
- そのうち上位20件

係り受け頻度解析で、フィルター条件に看護師を入れたのと同じ結果

注目！！
夜間-看護師、夜勤-看護師同じような内容では？

- 原文を検索

原文参照

原文参照

テキスト 属性

オリジナルテキスト --- 表示ページ: 1 / 全ページ数: 1

テーブル形式 レポート形式

テキスト検索条件

NOT テキスト

▶

夜間->看護師

夜勤->看護師

* [検索結果]

※項目内はANDで、項目間はORで結ば

完全一致で検索

全表示 分割表示 50

検索条件 分かち表示 複

検索

ファイルID	行ID	NO	性別	年齢	項目	半期	年月日	テキスト名	テキスト
1	94	263	男	73	看護師	上半期	2009年7月14日	自由記述	大変お世話を頂き、ありがとうございました。スタッフの方々に優しくしてもらい、良い入院生活を送ることが出来ました。 夜勤の看護師の方々 には特にご苦勞をおかけしました。御礼申し上げます。
1	160	162	男	72	入院	上半期	2009年9月1日	自由記述	夜間の看護師さん の仕事はハードすぎる。
1	163	161	女	69	入院	下半期	2009年9月8日	自由記述	特に 夜間の看護師 の負担が多すぎる。
1	221	234	女	44	看護師	下半期	2009年12月2日	自由記述	夜勤の看護師さん にカーテンを開けてとお願いしたら、ものすごい声で怒鳴られ一晩眠れなかった。一生忘れません。

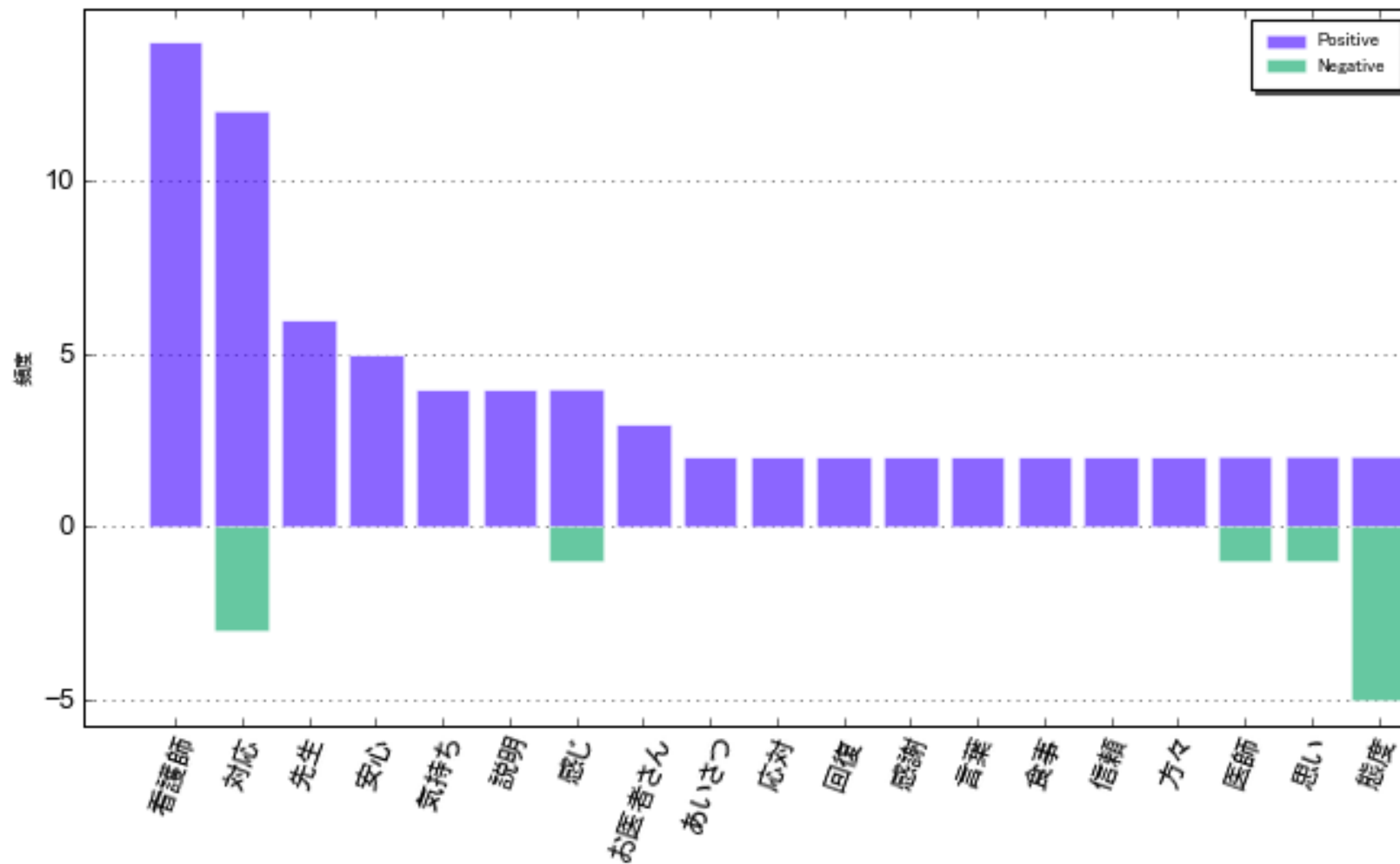
両方とも夜勤看護師についての文章

- “夜勤”と“夜間”を類義語辞典で一つにまとめて登録することもできる
- 夜間、夜勤のみの原文も参照することが重要！

評判分析

- 分析対象となる単語を、その単語に係り受けする好評語(ポジティブなイメージを持つ単語)と不評語(ネガティブなイメージを持つ単語)の頻度から単語の評価を分析する。
- 好評語と不評語はもともとTMSのなかで割り当てられている。
- ただ、文章が複雑になると、正確に評価されない場合があるため、原文に戻って確認、追加する必要がある。

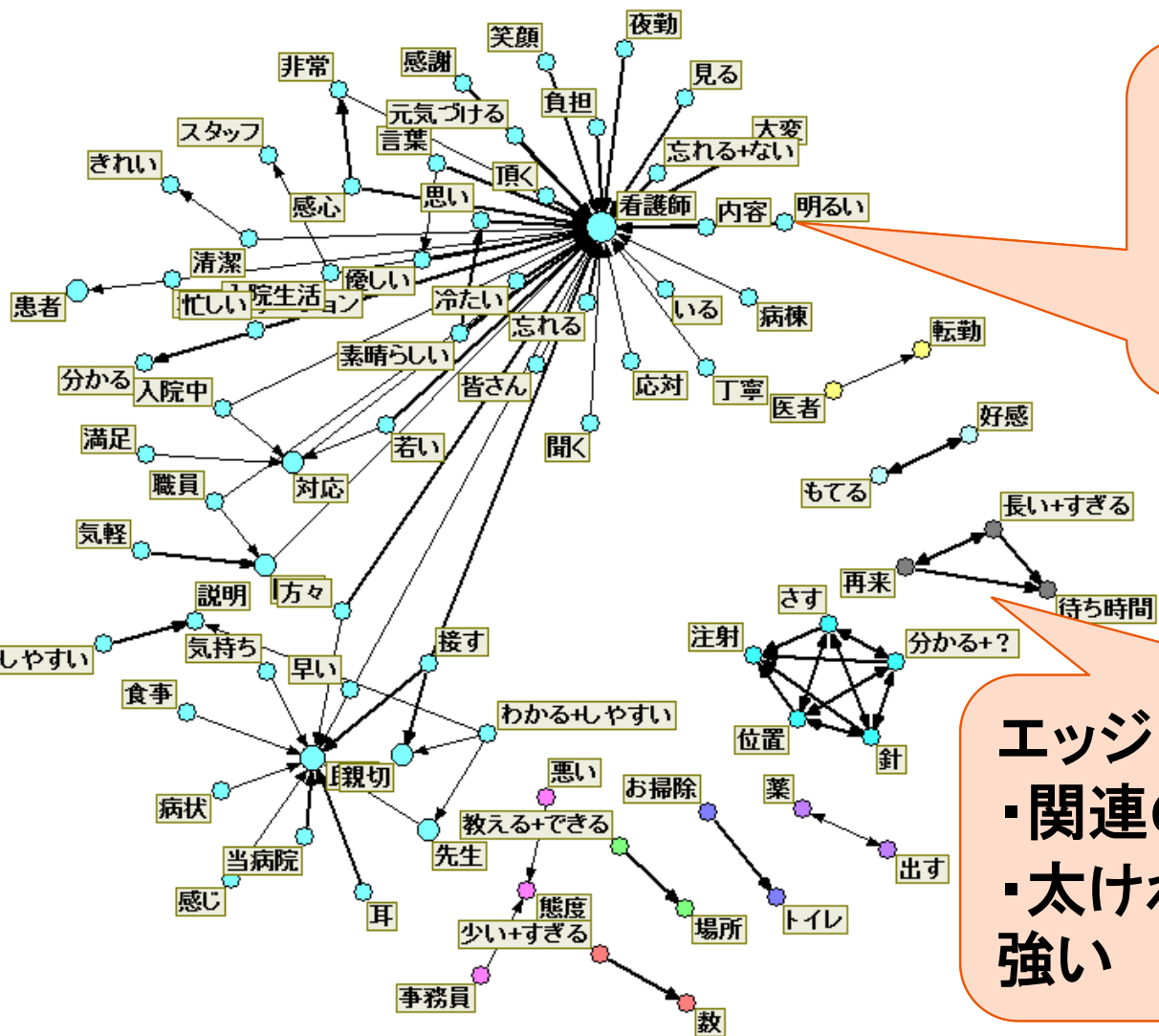
評判分析



言葉ネットワーク分析

- ことばとことばの関連をグラフを使用して可視化する分析方法
- 共起関係や係り受け関係をもとに、ことば同士の関連の強さを表す方法がある
- 頻度が高いほど関連が強いという考えをもとした分析

ことばネットワーク



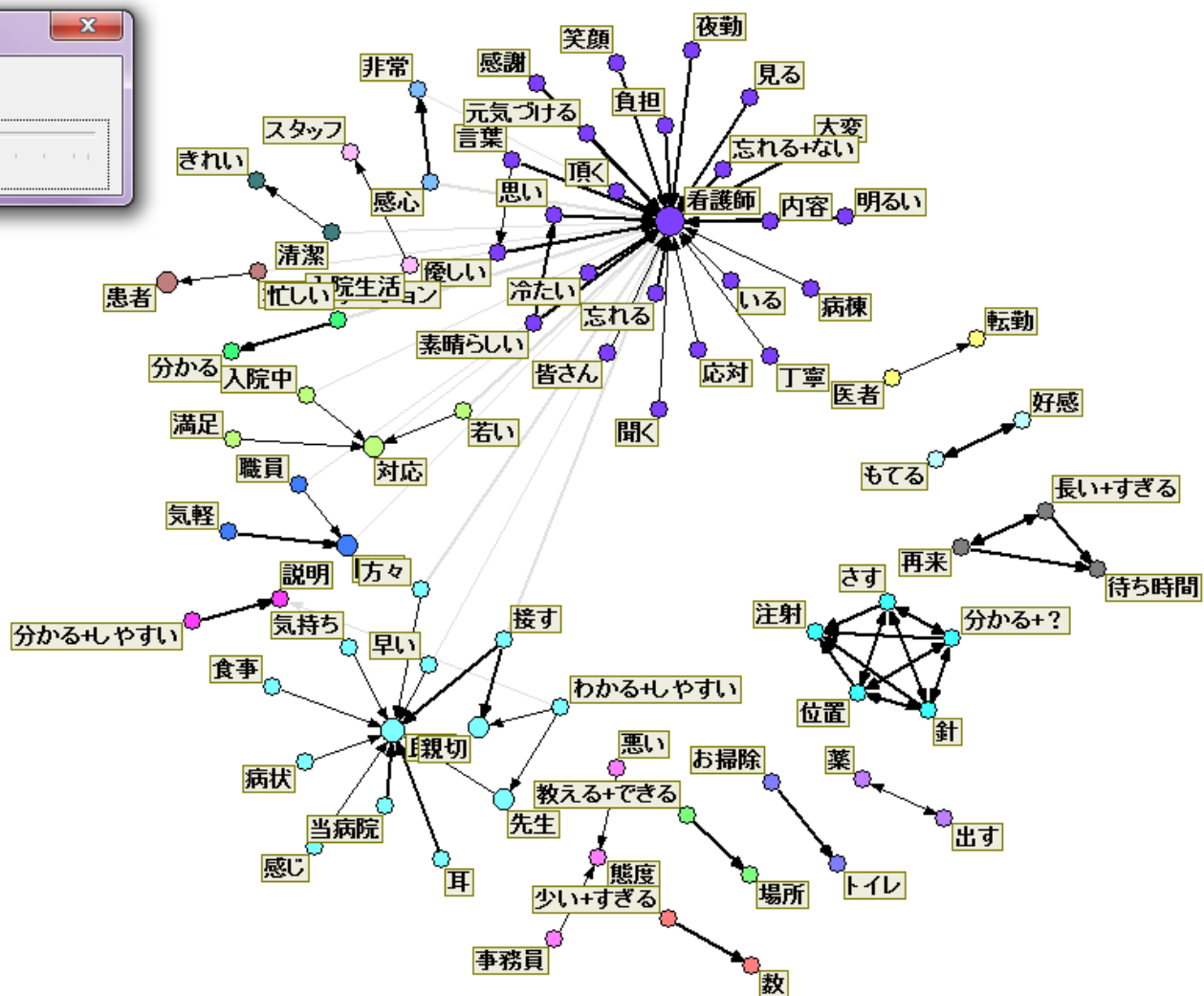
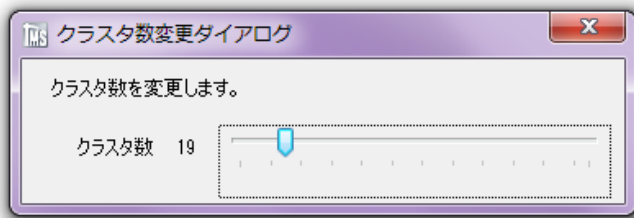
ノード

- ・頻度を表している
- ・大きければ大きいほど頻度が高い

エッジ

- ・関連の強さを表している
- ・太ければ太いほど関連が強い

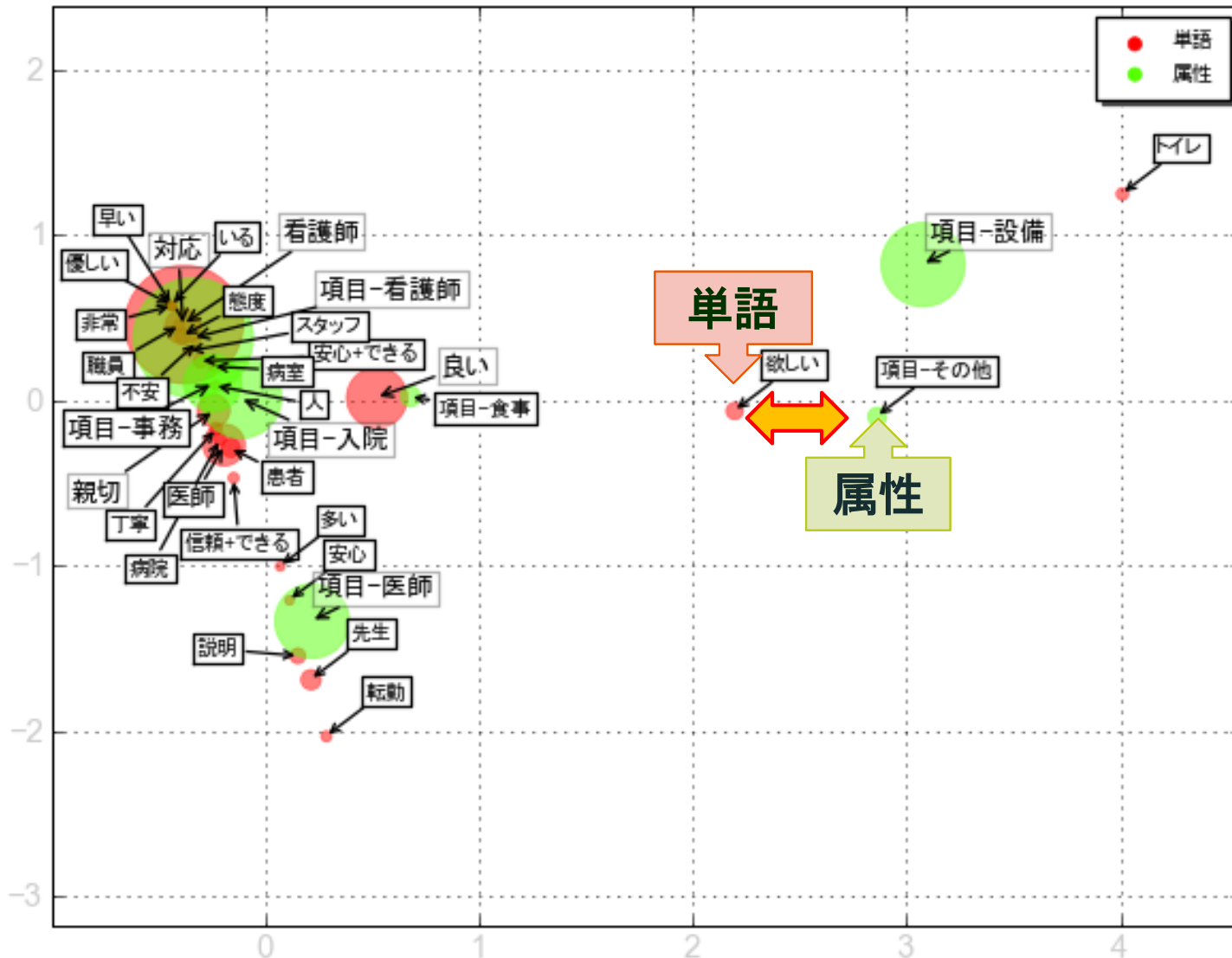
クラスタ数変更



対応バブル分析

- 対応(コレスポンディング)分析をもとにした分析手法。文中の単語と属性との関係を、単語頻度をもとに2次元のバブルで視覚化する。
- 単語の使われ方が似ている属性は近くに表示される。
- バブルの大きさは単語や属性の頻度に対応する。

対応バブル分析



注意して分析したい点

- 今まで困難であったテキストデータの分析がパソコンで容易にできる夢のツール・・・ではない
- コンピュータにすべてを任せるのではなく、研究者の目によるアナログな確認が重要
- 結果から何を読み取るかを常に考えて分析しないとことばの確率論に終始する

最後に・・・

- 壊れることはほぼありません！
 - まずは、触りながらいろいろな分析を試してみましょう!
- 分析のOKボタンを押す時に、自分で問いかけましょう。この分析は何の目的で、何を知りたくてするのか。ただ押すだけでは、カチカチクリックの中身のない結果になってしまいます。
- 結果を読むのは人間です。他の統計分析と同じで、出てきた結果をどう読み取るかは、研究者のセンスや知識、経験が重要になります。