



IPランドスケープによる経営戦略支援のための 共起語ネットワーク作成

2020032 平井遙斗

情報システム工学科 指導教員 奥原浩之

要約

近年人工知能(AI)などの情報処理技術は目覚ましく発展しており、記憶や認識、データ収集やその処理などは人間を凌駕しつつある。特許情報は過去の情報をアーカイブしたいわけは発明の保管庫的なデータであり、それを活用することで経営戦略・技術的発展等広く社会に役立つことができる。しかし、現状の特許プラットフォームは人手で少数の特許事例を調べるのには必要充分であるが、ビッグデータとして特許全体の分析を行いたい場合には整理されているとはいいがたい。また、コロナ化をはじめ、サステナビリティやESG(環境、会社、ガバナンス)の推進など、昨今の社会変化に手に対応していくためには多面的な視点から経営戦略を策定することが不可欠である。そこには攻めの知財情報活用であるIPランドスケープが有効である。

キーワード：自然言語処理、知的財産戦略、特許情報処理、テキストマイニング

1 はじめに

目的に作成された世界共通の特許分類です。特許文献の「Int.Cl.」の項に記載されています。2023年9月現在、IPC第8版(2006年1月発効)が最新の分類となっていますが、技術の進展に柔軟に対応するため、適宜改正が行われています。

我が国では独自の特許分類FI(File Index)も導入している。FIはIPCを細分化したものであり、わが国で特許出願が多い分野の細分化がIPCではなく、IPCで特許検索を行っても十分な絞り込みができることが多い。そのため、特許庁はIPCを我が国の技術に合わせて再展開したFIを用意している。

Fターム(File Forming Term)は、特許審査のための先行技術調査を迅速に行うために開発された検索インデックスで、関連先行技術を効率的に絞り込むことを目指している

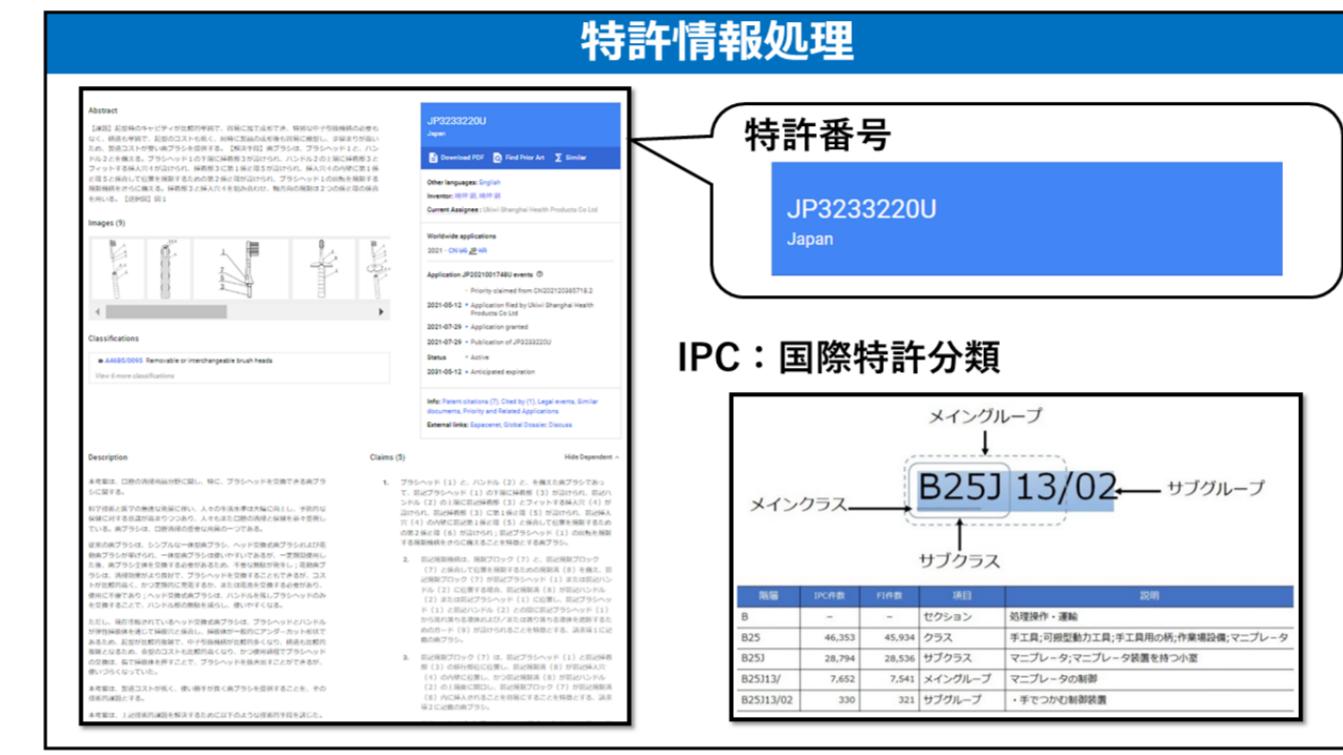


図1 多目的最適化の定式化

2.3 自然言語処理

自然言語処理とは、人が書いたり話したりする言葉をコンピュータで処理する技術です。人工知能(AI)の研究分野で中核を成す要素技術の一つといえます。自然言語処理技術は「言語理解」と「言語生成」に大きく二つに分けることができます。「言語理解」は人が書いた文章に対してなんらかの処理をする技術で、メールの自動分類、ウェブ検索などが典型的な応用になります。「言語生成」は、コンピュータに文章を生成させる技術で、文章の要約や機械翻訳などを含みます。これまで個別の用途ごとに技術開発が進んできましたが、ChatGPTをはじめとする最近のシステムがこの常識を変えました。高度な「言語理解」と「言語生成」が必要な質問応答もできるようになり、さまざまな作業(タスク)を一つのシステムでこなせるようになっています。実はChatGPTが世間を騒がせる前から、自然言語処理の研究者の間では技術の急速な進歩に驚きの声が上がっていました。2018年にGoogleが発表した「BERT(バート)」というシステムでは、開発者が少し手を加えるだけでさまざまなタスクに使えるようになりました。それだけでも驚きましたが、ChatGPTの前身である「GPT-2」や「GPT-3」では、システム自体を変更せずとも、人がシステムにあわせて入力を工夫するだけで多様なタスクの実行が可能になりました。さらに、ChatGPTは人との対話能力が強化され、人間が人間に頼むような言葉で指示をすることで、さまざまなタスクができるようになりました。形態素解析とは文を形態素ごとに分解する技術である[3]。自然言語処理の一つでテキストを品詞ごとに分解することである。形態素解析を行うことで取り出したい品詞を絞って分析できる。一般的に助詞や助動詞はよく使われるが、キーワードごとのテキストを分析したときに特徴が見られないと考える。

—3 特許情報の可視化—

3.1 特許情報のベクトル化

テキストデータは、「定性データ」の代表的なもので、この「定性データ」から付加価値の高い情報を収集することがテキストマイニングの目的である。アイデア発想において人間は自然言語から思考して発想することが一般的である。そこでサイバー空間にあるテキストデータを自然言語処理することを考える。現代社会においてインターネット上の情報量は莫大になっており、今後も増え続けることが予想される。このインターネット上の情報を収集して分析することで発想支援に生かせると考える。発想支援において重要なことはキーワードからより関連度の高い単語をより多く表示されることである。そこで、より良いデータを多く収集するためにインターネットからテキストデータを収集することとする。

今回、GooglePatents複数キーワードのand検索の結果を年代ごとに取得し、そのURLからテキストを抽出し、そのテキストに対し自然言語処理を行う。

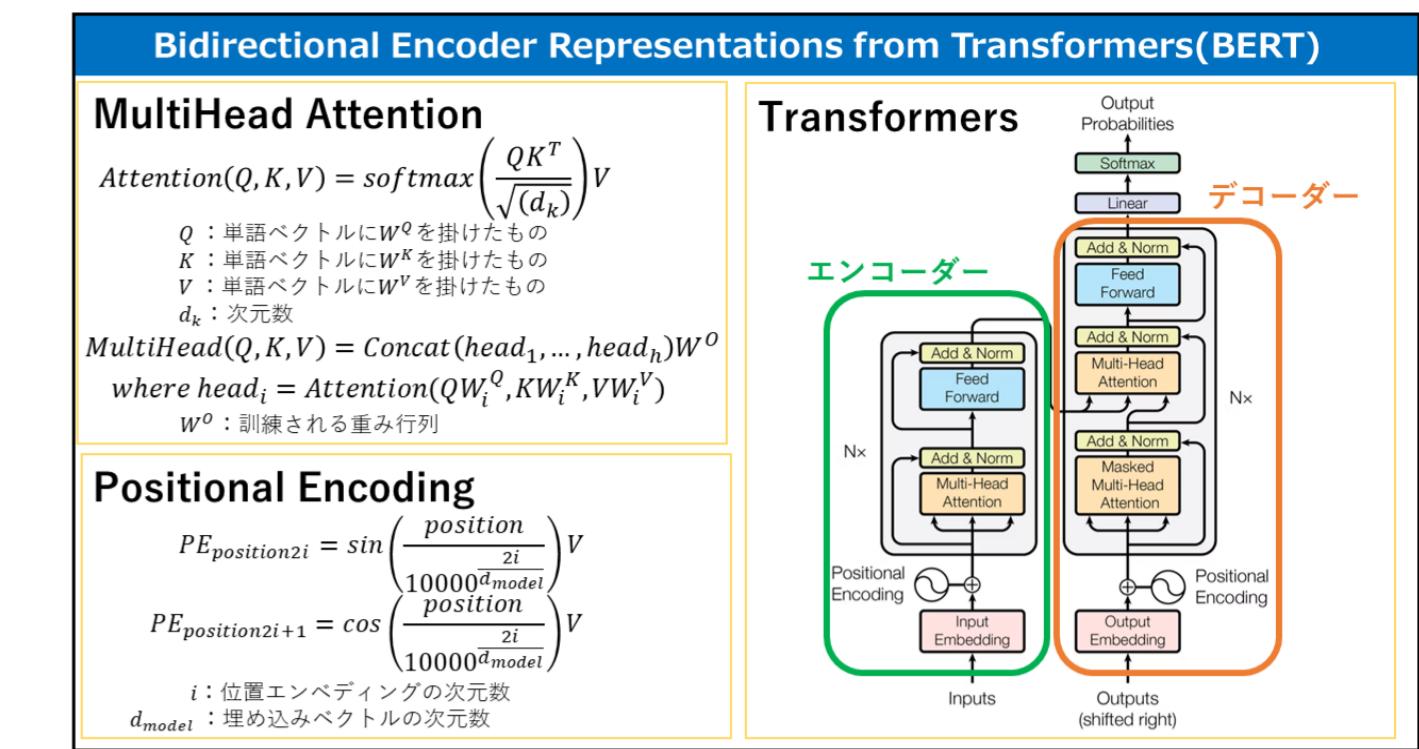


図2 BERT

3.2 次元削減とクラスタリング

次元削減技術は、データ解析の歴史と共に進化してきました。最も古典的な手法であるPCAは、データの分散が最大となる方向に線形変換を行うことで、次元を削減します。これにより、データの本質的な構造を抽出することができますが、非線形な関係性を捉えるのには限界がありました。

その後、2008年に提案されたt-SNEは、データ間の距離を確率分布に変換し、高次元空間と低次元空間の確率分布の違いを最小化することで次元削減を行います。この手法は、特にクラスタリングや可視化に優れた結果をもたらすことで広く使われるようになりました。ただし、計算コストが高いため、大規模データに対しては適用が難しいという課題がありました。

最新の次元削減手法であるUMAPは、t-SNEの持つ計算コストの高さを改善しつつ、データの局所的な構造と大域的な構造の両方を保持することができる特徴を持っています。これにより、より高速で精度の高い次元削減が実現され、機械学習やデータ解析の分野で急速に広まっている現在、UMAPは様々なアプリケーションで活用されています。

クラスタリング(clustering)とは、機械学習の1種でデータ間の類似度にもとづいて、データをグループ分ける手法である。この単語は機械学習や統計学の文脈以外でも使われるクラスタリングにはハードとソフトの2種類がある。それぞれのデータが单一のグループに所属するようにグルーピングするものをハードクラスタリング、それぞれが複数のグループに所属できることを許してグルーピングするものをソフトクラスタリングという。

「群平均法」は、2つのクラスターに属している対象の間のすべての組み合わせの距離を求めて、それらの平均値をクラスター間の距離として定める手法です。群平均法は鎖効果を防止できるメリットがあるためウォード法を実行した時に起こってしまう鎖効果(1つのクラスターに対象が1つずつ吸収されていき、新しいクラスターが作られる現象)を未然に防ぐことができます。

「ウォード法」は、凝集型のクラスター分析の手法の1つで「凝集型階層的クラスタリング」とも呼ばれています。ウォード法はすでにあるクラスターの中で、1番距離の近い2つのクラスターが選ばれ、1つのクラスターに結合されていく操作を、目標のクラスター数になるまで続ける方法です。

最短距離法は単連結法とも呼ばれる、2つのクラスター間で一番近いデータ同士の距離を、クラスター間の距離として採用する手法です。群平均法と同様に、クラスターを構成する要素同士の距離をすべて求め、その中で一番距離の短い組み合わせを選ぶことでクラスター間の距離として求めます。この方法のメリットはウォード法などと比較した場合に、計算量が少くなりますが、同時に外れ値に弱いというデメリットも抱えています。

最長距離法とは、最短距離法とは逆の方法で行う計算手法です。完全連結法と呼ばれることもあります。クラスターを構成している要素同士のすべての距離の中で、最も距離が長いものをクラスター間の距離として採用するという手法です。

非階層クラスタリングは、階層を作らずにデータをグループ化していく手法です。母集団の中で近いデータを収集し、指定された数のクラスターに分類します。この方法では階層クラスタリングとは対照的に、クラスターを形成した後で自由にクラスターを分けることができないため、事前にクラスター数を指定する必要があります。

k-means法とは、非階層クラスタリングを行うためのアルゴリズムのことです。「指定されたk個のクラスターに、平均(means)を用いて分類していく」という意味が込められています。そんなk-means法は、初めに指定したクラスターの数だけ「重心」をランダムに指定して、その重心をもとにクラスターをグルーピングしていくという手法です。k-means法を活用すれば、データ間の距離を計算する必要がなくなるというメリットがあります。

Uniform Manifold Approximation and Projection(UMAP)	
重み付きk近傍(高次元)	トポロジカル表現の最適化
$v_{j i} = \exp \left[\frac{-(r_{ij} - \rho_i)}{\sigma_i} \right]$ $\rho_i = \min_{j \in K_i} \{r_{ij}\}$ r_{ij} : 点 x_i と x_j のユークリッド距離 K_i : 点 x_i に対して、 k 近傍の集合 σ_i : 点の疎密に対応するための変数	$L = \sum_{i,j} \left[v_{ij} \log \frac{v_{ij}}{w_{ij}} + (1 - v_{ij}) \log \frac{1 - v_{ij}}{1 - w_{ij}} \right]$
対称化 $v_{ij} = (v_{j i} + v_{i j}) / 2$	確率的勾配法
重み付きk近傍(低次元) $w_{j i} = \exp(-\max\{0, d_{ij} - \rho'\}) = \tilde{w}_{ij}$ $w_{ij} = \frac{1}{1 + a \cdot d_{ij}^{2-b}}$	$\frac{\partial L^+}{\partial y_i} = \frac{\partial L^+}{\partial d^2} \frac{\partial d^2}{\partial y_i} = \frac{-2abd_{ij}^{2(b-1)}}{1+ad_{ij}^{2b}} v_{ij} (y_i - y_j)$ $\frac{\partial L^-}{\partial y_i} = \frac{\partial L^-}{\partial d^2} \frac{\partial d^2}{\partial y_i} = \frac{2b}{(e+ad_{ij})(1+ad_{ij})} (1 - v_{ij}) (y_i - y_j)$

図1 UMAP

3.3 共起関係と共起ネットワーク

ある単語とある単語が同時に出現することを共起するといい、文章において関係深い単語は共起することが多い。共起分析では単語同士の Jaccard 係数を比較したり、共起関係を持つ単語と単語を線で結んで描かれる共起ネットワークが利用される。文章または単語群に対して共起する単語をネットワークで表した共起ネットワークという。今回、キーワードごとに集めたテキストに対してそれぞれの共起ネットワークを作成する。

4 提案手法



図4 自動献立作成の流れ

5 数値実験並びに考察



図5 実験結果

6 おわりに

参考文献

- [1]
- [2]
- [3]
- [4]
- [5]