

Level4

文書を眺めてみよう 2。

・仮説

対象となる DataFrame には”love”と”others”の 2 カテゴリーあり、このカテゴリーは、前回の level3 で text の内容から後付けたものである。”love”には恋愛や結婚について書かれた文書、”others”にはそれ以外の文書という基準で分類した。従って、(5)のカテゴリー別の単語の出現数をカウントするという問題で、”love”というクラスにおいて、”恋愛”、”結婚”、”男性”、”デート”が上位にくると考える。また、(6)の scattertext では、2 カテゴリーとも独女について書かれた文書であるため、両文書で出現回数の多い単語(Frequency)は”女性”ではないかと考える。

(5) カテゴリー別に単語の出現回数をカウントし、積み上げ棒グラフにより描画せよ。グラフから分かることを述べよ。

カテゴリーの種類とカテゴリカル文書の数

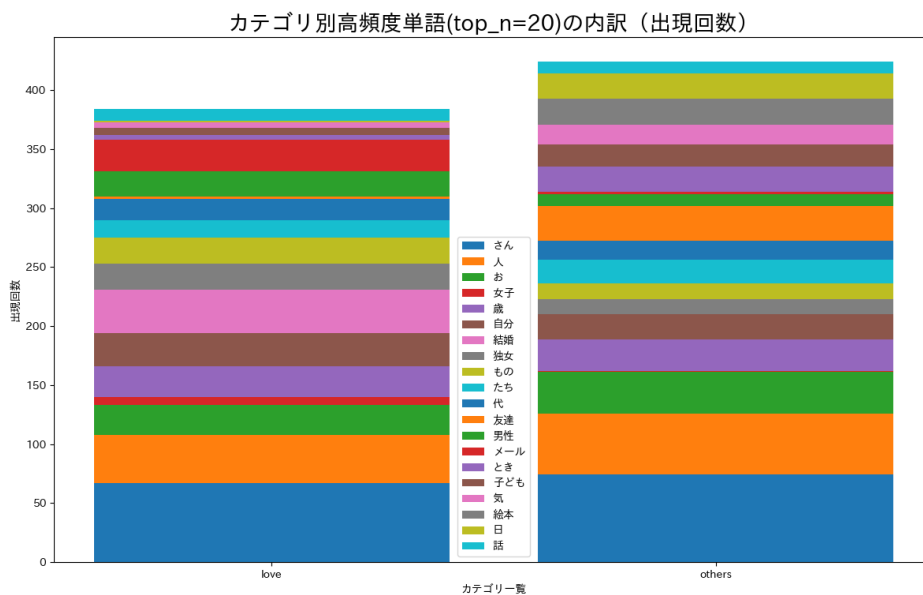
```
love      10
others    10
Name: class, dtype: int64
```

カテゴリ数は 2 つで、”love”, ”others”となっている。それぞれのカテゴリカルな文書のは数は 10 文ずつである。

品詞別カウント出力結果

```
pos = PROPN
[('お金', 12), ('日本', 10), ('オフィスエムツー', 8), ('紀世子', 6), ('セツコ', 6), ('エリ', 5), ('佐竹', 5), ('麻巳子', 5), ('カオリ', 5), ('フランス', 4)]
pos = NOUN
[('さん', 141), ('人', 93), ('お', 59), ('女子', 54), ('歳', 53), ('自分', 49), ('結婚', 37), ('独女', 35), ('もの', 35), ('たち', 35)]
pos = VERB
[('いう', 143), ('言う', 31), ('しまう', 30), ('くる', 25), ('行く', 25), ('ある', 23), ('聞く', 22), ('できる', 22), ('くれる', 20), ('とる', 20)]
pos = ADJ
[('ない', 60), ('多い', 45), ('いい', 21), ('少ない', 16), ('同じ', 9), ('若い', 9), ('好き', 8), ('楽しい', 7), ('大切', 7), ('良い', 7)]
pos = ADV
[('どう', 21), ('すぐ', 10), ('たくさん', 10), ('ちょっと', 10), ('もちろん', 9), ('もう', 8), ('そう', 8), ('もし', 7), ('少し', 6), ('とても', 6)]
```

積み上げ棒グラフ



グラフからわかること

“love”, “others”の2つの積み上げ棒グラフを見て、Y座標が0に近い単語(下に積み上がっている単語)ほど共通している word が多く、Y座標が400に近い単語(上に積み上がっている単語)ほど、それぞれ特有の単語が出てきていることがわかる。“love”において、仮説で高頻出単語に出てくるであろうと考えた、“恋愛”、“結婚”、“男性”、“デート”のうち“結婚”、“男性”の2単語が上位20単語の中に入っていた。また、独身女性の略語である、“独女”や“女子”が“others”と比べて多いことがわかった。

Scattertext 出力結果

