

プロジェクト 最終発表

目標

- Word2vecとgoogle検索からの連想語から作成されたネットワークの比較を行う。
- 二種類のネットワークと単語ごとの関係性に対し、理解を深める。

ネットワーク（グラフ）作成に使用したライブラリ

- Matplotlib

ネットワークの大きさの指定と、保存、表示に使用

- NetworkX

ネットワークにノード、エッジを追加するために使用。

ここで、出来るだけネットワークが重複しないように、下記を設定

```
pos = nx.spring_layout(G,k=1.5,iterations=10000)
```

引数kは大きくなるほど、ノード間の距離が広がる。

Iterationsは最適化を行う反復回数を指定。多いほど、レイアウトの精度が向上し、計算時間も増加する。（参考資料1）

Google検索 知識抽出

- 検索されるページの数を増やす。
上位10ページを対象にする
- 対象にするページを増やすと、スクレイピング制限にかかる時間が早まる。そのため、timeモジュールを用いて実行間時間を増やし、制限がかかるまでの時間を極力伸ばす。

```
import time
```

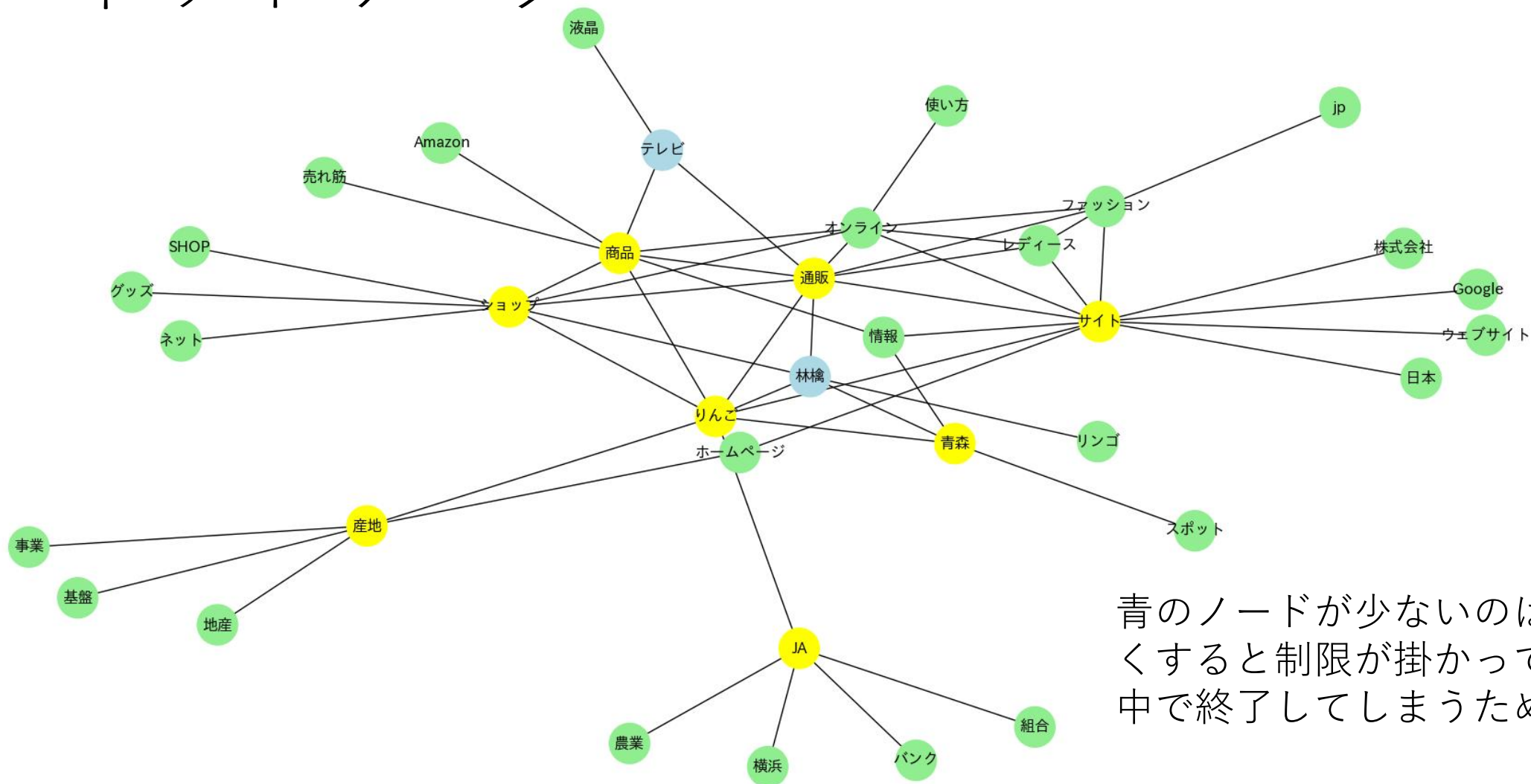
```
time.sleep(1) # 1秒停止
```

- リストに含まれる単語の数を増やし、足切りを行うことで関連性が低い単語を減らす。

```
c = collections.Counter(saihin)
```

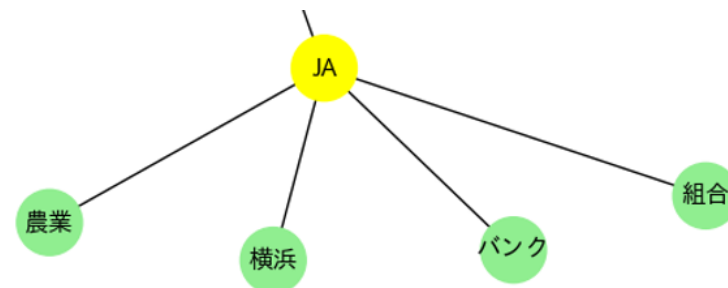
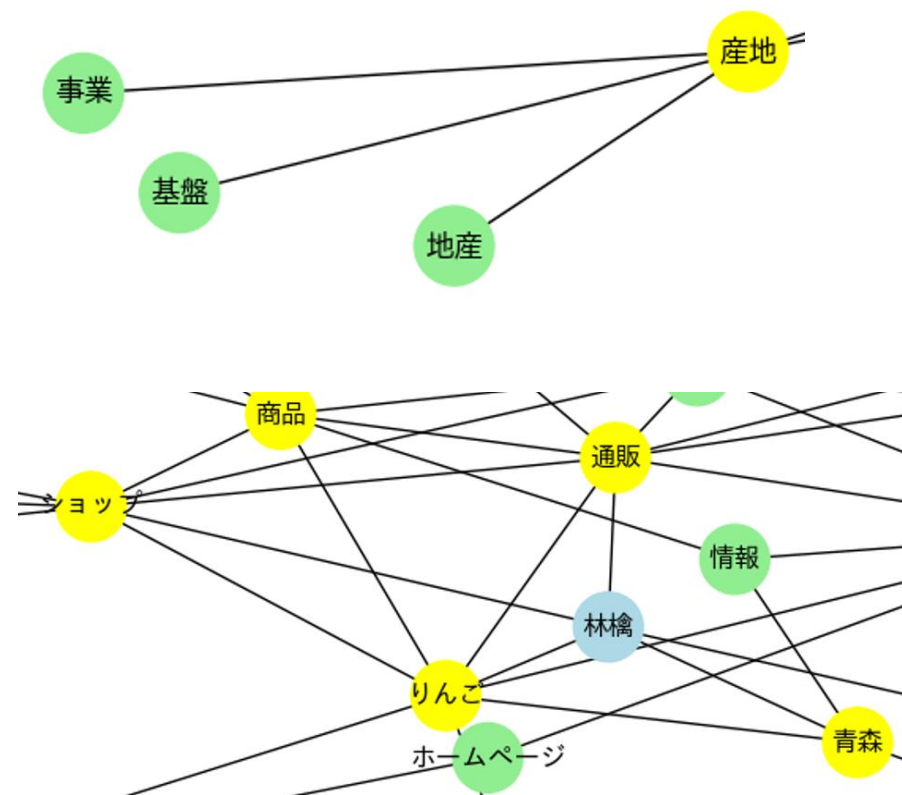
```
a = [i[0] for i in c.items() if i[1] >= 5] #5回以上出現のみ許容
```

Google検索での関連語句から作成したネットワーク



青のノードが少ないのは、検索を深くすると制限が掛かってしまい、途中で終了してしまうため。

Google検索での関連語句から作成したネットワーク



黄色のノードが最初に関連度の高い単語の集合。
青のノードは黄色のエッジから派生された単語の集合。
緑のノードは青のエッジから派生された単語の集合

Word2vec 知識抽出

- 日本語Wikipediaエンティティベクトルのバイナリ版を使用
txt版より動作が軽く、実行時間が短くなる。

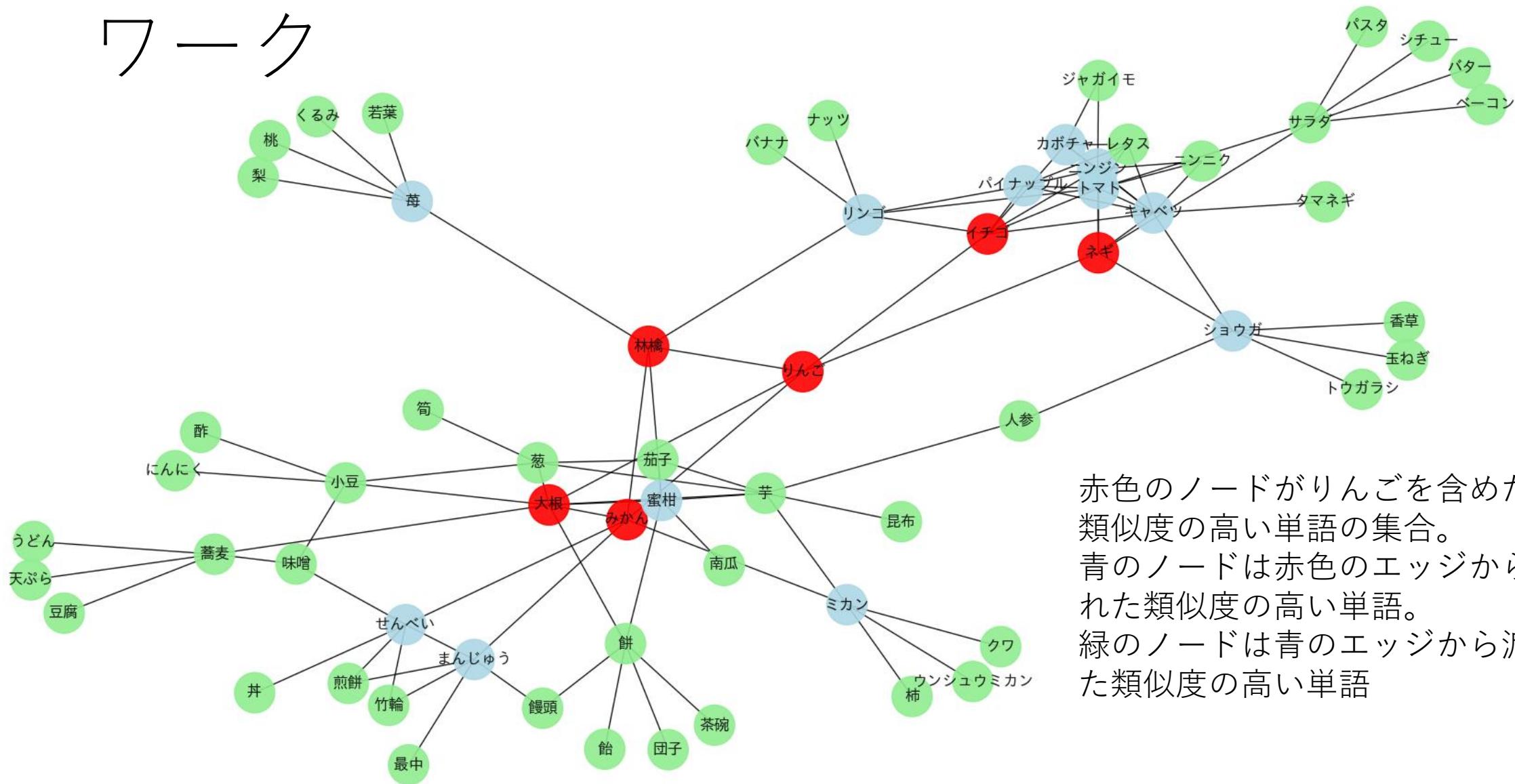
- 主な使用したライブラリ

gensim

KeyedVectorsを用いてバイナリファイルからモデルを作成。

類似度が高い単語を繰り返し調べ、ネットワークを作成する。

Word2vecでの類似度で作成したネットワーク



赤色のノードがりんごを含めた最初に類似度の高い単語の集合。
青のノードは赤色のエッジから派生された類似度の高い単語。
緑のノードは青のエッジから派生された類似度の高い単語

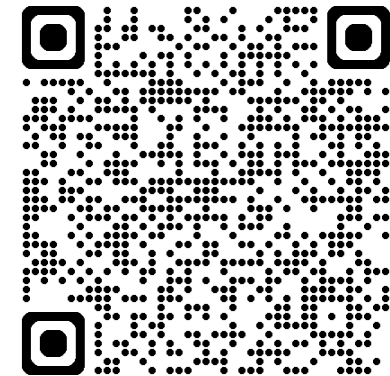
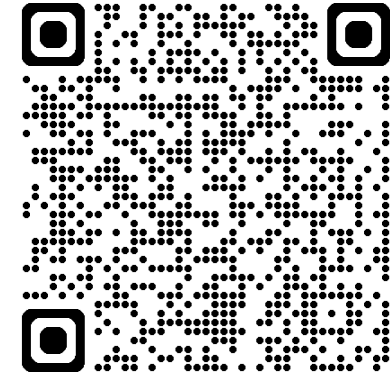
結果

- 人間が連想する単語はgoogle検索寄りだと考える。
- Word2vecではりんごから野菜や果物、食べ物に繋がる
これは、同じような文面の時にその単語の代わりとなる単語がそのようなものが多いためだと考える。
- 逆に、google検索は人間が実際に検索し、その検索結果でのクリック数により表示場所が変動するため、人間が想定する単語がネットワークに表示されやすくなる。

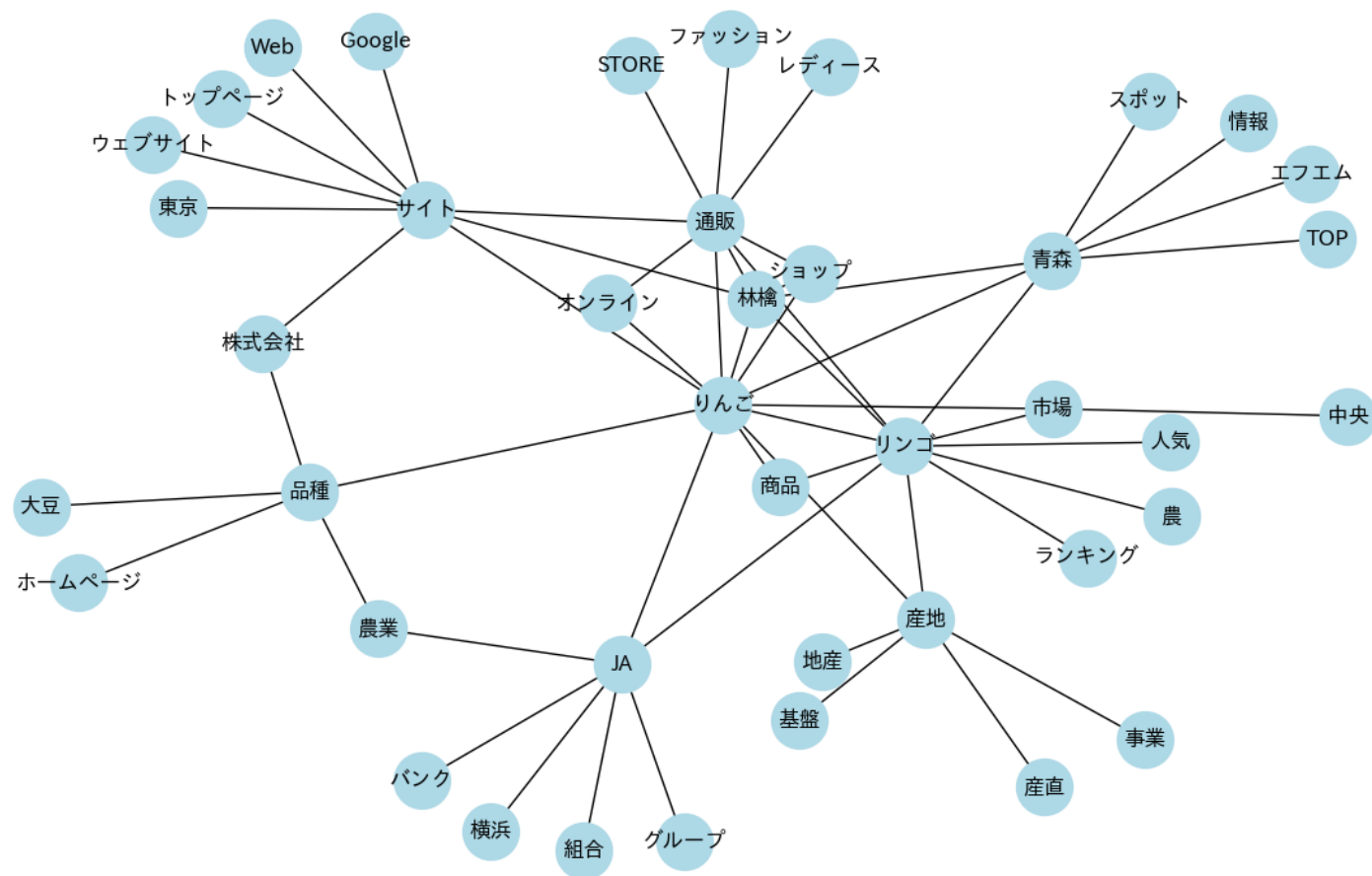
他にも複数の要因は存在する（参考資料2）

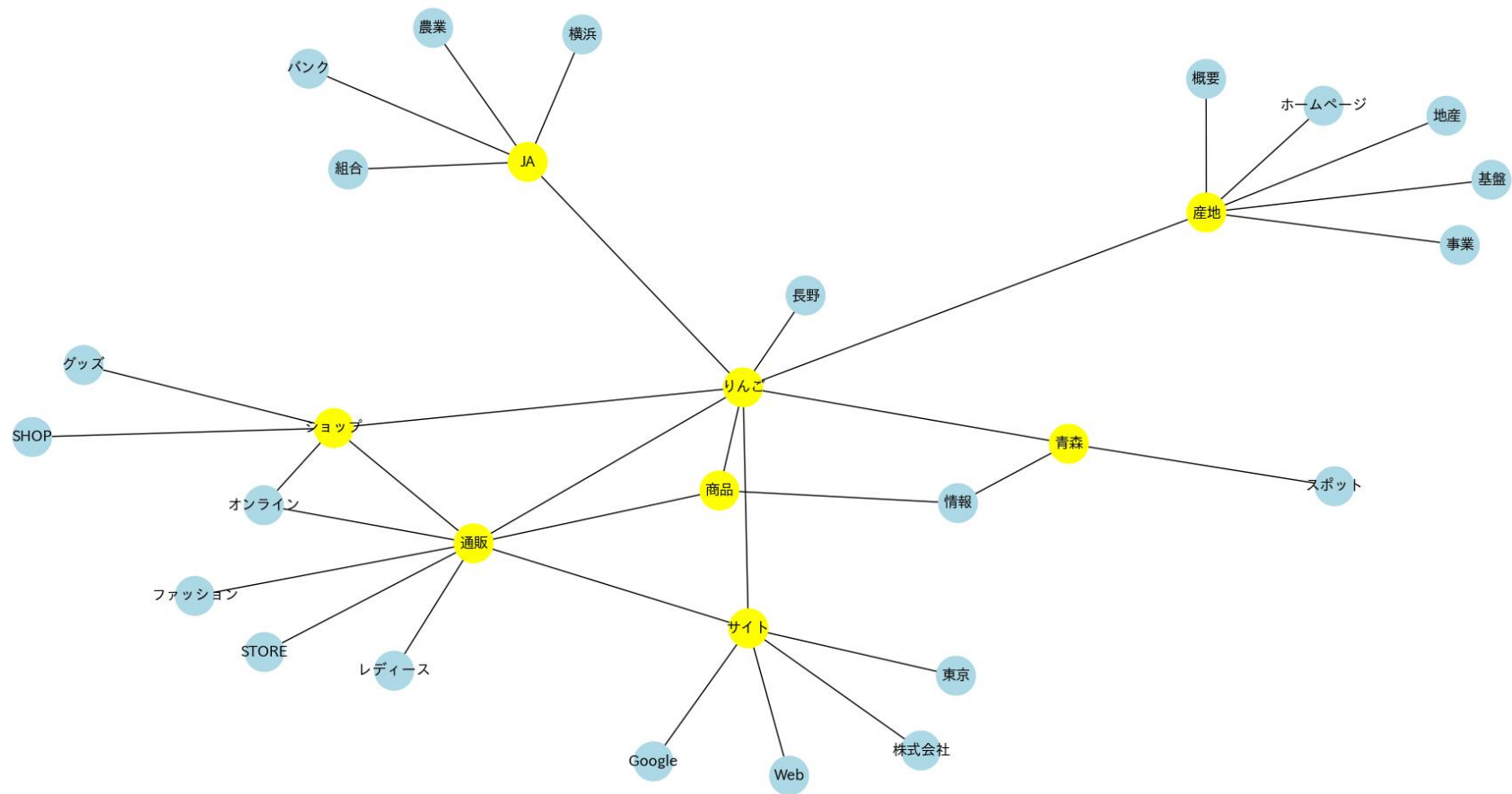
参考資料

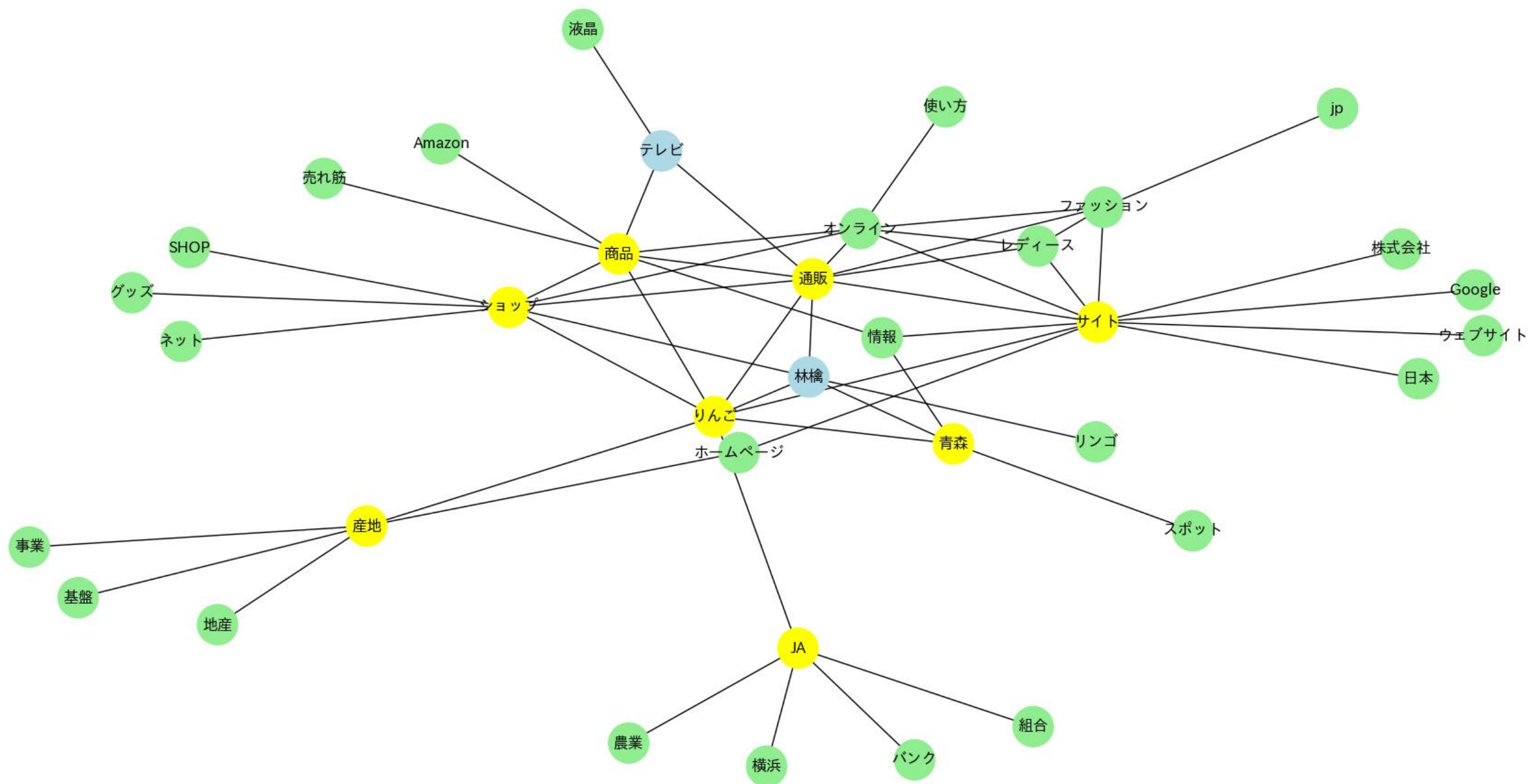
- 1 NetworkX Developers, [spring_layout](#)
[spring_layout — NetworkX 3.3 documentation](#)
- 2 Google, 結果を自動的に生成する仕組み
[ランキング結果 – Google 検索の仕組み](#)



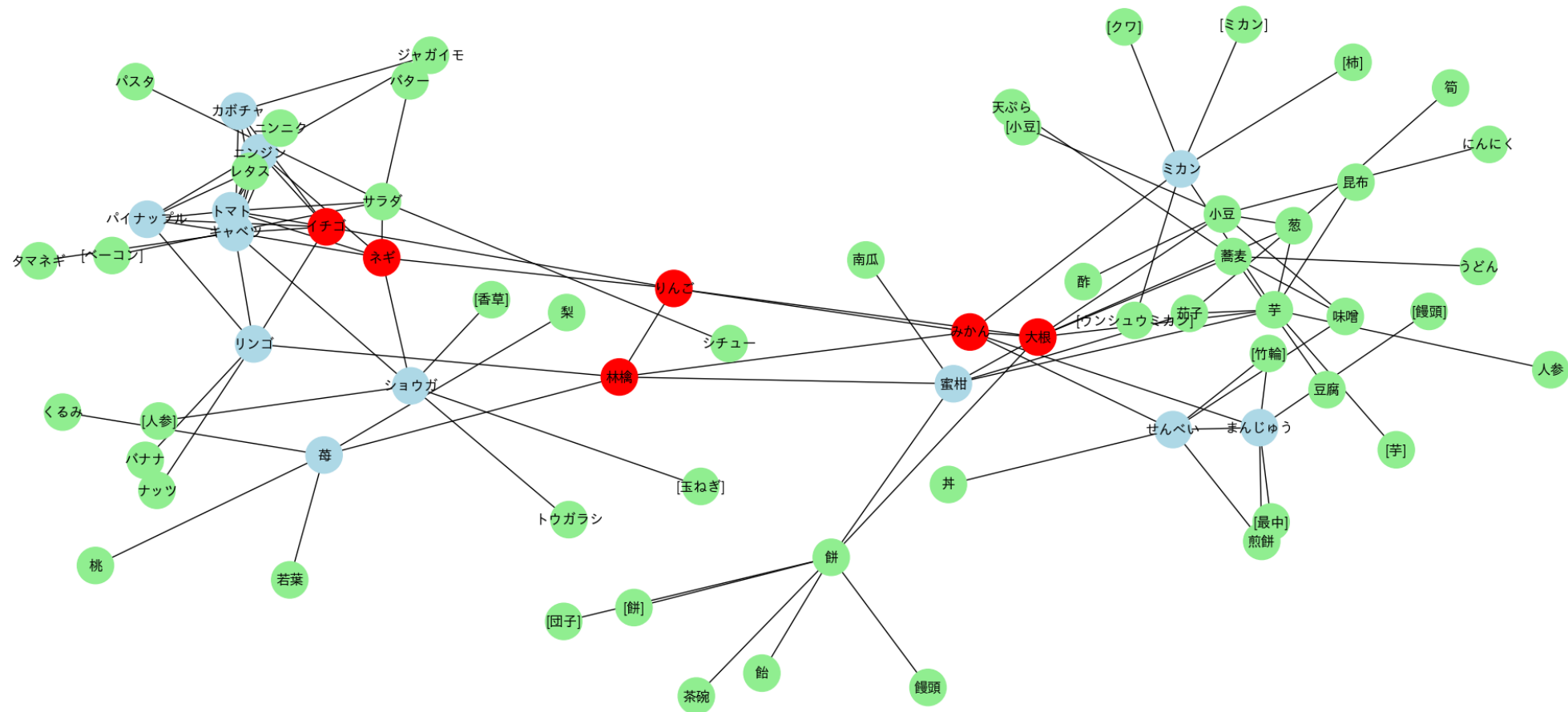
Google検索で単語の関係の表示



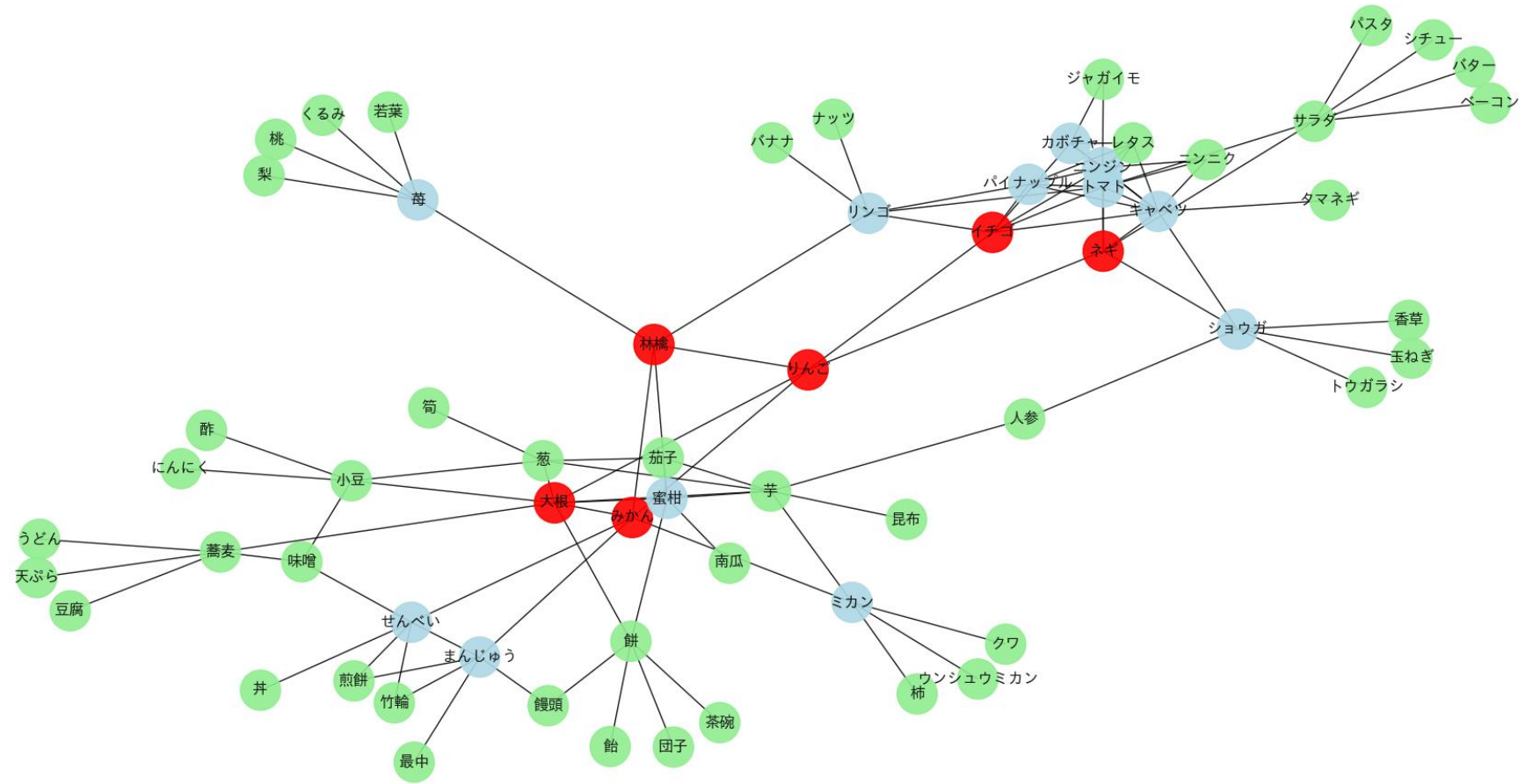




word2vec







最新 google

