

## 第 11 節 トピックモデルとベイジアンネットワークを応用した特許文書分析

(株) アナリティクスデザインラボ 野守 耕爾

はじめに

従来、特許分析というとパテントマップ<sup>1)</sup>と呼ばれる手法が代表的であり、特許調査において長く親しまれてきた。パテントマップは主に発願人や発願年、特許分類 (IPC, FI, F タームなど) を軸にして特許件数を集計し、技術の実態や傾向を可視化したものである。また、テキストマイニングを適用することで、特許の要約文や請求項、明細書といった文書情報も活用してパテントマップを形成することもある。テキストマイニングは非構造化データであるテキストデータを統計的に分析可能な形にする自然言語処理技術であり、文書情報に含まれる単語を抽出し、その出現頻度をカウントする。単語の出現頻度のデータに基づいて、頻出単語をランキングしたり、単語間の出現関係をネットワーク図やマップ図で可視化することで、テキストデータの全体像を単語ベースに把握する。特許文書にテキストマイニングを適用した分析には、たとえば、特許文書の記述内容の類似性に基づいたポジショニングマップを作成し、各特許の位置づけを把握する分析がある<sup>2)</sup>。また、課題と解決手段に相当するキーワード群から分類軸を設定し、そのクロス集計を取ることで課題解決マトリクスを作成し、用途と技術の関係性を把握する分析などもある<sup>3)</sup>。

近年では自然言語処理に AI を適用したテキストデータの分析が活発となっている。自然言語処理とは、人間が使う言葉を機械で処理できるようにする技術の集合体であるが、そのためにはその言葉を機械が認識できる数値のベクトルで表現する必要がある。テキストマイニングはその最も単純な方法を採用しており、単語の出現頻度という数値によって文書をベクトル表現している。近年における AI を活用した自然言語処理の取り組みとは、このベクトル表現を AI によってより高度で表現力の高いものを獲得し、より精度の高い分析を実現するというものである。AI と言ってもその技術領域の範囲は広く、さまざまな手法が存在するが、昨今の第 3 次 AI ブームに火をつけた深層学習 (ディープラーニング)<sup>4)</sup> がその代表格である。

深層学習を応用した自然言語処理は、主に 2010 年代以降に数多くの技術が登場した。ここではその技術の流れを簡単に紹介する。まず 2013 年に発表された word2vec<sup>5)</sup> は、大量のコーパスの穴埋め問題を深層学習で解くことで、単語の意味や類似性を捉えたベクトル表現を獲得した。この技術を文書全体に適用して文書のベクトル表現を獲得する doc2vec<sup>6)</sup> も提案された。word2vec は単語の順序情報がなく、文書の文脈を捉えられなかったが、文書を単語が並ぶ系列データと考え、系列データを対象とした深層学習モデルである RNN<sup>7)</sup> や、より長い系列を扱える LSTM<sup>8)</sup> を適用することで、文脈のベクトル表現の獲得が進められた。そして、その系列データを別の系列データに変換するエンコーダ-デコーダモデルと呼ばれる seq2seq<sup>9)</sup> が 2014 年に発表され、文脈を捉えた機械翻訳などに応用された。また、seq2seq において、注目すべき単語の重みを評価する Attention<sup>10)</sup> という仕組みが提案され、より効果的に文脈を捉えることができ精度が向上した。そして RNN などを使わず Attention だけで seq2seq を実現する Transformer<sup>11)</sup> という技術が 2017 年に発表され、精度と計算効率が格段に上がり、自然言語処理領域において大きなブレイクスルーとなった。その後は Transformer をベースに膨大なテキストデータを事前学習させることで汎用的な言語特徴を獲得した大規模言語モデルの開発が盛んになった。その代表として、文脈理解に優れた BERT<sup>12)</sup> や、文章生成に優れた GPT<sup>13)</sup> などが 2018 年に発表された。特に GPT を開発した OpenAI は 2022 年 11 月に AI チャットボットの ChatGPT を発表し、その賢すぎる AI に世界は騒然とし、生成 AI ブームが勃発した。

特許文書の分析においても、こうした深層学習による自然言語処理技術が特許の分類や類似性評価などに応用されている。たとえば、doc2vec を用いて文書をベクトル化したデータを用いて特許分類を検討したり<sup>14)</sup>、self-attention 機能を有する LSTM モデルで文書をベクトル化したデータで技術俯瞰マップを作成したり<sup>15)</sup>、BERT で文書をベクトル化したデータで二つの請求項の同一性を判定したり<sup>16)</sup>、ChatGPT で特許の要約から課題と解決手段の要約と分類を実現する取り組み<sup>17)</sup> などがある。また、特許庁では大量の特許公報で事前学習させた特許事前学習モデルという BERT モデルを自ら構築し、その有用性を評価する実証事業を 2022 年に実施している<sup>18)</sup>。

こうした深層学習を応用した自然言語処理技術は、文脈を捉えた表現力の高いベクトルを獲得することで、精度の高いテキストデータの分析ができることが特長である。ただし、そのベクトルとは大量の数字の羅列であり、それ自体の意味については解釈できない。また、深層学習はモデルの構造が複雑で、なぜそのような結果になったのかというプロセスはブラックボックスになってしまう。特に大規模言語モデルは、事前学習されたモデルの推論に基づくアウトプットとなるため、その信憑性には注意が必要であり、GPTなどの生成系AIでは、もっともらしいが誤っている情報を生成してしまうハルシネーション（幻覚）と呼ばれる問題がある。

一方、古典的な自然言語処理技術であるテキストマイニングは、単語の出現頻度のみでベクトル表現をしており、複雑な文脈は捉えられないが、そのベクトル自体の意味はどんな単語が何件出現したかという情報であり、シンプルで分かりやすい表現形式となっている。また、その単語頻度は推論ではなく事実であり、入力テキストデータには確実にその頻度の単語が出現しており、データに即した正しいアウトプットがされている。つまり、より最新のAIを使った分析の方が優れているということではなく、それぞれの技術の特徴を理解し、目的に応じて使い分けことが重要である。

特許文書分析のアプローチにおいても同様である。深層学習のAIを使った分析が活発化するなかで、従来のパテントマップのアプローチが陳腐化しているかということ、そうではない。そもそもパテントマップは分析者が母集団の特許に対して技術の全体像を把握するための手法であり、出願人や出願年、分類コード、キーワードなど、さまざまな軸から特許の傾向を可視化することで、関連技術の実態を人間が詳しく探索することを目的としている。可視化の仕方こそ単純ではあるものの、それ故にとってもシンプルで人間が容易に理解しやすく、誰もがその分かりやすい可視化結果から着想を得ることができる。パテントマップを用いることで、新しい技術の動向や未着手の技術の発見、競合他社との差異や関係性を把握でき、企業の技術戦略を人間が検討する上で有益な気づきを与えてくれる。深層学習のAIを使った分析は、自動化されたプロセスによって即座に結果が出力され、知財業務の効率化という面では優れているといえる。しかし、人間が考える余地は限定的で、プロセスが不明の結果だけでは重要な気づきを見落としてしまう可能性もある。一方、パテントマップはその分析プロセスを通じて、分析者が理解しながら考えるアプローチとなる。人間による納得感のある技術戦略を検討するという目的では、これまで長く親しまれてきたシンプルで分かりやすいパテントマップは大変有用といえる。

ただし、特許の文書情報を用いた従来のパテントマップの分析アプローチは、分析対象となる特許データの量が多くなると対応が難しくなる。データ量が少なければ人手を介した丁寧な処理もされるが、データ量が多くなるとこれを人間が一つひとつ実施することは難しく、分析結果に主観的な偏りも生じてしまう。また、テキストマイニングでは読み込むテキストの量が多くなれば、それだけ大量の単語が抽出されるため、その単語をベースに可視化をすればとても複雑な結果となってしまう、解釈の容易性が損なわれてしまう。

本節では、大量の特許文書データであっても、従来のパテントマップのように分かりやすくシンプルな分析が実行できるように、「トピックモデル」と「ベイジアンネットワーク」という深層学習とは異なる複数のAI技術を応用した分析手法を提案し、その適用事例を紹介する。

## 1. 特許の文書情報を用いた従来のパテントマップ分析とその課題

まずは特許文書にテキストマイニングを適用した従来のパテントマップの分析とその課題を確認する。

### 1.1 テキストマイニングという分析手法

テキストマイニングは、特許の要約や請求項や明細書といった文書情報の定性データを、定量的に統計分析可能にするデータマイニング手法である。大量のテキストデータからその文書に含まれる単語を抽出し、単語単位による集計や統計解析に基づいて、文書に記述されている傾向を可視化し現状把握をする。

テキストマイニングにおける自然言語処理には「形態素解析」と「構文解析」という2つの基本技術がある<sup>19)</sup>。形態素解析とは、文章を単語に分割し、その単語の品詞を割り当てる技術である。正確には単語ではなく、形態素と

呼ばれる意味を持つ最小単位の文字列として分割する手法だが、実用面では形態素ではなく単語の単位で分割されることが通例となっている。単語と形態素の違いの例として、たとえば「消費電力」という単語は「消費」「電力」という2つの形態素で、「電気掃除機」という単語は「電気」「掃除」「機」という3つの形態素で構成される。このように、日本語の意味情報を抽出するには最小単位の形態素よりも単語の方が適しているとされる。一方、構文解析とは、文法規則に基づいて文章の構造を解析し、単語間の関係性を識別する技術である。たとえば、主語と述語や、修飾語と被修飾語といった係り受け関係を抽出する。日本語の場合は単語ではなく文節を単位に係り受け関係を抽出するのが一般的である。文節とは、日本語を意味のわかる単位に区切ったものであり、助詞や助動詞は名詞や動詞とセットで括られる。形態素解析と構文解析の例を図1に示す。これらの処理により、文書に含まれる単語とその品詞、また単語間（文節間）の係り受け関係を抽出することができる。なお、実際に抽出される単語や係り受けは、登録されている辞書によってそれぞれ原形に変換された形で抽出されることが一般的である。形態素解析と構文解析には公開されたフリープログラムがあり、形態素解析ではJUMAN、ChaSen、MeCabなどが、構文解析ではKNPやCaboChaなどが広く知られている<sup>20)</sup>。

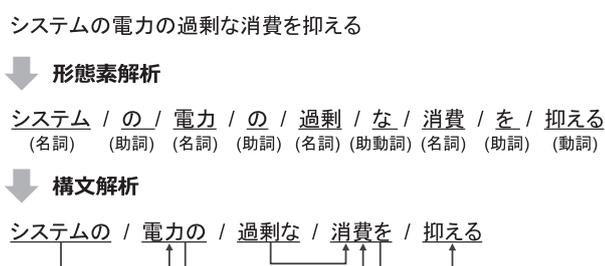


図1 形態素解析と構文解析の例

先述の通り、自然言語処理は文書情報を数値のベクトルで表現する必要がある。テキストマイニングで実行される形態素解析はそのベクトル化の最も単純な方法となる。形態素解析で文書に含まれる単語を抽出し、各文書にどの単語が何回出現したのかをカウントすることで、その出現頻度という数値によって文書一つひとつをベクトルとして表現できる。このように単語の出現頻度でもって文書をベクトル化する手法をBag-of-Words (BoW) と呼ぶ。Bag-of-Wordsのベクトル表現は「文書×単語」という行列形式のデータで考えると分かりやすい。一つの行に一件の文書を、一つの列に一つの単語を取ったとき、各セルの値はその文書内におけるその単語の出現回数を示す。Bag-of-Wordsによるデータ形式の例を図2に示す。これは最もシンプルな文書のベクトル表現だが、この処理だけでもテキストデータのさまざまな定量分析が可能であり、とても強力な処理手段である。テキストマイニングは現在多くのツールが存在するが、そこで実行される可視化や統計解析は基本的にこのデータ形式をベースにしている。

特許ID	要約の課題	空気清浄機	消費電力	抑制	コンパクト	加湿装置	菌	発生	増殖	空気	清浄化	効果的	可能	電気掃除機	除菌	消臭	抗菌	ドライヤー	実現	提供	⋮
1	消費電力量を抑制できる空気清浄機を提供する。	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	
2	コンパクトな空気清浄機を実現する。	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	
3	菌の発生や増殖を抑制することができる加湿装置を提供する。	0	0	1	0	1	1	1	1	0	0	0	0	0	0	0	0	0	0	1	
4	効果的に空気を清浄化することが可能な加湿装置を提供する。	0	0	0	0	1	0	0	0	1	1	1	1	0	0	0	0	0	0	1	
5	効果的に除菌および消臭できる電気掃除機を提供する。	0	0	0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	1	
6	除菌、抗菌及び消臭ができるドライヤーを提供すること。	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	0	1	
...	.....																				

図2 Bag-of-Wordsの例

## 1.2 テキストマイニングを適用したパテントマップの可視化

特許の文書情報を用いたパテントマップの分析目的として、よく取り上げられるものには、(1) 全体像の把握、(2) トレンドの把握、(3) 競合他社の動向の把握、(4) 用途と技術の関係の把握といったものが挙げられる。テキストマイニングの適用によってよくアウトプットされる可視化の例を図3に示しながら、それぞれの分析の概要について以下に述べる。

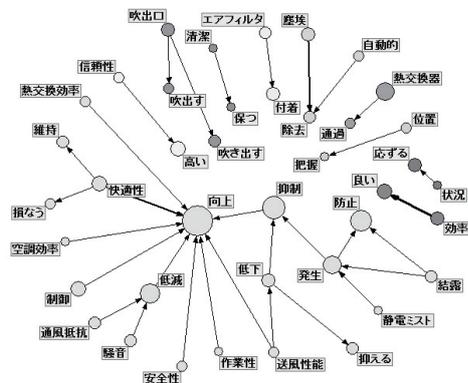
### (1) 全体像の把握

対象となる技術領域の全体像を俯瞰して把握するための分析である。最も基本的なものでは、図3(A)のように特許文書に含まれる単語や文法的な単語のペアとなる係り受けの出現頻度を集計し、頻度の多い言葉から全体像を把握する。また図3(B)のように各単語が同じ特許文書内で出現する共起関係をネットワークで可視化し、その単語のかたまり状況から、形成されている話題について定性的に考察することもある。

### (A) 単語や係り受けの頻度集計



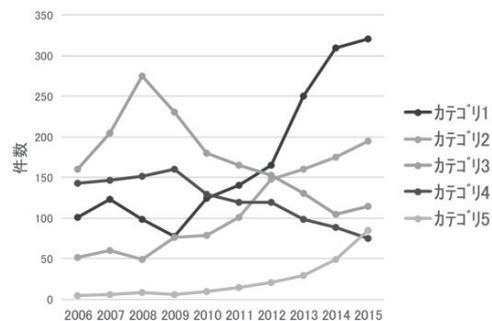
### (B) 単語の共起ネットワーク



### (C) カテゴリリストの作成

掃除機カテゴリのリスト	
掃除機	塵埃⇒分離
集塵	塵埃⇒吸い込む
集塵容器	塵埃⇒収容
集塵室	塵埃⇒遠心分離
吸引力	含塵空気⇒分離
サイクロン	塵埃⇒溜める

### (D) 各カテゴリの出願数の経年変化



### (E) 単語と出願人の対応マップ



### (F) 課題と解決手段のクロス集計

課題	解決手段カテゴリ					
	カテゴリ01	カテゴリ02	カテゴリ03	カテゴリ04	カテゴリ05	カテゴリ06
課題01	206	80	71	184	26	47
課題02	208	76	87	182	23	48
課題03	172	74	53	57	31	35
課題04	176	54	37	59	26	46
課題05	85	39	13	23	14	16
課題06	87	53	31	33	59	37

図3 テキストマイニングを適用したパテントマップの可視化例

## (2) トレンドの把握

今後成長が見込まれる技術や、逆に衰退している技術を把握し、研究開発戦略を検討していくための分析である。抽出した単語の出現頻度を出願年で集計することもあるが、一つひとつの単語で集計すると複雑になりすぎるため、しばしば図3(C)のように複数の単語や係り受けを人がグルーピングして意味性のあるカテゴリを形成し、図3(D)のようにカテゴリ別に該当特許数を出願年で集計し、そのトレンドを把握する。

## (3) 競合他社の動向の把握

他社と差別化するような研究開発戦略を検討したり、自社技術のライセンス先候補や技術提携先候補となるような企業を検討するうえでニーズがある分析である。コレスポネンス分析や数量化Ⅲ類と呼ばれる手法がよく用いられ、図3(E)のように単語と出願人を同じ平面上にマッピングし、出願人と近くに位置する単語から各出願人の技術開発動向を把握する。

## (4) 用途と技術の関係の把握

特に自社技術の新たな用途展開を探索するうえでニーズがある分析である。国内特許の要約文で記述されていることの多い「課題」と「解決手段」という2つの項目に着目し、それぞれの文章をテキストマイニングして図3(C)のようにカテゴリを形成し、図3(F)のように課題のカテゴリと解決手段のカテゴリに該当する特許件数をクロス集計することで、用途と技術の対応関係を把握する。

### 1.3 テキストマイニングを適用した従来の分析の課題

こうした分析は、人間ではなかなか読み切れない特許文書の全体像を把握するうえで有効な手段である一方、①単語をベースにした分析であるため単語数が増えると結果が複雑で考察しにくい、②単語のカテゴリの設定が主観的で作業負荷も大きい、③用途と技術の関係は単純なクロス集計で統計的な関係が分析できていない、といった課題もある。①の課題については、特にビッグデータのような件数がとても多い特許文書を分析対象とする場合、テキストマイニングで抽出される単語も膨大となるため、そうした大量の単語をベースにした可視化結果は複雑で解釈困難となることが多い。②の課題については、そのカテゴリルールが属人的となるため、作業者が変わればルールも変わってしまう曖昧なものであり、知識継承もされにくい。特に分析対象がビッグデータの場合、人間がその結果を整理してカテゴリルールを作成するにはあまりにも負荷が大きい。③の課題については、単純な頻度の大きさでは一見関係がありそうな用途と技術でも、それが統計的に意味のある関係であるとは限らない。つまりその用途や技術に該当する特許の元々の件数が多ければ当然クロス集計の頻度も大きくなるため、単純なクロス集計では関係性の考察を誤る可能性がある。

本節では、上記の①②の課題を解決する方法として、使われ方の似ている単語群をトピックという単位に集約できるトピックモデルというAI技術を紹介する。これにより、図3(C)のように従来人手で作成していた単語のカテゴリを機械的に実行できる。また、特許文書全体の傾向を従来のように大量の単語を単位に可視化するのではなく、いくつかのトピックを単位に可視化することによってシンプルで分かりやすい結果を得ることができる。こうしたトピックはパテントマップにおける新たな探索軸と捉えることもできる。一方、③の課題を解決する方法としては、変数間の要因関係を確率統計的にモデリングできるベイジアンネットワークというAI技術を紹介する。ベイジアンネットワークでは、用途と技術の統計的な関係をモデル化することに適用する。事前にトピックモデルで課題に関する用途のトピックと解決手段に関する技術のトピックを抽出し、それらの間に存在する統計的な関係性をベイジアンネットワークで分析する。次項ではトピックモデルとベイジアンネットワークについて解説する。

## 2. トピックモデル

ここではトピックモデルを適用することの有効性について述べ、トピックモデルの各手法の概要を解説する。その中でもPLSAという手法の有効性を考察し、PLSAの理論の詳細とテキストマイニングへの適用で得られるアウトプットについて説明する。

## 2.1 トピックモデルの適用効果

トピックモデルは主に 1990 年代から 2000 年代初頭にかけて開発が盛んになった自然言語処理技術である。Bag-of-Words による「文書×単語」の行列データをインプットに、教師なし学習で文書に潜むトピックをアウトプットする。そのトピックは文書の集合と単語の集合で構成され、使われ方の似ている単語群は同じトピックに集約されるようになる。トピックモデルを適用することで、単語単位ではなく集約されたトピックを単位に文書全体の傾向を把握することができる。

教師なし学習でデータを機械的に集約する方法としては、同じく Bag-of-Words のデータに対して、階層クラスタ分析の Ward 法や非階層クラスタ分析の k-means 法といった従来のクラスタ分析を適用する方法も考えられる。しかし、以下の点で挙げられるように、トピックモデルは従来のクラスタ分析より Bag-of-Words のデータクラスタリングという観点において優れている。

### (1) 高次元データに対応できる

従来のクラスタ分析では、データ間の距離に基づいて類似度を計算し、距離の近いデータをまとめていく。しかし、変数の数が大量にある高次元データになると、データ間の距離がどれも大きく離れ、妥当な結果が得られにくくなる「次元の呪い」と呼ばれる問題が起きてしまう。Bag-of-Words は大量の単語を列に取った超高次元のデータとなるため、こうした従来のクラスタ分析手法は適していない。これに対してトピックモデルでは、距離に基づかず、高次元の情報をできるだけ保持した形でその次元数を削減し、低次元に変換できる手法である。そのためトピックモデルはしばしばデータの次元削減でも用いられ、次元削減法と呼ばれることもある。トピックモデルを適用すれば、膨大な単語で構成される Bag-of-Words の高次元データでもトピックに集約することができる。

### (2) 行の要素と列の要素を同時にクラスタリングできる

上記のような従来のクラスタ分析手法は、列をベースに行をクラスタリングする、あるいは行をベースに列をクラスタリングするため、どちらか一方のみがクラスタリングの対象となる。トピックモデルでは、トピックは行の要素と列の要素の 2 つの軸に基づいて抽出され、各トピックに対する行の要素と列の要素のそれぞれの関連度が同時に計算される。つまり、Bag-of-Words による「文書×単語」の行列データに適用すれば、一つのトピックに対して文書と単語が同時にクラスタリングされていると解釈できる。このように、抽出されるトピックは行と列の 2 つの軸の情報を持つことができるため、従来のクラスタ分析よりも情報量が多い結果となり解釈がしやすくなる。

### (3) ソフトクラスタリングできる

上記のような従来のクラスタ分析手法はハードクラスタリングとも呼ばれ、ある要素が一つのグループに所属してしまうと他のグループには重複して所属が許されない。一方、トピックモデルはソフトクラスタリングとも呼ばれ、すべての要素がすべてのトピックにまたがって所属し、その関連度合いが計算される。これにより、複数の意味を持つ要素がある場合でも柔軟なクラスタリングを実現できる。

## 2.2 トピックモデルの各手法

トピックモデルの代表的な手法には、「LSA」「NMF」「PLSA」「LDA」などがある。ここでは各手法の概要を解説する。

### 2.2.1 LSA

LSA (Latent Semantic Analysis, 潜在意味解析) は、Bag-of-Words による「文書×単語」の行列を SVD (Singular Value Decomposition, 特異値分解) によって次元削減することでトピックを抽出する手法である。1990 年に発表された手法であるが<sup>21)</sup>、その基礎となる考え方の研究はそれ以前から存在している。情報検索の分野では LSI (Latent Semantic Indexing) とも呼ばれている。

LSA による特異値分解のイメージを図 4 に示す。LSA の具体的な処理は、「文書×単語」行列を①「文書×トピック」(左特異ベクトル)、②「トピック×トピック」(特異値)、③「トピック×単語」(右特異ベクトル) の 3 つの行列に分解することである。特異値「トピック×トピック」は対角行列であり、その対角成分には特異値が大きい順に配列される。左特異ベクトルの「文書×トピック」と右特異ベクトルの「トピック×単語」は、トピックを示すベクトル(左

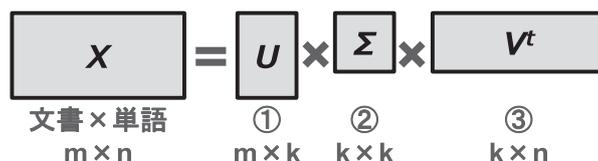
特異ベクトルでは各列, 右特異ベクトルでは各行) が直交しており, 各トピックの軸は数学的に互いに独立している。SVD による数学的な処理は, 元の行列を最もよく近似する行列に分解すること, つまり元の行列と分解後の行列の誤差を最小化する最適化問題を解くことである。特に, SVD はこの最適解の存在が数学的に保証されているという点で強力であり, 計算過程のシンプルさから計算効率も高いメリットがある。

SVD に類似する手法に PCA (Principal Component Analysis, 主成分分析) があり, PCA は SVD を用いた分析手法と捉えることができるが, 両者の関係について簡単に説明する。PCA は元のデータセットの分散を最大化させる方向 (主成分) を見つけ出し, 最も情報量の多い特徴を抽出することを目指している。その主成分の方向に元のデータを射影することで, 元のデータの重要な情報を保持しながら次元を削減する。SVD は分散を最大化させるという目的で用いられるものではないが, PCA を効率的に実行できる手法として用いられるため, PCA は特定の条件下における SVD の応用の一例と見ることができる。つまり, SVD は PCA よりも一般的な手法といえるが, 実用面でいえば両者は同様の手法のように扱われることがある。

SVD は大きな値に引っ張られてトピックが抽出される傾向があるため, LSA を実行する際には, Bag-of-Words に前処理として単語の重要度を評価する TF-IDF など重み付けした行列を用いることが多い。また LSA の課題として, 分解する行列の要素に負の値を許容しているため, 結果の解釈が難しくなることや, 結果が学習データに完全に依存するため, 過学習を起こしやすく, 新しい文書のトピックは推定できないことなどが挙げられる。

### LSAの特異値分解

$$X = U\Sigma V^t$$



### LSAでアウトプットされる各行列の例

①文書×トピック					②トピック×トピック					③トピック×単語					
U	トピック1	トピック2	トピック3	トピック4	Σ	トピック1	トピック2	トピック3	トピック4	Vt	単語1	単語2	単語3	単語4	単語5
文書1	-0.47	-0.75	-0.46	0.11	トピック1	7.7	0	0	0	トピック1	-0.51	-0.49	-0.38	-0.37	-0.47
文書2	-0.61	-0.10	0.69	-0.37	トピック2	0	2.8	0	0	トピック2	-0.27	0.14	0.80	-0.50	-0.09
文書3	-0.27	0.41	-0.54	-0.68	トピック3	0	0	2.0	0	トピック3	0.47	0.43	-0.35	-0.67	-0.14
文書4	-0.58	0.52	-0.10	0.62	トピック4	0	0	0	0.6	トピック4	0.61	-0.74	0.23	-0.18	0.07

図4 LSAの特異値分解

## 2.2.2 NMF

NMF (Non-negative Matrix Factorization, 非負行列因子分解) は, Bag-of-Words の「文書×単語」の行列を, 2つの非負の行列「文書×トピック」と「トピック×単語」に分解することでトピックを抽出する手法で, 1999年に発表された<sup>22)</sup>。NMFによる行列分解のイメージを図5に示す。

計算アルゴリズムは, 元の行列と分解後の行列の積との誤差を最小化することを目的関数に, 初期値を与えた反復計算により最適解を得る。誤差の定義の仕方は, 平方ユークリッド距離や KL ダイバージェンスなどが使われる。反復計算は乗法更新式 (Multiplicative update rule) が代表的であり, これは得られる行列が必ず非負であるという制約条件下で適用される勾配降下法である。LSAと比較して, 分解後の行列の要素がすべて非負であるため, 結果の解釈がしやすいメリットがある。

LSAと異なり, 分解された「文書×トピック」と「トピック×単語」の行列は, 各トピックのベクトル間で直交し

ておらず、各トピックは数学的に互いに独立していない。これにより、トピック間に意味的な重複が生じる可能性がある。また、LSAと同様に、NMFでも結果が学習データに完全に依存するため、過学習を起しやすく、新しい文書のトピックは推定できないという課題がある。

### NMFの行列分解

$$X = WH^t$$

### NMFでアウトプットされる各行列の例

W	トピック1	トピック2	トピック3	トピック4
文書1	0.22	0.48	0.16	0.39
文書2	0.35	0.12	0.07	0.88
文書3	0.11	0.58	0.24	0.14
文書4	0.29	0.27	0.60	0.20

H <sup>t</sup>	単語1	単語2	単語3	単語4	単語5
トピック1	4.34	0.03	1.21	2.24	1.87
トピック2	3.41	1.89	2.81	0.46	1.03
トピック3	0.38	2.24	1.09	5.62	0.06
トピック4	1.85	3.66	0.08	2.11	4.16

図5 NMFの行列分解

### 2.2.3 PLSA

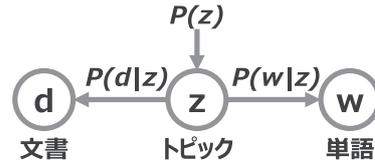
PLSA (Probabilistic Latent Semantic Analysis, 確率的潜在意味解析) は、確率モデルを導入してトピックを抽出する手法である。特異値分解によってトピックを抽出するLSAの処理を、確率的な枠組みによって発展させたもので、1999年に発表された<sup>23)</sup>。情報検索の分野ではPLSI (Probabilistic Latent Semantic Indexing) と呼ばれることもある。

PLSAでは、Bag-of-Wordsの「文書×単語」の行列を確率的に分解することでトピックを抽出する。PLSAの確率モデルのイメージを図6に示す。PLSAのモデルは文書 $d$ と単語 $w$ の同時確率 $P(d,w)$ を、トピック $z$ を使った3つの確率分布① $P(d|z)$ 、② $P(w|z)$ 、③ $P(z)$ に分解して表現する。この3つの確率分布をBag-of-Wordsの「文書×単語」行列データに基づき推定する。具体的には、 $P(d,w)$ の対数尤度関数を最大化するEMアルゴリズムを実行し、初期値を与えた反復計算により最適解を得る。

PLSAのメリットとして、LSAでは「文書×単語」の行列に対して事前にTF-IDFなどで重みづけする必要があったが、PLSAでは確率的な処理によりそうした重みづけをせずに実行できる。また、PLSAは推定される $P(d|z)$ と $P(w|z)$ によって、文書および単語のトピックに対する関連度が所属確率として出力されるため、結果が解釈しやすいメリットがある。なお、各トピックは互いに独立しているという数学的な仮定を置いている。

PLSAのデメリットとしては、入力する「文書×単語」のデータセットが大規模になると計算コストが高くなることが挙げられる。また、PLSAもLSAやNMFと同様に、結果が観測データに完全に依存するため、過学習を起しやすく、新しい文書のトピックは推定できないという課題があるが、一方で、観測データに対する再現度は高く、その観測データ固有の特徴をそのまま反映できるモデルと捉えることもできる。

### PLSAの確率モデル



$$P(d, w) = \sum_z P(d|z)P(w|z)P(z)$$

### PLSAの出力結果例

P(文書d   トピックz)					P(単語w   トピックz)					P(トピックz)				
P(d z)	トピック1	トピック2	トピック3	トピック4	P(w z)	トピック1	トピック2	トピック3	トピック4	P(z)	トピック1	トピック2	トピック3	トピック4
文書1	0.41	0.16	0.06	0.27	単語1	0.11	0.09	0.44	0.20		0.31	0.27	0.23	0.19
文書2	0.29	0.54	0.11	0.06	単語2	0.09	0.38	0.04	0.18	$\sum_z P(z) = 1$				
文書3	0.22	0.13	0.04	0.58	単語3	0.50	0.11	0.29	0.09					
文書4	0.08	0.17	0.79	0.09	単語4	0.08	0.14	0.17	0.31					
$\sum_d P(d z) = 1$					単語5	0.22	0.28	0.06	0.22					
					$\sum_w P(w z) = 1$									

図6 PLSAの確率モデル

### 2.2.4 LDA

LDA (Latent Dirichlet Allocation, 潜在ディリクレ配分法) は、PLSA をベイズ的に拡張させた手法で、ディレクレ分布を事前分布に導入しており、2003年に発表された<sup>24)</sup>。

LDAのグラフィカルモデルを図7に示す。LDAは確率的な生成モデルともいわれ、文書d内の各単語wが特定のトピックzから生成されたと仮定する。このトピックzの分布は文書dごとに異なり、その確率分布P(z|d)はθで表される。一方、特定のトピックzが各単語wを生成する確率分布P(w|z)はφで表し、そのトピックzの下で単語wを生成するプロセスではφが参照される。つまり、LDAではθ→z→wという生成プロセスが仮定され、その中でz→wというアローの確率過程はφ→wというアローによる確率分布φが参照される。なお、θとφはそれぞれハイパーパラメータαとβを持つディレクレ分布に従う。LDAの推定対象はθ=P(z|d)とφ=P(w|z)であるが、これらの生成過程を逆にたどるアルゴリズムにより、単語w(観測データ)からθとφを推定する。具体的には、各文書d中の各単語w(観測データ)に対して、まずその単語wがどのトピックzから生成されたと考えられるか、θとφの暫定値(初期値)から推定する(zの推定)。次に、その結果を基にして、各文書dがどのトピックzをどれだけ含むか推定する(θの推定)。そして、各トピックzがどの単語wを生成する可能性が高いか推定する(φの推定)。これを繰り返すことでθとφを逐次的に更新していく。推定アルゴリズムは、ギブスサンプリングや変分ベイズ法などが適用され、反復的に計算される。

PLSAはパラメータが固定的で、観測データのみから直接推定するため、結果が完全に観測データに依存するが、LDAはパラメータが事前分布に従って変動する確率分布とし、観測データと事前分布から推定する。事前分布を導入することで、確率的なスムージング効果があり、過学習を抑制できる。また、LDAは新しい文書についても推定ができ、新しい文書に対応するθを未知とし、新しい文書に含まれる単語w(観測データ)と学習済みのφとαからθを推定する。

一方、LDAはハイパーパラメータαとβによって結果が変動しやすく、この値の設定の仕方、推定の仕方が難しい。ハイパーパラメータα,βの大きさとディレクレ分布の特徴は、α,β>1では、値が大きいほど確率が均一化し、どの文書も同じようなトピックの分布を持つようになり、トピックは多様な単語に分散して関連づけられる傾向がある。

$\alpha, \beta = 1$  では分布は一様分布となり、さまざまなトピックにランダムに確率が割り当てられる。これは事前分布を仮定していない PLSA と近い振る舞いをする。 $1 > \alpha, \beta > 0$  では、値が小さいほど確率が局所的になり、文書は特定のトピックに確率が集中し、トピックは一部の単語に偏って関連づけられる傾向がある。また、LDA は新しい文書の推定ができる高い汎化性能を持つ一方で、トピックの結果が一般的で抽象度が高くなることもある。

なお、PLSA と LDA はどちらも文書  $d$  と単語  $w$  の関係をモデル化している手法だが、文書  $d$  と単語  $w$  の関係の扱い方の違いには注意が必要である。PLSA の推定対象は  $P(d|z)$  と  $P(w|z)$  であるが、「文書  $d \times$  単語  $w$ 」の行列データに対して文書  $d$  と単語  $w$  の役割には対称性があり、行と列を入れ替えても推定結果に影響はない。一方、LDA の推定対象は  $P(z|d)$  と  $P(w|z)$  であるが、文書  $d$  と単語  $w$  の役割は対称的ではなく、それぞれの役割は明確に区別されたモデルとなっている。そのため、行と列を入れ替えて適用すれば異なる結果となる。

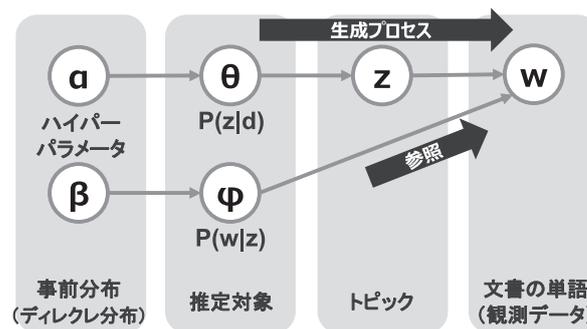


図7 LDAのグラフィカルモデル

### 2.3 テキストマイニングに適用するトピックモデルの考察

トピックモデルの中ではLDAを適用することが一般的である。しかし、LDAがあらゆる課題で優れたトピックモデルということではない。トピックモデルの各手法の特徴を俯瞰的に理解し、解決したい課題に応じて使い分けることが重要である。

LSA, NMF, PLSAは、過学習を起こしやすいことがデメリットとして挙げられ、そうした背景からLDAが開発されているが、これは新しい文書の推定を課題としたときのデメリットである。たとえば、情報検索などへの適用を想定した場合、学習データにないテキストデータのトピックを推定できることは確かに重要であり、LDAを適用することが相応しいと考えられる。では、テキストマイニングに適用する場合はどうか考察する。先述したように、テキストマイニングによる単語ベースの可視化は、単語数が増えると複雑性が増し解釈性が損なわれる課題が生じる。その課題を解決する手段としてトピックモデルを適用する場合、テキストマイニングが現状把握をするための手法であることを考えれば、新しいデータのトピックを推定することよりも、観測データのありのままの特徴をトピックで反映できることが重要といえる。LSA, NMF, PLSAは観測データに依存した過学習でトピックを抽出するモデルだが、観測データに対する再現度は高く、その観測データ固有の特徴をそのまま反映できるモデルと捉えることができる。

この点において、LDAは事前分布にディレクレ分布を仮定していることで、観測データをそのまま再現するモデルではなくなっている。またLDAは新しい文書の推定ができる高い汎化性能を持つ一方で、抽出されるトピックの内容が一般的で抽象度が高くなることもある。さらに、LDAはディレクレ分布のハイパーパラメータによって結果が変動しやすく、その設定の仕方、推定の仕方が難しい。テキストマイニングは専門知識がなくても直感的に操作できる使いやすさが特長の一つであり、ビジネス現場では文系理系問わず人気の分析手法となっている。こうした点においても、より高度なLDAを適用することは複雑で扱いにくさが生じてしまう懸念がある。

したがって、テキストマイニングによる単語ベースの可視化の複雑性を解決する手段としては、LDAよりもLSA, NMF, PLSAの方が適していると考えられる。その中でもPLSAは、LSAのようなTF-IDFなどの前処理がいらす、トピックの結果が確率で解釈しやすく、かつ各トピックが独立しているという特徴を有していることから、筆者は

PLSA がテキストマイニングに適用するトピックモデルの手法として相応しいと考えている。

## 2.4 PLSA の適用によるトピック抽出

ここでは PLSA という手法についてより詳細に解説する。PLSA は自然言語処理に限らずデータクラスタリングや次元削減法として用いられる手法である。改めて PLSA によるクラスタリングのイメージを図 8 に示す。自然言語処理に限定せず、より一般的に PLSA という手法を表現すると、共起行列と呼ばれる行列データをインプットに教師なし学習をし、行の要素  $x$  と列の要素  $y$  の背後にある共通する特徴となる潜在クラス  $z$  を抽出する手法となる。自然言語処理の分野では、共起行列を Bag-of-Words の「文書×単語」行列とし、潜在クラスをトピックと呼んでいる。つまり、文書  $x$  とそこに出現する単語  $y$  の間には潜在的な意味クラス  $z$  があることを想定し、文書と単語に共通するトピックを抽出している。これにより文書全体をいくつかのトピックを軸に解釈できるとともに、使われ方の似ている単語群をそのトピックの構成要素として集約できる。自然言語処理の分野以外では、たとえば「画像×画像特徴量」の行列データ (Bag-of-visual-words と呼ばれるデータ) に PLSA を適用して画像分類を行う取り組みや<sup>25)</sup>、「顧客 ID × 購買商品」の行列データ (ID-POS データと呼ばれるデータ) に PLSA を適用して顧客と商品を同時に分類する取り組みなどがある<sup>26)</sup>。

以下に PLSA によるトピック抽出までの計算過程を説明する。PLSA では文書  $x$  における単語  $y$  の出現確率  $P(y|x)$  を、潜在的なトピック  $z$  を介した式 (1) でモデル化する。これは Aspect モデルと呼ばれるが、Aspect とは文書内に潜むトピックを指している。この Aspect モデルに、確率の連鎖律により成り立つ同時確率と条件付確率の関係式 (2) と、ベイズの定理により成り立つ式 (3) を適用する。これによって式 (4) で示すように、文書  $x$  と単語  $y$  の同時確率  $P(x,y)$  を 3 つの確率分布①  $P(x|z)$ 、②  $P(y|z)$ 、③  $P(z)$  に分解して表現する。ここで、文章  $x$  における単語  $y$  の出現回数を  $N(x,y)$  とすると、式 (5) の対数尤度を最大にする  $P(x|z)$ 、 $P(y|z)$ 、 $P(z)$  を、EM アルゴリズムを用いて計算する。つまり式 (6) の E ステップと式 (7) ~ (9) の M ステップを計算することで最尤推定する。以上から PLSA の実行によって得られるアウトプットは 3 種類の確率変数  $P(x|z)$ 、 $P(y|z)$ 、 $P(z)$  の値となる。 $P(x|z)$  はそれぞれのトピック  $z$  においてすべての文書  $x$  の所属確率を合計すると 100% となり、 $P(y|z)$  はそれぞれのトピック  $z$  においてすべての単語  $y$  の所属確率を合計すると 100% となる。またトピックの周辺確率である  $P(z)$  も、すべてのトピック  $z$  の確率を合計すると 100% となる。こうした確率分布として得られるアウトプットは解釈性に優れており、PLSA を適用することで、どのような単語群で構成されるトピックが抽出され、それぞれの文書はどのトピックの要素が多いのか把握できる。

このようにテキストマイニング× PLSA の適用により、大量のテキストデータを複雑な単語単位ではなく、集約されたトピックを単位としてシンプルに傾向を分析できる。

$$P(y|x) = \sum_z P(y|z)P(z|x) \quad (1)$$

$$P(x,y) = P(y|x)P(x) \quad (2)$$

$$P(z|x) = \frac{P(x|z)P(z)}{P(x)} \quad (3)$$

$$P(x,y) = \sum_z P(x|z)P(y|z)P(z) \quad (4)$$

$$L = \sum_x \sum_y N(x,y) \log P(x,y) \quad (5)$$

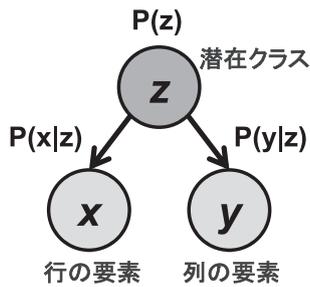
$$P(z|x,y) = \frac{P(x|z)P(y|z)P(z)}{\sum_z P(x|z)P(y|z)P(z)} \quad (6)$$

$$P(x|z) = \frac{\sum_y N(x,y)P(z|x,y)}{\sum_x \sum_y N(x,y)P(z|x,y)} \quad (7)$$

$$P(y|z) = \frac{\sum_x N(x,y)P(z|x,y)}{\sum_x \sum_y N(x,y)P(z|x,y)} \quad (8)$$

$$P(z) = \frac{\sum_x \sum_y N(x,y)P(z|x,y)}{\sum_x \sum_y \sum_z N(x,y)P(z|x,y)} \quad (9)$$

### PLSAのグラフィカルモデル



xとyの共起確率を潜在クラスzを使って表現する

$$P(x,y) = \sum_z P(x|z)P(y|z)P(z)$$

※条件付確率P(A|B)  
事象Bが起こる条件の下で事象Aの起こる確率

### PLSAのインプットとアウトプット

共起行列  
N(x,y)

	y1	y2	y3	y4
x1	0	2	0	1
x2	1	0	0	1
x3	1	1	0	0
x4	0	0	2	0

P(x|z)  
 $\sum_x P(x|z)=1$

	z1	z2	z3
x1	0.45	0.14	0.15
x2	0.22	0.41	0.07
x3	0.28	0.34	0.18
x4	0.05	0.11	0.60

P(y|z)  
 $\sum_y P(y|z)=1$

	z1	z2	z3
y1	0.16	0.58	0.13
y2	0.44	0.17	0.05
y3	0.30	0.09	0.65
y4	0.10	0.16	0.17

P(z)  
 $\sum_z P(z)=1$

	z1	z2	z3
	0.47	0.32	0.21

図8 自然言語処理に限定しないPLSAによるクラスタリング

## 3. ベイジアンネットワーク

ここではベイジアンネットワークの概要を解説し、テキストマイニングに適用することの有効性について、特にトピックモデルの適用を介してベイジアンネットワークを適用することの有効性について述べる。

### 3.1 ベイジアンネットワークの概要

ベイジアンネットワークは、複数の変数の確率的な因果関係を有向リンクのネットワーク構造で表わし、その関係の強さを条件付確率で表現した確率モデルである<sup>27)</sup>。このモデルを用いることで、ある変数の状態を条件として与えたときの他の変数の起こりうる確率を推論することができる。ベイジアンネットワークの概要図を図9に示す。ベイジアンネットワークは、①確率変数、②確率変数間のリンク構造、③各リンクの条件付確率表の3つによって定義される。ベイジアンネットワークはすべての確率変数の同時確率を、各変数間のリンク関係が示す条件付確率の連鎖律で表現している。そのリンク構造は、観測データと変数の定義に基づき、それを数学的に最もよく説明する（尤度を最大化する）モデルを学習して獲得される。これにより各変数間の確率統計的な関係性を把握することができる。また、構築されたモデルを用いることで、観測した変数群から未観測の変数の確率分布を条件付確率表に基づいて推論することができる。なお、ベイジアンネットワークで用いる確率変数は質的変数（カテゴリカル変数）となるため、量的変数の場合は閾値を設けて事前にカテゴリに分割する必要がある。ベイジアンネットワークの特長には以下の点が挙げられる。

(1) 要因の関係構造を理解できる

どの変数がどの変数に影響しているのか、構築されたモデルの構造によってデータ全体に潜む要因関係を理解することができる。なお、ベイジアンネットワークでモデル化される有向リンクの関係性は、本当の因果関係ではなく、あくまでも確率的な因果関係である。本当の因果関係は原因と結果の関係であるが、確率的な因果関係では、その変数を条件に与えた条件付確率を形成する方がモデルの尤度（データへの当てはまりの良さ）が向上するというものである。

(2) モデルの構造を指定できる

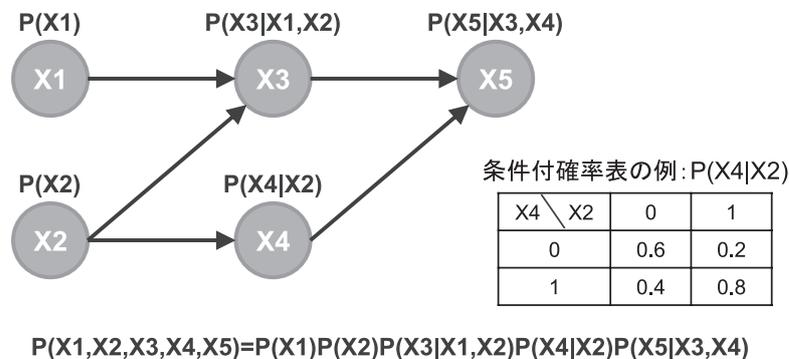
各変数の関係構造は、観測データに基づいて数学的な基準のみで探索することもできるが、人間が構造を指定することもできる。たとえば、この変数とこの変数は関係があることは分かっているといった経験則があったり、この変数群とこの変数群の関係にフォーカスして確認したいというような目的が定まっていれば、それをモデル構築の条件に採用することができ、それ以外の部分を観測データから数学的に探索するということもできる。たとえば、特許文書情報から用途と技術の関係を分析するという点においても、用途を実現する上での重要な要素技術を把握したいときは用途群⇒技術群という向きのリンク構造を指定し、技術を応用する用途の展開を把握したいときは技術群⇒用途群という向きのリンク構造を指定してモデルを構築すると効果的である。つまり業務における経験則や分析目的といった事前知識と数学のハイブリッドで、より実務に即したモデルを構築できる。

(3) 互いの変数間の関係を把握しさまざまな方向からの推論を実行できる

ベイジアンネットワークでは、目的変数と説明変数の区別なく変数の関係をモデル化し、ある変数を条件に与えたときの他の変数の確率を推論することができる。回帰分析や決定木分析、ニューラルネットワークなど、通常のモデリング手法では、目的変数と説明変数を設定し、一つの目的変数ごとにモデルを構築する必要があり、一つの目的変数に対する説明変数の関係しか把握することができない。また構築されたモデルを用いた推論の実行では、複数の説明変数群から一つの目的変数を推論するという一方の推論に限定される。一方、ベイジアンネットワークでは目的変数と説明変数の区別がないため、それぞれの変数が互いにどのような関係性を持っているのかという構造を把握できる。また、モデルを推論で用いる場合も、その推論対象と推論条件とする変数は自由に設定でき、さまざまな方向から各変数の起こり得る確率を推論できる。

(4) 非線形の関係や交互作用の効果も表現できる

ベイジアンネットワークは回帰分析のように線形処理によってモデルを構築するのではなく、確率論による非線形処理のモデルであるため、非線形の関係がある複雑な現象でもモデルで表現できる。また、ある条件が揃うときにだけ効果が発揮されるというものや、ある条件とある条件が組み合わせると逆の効果に転じてしまうといった交互作用がある場合でも、確率的に意味のある関係としてモデル化することができる。



※条件付確率P(A|B)  
事象Bが起こる条件の下で事象Aの起こる確率

図9 ベイジアンネットワークの概要図

### 3.2 テキストマイニングにベイジアンネットワークを適用するメリット

本来ベイジアンネットワークはテキストデータを分析対象として開発された技術ではないが、これを応用することで、テキストデータの中に潜む要因関係を構造化することが可能となる。たとえば、テキストマイニングによって構築されたBag-of-Wordsのデータの単語一つひとつを確率変数に設定し、その変数間の関係をベイジアンネットワークでモデル化する取り組みが報告されている<sup>28)</sup>。この取り組みでは病院で収集された子どもの傷害事故の診療記録データ4,238件を対象に、事故状況がテキストで記されたデータにテキストマイニングを実行することで事故に関わる製品や行動の単語を抽出し、それら一つひとつを確率変数としている。また子どもの性別・年齢の情報や事故の情報も確率変数とし、これらの関係をベイジアンネットワークでモデル化している。構築されたモデルを図10にて紹介する。

このようにテキストマイニングにベイジアンネットワークを適用することのメリットをまとめると、主に以下の2つが挙げられる。

#### (1) テキストデータに潜む要因関係を構造的にモデル化できる

テキストマイニングでは、文書でどのような単語が多く使われる傾向があるのか、またその傾向が属性によってどのような違いがあるのか把握できる。ここにベイジアンネットワークをさらに適用することで、そうした傾向がどのような要因によって影響を受けているのか、あるいは他の事象に影響を与えているのかという、統計的な関係性を構造的なモデルとして可視化できる。たとえば図10のモデルでは、子どもの製品に対する行動はどのような発達段階を要因としており、またどのような事故の要因となり得るのか、その統計的な関係性を分析できている。

#### (2) 現状把握だけでなく状況の変化に伴うシミュレーションができる

先述したとおり、テキストマイニングは基本的に文書の記述傾向の現状把握をする手法であるが、その現状から状況が変化したときに、それに伴って結果がどう変化・影響を受けるのかシミュレーションすることはしない。ここにベイジアンネットワークをさらに適用することで、その確率的なシミュレーションが可能になる。上記(1)で述べたように、ベイジアンネットワークを適用することで、テキストデータに潜む要因関係を構造的にモデル化できるようになる。そのモデルを利用することで、たとえば図10のモデルでは、子どもの発達段階の変化によって行動の確率がどの程度変化するのか、同じ行動でも製品が異なれば起こり得る事故の確率がどの程度変化するのかというように、与えた条件に対する結果の挙動を確率的にシミュレーションすることができている。たとえば、「コイン」という製品では「飲む」という行動が57.3%で、「誤飲」という事故が32.6%と確率が高いことがシミュレーションで計算されている。また「ミニカー」という製品では「押す」という行動が50.0%、「乗る」という行動が23.8%で、「転倒」という事故が31.4%、「転落」という事故が27.1%と確率が高く、さらにこれらは1歳児で確率がさらに高くなることがシミュレーションで計算されている。こうした分析によって、子どもの事故の発生を抑制するための具体的な施策をデータドリブンで検討できる。

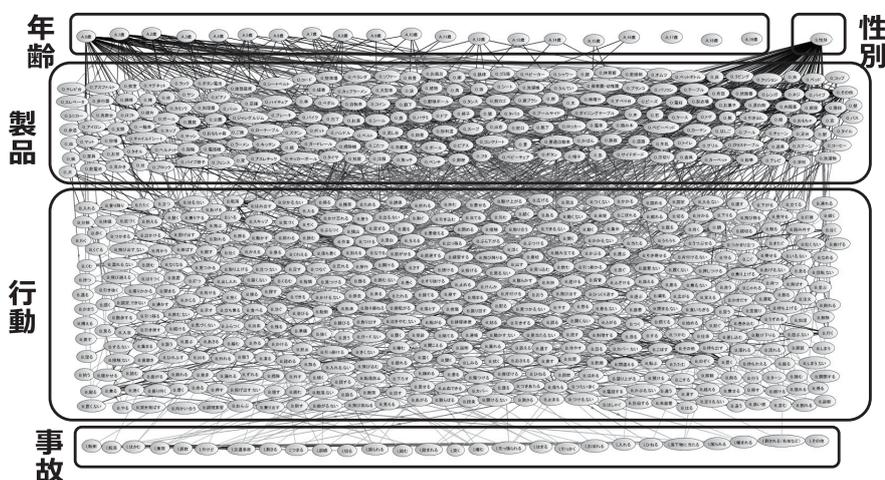


図10 Bag-of-Words にベイジアンネットワークを適用したモデルの例 (※文献<sup>28)</sup> より引用)

### 3.3 テキストマイニング×ベイジアンネットワークの課題

図 10 のモデルで取り上げた分析例では、テキストマイニングで抽出された単語一つひとつを確率変数に採用しているため、構築されたモデルはとても複雑で解釈が難しいものとなっている。つまり、Bag-of-Words のデータにそのままベイジアンネットワークを適用し、単語一つひとつを確率変数にすると、重要な傾向や気づきが埋もれてしまう懸念がある。本来モデルとはデータに潜む特徴や傾向を抽象化したものであり、情報が集約されて全体としてシンプルに表現されていることが望ましい。

そこでテキストマイニングに直接ベイジアンネットワークを適用するのではなく、まずトピックモデルを適用することを考える。この処理を挟むことによって、単語一つひとつではなく、単語群が集約されたトピックを確率変数として定義し、その上でベイジアンネットワークを適用する。つまり、Bag-of-Words ではなく Bag-of-Topics のデータにベイジアンネットワークを適用する。これによって、テキストデータ全体に存在する要因関係をよりシンプルに把握できる。このようにテキストマイニング、トピックモデル、ベイジアンネットワークの 3 つの手法を組み合わせることでテキストデータを分析する手法を次項で解説する。

## 4. テキストマイニング×PLSA×ベイジアンネットワークの分析手法：Nomolytics

ここでは従来のテキストマイニングに加え、トピックモデル (PLSA) とベイジアンネットワークという 2 つの AI 技術を連携して適用し、テキストデータに潜む特徴や要因関係を構造的にモデル化する手法として、筆者が開発した Nomolytics (Narrative Orchestration Modeling Analytics)<sup>29)</sup> を紹介する。なお本手法は筆者が有限責任監査法人トーマツに所属していたときに特許として出願し登録されたものであり (特許第 6085888 号)<sup>30)</sup>、令和 7 年 3 月 10 日現在で有限責任監査法人トーマツと株式会社アナリティクスデザインラボが権利を保有している。

### 4.1 Nomolytics という分析手法

Nomolytics の概要を図 11 に示す。本手法では、まずテキストマイニングによりテキストデータの文書情報から単語を抽出し、各単語の共起頻度をデータ化した共起行列を作成する。次に、その共起行列をインプットとして PLSA を適用し、使われ方の似ている単語群を複数のトピックにまとめ上げる。最後に、すべてのテキストデータに対して各トピックの該当度 (該当有無) を計算することでトピックを確率変数として定義し、ベイジアンネットワークによってトピック間あるいは他の属性情報との間の確率的な要因関係を構造的にモデル化する。

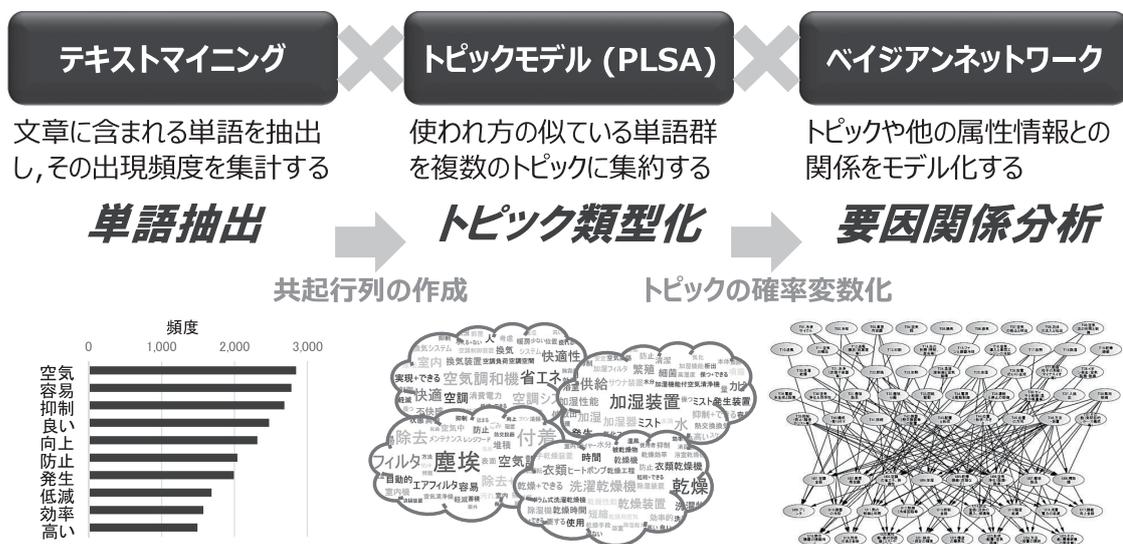


図 11 Nomolytics の概要図

こうした3つの手法を組み合わせることで、膨大なテキストデータをいくつかのトピックという人間が理解しやすい形に整理でき、そのトピックを新たな分析軸としてさまざまな特徴の探索が可能となる。さらにベイジアンネットワークによってそのトピック周辺に潜む要因関係を構造的にモデル化できる。そしてそのベイジアンネットワークのモデルを用いることで、ある変数の条件を変化させたときに、それに伴って他の変数がどのように変化するかという確率的なシミュレーションを実行することができる。

本技術はテキストデータであればあらゆる分野で適用でき、例えば旅行の口コミデータに適用して地域観光のマーケティングを検討する事例もある<sup>31)</sup>。本節ではこれを特許文書に適用した事例について紹介する。

## 4.2 Nomolytics における各手法の連携の工夫

Nomolytics では、「テキストマイニング×PLSA」と「PLSA×ベイジアンネットワーク」というそれぞれの手法を連携する際の処理に独自の工夫を施しており、それが Nomolytics の技術的な特徴にもなっている。以下にその連携の工夫について紹介する。

### 4.2.1 テキストマイニング×PLSA の工夫

先述したトピックモデルの解説では、トピックモデルは従来のクラスター分析と異なり行と列を同時にクラスタリングできると説明した。つまり、PLSA の共起行列の行と列は、双方が十分意味を持つ情報で構成すれば、抽出されたトピックの意味は2つの軸から解釈することができる。図12（左）に示すように、一般的な PLSA の適用では、「文書×単語」という構成の Bag-of-Words を共起行列としてインプットする。しかし、行に設定された「文書」とは文書 ID の情報であり、それ自体に意味は持たないため、抽出された潜在クラスの意味解釈には使用しにくい情報である。また、「文書×単語」の Bag-of-Words は、ほとんどは0となる疎なデータであるため、文書間・単語間で違いが現れにくく、特徴がクリアなトピックが得られにくい。

そこで Nomolytics では、共起行列の構成によって解釈のしやすいトピックを抽出する工夫を施している。たとえば、「名詞の単語×動詞の単語」というように、行と列をそれぞれ異なる品詞の単語で構成する方法や、図12（右）のように「単語×係り受け」という構成で、軸の一方を係り受け表現（文法的なつながりのある単語のペア）とする方法を取る。こうした「単語×係り受け」の共起行列に PLSA を適用することで、単語という話題の観点となる軸と、その観点の具体的な内容となる係り受け表現を軸にクラスタリングができ、より文脈上近い言葉・表現でまとめられた、解釈のしやすいトピックを抽出できることが期待できる。また、この共起行列は各単語と各係り受けが同時に出現する文章数（共起頻度）が値として入るクロス集計型の行列であり、一般的な PLSA の共起行列のような0ばかりの疎なデータとは異なり、密なデータとなりやすくなる。これにより、スパース性の問題の影響を受けにくく、よりまとまりのあるクリアなトピックを抽出できることが期待できる。さらにその共起行列のサイズは一般的な PLSA で用いる共起行列に比べて、特に行数においてとても小さいものになり、計算時間も大幅に削減できる効果がある。

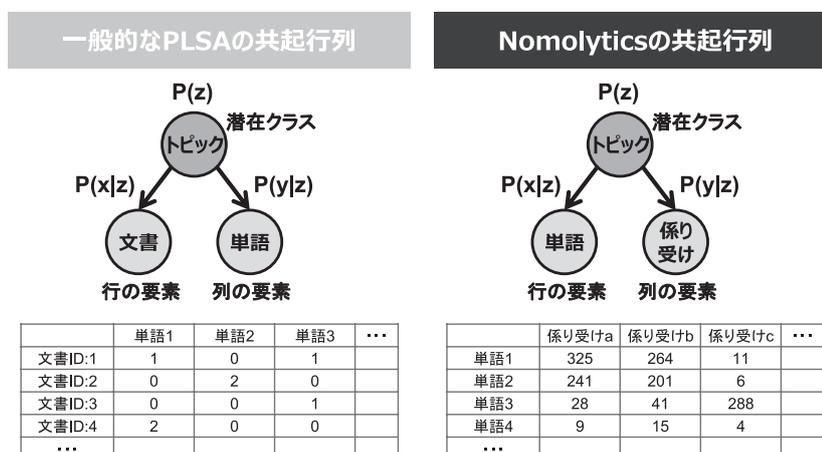


図12 Nomolytics における PLSA の共起行列構成の工夫

## 4.2.2 PLSA ×ベイジアンネットワークの工夫

Nomolytics では抽出されたトピックを確率変数として扱い、ベイジアンネットワークを適用することでトピック周辺に存在する要因関係を構造的にモデル化するが、トピックを確率変数として扱うための変換処理にも工夫がある。処理の詳細は後述する実際の分析事例で解説するが、ここではその考え方を紹介する。まず、「文書×単語」という共起行列をインプットとする一般的な PLSA で抽出されるトピックは、図 6 で示したような所属確率  $P(d|z)$  という値によって各文書 ID にトピックが紐づいた形で結果が得られる。そのため、この所属確率という連続値をカテゴリ化処理すれば、各トピックは元のテキストデータに対応する確率変数としてそのまま扱うことができる。一方、Nomolytics で採用する「品詞 A の単語×品詞 B の単語」という構成の共起行列や、「単語×係り受け」という構成の共起行列で抽出されたトピックでは、各単語や各係り受けはトピックに紐づいて所属確率が計算されているが、元のテキストデータには紐づきがされていない結果となる。つまり、元のテキストデータと抽出されたトピックの紐づきを計算する処理が必要になる。そこで、各テキストデータの文書情報に含まれる単語・係り受けと、その単語・係り受けが各トピックに対して持つ所属確率から、そのテキストデータに対する各トピックの該当度を示すスコアを確率的に計算する。最終的にはそのスコアに閾値を設け、各トピックに該当するか否かを 1.0 のフラグに変換し、トピックの確率変数を作成する。

## 5. Nomolytics を特許文書データに適用した分析事例（その 1）

ここでは Nomolytics を実際の特許文書データに適用した分析事例について紹介する。なお、本節で紹介する事例は筆者のコンサルティング業務で実際にクライアントから依頼のあった分析の事例ではない。あくまでも対外紹介用に作成した分析事例となるが、ここで紹介している分析のアプローチ自体は実際に筆者のコンサルティング業務でよく提供しているものとなる。

### 5.1 分析で用いるデータ

本分析事例では、特許の要約と請求項に「風」「空気」という 2 つのキーワードを含む国内の特許公報データ 30,039 件を分析対象としている。出願期間は 2006 年 1 月 1 日から 2015 年 12 月 31 日までのちょうど 10 年分の特許公報を対象に抽出した。特許は出願してから公開されるまで 1 年半の期間を要するため、それを考慮して、データの抽出は 2017 年 7 月 25 日に実施した。ここでは特許の要約文の記述情報を主な分析対象としている。国内特許の要約文は 400 字以内という字数制限が設けられており、また厳密なルールではないが、【課題】や【解決手段】などの見出しをつけてから、それに該当する内容を記述することが慣例となっている。要約文の例を図 13 に示す。「風」「空気」をキーワードに含む特許のデータということで、たとえば、エアコンや扇風機、空気清浄機、加湿器、掃除機、洗濯乾燥機など、さまざまな生活家電が関連した特許が含まれる。

【要約】【課題】ユーザーの快適性を維持しつつ、省エネ運転を行うことができる空気調和機を提供すること。【解決手段】本発明の空気調和機は、室内温度を検出する室内温度検出手段と、人体の活動量を検出する人体検出手段と、基準室内設定温度を設定するリモコン装置を備え、室内温度が基準室内設定温度となるように空調制御を行う空気調和機であって、人体検出手段で検出する活動量が所定の活動量以内であるときは、室内温度が、基準室内設定温度を補正した補正室内設定温度となるように空調を行い、補正室内設定温度よりも低い状態を継続すると、圧縮機を停止させ、圧縮機の復帰は、基準室内設定温度に基づいて行う。

※例示のための要約文であり、一部内容は筆者が加工している

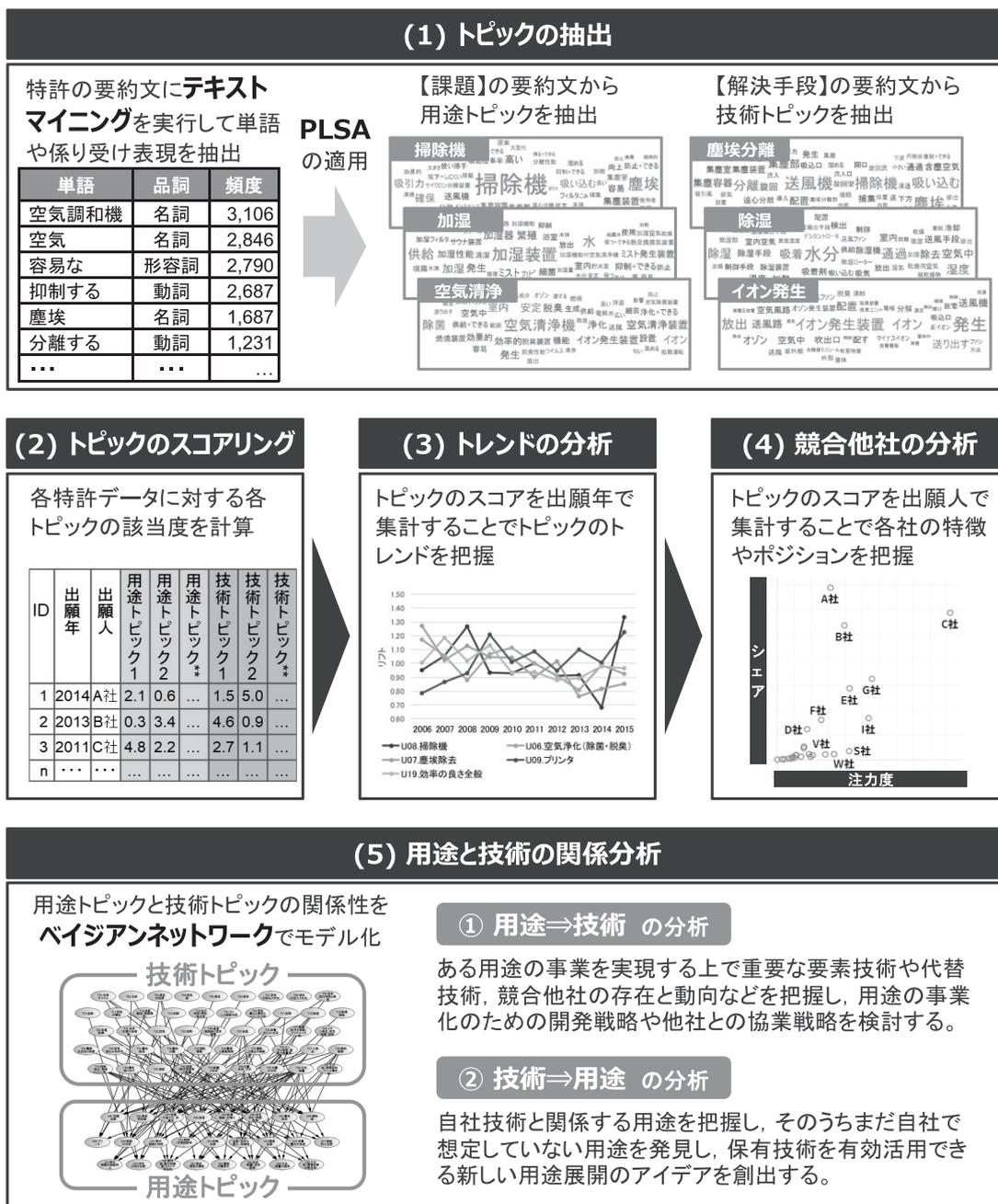
図 13 分析対象とする特許の要約文の例

## 5.2 分析の全体像

本分析事例では、特許の要約文の情報を対象に、そこに記載されている「課題」の項目の文章と「解決手段」の項目の文章にテキストマイニングとPLSAを適用することで、それぞれ用途に関するトピックと技術に関するトピックを抽出し、そのトピックを軸に特許データ全体の傾向を把握する。さらに、用途トピックと技術トピックの関係をベイジアンネットワークでモデル化する。分析のプロセスの概要を図14に示す。ここでは(1)トピックの抽出、(2)トピックのスコアリング、(3)トレンドの分析、(4)競合他社の分析、(5)用途と技術の関係分析、という5つのステップで特許文書データを分析する。それぞれの概要は以下のとおりとなる。

### (1) トピックの抽出

特許の要約文に記述されている「課題」と「解決手段」という項目の文章を対象に、テキストマイニングとPLSAを適用し、それぞれ用途に関するトピックと技術に関するトピックを抽出する。ここで得られた結果により、分析対象の特許に記されている用途と技術の全体像を把握する。



## (2) トピックのスコアリング

分析対象とした全特許データに対して抽出したトピックのスコア（該当度）を計算する。ここで得られた結果により、抽出したトピックを新たな分析軸とする。

## (3) トレンドの分析

各トピックのスコアを「出願年」の属性で集計してトレンドを分析することで、用途や技術のトレンドを把握する。ここで得られた結果により、有望なニーズやシーズを探る。

## (4) 競合他社の分析

各トピックのスコアを「出願人」の属性で集計することで、各出願人の傾向やポジショニングを分析する。ここで得られた結果により、自社の技術の開発戦略や差別化戦略、他社との提携戦略、自社技術の売却先などを検討する。

## (5) 用途と技術の関係分析

用途のトピックと技術のトピックの確率的な因果関係をベイジアンネットワークでモデル化し、用途と技術の関係性を分析する。この関係分析は、①用途⇒技術の分析（用途に対する技術の関係分析）と②技術⇒用途（技術に対する用途の関係分析）という2つのパターンがある。

①用途⇒技術の分析では、ある検討中の用途を実現する際に重要となる要素技術を把握するための分析である。ここで得られた結果により、その用途を達成するためにどの技術開発に注力すべきか、またその技術領域で競合となりそうな他社はどこか、そのなかで他社が牛耳る技術の代替技術は存在するか、どの会社と技術提携すると効果的かなど、自社の開発戦略や他社との協業戦略を検討する。

②技術⇒用途の分析では、自社で保有している技術と関係のある用途を把握するための分析である。ここで得られた結果により、自社技術と関係のある用途のうち自社でまだ想定していない用途を見つけ、自社の技術をさらに有効活用できる新しい用途展開のアイデアを創出する。

## 5.3 PLSA の適用による用途と技術のトピック抽出

分析ではまずテキストマイニングと PLSA を用いて特許の要約文の内容をいくつかのトピックに集約する。先述のとおり、国内特許の要約文では、【課題】と【解決手段】という2つの項目が記載されていることが多いが、【課題】の項目の文章と【解決手段】の項目の文章をそれぞれ抽出し、課題からは用途に関するトピックを、解決手段からは技術に関するトピックを抽出する。以下にトピック抽出の手順を述べる。

### (1) テキストマイニング

課題の項目で記述された文章と解決手段の項目で記述された文章を切り出し、それぞれにテキストマイニングを適用し、単語とその文法的なペアとなる係り受け表現を抽出する。課題と解決手段という項目のない特許は、要約文全体を対象とする。単語は名詞、動詞、形容詞、形容動詞を、係り受けは名詞に対する動詞・形容詞・形容動詞の単語ペアおよびその逆の係り受け（動詞・形容詞・形容動詞に対する名詞）の単語ペアを抽出する。なおテキストマイニングの実行には Text Mining Studio(株式会社 NTT データ数理システム)を使用した。本ツールの形態素解析において、文章の区切りを定義する文字列には「句点 (。)」の他に、特許要約文の中で文章のラベル付けによく使用される「隅付括弧 (【 》)」も設定した。

また、テキストマイニングの形態素解析で抽出された単語には複数の類義語が存在する。PLSA で抽出するトピックがよりクリアな内容となるように、目視で類義語辞書を作成し、テキストマイニングのツールに反映した。作成した類義語辞書は、7,240 語の類義語に対して 2,204 語の代表語を登録した。類義語辞書の例を表 1 に示す。表 1 では、背景が灰色の単語が代表語となっている。

表1 風・空気に関連する特許の要約文から作成した類義語辞書の例

<b>温度調節</b>	<b>給気風路</b>	<b>螺旋状</b>	<b>隣接</b>
温度調節	給気風路	螺旋状	隣接
温度調整	給気路	渦巻き状	隣り合う
温調	給気通路	渦巻状	隣合う
調温	給気経路	<b>発生</b>	<b>除霜運転</b>
温度設定	<b>切り替える</b>	発生	除霜運転
<b>中心</b>	切り替える	起こる	デフロスト運転
中心	切り換える	起きる	デフロスタモード
中心部	切替える	<b>使用</b>	<b>行う</b>
中心側	切換える	使用	行う
中心位置	<b>方法</b>	利用	行なう
中心部分	方法	使う	おこなう
<b>筒状</b>	手段	<b>熱伝導</b>	<b>濾過</b>
筒状	本方法	熱伝導	濾過
筒形状	方式	熱伝達	ろ過
筒体	<b>屋外</b>	熱移動	<b>マイナスイオン</b>
筒形	屋外	<b>排ガス</b>	マイナスイオン
筒状体	室外	排ガス	負イオン
<b>排気風路</b>	室外側	排気ガス	<b>攪拌</b>
排気風路	屋外側	排出ガス	攪拌
排気通路	<b>吸込空気</b>	<b>噴出</b>	攪拌
排気路	吸込空気	噴出	<b>測定結果</b>
排気流路	吸込み空気	噴き出す	測定結果
排気経路	吸い込み空気	噴出す	計測結果

(2) 共起行列の作成

続いてPLSAでトピックを抽出する際のインプットとする共起行列を作成する。先述したとおり、一般的なPLSAでは、「文書×単語」という構成の共起行列を用いるが、本分析事例で適用するNomolyticsでは、「単語×係り受け」という構成をとり、それぞれの単語と係り受けが同時に出現する共起頻度（同時に出現する文章数）を集計したデータを用いる。その構成に採用する単語と係り受けは、特許単位でカウントした頻度が10件以上のものを対象とした。「課題」の文章からは「単語(3,256語)×係り受け(2,084表現)」の共起行列を、「解決手段」の文章からは「単語(5,187語)×係り受け(7,174表現)」の共起行列を作成した。それぞれの共起行列の例を表2,3に示す。表2,3は文章単位にカウントした頻度が上位10個の単語と係り受けに限定した共起行列を例として掲載している。なお、共起行列を作成するにあたり、多くの特許で共通して使用される単語で、かつ頻度が高すぎる単語（「提供」や「備える」など）や、重要な意味を持たない単語（「前記」や「本発明」など）は、トピック抽出においてノイズになり得るため、ストップワードとして共起行列を構成する単語の対象から除外した。ただし、「提供」や「備える」といった単語は係り受けを構成する単語としては含めることとした。

表2 課題の文章から作成した共起行列の例

	全体の頻度	1,578	1,331	578	548	540	335	296	279	272	264
	係り受け	空気調和機⇒提供	効率⇒良い	車両用空調装置⇒提供	掃除機⇒提供	容易⇒構成	画像形成装置⇒提供	抑制⇒提供	向上⇒図る	方法⇒提供	装置⇒提供
全体の頻度	単語										
3,163	空気調和機	1,578	100	4	1	55	0	39	27	2	1
2,802	容易	190	105	51	67	540	28	14	14	32	25
2,700	抑制	142	95	64	63	36	55	296	27	15	7
2,525	良い	113	1,331	12	56	43	31	27	9	19	24
2,339	向上	122	33	24	51	22	8	15	279	11	4
2,062	防止	83	51	23	21	35	34	14	30	6	18
2,042	発生	79	108	40	13	30	43	43	12	13	22
1,745	使用	34	83	14	19	25	3	11	6	23	28
1,717	低減	80	49	44	18	19	10	14	27	16	8
1,582	効率	67	1,331	8	31	26	18	21	7	15	14

表3 解決手段の文章から作成した共起行列の例

		1,146	917	776	709	701	695	675	660	625	535
全体の頻度		1,146	917	776	709	701	695	675	660	625	535
係り受け		空気⇒吸い込む	吸い込む⇒空気	連⇒通す	空気⇒吹き出す	備える⇒構成	吹出口⇒吹き出す	空気⇒供給	空気⇒送風	空気⇒排出	制御部⇒備える
全体の頻度	単語										
7,418	配置	210	191	150	123	145	143	100	127	100	59
4,119	送風機	225	197	112	113	95	73	125	109	58	72
4,060	供給	127	101	67	54	108	53	675	73	81	65
3,830	内部	134	107	126	76	64	44	84	71	111	39
3,738	制御	81	76	29	72	64	103	68	92	47	302
3,660	吹出口	267	301	107	370	108	695	42	63	54	60
3,458	位置	82	72	70	68	72	72	39	44	51	28
3,291	吸い込む	1,146	917	99	313	75	244	68	70	112	58
3,170	発生	118	102	63	80	67	64	49	42	40	47
3,138	外部	142	99	93	70	74	39	62	65	169	36

### (3) PLSA の実行と評価

作成した共起行列に PLSA を適用することで、使われ方の似ている単語と係り受けでまとめられたトピックを抽出する。「課題」の共起行列からは用途のトピックを、「解決手段」の共起行列からは技術のトピックを抽出する。なお PLSA は予めトピック数を設定する必要がある、また与える初期値により解が異なる初期値依存性がある。そこでトピック数を 1 刻みで変化させ、それぞれのトピック数に対して初期値をランダムに変えて PLSA を 5 回ずつ実行し、それぞれの解を情報量基準 AIC (Akaike's Information Criterion, 赤池情報量基準)<sup>32)</sup> で評価して最も評価の良い解を採用することとした。「課題」と「解決手段」のそれぞれの共起行列に対する PLSA の実行解の AIC の評価の結果を図 15,16 に示す。図 15 の「課題」では、トピック数を 15 から 35 まで 1 刻みで変化させて実行した結果であり、図 16 の「解決手段」では、トピック数を 35 から 60 まで 1 刻みで変化させて実行した結果となっている。各トピック数に対して 5 つの黒いプロットがあるが、この一つひとつが毎回初期値を変えて PLSA を実行した結果であり、灰色の折れ線グラフはその 5 つの実行解における AIC の平均値を示している。なお PLSA の実行には Visual Mining Studio (株式会社 NTT データ数理システム) の二項ソフトクラスタリング<sup>33)</sup> という PLSA と同様の分析機能を使用した。

図 15, 16 の縦軸の AIC は値が小さいほど良いモデルと評価する指標であるが、その計算式は以下の式 (10) で示される。式 (10) の第一項は元のデータに対するモデルの当てはまりの良さを示し、第二項はモデルのパラメータの自由度であり複雑度を示す。つまり、モデルの良さを当てはまりの良さだけで評価するのではなく、モデルの複雑度をペナルティとしている。これは、モデルは複雑にするほど (パラメータ数を多くすればするほど)、元のデータへの当てはまりは良くなるが、当てはまり過ぎるモデルは一般性を損なう可能性があり、それを第二項のペナルティで抑制しているということである。AIC で評価することで、程良くシンプルで元のデータへの当てはまりも良いモデルを選択することができる。図 15,16 から AIC の評価はトピック数に対して下に凸のカーブを描いていることがわかる。つまりトピック数を多くするとモデルの当てはまりが良くなるため AIC が改善していくが、あるときからトピック数の多さにより複雑度のペナルティが効き始め、AIC が悪化していくという現象である。なお、「行 X × 列 Y」の共起行列に PLSA を実行してトピック Z を抽出するモデルの場合、求めるパラメータは  $P(X|Z)$ ,  $P(Y|Z)$ ,  $P(Z)$  の 3 つの確率変数であるが、共起行列の行数を  $N_x$ 、列数を  $N_y$ 、トピック数を  $k$  とすると、各パラメータの自由度はそれぞれ  $k(N_x-1)$ ,  $k(N_y-1)$ ,  $k-1$  となる。これらの総和が AIC における第二項の自由度となる。

$$AIC = -2 (\text{モデルの最大対数尤度}) + 2 (\text{モデルの自由度}) \quad (10)$$

採用するトピックの選定の仕方は、まず AIC の平均が最小となるトピック数を定め、その中で個別の AIC が最小

となる実行結果を採用した。図 15, 16 より, 灰色の折れ線が示す各トピック数における AIC の平均値を比較したとき, それが最小となるトピック数は, 「課題」から得られる用途トピックについては 25 個, 「解決手段」から得られる技術トピックについては 47 個となった。そしてそれぞれのトピック数における 5 つの実行解の中で AIC が最小となる結果を採用した。

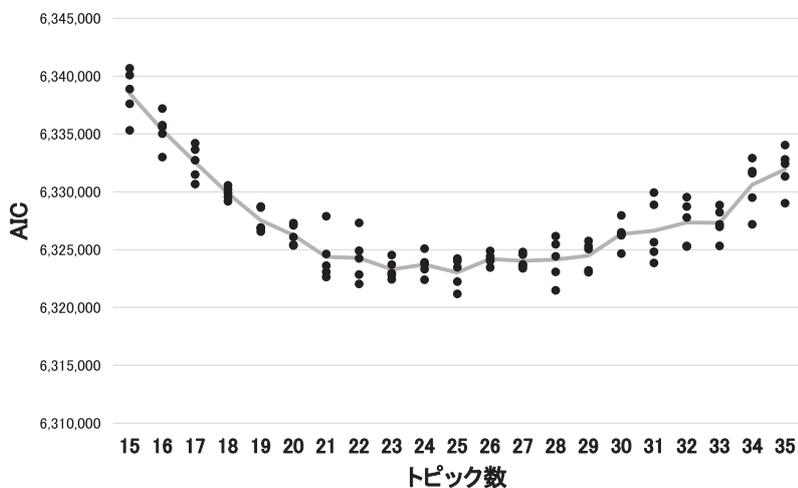


図 15 用途トピック選定のための各トピック数における PLSA の実行解の AIC 評価

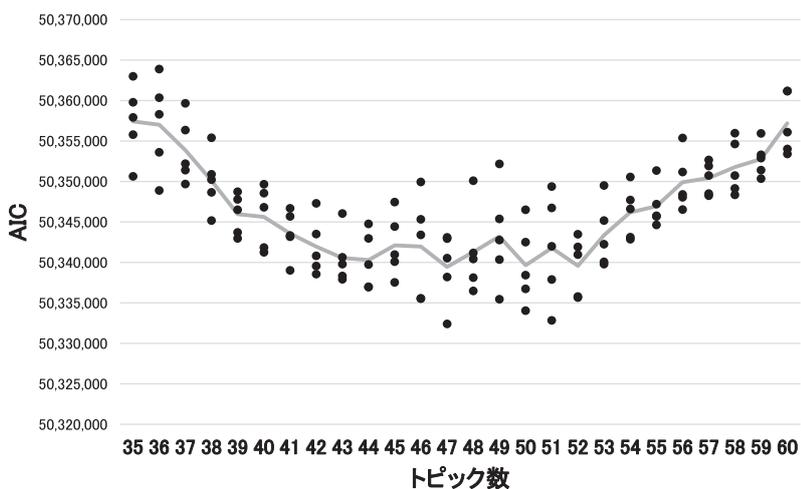


図 16 技術トピック選定のための各トピック数における PLSA の実行解の AIC 評価

#### (4) トピックの解釈

PLSA のアウトプットは, ①各トピックにおける行要素(単語)の所属確率, ②各トピックにおける列要素(係り受け)の所属確率, ③各トピックの周辺確率, という 3 つの確率が計算される。抽出された用途と技術のトピックの内容の例を表 4 に示す。なお, 単語と係り受けは所属確率の高い順に並べている。表 4 (左) の用途トピック U04 では, 単語は, 加湿装置, 水, 供給, 加湿, カビなどで所属確率が高く, 係り受けは, 加湿装置⇒提供, 加湿器⇒提供, ミスト発生装置⇒提供, 水⇒供給, 細菌⇒繁殖といった表現で所属確率が高い。つまり, この結果は加湿に関するトピックであると解釈できる。表 4 (右) の技術トピック T32 では, 単語は, 送風機, 塵埃, 掃除機, 分離, 吸い込む, 集塵部などで所属確率が高く, 係り受けは, 塵埃⇒分離, 分離⇒塵埃, 塵埃⇒含む, 吸い込む⇒塵埃, 含む⇒空気, 空気⇒分離といった表現で所属確率が高い。つまり, この結果は塵埃の分離に関するトピックであると解釈できる。このように「単語×係り受け」という構成の共起行列をインプットに PLSA でトピックを抽出することで, 行と列の両方の情報軸が十分な意味を持ち, 特に「単語」という話題の観点となる軸と, その観点の具体的な内容となる「係

り受け」という軸でトピックが構成され、解釈のしやすい結果を得ることができている。解釈をつけた 25 個の用途トピックと 47 個の技術トピックの一覧をそれぞれ表 5、表 6 に示す。

表 4 用途トピックと技術トピックの例

用途トピックU04				技術トピックT32			
確率	単語	確率	係り受け	確率	単語	確率	係り受け
5.5%	加湿装置	6.8%	加湿装置⇒提供	5.5%	送風機	2.1%	塵埃⇒分離
3.7%	水	3.1%	加湿器⇒提供	5.2%	塵埃	1.7%	分離⇒塵埃
3.3%	供給	2.9%	ミスト発生装置⇒提供	4.1%	掃除機	1.7%	塵埃⇒含む
2.4%	加湿	1.9%	水⇒供給	3.6%	分離	1.5%	吸い込む⇒塵埃
2.3%	カビ	1.7%	細菌⇒繁殖	3.5%	吸い込む	1.3%	含む⇒空気
2.1%	加湿器	1.5%	加湿⇒行う	2.3%	集塵部	1.0%	空気⇒分離
2.1%	発生	1.4%	加湿機能付空気清浄機⇒提供	1.9%	配置	1.0%	送風機⇒吸い込む
2.0%	繁殖	1.3%	ミスト⇒噴霧	1.9%	集塵容器	1.0%	発生⇒送風機
1.9%	ミスト	1.3%	繁殖⇒抑制	1.6%	旋回	0.9%	含塵空気⇒分離
1.7%	加湿性能	1.2%	十分⇒量	1.5%	含塵空気	0.9%	備える⇒掃除機
1.5%	ミスト発生装置	1.2%	カビ⇒発生	1.4%	捕集	0.8%	掃除機⇒設ける
1.4%	細菌	1.2%	効率⇒良い	1.3%	集塵室	0.8%	空気⇒吸い込む
1.3%	室内	1.2%	空気調和機⇒提供	1.3%	通過	0.8%	送風機⇒備える
1.3%	抑制+できる	1.1%	加湿⇒加湿装置	1.3%	発生	0.8%	掃除機⇒備える
1.1%	浴室	1.1%	空気⇒加湿	1.2%	集塵装置	0.7%	吸い込む⇒空気
...	...	...	...	...	...	...	...

表 5 用途トピックの一覧

No.	トピック名	No.	トピック名
U01	空調全般	U14	防止全般(流体の侵入、破損等)
U02	車両用空調	U15	騒音低減
U03	空調の省エネ、快適性	U16	消費電力の低減
U04	加湿	U17	機能向上全般
U05	乾燥機能(衣類等)	U18	熱交換器の機能向上
U06	空気浄化(除菌・脱臭)	U19	効率の良さ全般
U07	塵埃除去	U20	高性能・高付加価値(コストや安全性等)
U08	掃除機	U21	検出・測定の精度
U09	プリンタ	U22	構造の簡素化
U10	機器の冷却	U23	形成・配置(空気路等)
U11	熱の制御と利用	U24	方法・装置の提供
U12	制御(冷媒回路等)	U25	その他(環境破壊の懸念等)
U13	抑制全般		

表 6 技術トピックの一覧

No.	トピック名	No.	トピック名
T01	冷凍サイクル	T25	加湿
T02	冷却	T26	放電式ミスト生成
T03	車室内空調	T27	微細粒子の飛散(マイナスイオン等)
T04	空気路	T28	イオン発生・空気除菌・脱臭
T05	換気	T29	電解水生成と除菌
T06	排気	T30	空気浄化&効率性
T07	空気の吸込と吹出	T31	塵埃除去
T08	流体の流入と吐出	T32	塵埃分離
T09	空気流の利用と制御	T33	回転駆動
T10	送風	T34	電源と駆動制御
T11	空気の噴出	T35	運転と停止の制御
T12	送風搬送(紙葉類等)	T36	センサと制御(温度や風量等)
T13	印刷	T37	人検出
T14	光の利用(照射、発光等)	T38	風向制御
T15	ファンと機器冷却	T39	抑制・防止(騒音やコスト等)
T16	空気導入と車両エンジンの冷却	T40	構成・取り付け
T17	放熱	T41	接続
T18	除湿	T42	機器(熱交換器等)の配置
T19	乾燥機能	T43	配置と形成
T20	洗濯乾燥	T44	位置・形状・大きさ
T21	洗浄(衣類や食器等)	T45	位置の方向
T22	燃焼	T46	方法・装置
T23	加熱	T47	その他(発明目的、ケース構成等)
T24	温湿度制御と空気循環		

#### 5.4 トピックのスコアリング

続いて分析対象とした 30,039 件の特許データに対して、今回抽出された 25 個の用途トピックと 47 個の技術トピックのスコア（該当度）を計算する。スコアの計算では、1 件の特許の要約文章には複数の文が存在するため、まず文単位（句点「。」で区切られたセンテンス単位）に各トピックのスコアを計算し、それを特許単位に集約する。なお、文 S におけるトピック T のスコアは  $P(S|T)/P(S)$  で定義する。これは事後確率と事前確率の比率を示し、リフト値と呼ばれることもある指標である。トピック T を条件とすることでその文 S の発生確率が何倍になるのかを示すため、そのトピックをよく話題にしている文ほど値は高くなる。以下にこの指標の中の  $P(S|T)$  と  $P(S)$  の計算方法について説明する。

$P(S|T)$  については、文 S を単語 X で定義される文  $S_x$  と係り受け Y で定義される文  $S_y$  に分解し、それぞれについて  $P(S_x|T)$  と  $P(S_y|T)$  を計算し、それらを一つに統合して  $P(S|T)$  を計算する。 $P(S_x|T)$  と  $P(S_y|T)$  はそれぞれ式 (11) と式 (12) で計算される。単語 X と係り受け Y が含まれる文の数をそれぞれ  $N(X)$  と  $N(Y)$  とすると、式 (11) の  $P(S_x|X)$  は  $N(X)$  の逆数、式 (12) の  $P(S_y|Y)$  は  $N(Y)$  の逆数として計算される。式 (11) の  $P(X|T)$  と式 (12) の  $P(Y|T)$  はそれぞれ PLSA の実行結果によって得られている単語と係り受けの所属確率に該当する。そして  $P(S|T)$  は式 (13) で計算されるが、 $S_x$  と  $S_y$  は文 S において重みは同じであるため、 $P(S|S_x)$  と  $P(S|S_y)$  それぞれ 0.5 とする。また  $P(S)$  は式 (14) で計算され、 $P(T)$  は PLSA の実行結果によって得られているトピックの周辺確率に該当する。

$$P(S_x | T) = \sum_x P(S_x | X) P(X | T) \tag{11}$$

$$P(S_y | T) = \sum_y P(S_y | Y) P(Y | T) \tag{12}$$

$$P(S | T) = P(S | S_x) P(S_x | T) + P(S | S_y) P(S_y | T) \tag{13}$$

$$P(S) = \sum_T P(S | T) P(T) \tag{14}$$

以上から  $P(S|T)/P(S)$  で定義されるスコアを文単位に計算し、それを特許単位に見たとき、各トピックのスコアの最大値をその特許のトピックスコアとして採用する。さらにこのスコアの閾値を 3 に設定し、各特許データに対してそのトピックの該当有無を示す 1.0 のフラグ情報を付与する。なお、 $P(S|T)/P(S)$  で定義したスコアは本来 1 が基準の目安となる。つまりスコアが 1 より大きいということはトピック T を条件にすることで文 S の確率が上昇することであり、トピック T と文 S の間に関係があると判定できる。本分析事例では各トピックの特徴を抽出するため、特に関連の強い特許に対して該当ありのフラグを立てることを考え、またこのスコアの分布や実際の文章の内容も確認しながら、その閾値を基準の 3 倍と厳しく設定した。

以上の計算処理により、表 7 に示すようなデータが作成された。30,039 件の特許データには、出願年、出願人、要約文という情報が元々あるが、そこに加え、用途トピック 25 個、技術トピック 47 個の 1.0 のフラグ情報が付加されたデータとなる。つまりこれは Bag-of-Topics という形式のデータであり、PLSA でトピックを抽出しそのスコ

表 7 トピックのスコア（フラグ情報）を紐づけた特許データ

特許ID	出願番号	出願年	出願人	要約文		用途トピック	用途トピック	...	用途トピック	技術トピック	技術トピック	...	技術トピック
				【課題】	【解決手段】	U01	U02	U25	T01	T02	T47		
1	特願2006-XX	2006	A社	空気調和機の高	吸気口から導入	1	0	0	0	1	0	0	
2	特願2009-XX	2009	B社	短時間で除霜を	着霜検出手段が	0	1	0	1	0	0	0	
3	特願2011-XX	2011	C社	乾燥運転が中断	通風路を通して	0	0	1	1	0	0	0	
4	特願2013-XX	2013	D社	ウインドシールド	車両用空調装置	0	1	0	0	1	1	0	
...	...	...	...			...	...	...	...	...	...	...	
30,039	特願2012-XX	2012	Z社	プリ空調時に、除	冷暖房空調ユニ	0	1	0	1	1	0	0	

アを計算することで、特許文書の記述情報が集約された新たな分析軸が追加された。このデータセットを用いることでトピックを軸にしたさまざまな分析を実行することができる。なお、この先の各分析はすべてこのデータセットをベースとしている。

### 5.5 トピックのトレンド分析

表7のトピックのフラグデータを用いて、出願年の情報と、トピックのフラグ情報から、各トピックのトレンドを分析する。ここで得られた結果により、有望なニーズやシーズの探索に活用することができる。

具体的には出願年  $Y_r$  とトピック  $T$  の関連度を示す指標として  $P(Y_r|T=1)/P(Y_r)$  を計算し、この値の経年変化を可視化する。この指標はトピックのスコアリングで用いた事後確率と事前確率の比率を示すリフト値となっている。このリフト値はトピック  $T$  を条件とすることで出願年  $Y_r$  の元々の出願件数の割合が何倍になるのかを示している。各トピックで見れば、ある出願年の該当件数の割合が全体における該当件数の割合よりも大きいほど値は高くなる。つまりこのリフト値は、各年における元々の出願件数の大小の影響を除いた形で、トピックと出願年の関係の強さを表している。用途トピックと技術トピックにおいて、2013年からの直近3年でこのリフト値の上昇率が高い上位5つのトピックのトレンドをそれぞれ図17, 18に示す。

図17より、用途トピックでは、特に「U08. 掃除機」が上昇しており、それに関連してなのか「U06. 空気浄化（除菌・脱臭）」や「U07. 塵埃除去」も上昇している。一方図18より、技術トピックでは、「T32. 塵埃分離」や「T14. 光の利用（照射、発光等）」、「T19. 乾燥機能」、「T16. 空気導入・車両エンジンの冷却」に関する技術が上昇している。用途で掃除機のトピックが大きく上昇しており、技術では塵埃分離のトピックが上昇しているということは、直近ではサイクロン掃除機の需要と開発がホットであったと考えられる。

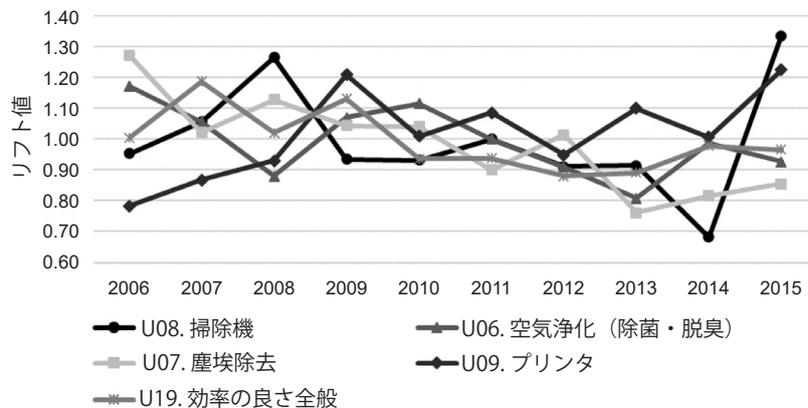


図17 2013年からの上昇率トップ5の用途トピックのトレンド

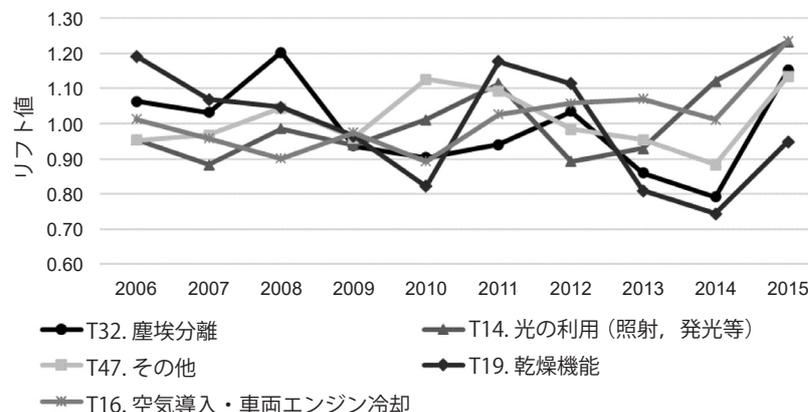


図18 2013年からの上昇率トップ5の技術トピックのトレンド

## 5.6 トピックを用いた競合他社の分析

表7のトピックのフラグデータを用いて、出願人の情報と、トピックのフラグ情報から、各トピックにおける出願人の特徴を分析する。ここで得られた結果により自社の技術開発戦略や差別化戦略、他社との協業戦略、自社技術の売却先候補などの検討に活用することができる。

### (1) 静的分析：出願人のポジショニングマップ

本分析事例では、まずトピックにおける各出願人のポジショニングを可視化する。具体的には出願人AとトピックTの関連度を示す「シェア」と「注力度」という2つの指標を計算し、縦軸にシェア、横軸に注力度を設定することで、トピックごとに各出願人をプロットしたポジショニングマップを作成する。シェアとは  $P(A|T=1)$  で定義され、そのトピックTが該当する全特許の中での出願人Aの出願割合を示す。出願件数が多いほど値が高く、そのトピックTにおけるシェアが高いということの意味する。注力度とは  $P(T=1|A)$  で定義され、その出願人Aが出願した特許の中におけるトピックTの該当割合を示す。この値が高ければそれだけトピックTに全社的に注力しているということであり、独自の特有な技術を保有している可能性がある。

ここでは先のトレンド分析で上昇していた技術トピック「T32. 塵埃分離」を例とした結果を図19に示す。塵埃分離の技術とはサイクロン掃除機に代表される遠心分離などの技術になる。図19より、まずシェアの高さを確認するとA社、B社、C社が高く、特にC社は注力度が非常に高いため、C社は特有の技術力があると考えられる。一方、E社、G社、I社はシェアは中程度だが、注力度は比較的高いため、技術力もあると考えられる。たとえば、高いシェアの企業は、中程度のシェアの企業と連携することで、より技術力を高めながらシェアを伸ばすことが期待できる。あるいは中程度のシェアの企業の間で連携をすることで、業界大手に対抗するということも考えられる。

このように塵埃分離に関する技術は、3社のシェアが高いものの、他にもある程度のシェア・注力度を持つ企業が何社か存在し、また先のトレンド分析でも近年ホットであったため、今後企業連携などの動きも十分考えられる領域と考察できる。

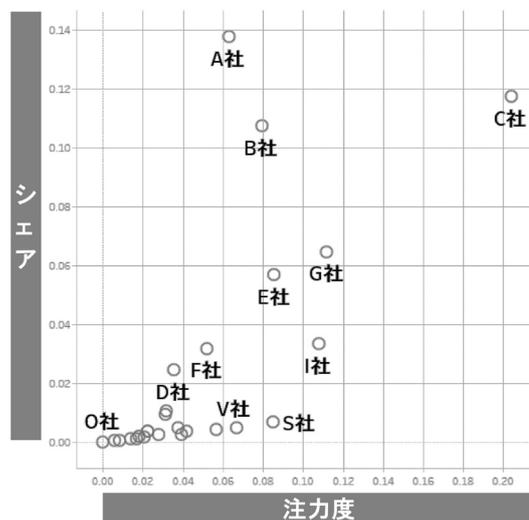


図19 技術トピック「T32. 塵埃分離」における出願人のポジショニングマップ

### (2) 動的分析：出願数の推移

こうしたポジショニングマップで可視化をすると各出願人の位置づけが分かりやすいが、このマップの結果はあくまでも今回の10年分の特許データをまとめた静的な結果である。つまり、このマップからはその時系列の変化といった動的な傾向までは把握できない。そこで、今回注目の対象となったA社、B社、C社、E社、G社、I社について、この「T32. 塵埃分離」という技術トピックに該当する特許の出願件数の推移を可視化した。その結果を図20に示す。

図20より、シェア1位のA社は、近年は出願が少なく、今はあまり力を入れて開発していない可能性が考えられる。シェア3位のB社は、ここ10年で徐々に出願が増えており、今特に力を入れている技術である可能性がある。

注力度 1 位シェア 2 位の C 社は、10 年ほど前には出願件数が多く、その後一度落ち着き、直近でまた出願が急激に増えているため、再び力を入れ始めている可能性がある。E 社、G 社、I 社は、A 社、B 社、C 社と比べて全体の件数は少ないが、たとえば G 社は近年出願が増えており、徐々に開発に力を入れている可能性があり、このトピックの領域において要注目と思われる。E 社は、近年では出願件数の変化が少なく、I 社は、すでに他社に買収されている会社でもあるため、2013 年以降の出願は存在していない。ここから、「T32. 塵埃分離」という技術領域では、高いシェアを誇る企業で近年競合関係にあると考えられるのは特に B 社と C 社である。またシェアは低いものの徐々に出願を増やしている G 社の動向も今後注目すべきといえる。

このように全体でのポジショニングマップでは静的な出願人の位置づけを把握できるが、時系列での出願件数の推移も組み合わせて確認することで、動的な出願人の動向も把握することができる。こうした結果からさまざまな技術戦略を検討するヒントが得られると期待できる。

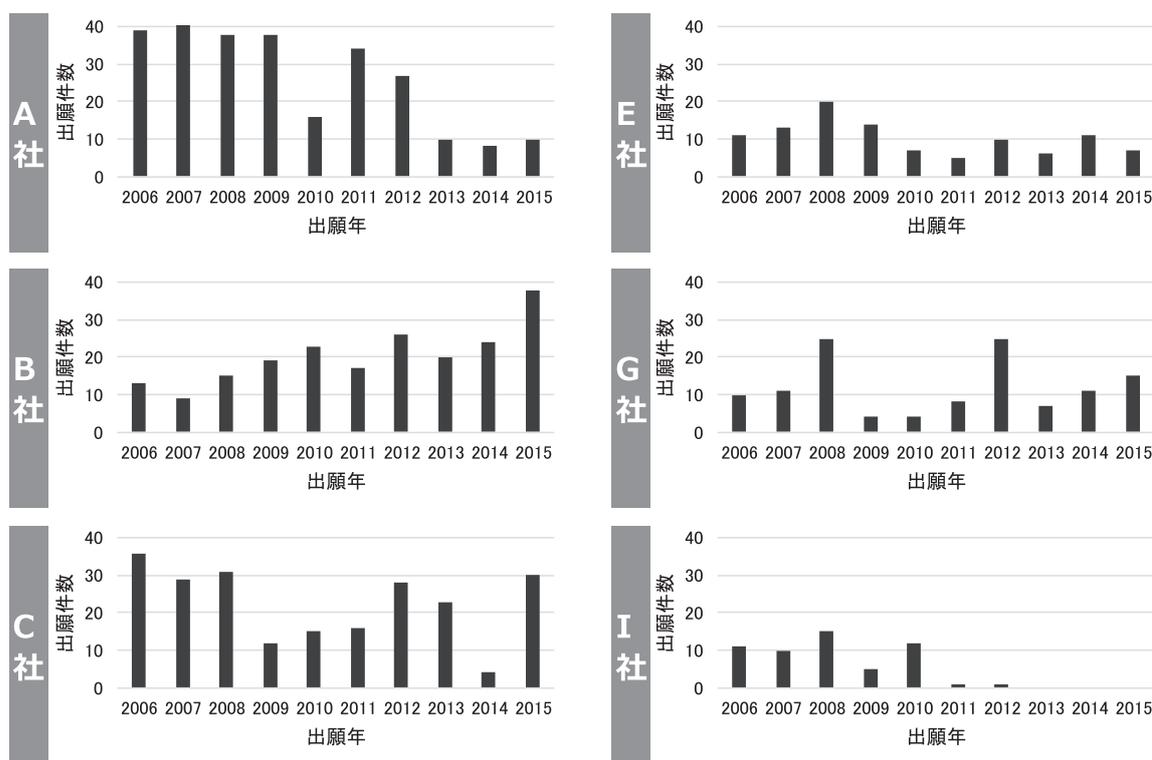


図 20 技術トピック「T32. 塵埃分離」における出願人の出願件数の推移

### 5.7 ベイジアンネットワークの適用による用途と技術の関係分析

ここからはベイジアンネットワークを適用した用途と技術の関係分析について解説する。表 7 のトピックのフラグデータを用いて、用途トピックのフラグ情報と、技術トピックのフラグ情報を確率変数とし、ベイジアンネットワークを適用することで用途トピックと技術トピックの関係構造を分析する。

この関係分析は、①用途⇒技術の分析（用途に対して関係する技術の分析）と②技術⇒用途（技術に対して関係する用途の分析）という 2 つのパターンがあり、それぞれはベイジアンネットワークのリンク構造が逆転する。①用途⇒技術の分析では、自社で検討しているある用途の事業を実現する際に、重要となる要素技術を把握するための分析である。ここで得られた結果により、その用途を達成するためにどのような技術開発に注力すべきか、またその技術領域で競合となりそうな他社はどこか、そのなかで他社が牛耳る技術の代替技術は存在するか、どの会社と技術提携すると効果的かなど、自社の開発戦略や他社との協業戦略の検討に活用することができる。一方②技術⇒用途の分析では、自社で保有している技術と関係のある用途を把握するための分析である。ここで得られた結果により、自社

技術と関係のある用途のうち自社でまだ想定していない用途を見つけ、自社の技術をさらに有効活用できる新しい用途展開のアイデアを考えることができる。以下にそれぞれのパターンの分析結果と考察の例を解説する。

### 5.7.1 用途⇒技術の関係分析

用途に対して関係する技術の分析では、用途トピック 25 個をリンク元に、技術トピック 47 個をリンク先に指定してベイジアンネットワークのモデルを構築し、用途に対する技術の関係構造を可視化する。構築されたモデルの結果を図 21 に示す。なおベイジアンネットワークのモデル構築には BayoLink（株式会社 NTT データ数理システム）を使用した。

#### (1) 各用途と関係する技術の把握

図 21 のモデルを用いることで、各用途トピックと確率統計的に関係があると判定された技術トピックを把握できる。特にベイジアンネットワークの確率シミュレーションにより、ある用途トピックを条件に与えたときの各技術トピックの確率分布を推論できるため、その用途条件下で確率が上昇するような関係性の強い技術トピックを把握できる。ここでは先のトレンド分析で上昇していた用途トピック「U06. 空気浄化（除菌・脱臭）」を対象に関係の強い技術トピックを確認した例を紹介する。図 21 のモデルを用いて、用途トピック「U06. 空気浄化（除菌・脱臭）」を条件に与えたときに、これと関係構造を持つ技術トピックの確率を推論した結果を図 22 に示す。図 22 のグラフでは、用途トピック U06 と関係が見られた各技術トピックの元々の確率（事前確率）と、用途トピック U06 を条件に与えたときの条件付確率（事後確率）を示している。どの技術トピックも事後確率の方が高くなっているため、用途トピック U06 との関係が強いことがわかる。したがって、「U06. 空気浄化（除菌・脱臭）」の用途と関係の強い技術トピックは、「T26. 放電式ミスト生成」、「T28. イオン発生・空気除菌・脱臭」、「T29. 電解水生成と除菌」、「T30. 塵埃吸込 & 効率性」、「T47. その他」と確認できる。

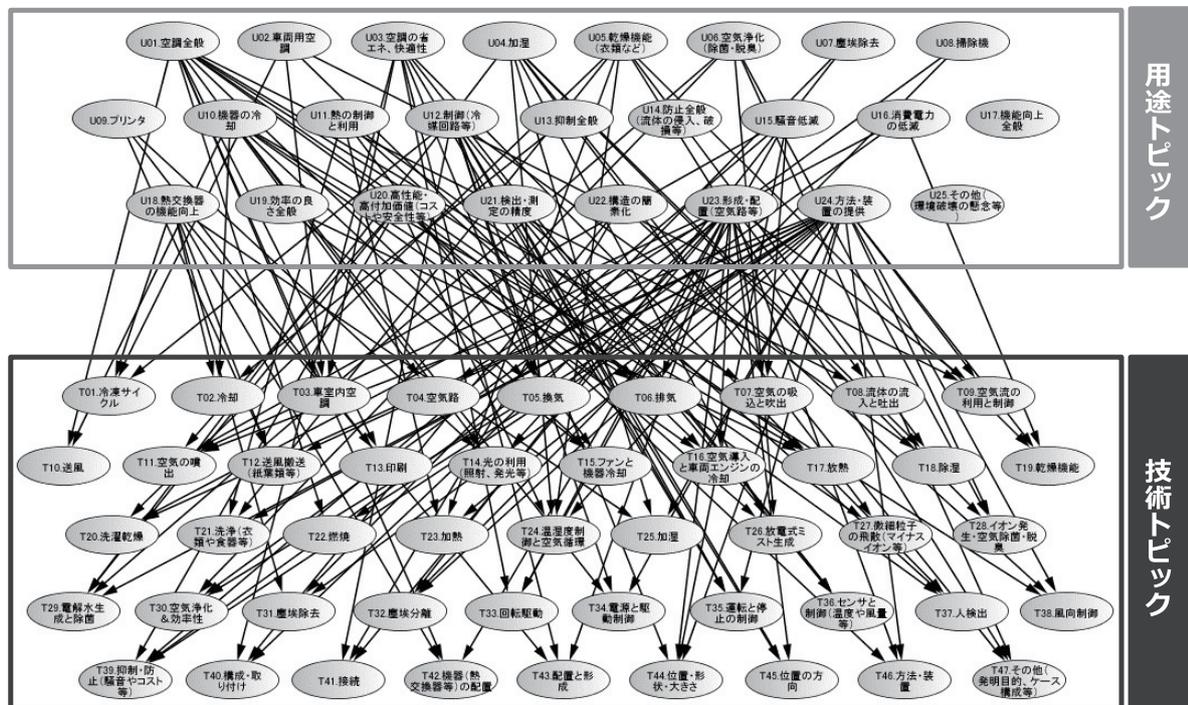


図 21 ベイジアンネットワークを適用した用途⇒技術の関係モデル

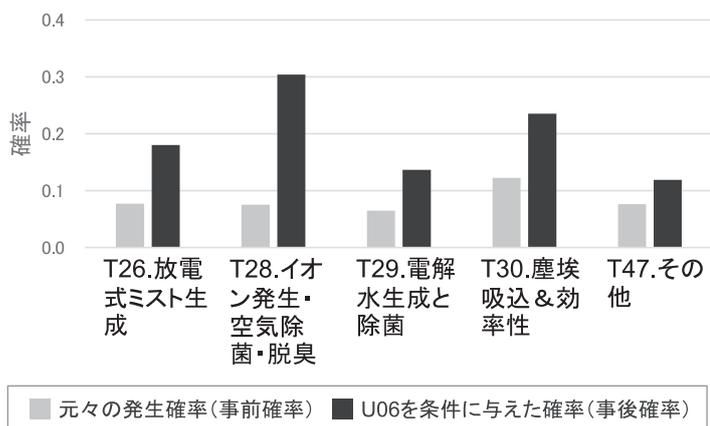


図 22 用途トピック「U06. 空気浄化」を条件に与えたときに確率が上昇する技術トピック

(2) 用途×技術における競合他社の動向の把握

続いて、用途トピック「U06. 空気浄化（除菌・脱臭）」と関係が見られた「T.47 その他」を除く 4 つの技術トピックについて、各技術を保有している出願人を確認するため、先ほどと同様の競合他社の分析を実施した。ここでは、用途トピック「U06. 空気浄化（除菌・脱臭）」が該当する特許データを対象に、「T26. 放電式ミスト生成」、「T28. イオン発生・空気除菌・脱臭」、「T29. 電解水生成と除菌」、「T30. 塵埃吸込 & 効率性」について、それぞれ各出願人のシェアと注力度を計算してポジショニングマップを作成した。その結果を図 23 に示す。たとえば「T26. 放電式ミスト生成」は、シェアは A 社と G 社が高いが、高シェア高注力度のポジションは空いていることがわかる。「T28. イオン発生・空気除菌・脱臭」と「T29. 電解水生成と除菌」は一社が高シェア高注力度のポジションを確立した一強状態にある技術領域であり、T28 は G 社が、T29 は I 社が牛耳っている技術であることがわかる。「T30. 塵埃吸込 & 効率性」は、シェアは A 社が高いが、T26 と同じく高シェア高注力度のポジションは空いている。

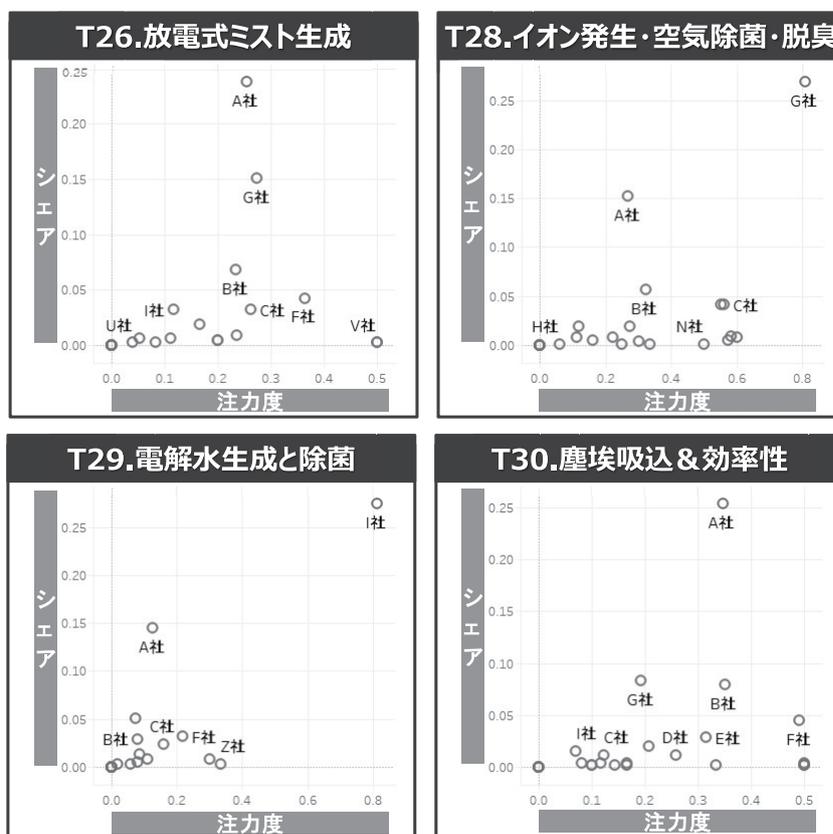


図 23 用途トピック「U06. 空気浄化」と関係のある技術トピックの出願人ポジショニングマップ

この結果より、「U06. 空気浄化（除菌・脱臭）」の用途を実現しようとしたとき、たとえば一強状態の技術を避けるのであれば、T26 や T30 の技術が狙い目と考えることができるかもしれない。あるいは、逆に一強状態にある T28 や T29 の技術においては、その一強企業と提携したり、M&A で買収を実現すればその技術領域ごと獲得できることになる。実際に「T29. 電解水生成と除菌」を牛耳る I 社はすでに買収されており、その買収した会社からは電解水（次亜塩素酸）で空気を洗うという全く新しい空気浄化家電が、2013 年に業務用として、2017 年には家庭用として発売されている。つまり、その家電のコア技術となる電解水を生成する技術は I 社で培われた技術であったと考察できる。

このように自社で事業化を検討している用途について、それと関係する重要な要素技術を分析することで、その用途を達成するためにどのような技術開発に注力すべきか、また競合となりそうな他社はどこか、他社が牛耳る技術を回避するような代替技術は存在するか、あるいはどの会社と連携すると効率的にその技術を獲得できるかといった、開発戦略や協業戦略を検討することができる。

### 5.7.2 技術⇒用途の関係分析

技術に対する用途の関係分析では、先ほどの用途⇒技術の関係分析におけるモデルのリンク構造を逆転させ、技術トピック 47 個をリンク元に、用途トピック 25 個をリンク先に指定してベイジアンネットワークのモデルを構築し、技術に対する用途の関係構造を可視化する。図 24 に構築されたモデルの結果を示す。

#### (1) 各技術と関係する用途の把握

図 24 のモデルを用いることで、各技術トピックと確率統計的に関係があると判定された用途トピックを把握できる。また、ある技術トピックを条件に与えたときの各用途トピックの確率分布を推論することができるため、特にその技術条件下で確率が上昇するような関係性の強い用途トピックを把握することができる。ここでは技術トピック「T18. 除湿」と「T32. 塵埃分離」を対象に関係の強い用途トピックを確認した例を紹介する。図 24 のモデルの結果より、技術トピック「T18. 除湿」と「T32. 塵埃分離」を条件に与えたときに、これと関係構造を持つ用途トピックの確率を推論した結果をそれぞれ図 25, 26 に示す。図 25, 26 のグラフでは、技術トピック T18 あるいは T32 と関係が見られた各用途トピックの元々の確率（事前確率）と、技術トピック T18 あるいは T32 を条件に与えたときの条件付確率（事後確率）を示している。どの用途トピックも事後確率の方が高くなっているため、これらの用途と技術の関係が強いことがわかる。特に「T32. 塵埃分離」に対して関係が見られた用途トピックはただ一つ「U08. 掃除機」のみであり、事後確率の上昇の仕方も顕著である。

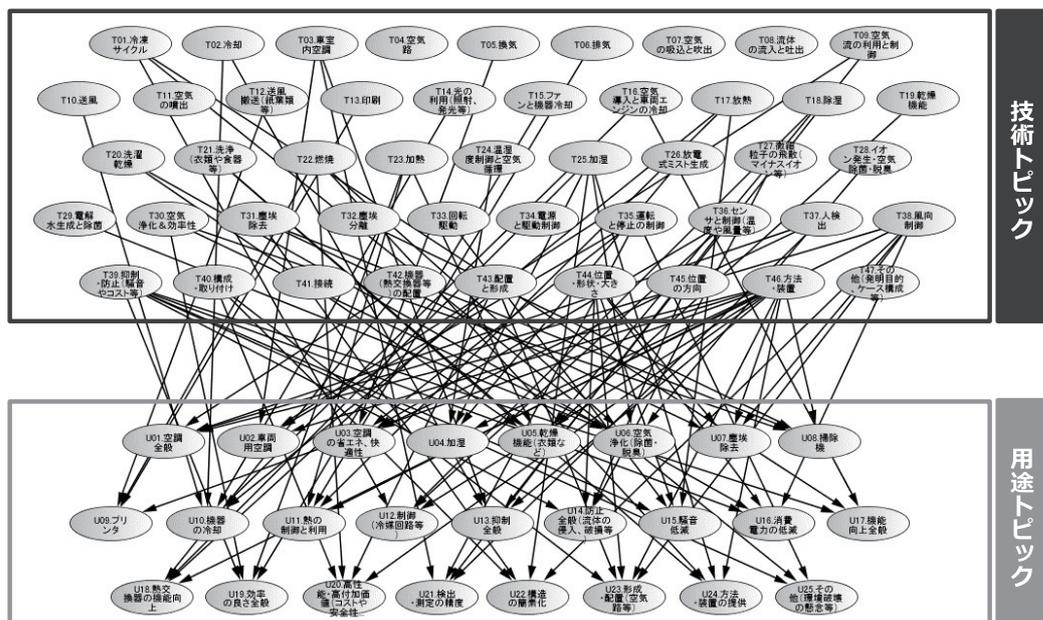


図 24 ベイジアンネットワークを適用した技術⇒用途の関係モデル

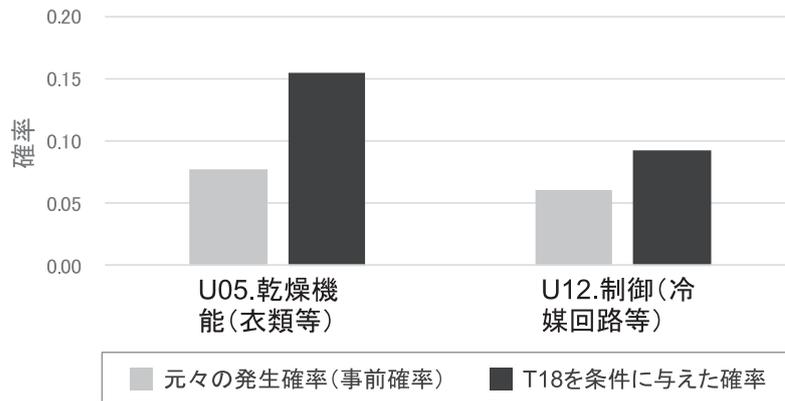


図 25 技術トピック「T18.除湿」を条件に与えたときに確率が上昇する用途トピック

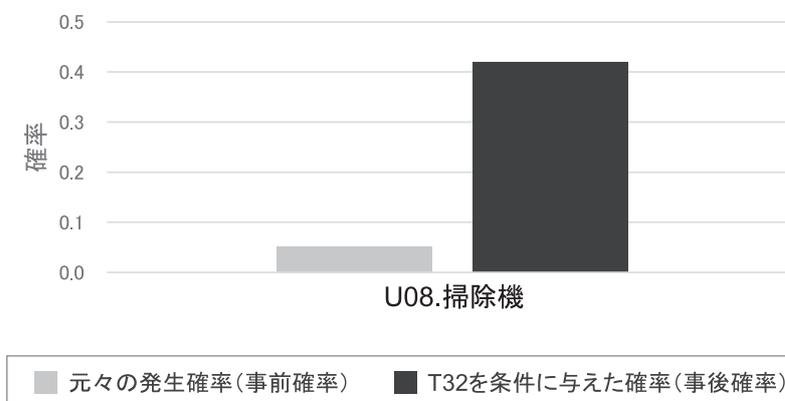


図 26 技術トピック「T32.塵埃分離」を条件に与えたときに確率が上昇する用途トピック

さらにここから、技術「T18.除湿」×用途「U05.乾燥機能(衣類等)」の関係と、技術「T32.塵埃分離」×用途「U08.掃除機」の関係を取り上げ、技術の新しい用途展開を探索する考え方について紹介する。図 25,26 の結果から、当然 T18 の除湿の技術は U05 の衣類などの乾燥機能の用途を想定して出願されている特許が多く、T32 の塵埃分離の技術は U08 の掃除機の用途を想定して出願されている特許が多いということになる。しかし、この傾向に当てはまらない特許も存在している。技術の新しい用途展開を探索する上ではそうした特許がヒントになり得る。以下に具体的な探索例を紹介する。

#### (2) 技術「T18.除湿」×用途「U05.乾燥機能(衣類等)」から探索する新規用途

まず「T18.除湿」の技術が「U05.乾燥機能(衣類等)」の用途を想定して出願されている特許の要約文の典型例を図 27 (左) に示す。これはドラム式洗濯乾燥機の特許である。特許の要約の内容を確認すると、洗濯物を短い時間でムラなく乾燥させ、乾燥工程の時間を短くするための除湿技術として出願されている。一方、「T18.除湿」の技術が「U05.乾燥機能(衣類等)」の用途を想定しないで出願されている特許の要約文の例を図 27 (右) に示す。これはインクジェットプリンタに関する特許である。特許の要約の内容を確認すると、インク液を吸収した紙の湿気をムラなく取り除いて、コックリング(紙の波打ち)を防ぐ除湿の技術として出願されている。この 2 つの特許の比較から、たとえば、プリンタという空間の中で、インク液を吸収した用紙の湿気をムラなく取り除いて紙の波打ちを防ぐ除湿技術は、洗濯乾燥機という空間の中で、洗濯物をムラなく効率的に乾燥させることにも応用できる可能性を考えることができる。

T18の技術がU05の用途を想定して出願された特許例	T18の技術がU05の用途を想定せず出願された特許例
<p><b>【発明の名称】</b> ドラム式洗濯乾燥機</p> <p><b>【課題】</b> 洗濯物を短い時間でムラ無く乾燥させ、乾燥工程の時間を短くすることができるドラム式洗濯乾燥機を提供する。</p> <p><b>【解決手段】</b> 送風機に吸い込まれた空気は、風路切替弁の切り替えにより、ドラム開口部に対向する前側吹出口へ流れたり、回転ドラムの後部に設けられた後側吹出口へ流れたりする。制御装置が風路切替弁の切り替えを制御することによって、恒率乾燥過程時、前側吹出口から乾燥用空気が吹き出し、かつ、減率乾燥過程時、後側吹出口から乾燥用空気が吹き出す。これにより、恒率乾燥過程において乾燥用空気が効果的に当たらなかった、回転ドラムの後端壁側の洗濯物に、乾燥用空気が減率乾燥過程で効果的に当たる。</p>	<p><b>【発明の名称】</b> インクジェット記録装置及び画像記録方法</p> <p><b>【課題】</b> 処理液の厚みムラを低減するとともに処理液による用紙のコックリングを低減することで、高品質かつ高速の画像記録を可能とするインクジェット記録装置及び画像記録方法を提供する。</p> <p><b>【解決手段】</b> 記録媒体に処理液を付与する処理液付与部の後段には、記録媒体表面に残存する溶媒を蒸発させるプレ加熱部が設けられている。プレ加熱部はIRプレヒータにより記録媒体表面を輻射加熱するとともに、吸引ファンにより記録媒体表面の湿り空気を置換する。液状の処理液が不均一にならないように乾燥処理を施すことで、均一な膜厚を持つ固体状の凝集処理層が形成される。その後、本加熱部による熱風噴射加熱により、コックリング量が所定量以下になるように本加熱処理が施される。</p>

※例示のための要約文であり、一部内容は筆者が加工している

図 27 技術「T18. 除湿」が応用されている 2 つの用途の特許例

(3) 技術「T32. 塵埃分離」×用途「U08. 掃除機」から探索する新規用途

同様に技術「T32. 塵埃分離」と用途「U08. 掃除機」の関係についても探索していく。まず「T32. 塵埃分離」の技術が「U08. 掃除機」の用途を想定して出願されている特許の要約文の典型例を図 28（左）に示す。これはサイクロン掃除機の特許である。特許の要約の内容を確認すると、サイクロン掃除機の排気筒の詰まりを防止して、集塵性能を向上させる技術として出願されている。一方、「T32. 塵埃分離」の技術が「U08. 掃除機」の用途を想定しないで出願されている特許の要約文の例を図 28（右）に示す。これは画像形成装置、つまりプリンタに関する特許である。このプリンタでは、トナーが含まれる空気をサイクロンで分離して回収しており、そのサイクロン部の清掃時期をセンサで判断し、自動で清掃モードという処理を実行することで、トナーの分離効率の低下を抑制するという技術として出願されている。この 2 つの特許の比較から、たとえば、プリンタでトナーを分離・回収するサイクロン部の清掃時期を、センサで判断して分離効率を維持する技術は、サイクロン掃除機の集塵性能の向上にも応用できる可能性を考えることができる。

T32の技術がU08の用途を想定して出願された特許例	T32の技術がU08の用途を想定せず出願された特許例
<p><b>【発明の名称】</b> 電気掃除機</p> <p><b>【課題】</b> 集塵性能が向上しメンテナンスの軽減が図れる電気掃除機を提供すること。</p> <p><b>【解決手段】</b> 塵埃を含む空気を巡回させ塵埃分離する略円筒状の 1 次旋回室と、1 次旋回室に連通した 2 次旋回室と、1 次旋回室の下方に位置し塵埃を溜める集塵室と、塵埃を圧縮する圧縮板と、塵埃が流入する流入口を有し、圧縮板の底面の一部に突出部を流入口から見て集塵室の奥側に配設する構成としたことより、集塵室内に入った塵埃は、圧縮板の突出部に引っかかり動きが止められ、流れに乗って 2 次旋回室や 1 次旋回室側に戻ることが無い。そのため集塵性能が向上し、排気筒の詰まり防止によるメンテナンスの軽減を図ることができる。</p>	<p><b>【発明の名称】</b> 画像形成装置</p> <p><b>【課題】</b> サイクロン部の清掃時期を適正に判断して、トナーの分離効率の低下を抑制することが可能な画像形成装置を提供する。</p> <p><b>【解決手段】</b> 画像形成装置は、トナー含有空気からトナーを遠心分離するサイクロン部と、サイクロン部によって分離されたトナーを回収する回収部と、サイクロン部によってトナーが分離された空気を通過させ、残留トナーを捕集するフィルタ部と、空気を吸引する送風部と、フィルタの汚れを検知する汚れ検知センサが設けられたトナー捕集部を備え、汚れ検知センサで検知されたフィルタの汚れから推定した風量と、風速センサで取得した風量の実測値の差分が、サイクロン清掃閾値を超えたと判断すると、サイクロン部の清掃モードを実行する。</p>

※例示のための要約文であり、一部内容は筆者が加工している

図 28 技術「T32. 塵埃分離」が応用されている 2 つの用途の特許例

#### (4) イノベーションの鍵となるアイデアの創出

このように自社で保有している技術と関係のある用途を把握し、そのうちまだ自社で想定していない用途を見つけることで、自社の技術をさらに有効活用できる新しい用途展開のアイデアを創出できる。自社で注力してきた技術を、そっくりそのまま別の用途に転用できるようなことはなかなかないかもしれないが、これまで自社で培ってきた技術や経験と関連のある用途をいかに発想できるかということがイノベーションの鍵となる。ここで紹介した例はあくまでも分析結果から筆者が発想したアイデアであり、現実性は検討していないが、こうした分析を業界の知識・経験が豊富な技術担当者が実施していくことで、これまで発想していなかった新しい用途展開の気づきが得られるものと期待できる。

#### 5.8 Nomolytics を適用した特許文書データ分析のまとめ

分析事例の最後として、ここでは Nomolytics を適用した特許文書データ分析のメリットを二つにまとめる。

まず一つ目のメリットは、PLSA を適用することで単語ではなく集約されたトピックを軸に分析を実行でき、膨大な特許情報に潜む傾向をわかりやすく理解できることである。従来のテキストマイニングのみを適用した特許文書データの分析では、大量の単語をベースにした複雑な可視化アウトプットとなるため、それを解釈することが難しいという課題があった。また、その大量の単語を人がグルーピングしていくつかのカテゴリを作成し、カテゴリベースに分析することもあるが、そのカテゴリの作成が属人的で作業負荷も大きいという課題もあった。これに対して Nomolytics の分析では、特許文書全体に存在するトピックを PLSA で機械的に抽出して分類・整理でき、単語ではなくそのトピックをベースにトレンドや出願人の動向を可視化することができる。これにより、膨大な特許情報に潜む特徴をシンプルに分かりやすく把握できる。

二つ目のメリットは、ベイジアンネットワークを適用することで、用途と技術の統計的な関係を把握でき、各用途を実現するための重要技術を確認して技術戦略を検討したり、自社技術を有効活用できる新規用途のアイデアを創出できることである。従来の特許分析でも、図 3 に示したように用途と技術の関係を分析するアプローチはあったが、「課題」と「解決手段」それぞれに対して人がグルーピングして作成したカテゴリのクロス集計をし、その対応関係を考察するというもので、統計的な関係までは分析できていなかった。これに対して Nomolytics の分析では、PLSA によって客観的に抽出されたトピックをベースに課題と解決手段の統計的な関係性をベイジアンネットワークで把握できる。そしてその構築された関係モデルを用いることで、たとえば、事業化を検討しているある用途に対して関係の強い技術を確認し、その技術の出願人の動向から自社の技術戦略を検討したり、あるいは自社の保有技術と関係の強い用途を見つけ、そこでまだ想定していない用途を確認することで、自社技術の新しい用途展開のアイデアを創出することなどに活用できる。

このように、従来のテキストマイニングに PLSA とベイジアンネットワークという 2 つの AI 技術を組み合わせた Nomolytics で特許文書データを分析することで、人間では読み切れない膨大な特許文書に潜む傾向や要因関係を把握でき、企業の技術戦略の検討において有益な気づきを得る新たな切り口を提供することができる。今回は特許の要約文を分析対象としており、最終的には特許文書の明細まで詳細に確認することが求められるが、こうした要約文のトピックをベースにした分析を進めることで、詳細まで確認すべき重要な特許を効率的に絞り込んでいくことができる。

### 6. Nomolytics を特許文書データに適用した分析事例（その 2）

Nomolytics を実際の特許文書データに適用した分析事例をもう一つ紹介する。なお、二つ目の事例も特許の要約文を分析対象としているが、一つ目の事例のように要約文を「課題」の文章と「解決手段」の文章に分けてトピックを抽出するのではなく、今回は要約文の全文を対象にトピックを抽出する。また、一つ目の事例で適用したベイジアンネットワークはここでは適用せず、Nomolytics のアプローチの中でも PLSA でトピックを抽出し、そのトピックを軸に傾向を分析するまでの事例を紹介する。

## 6.1 分析で用いるデータ

本分析事例では、特許の要約と請求項に「車」「電気」という2つのキーワードを含む国内の特許公報データ26,419件を分析対象としている。出願期間は2007年1月1日から2016年12月31日までのちょうど10年分の特許公報を対象に抽出した。今後開発と普及の拡大が見込まれる「電気自動車」に関連した技術を分析する狙いから抽出したデータである。

## 6.2 分析の全体像

今回は特許の要約文全文を対象にトピックを抽出し、その傾向を分析するが、分析プロセスの全体像は一つ目の分析事例の図14で示した(4)のステップまでが該当する。すなわち(1)テキストマイニング×PLSAによるトピックの抽出、(2)トピックのスコアリング、(3)トピックのトレンドの分析、(4)トピックを用いた競合他社の分析というステップである。要約文を「課題」と「解決手段」に分割してトピックを抽出すること以外は、基本的に一つ目の分析事例と同様のアプローチで分析を進める。以下に分析ステップごとに説明する。

## 6.3 テキストマイニング×PLSAによる要約トピックの抽出

テキストマイニングとPLSAを用いて特許の要約文全体の内容をいくつかのトピックに集約する。そのアプローチは基本的に一つ目の分析事例と同様であるが、以下にトピック抽出の手順を述べる。

### (1) テキストマイニング

要約文にテキストマイニングを適用し、単語とその文法的なペアとなる係り受け表現を抽出する。単語は今回は名詞のみを抽出し、係り受けは名詞に対する動詞・形容詞・形容動詞の単語ペアおよびその逆の係り受け(動詞・形容詞・形容動詞に対する名詞)の単語ペアを抽出する。テキストマイニングの実行にはText Mining Studio(株式会社NTTデータ数理システム)を使用した。本ツールの形態素解析において、文章の区切りを定義する文字列には「句点(.)」の他に、特許要約文の中で文章のラベル付けによく使用される「隅付括弧(【】)」も設定した。

一つ目の分析事例と同様に、テキストマイニングの形態素解析で抽出された単語に対して、目視で類義語辞書を作成し、テキストマイニングのツールに反映した。作成した類義語辞書は、3,278語の類義語に対して1,343語の代表語を登録した。類義語辞書の例を表8に示す。表8において、背景が灰色の単語が代表語となっている。

表8 電気自動車に関連する特許の要約文から作成した類義語辞書の例

<b>充電スタンド</b>	<b>簡素化</b>	<b>フランジ</b>	<b>Li</b>
充電スタンド	簡素化	フランジ	Li
充電ステーション	簡略化	フランジ部	リチウム
充電設備	単純化	鏝部	<b>SMR</b>
充電施設	簡易化	<b>誤検出</b>	SMR
給電設備	<b>先端</b>	誤検出	システムメインリレー
給電スタンド	先端	誤検知	<b>車軸</b>
<b>給電装置</b>	先端部	誤判定	車軸
給電装置	先端側	<b>小型化</b>	アクスル
送電装置	先端部分	小型化	<b>共重合体</b>
電力供給装置	<b>被覆</b>	小型軽量化	共重合体
電源供給装置	被覆	コンパクト化	コポリマー
給電機器	コーティング	<b>送風機</b>	<b>羽根車</b>
<b>省エネ</b>	被膜	送風機	羽根車
省エネ	皮膜	送風装置	インペラ
省エネルギー	<b>電子機器</b>	プロフ	<b>接地線</b>
省エネルギー化	電子機器	<b>漏洩</b>	接地線
省電力	電子装置	漏洩	アース線
省電力化	電子機器用	漏れ	<b>緊急時</b>
<b>サスペンション</b>	電子デバイス	漏出	緊急時
サスペンション	<b>固定子巻線</b>	<b>熱伝導</b>	非常時
サスペンション装置	固定子巻線	熱伝導	<b>僅か</b>
懸架	ステータコイル	熱伝達	僅か
懸架装置	ステータ巻線	伝熱	わずか

## (2) 共起行列の作成

PLSA でトピックを抽出する際のインプットとする共起行列を作成する。共起行列は「単語×係り受け」という構成で、それぞれの単語と係り受けが同時に出現する共起頻度（同時に出現する文章数）を集計したデータを用いる。なお共起行列の構成に採用する単語と係り受けは、特許単位でカウントした頻度が 20 件以上のものを対象とし、「名詞単語（3,020 語）×係り受け（2,128 表現）」の共起行列を作成した。共起行列の例を表 9 に示す。表 9 は文章単位にカウントした頻度が上位 10 個の単語と係り受けに限定した共起行列を例として掲載している。なお、共起行列を作成するにあたり、多くの特許で共通して使用される単語で、かつ頻度が高すぎる単語（「課題」「解決手段」「選択図」「提供」など）や、重要な意味を持たない単語（「前記」「本発明」「当該」「個」など）は、トピック抽出においてノイズになり得るため、ストップワードとして対象から除外した。

表 9 電気自動車関連の要約文から作成した共起行列の例

		全体の 頻度	1,350	539	529	494	444	368	306	289	282	280
係り受け		電力⇒ 供給	否⇒判 定	バッテ リ⇒充 電	モー ター ⇒駆 動	効率⇒ 良い	供給⇒ 電力	充電⇒ 行う	電気自 動車⇒ 提供	並列⇒ 接続	モー ター ⇒供給	
全体の 頻度	単語											
5,880	構成	118	33	33	36	24	32	25	10	46	30	
5,188	モータ	239	61	58	494	31	54	19	33	27	280	
5,092	制御	268	73	85	108	12	115	41	2	40	74	
4,655	電気自動車	193	73	129	79	59	53	114	289	13	54	
4,604	配置	69	2	8	29	6	15	9	5	12	18	
4,602	バッテリー	337	87	529	60	44	81	72	20	70	106	
4,006	形成	31	4	8	20	1	7	5	0	12	4	
3,978	供給	1,350	43	85	53	10	368	43	7	43	280	
3,960	検出	134	99	36	56	4	44	19	8	27	33	
3,629	電力	1,350	50	127	75	35	368	57	20	36	189	

## (3) PLSA の実行と評価

作成した共起行列に PLSA を適用することで、使われ方の似ている単語と係り受けでまとめられたトピックを抽出する。PLSA の実行は一つ目の分析事例と同様にトピック数を 1 刻みで変化させ、それぞれのトピック数に対して初期値をランダムに変えて PLSA を 5 回ずつ実行し、それぞれの解を情報量基準 AIC で評価して最も評価の良い解を採用する。PLSA の実行解の AIC の評価の結果を図 29 に示す。図 29 では、トピック数を 20 から 45 まで 1 刻みで変化させて実行した結果となっている。各トピック数に対して 5 つの黒いプロット一つひとつが初期値を変えて PLSA を実行した結果であり、灰色の折れ線グラフはその 5 つの実行解における AIC の平均値を示している。灰色の折れ線が示す AIC の平均値を比較したとき、それが最小となるトピック数は 34 個となり、このトピック数 34 個における 5 つの実行解の中で AIC が最小となる結果を採用した。なお、PLSA の実行には Visual Mining Studio（株式会社 NTT データ数理システム）の二項ソフトクラスタリングを使用した。

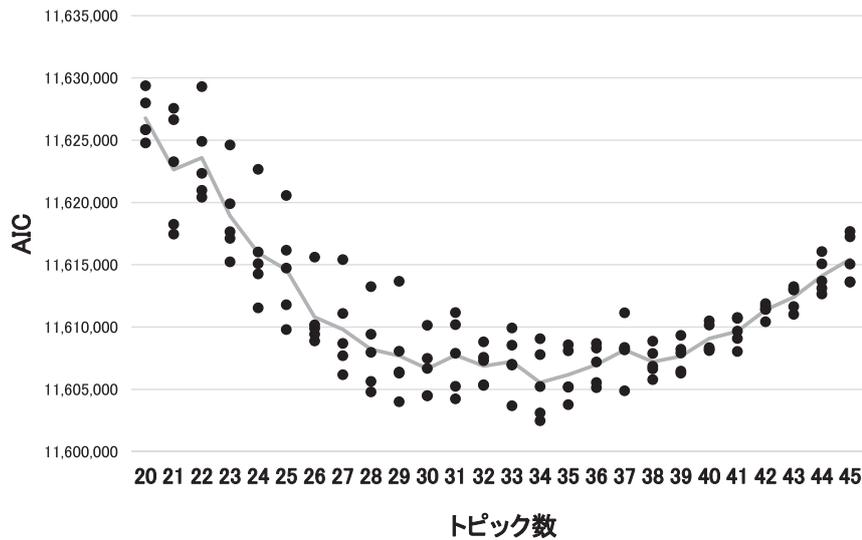


図 29 要約トピック選定のための各トピック数における PLSA の実行解の AIC 評価

(4) トピックの解釈

抽出されたトピックの内容の例を表 10 に示す。単語と係り受けは所属確率の高い順に並べている。表 10 (左) のトピック Z02 では、単語は、連結、伝達、モータ、出力軸、トランスミッション、クラッチなどで所属確率が高く、係り受けは、介する⇒伝達、介する⇒連結、動力⇒伝達、出力軸⇒連結、駆動力⇒伝達、トルク⇒伝達といった表現で所属確率が高い。つまり、この結果は動力の伝達に関するトピックであると解釈できる。表 10 (右) のトピック Z13 では、単語は、充電、電気自動車、蓄電装置、バッテリー、充電システム、蓄電池などで所属確率が高く、係り受けは、バッテリー⇒充電、充電⇒行う、電気自動車⇒充電、蓄電池⇒充電、蓄電装置⇒充電、電力⇒供給といった表現で所属確率が高い。つまり、この結果は電気自動車の蓄電池充電に関するトピックであると解釈できる。このように解釈をつけた全 34 個の要約トピックの一覧を表 11 に示す。

表 10 電気自動車関連の要約トピックの例

トピックZ02		トピックZ13	
確率	単語	確率	係り受け
6.8%	連結	2.2%	介する⇒伝達
4.7%	伝達	2.0%	介する⇒連結
3.6%	モータ	1.9%	動力⇒伝達
3.0%	出力軸	1.8%	出力軸⇒連結
2.8%	トランスミッション	1.8%	駆動力⇒伝達
2.6%	クラッチ	1.4%	トルク⇒伝達
2.1%	内燃機関	1.3%	クラッチ⇒介する
2.0%	駆動部	1.2%	駆動輪⇒伝達
2.0%	入力軸	1.1%	クラッチ⇒備える
1.7%	駆動輪	1.1%	出力軸⇒接続
1.6%	駆動軸	1.1%	機械的⇒連結
1.6%	動力	1.0%	モータ⇒備える
1.6%	駆動	1.0%	トランスミッション⇒備える
1.4%	トルク	0.9%	モータ⇒連結
1.3%	駆動装置	0.9%	駆動軸⇒連結
...	...	...	...

確率	単語	確率	係り受け
12.6%	充電	5.1%	バッテリー⇒充電
8.9%	電気自動車	4.0%	充電⇒行う
6.5%	蓄電装置	3.9%	電気自動車⇒充電
3.0%	バッテリー	1.9%	蓄電池⇒充電
2.0%	充電システム	1.6%	蓄電装置⇒充電
2.0%	蓄電池	1.6%	電力⇒供給
1.9%	電力	1.3%	充電⇒開始
1.7%	制御	1.2%	電気自動車⇒接続
1.5%	充電スタンド	1.2%	充電⇒蓄電装置
1.5%	放電	1.1%	用いる⇒充電
1.3%	外部電源	1.1%	充電⇒制御
1.2%	充電+できる	1.0%	供給⇒電力
1.0%	充電ケーブル	0.9%	蓄電装置⇒備える
0.8%	検出	0.9%	電力⇒充電
0.7%	情報	0.8%	蓄電装置⇒提供
...	...	...	...

表 11 電気自動車関連の要約トピックの一覧

No.	トピック名	No.	トピック名
Z01	エンジンの始動と停止	Z18	演算・推定
Z02	動力の伝達	Z19	機器の異常検出
Z03	モータ駆動	Z20	操作スイッチ
Z04	ロータ・ステータなど回転部品の構成	Z21	筐体
Z05	ブレーキ装置	Z22	表面の形成
Z06	動作制御	Z23	位置とその移動
Z07	動力伝達の制御	Z24	配置・位置・方向
Z08	スイッチの切り替え	Z25	構成の方位
Z09	交流・直流の変換	Z26	構成
Z10	エネルギーの変換	Z27	接続
Z11	電池モジュールの提供	Z28	方法の提供
Z12	二次電池の構成	Z29	損傷や浸水など不具合の防止
Z13	電気自動車の蓄電池充電	Z30	小型化・簡素化・低コスト化など付加価値
Z14	非接触受電など給電装置	Z31	効率性・安全性の向上
Z15	外部への電力供給	Z32	既存エンジンへの警鐘・樹脂組成物の提供
Z16	空調などの冷却・加熱	Z33	重力発電の活用による地球温暖化防止
Z17	情報通信	Z34	タービン発電の出力向上・燃費低減

#### 6.4 トピックのスコアリング

続いて分析対象とした 26,419 件の特許データに対して、今回抽出された 34 個の要約トピックのスコア（該当度）を計算した。スコアの計算方法は一つ目の分析事例と同様であるため割愛する。連続値であるスコアから該当有無を示す 1,0 のフラグに変換する閾値は、スコアの分布や実際の文章の内容も確認しながら 5 に設定した。トピックのスコアリングの計算処理により、表 12 に示すようなデータが作成された。26,419 件の特許データは、出願年、出願人、要約文という元々の情報に加え、要約のトピック 34 個の 1,0 のフラグ情報が付加されたデータとなり、このデータセットを用いることでトピックを軸にしたさまざまな分析を実行することができる。この先の各分析はすべてこのデータセットをベースとしている。

表 12 要約トピックのフラグ情報を紐づけた特許データ

特許ID	出願番号	出願年	出願人	要約	トピック Z01	トピック Z02	...	トピック Z34
1	特願2007-XXXX	2007	A社	【課題】電気式変速操作装	1	1		0
2	特願2009-XXXX	2009	B社	【課題】従来の電気自動車	0	1		1
3	特願2012-XXXX	2012	C社	エンジンのための方法及び	0	1		1
4	特願2013-XXXX	2013	D社	【課題】駐車場に設置され	1	0		0
...	...	...	...		...	...		...
26,419	特願2016-XXXX	2016	X社	充電ステーションが電気工	1	0		1

#### 6.5 トピックのトレンド分析

表 12 のトピックのフラグデータを用いて、出願年の情報と、トピックのフラグ情報から、各トピックのトレンドを分析する。一つ目の分析事例と同様に、出願年とトピックの関係の強さを示すリフト値の経年変化を可視化する。たとえば、2014 年からの直近 3 年でリフト値の上昇率が高い上位 4 つのトピックのトレンドを図 30 に示す。

図 30 より、特に「Z08. スwitchの切り替え」が上昇しており、他には「Z01. エンジンの始動と停止」や「Z19. 機器の異常検出」に関するトピックが上昇している。

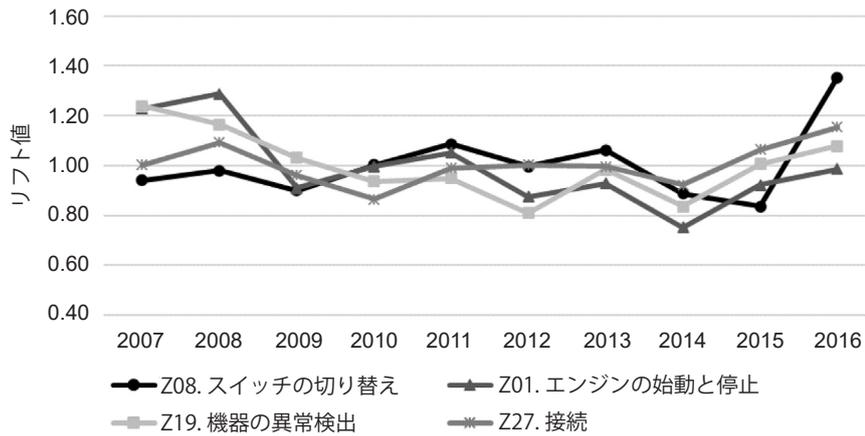


図30 2014年からの上昇率トップ4のトピックのトレンド

## 6.6 トピックを用いた競合他社の分析

表12のトピックのフラグデータを用いて、各トピックにおける出願人の特徴を分析する。

### (1) 静的分析: 出願人のポジショニングマップ

一つ目の分析事例と同様に、「シェア」と「注力度」という2つの指標を縦軸と横軸に設定し、トピックごとに出願人をプロットしたポジショニングマップを作成する。

ここでは先のトレンド分析で最も上昇していたトピック「Z08. スイッチの切り替え」を例とした結果を図31に示す。図31より、まずシェアの高さを確認すると、A社が圧倒的に高く、この領域で多くの特許を出願している。一方、注力度を確認すると、O社が最も高く他社とは大きなギャップがある。O社はこの領域に全社的に関心が高く、特有の技術力を保有している可能性が考えられる。また、中程度のシェア・注力度にはB社、D社、H社などがあるが、こうした企業の間で連携することで、より技術力を高めながらシェアを伸ばし、A社という巨人に対抗することも考えられる。

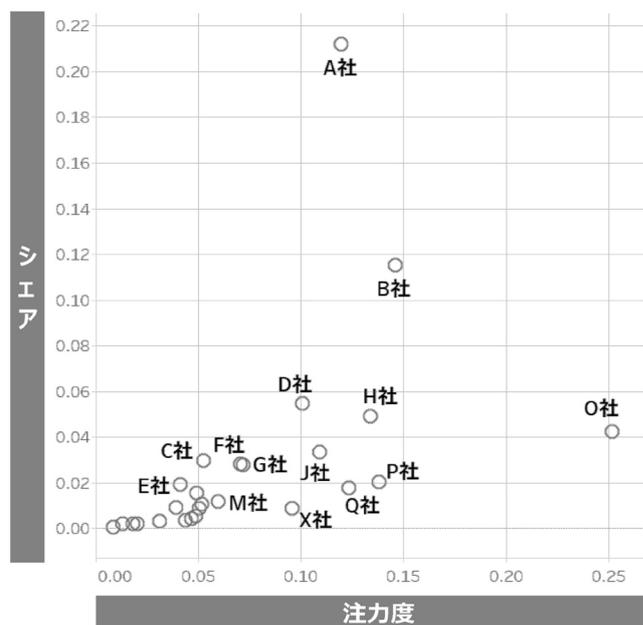


図31 トピック「Z08. スイッチの切り替え」における出願人のポジショニングマップ

## (2) 動的分析：出願人の出願件数の推移

ここで、注目の対象とした A 社、B 社、D 社、H 社、O 社について、一つ目の分析事例と同様に、この「Z08. スイッチの切り替え」というトピックに該当する特許の出願件数の推移も可視化した。その結果を図 32 に示す。図 32 より、シェア 1 位の A 社は、ここ最近では減少傾向にあったが、直近では出願が増加している。シェア 2 位の B 社は、最近では減少傾向にあり、シェア 3 位の D 社も最近では出願件数が少なくなっている。H 社と、注力度 1 位の O 社は、他と比べると件数は全体的に少ないが、どちらも直近で出願が増えている。ここから、直近で出願件数を急に増やしている A 社、H 社、O 社は、技術戦略の転換の可能性も考えられるため、直近での出願特許の内容や今後の出願動向には要注目といえる。

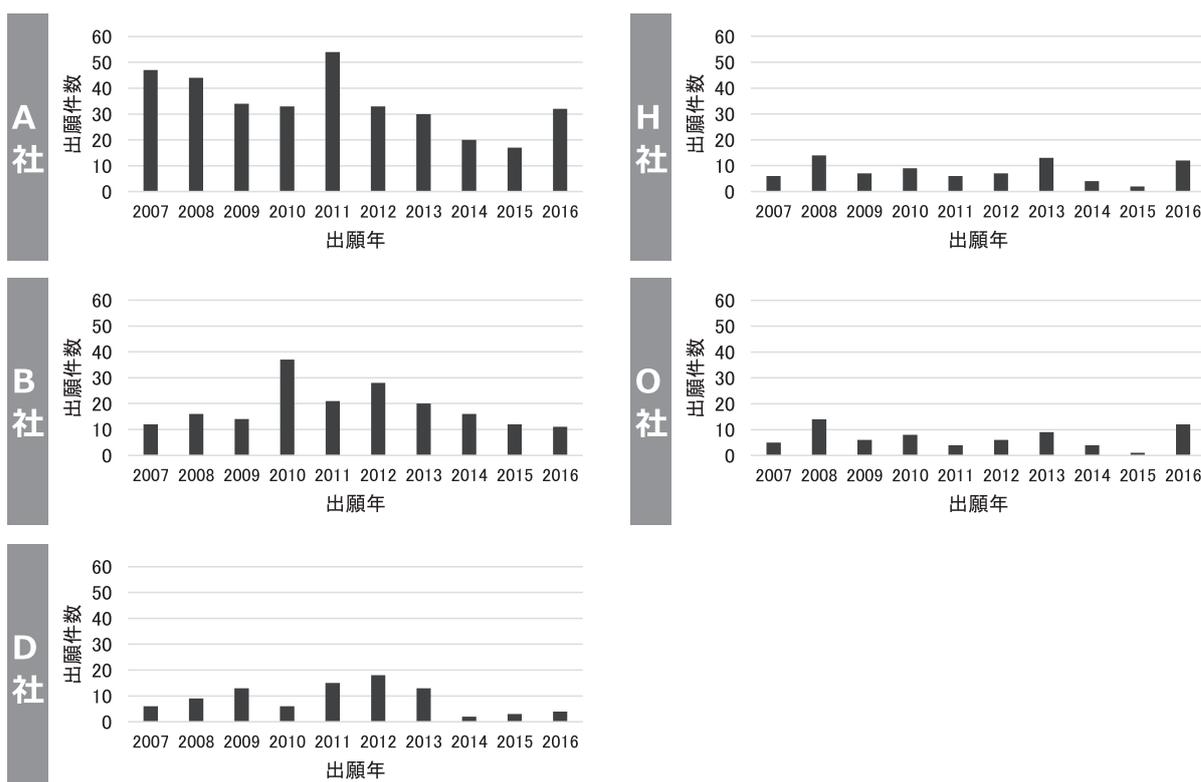


図 32 トピック「Z08. スイッチの切り替え」における出願人の出願件数の推移

## (3) 出願数が増加した特許内容の確認

直近で出願件数が増えている A 社、H 社、O 社の 3 社について、その直近一年（2016 年）に出願している特許の内容を確認した。

まず A 社が 2016 年に Z08 のトピックの特許を出願した件数は 32 件あったが、そのうち 8 件の特許を例に、「発明の名称」と「要約の課題」の部分抜粋したものを図 33 に示す。今回確認した特許は、どれも「Z08. スイッチの切り替え」というトピックに該当している特許であるが、たとえば、図 33 の左側の特許①～④は、電気自動車の衝突時にコンデンサを放電に切り替える技術に関する特許となる。これは衝突時の感電を防止する重要な技術となる。また、図 33 の右側の特許⑤～⑧は、何か異常が生じたときにそれを検知して回避する切り替え技術に関する特許となる。

続いて、H 社と O 社が 2016 年に Z08 のトピックの特許も確認したが、実は H 社と O 社の 2016 年の出願特許はすべて H 社と O 社の共同出願であった。この共同出願の件数は 12 件あったが、そのうち 6 件の特許を例に、「発明の名称」と「要約の課題」の部分抜粋したものを図 34 に示す。たとえば、図 34 の左側の特許①～③は、電力の供給が途絶えた場合でも給電を維持できる切り替え技術に関する特許となる。また、図 34 の右側の特許④～⑥は、

蓄電部などに異常が発生したときに、他にその影響が及ぶことを回避する切り替え技術に関する特許となる。このように、H社とO社の2社はすでに連携することでこの切り替え技術の領域において開発に力を入れていることが確認できる。

出願特許 ①		出願特許 ⑤	
発明の名称	電気自動車	発明の名称	電気自動車
要約【課題】	車両の衝突時に平滑化コンデンサを放電する確実性を向上させる。	要約【課題】	補機バッテリーが短絡しても、重要な補機への電力供給を継続できる電気自動車を提供する。
出願特許 ②		出願特許 ⑥	
発明の名称	電気自動車用の電源システム	発明の名称	電気自動車用の電源システム
要約【課題】	車両が衝突したときにより確実に平滑化コンデンサを放電する。	要約【課題】	通信不良が生じた場合であってもシステムスイッチを安全に開放することのできる電気自動車用の電源システムを提供する。
出願特許 ③		出願特許 ⑦	
発明の名称	電気自動車	発明の名称	電気自動車
要約【課題】	電気自動車の衝突時に、できるだけパーキングロックを使わずにモータを停止させて平滑化コンデンサを放電させる。	要約【課題】	第1インバータ回路と第2インバータ回路のうち一方で異常が生じたときに、両インバータ回路の複数のスイッチング素子を同時にオフするモータ制御ユニットを提供する。
出願特許 ④		出願特許 ⑧	
発明の名称	ハイブリッド車両	発明の名称	ハイブリッド車
要約【課題】	車両の衝突時に、モータの回転数を取得できない場合にも、インバータに接続されるコンデンサの電荷の放電を速やかに完了する。	要約【課題】	ハイブリッド車の第1コントローラに異常が生じたときでも、エンジンを始動し得る技術を提供する。

図33 A社が2016年に申請した「Z08.スイッチの切り替え」の特許例

出願特許 ①		出願特許 ④	
発明の名称	給電中継回路、副電池モジュール、電源システム	発明の名称	車両用電源装置
要約【課題】	一つの車載電源が失陥した場合であっても、車両の電氣的負荷に給電する。	要約【課題】	第1電源部の電力不足によって特定負荷が駆動できなくなる事態を、第2電源部に及ぼす影響を抑えて回避し得る車両用電源装置を提供することを目的とする。
出願特許 ②		出願特許 ⑤	
発明の名称	車載用のバックアップ装置	発明の名称	リレー装置及び車載システム
要約【課題】	電源部からの電力供給が途絶えた場合であっても電力供給対象への電力供給を途切れさせることなく供給源を蓄電部に切り替えることが可能な装置を、より簡易な構成で実現する。	要約【課題】	少なくとも2つの蓄電部に対して発電機から充電電流を供給することができ、且つ一方の蓄電部側に異常が発生した場合に、その異常が他方の蓄電部側に及ぶことを抑え得るリレー装置を提供する。
出願特許 ③		出願特許 ⑥	
発明の名称	車両用電源装置	発明の名称	車両用電源装置
要約【課題】	第1電源部からの電力に基づく急速充電動作と、充電中に第1電源部の電力供給が遮断されても遮断前後で放電状態を維持し得る充放電動作とを行い得る車両用電源装置を提供する。	要約【課題】	アイドリングストップ状態からの復帰の際に発電機側の蓄電部の電圧が低下しても負荷にその影響が及びにくい車両用電源装置を、電流集中を抑制し得る構成で実現する。

図34 H社とO社が2016年に共同申請した「Z08.スイッチの切り替え」の特許例

Nomolytics を特許文書データに適用した二つ目の分析事例は以上となる。より実践的な特許分析では、さらにここから特許の明細まで深く確認していくことが重要となるが、このようにトピックをベースにしたシンプルで分かりやすい分析をドリルダウンで進めていくことで、注目すべき特許を効率的に発見することができる。

## 7. インサイト獲得のためのPLSAの展開技術

ここからは、より深いインサイトを獲得できるようなトピック抽出の方法として、PLSAを応用して筆者が開発した"PCSA"と"differential PLSA"という2つの手法について紹介する。

これはテキストデータの分析に限った話ではないが、ビジネスにおけるデータ分析では、新たな気づきとなるインサイトを抽出できることがしばしば求められる。しかし、いざデータ分析に取り組むものの、結果が期待外れということは少なくない。それは、その業務担当者からすると経験的によく知っている結果であり、目新しさが無いということが大半である。この期待と現実のギャップから、データを分析してもその結果が業務に活かされないということも起きてしまい、データ分析の取り組み自体の有用性に対して現場で疑問が生じることもある。

しかし、データ分析の結果に目新しさが無いということは、案外自然なことである。そもそもデータをモデル化するということは、データに潜む普遍的な傾向やルールを抽象化するということである。たとえば観測データの平均値を算出するというありふれた集計も一種のモデル化であり、このモデルは統計学的には母集団の代表値を推定している。元のデータに対して説明力の高いモデルを構築しようとするれば、当然頻度の多い事象がよく反映され、結果的に経験的によく知っていることが優先的に表現される。ただ、経験的に分かっていることでも、その仮説をデータで示せたり、定性的な仮説を定量的な知識として獲得できることには十分価値があるはずであり、それが科学というものである。しかし残念ながらビジネスの現場ではそれだけでは満足されないというのも事実である。

こうした問題はPLSAでトピックを抽出するときにも該当する。PLSAではテキストデータから観測された共起行列を教師なし学習することでトピックを抽出するが、その際の評価基準は尤度という数学的な基準であり、いかに元のデータ全体を再現できているのかということが評価される。そうして抽出されるトピックはどうしても代表的で典型的なトピックが抽出される傾向となる。そこでより特徴的で個性的なトピックを抽出し、ビジネスにおけるインサイトの獲得を狙った新しい手法として、PLSAを応用した"PCSA"と"differential PLSA"という手法を開発した。

ここからはそのPCSAとdifferential PLSAという手法の具体的な分析アプローチと、その効果を示す分析事例を紹介する。

## 8. 課題のターゲットに特化したトピックの抽出手法:PCSA

PCSA (Probabilistic Causal Semantic Analysis, 確率的因果意味解析)<sup>34)</sup>は、分析において特徴を見たいターゲットを定め、そのターゲットに影響を与える要因となり得るトピックを優先的に抽出する手法である。以下に手法の内容とその適用事例について説明する。

### 8.1 課題となるターゲットの要因の探索

ビジネスの課題解決において有効なアクションを検討するためには、その課題となるターゲットに対して要因を探ることは重要である。たとえば、出願件数が最近増えている特許の技術的要因を探ることで今後投資をすべき研究開発テーマを検討できる。また、ある用途の要因となる要素技術を探ることで、その用途に関する事業を実現するための技術の候補を把握できる。知財業務以外のビジネス課題で考えると、たとえば、商品の顧客満足度の要因を探ることで満足度の高い新商品の企画やプロモーションを検討できる。あるいは、サービスの解約や会員の退会などの要因を探ることで顧客を維持する対応やサービスの内容を検討できる。

データ分析において、こうしたターゲットの要因を探るアプローチは、回帰分析や決定木分析など、しばしば教師あり学習の手法が適用される。つまりターゲットとなる目的変数とその要因となる説明変数との関係をモデル化する。これが説明変数の候補が多量となる高次元のデータになると、教師なし学習と教師あり学習を組み合わせるアプローチがよく適用される。多量な変数をいくつかの特徴に教師なし学習で次元圧縮し、その集約された特徴を説明変数として設定し、目的変数との関係を教師あり学習でモデル化するというアプローチである。従来よく用いられてきたものとして、まず主成分分析や因子分析を実行し、その因子を説明変数に回帰分析を実行するアプローチがある。たと

例えば、アンケートの多量な設問項目を因子分析でいくつかの因子にまとめ、その因子と顧客満足度などのターゲット設問との関係を回帰分析して顧客の価値観を理解し、マーケティングの検討に適用されるケースが挙げられる。

先述した Nomolytics も教師なし学習と教師あり学習を組み合わせた分析手法と解釈できる。Nomolytics では、テキストマイニングで抽出された単語群に教師なし学習の PLSA を適用してトピックを抽出し、そのトピックを確率変数として教師あり学習のベイジアンネットワークを適用してモデルを構築している。第 5 項で紹介した Nomolytics の分析事例では、技術の要因となる用途、あるいは用途の要因となる技術を探るモデルを構築した。

このように教師なし学習と教師あり学習を組み合わせることで、多量の変数を持つ高次元なビッグデータでも、次元圧縮された特徴量の変数を用いてターゲット変数との関係をシンプルに理解することができる。しかし、こうしたアプローチでは教師なし学習と教師あり学習はそれぞれ独立している。つまり、まず教師なし学習を完了してからその結果を教師あり学習に引き継ぐものであり、それは Nomolytics も同様である。そのため、教師なし学習で抽出された特徴量はどれもターゲット変数と強い関係を示すとは限らず、教師あり学習のモデリングの過程でターゲット変数に有効な特徴量のみが選択される。一方、ビジネスの課題解決の場面では、初めからターゲットに影響を与える要因に特化して特徴量が抽出される方が、効率的で望ましいとされることもある。

そこで課題となるターゲットを設定したときに、テキストデータからそのターゲットに影響を与えるトピックを優先的に抽出する手法として PCSA を開発した。以下に手法の内容を説明する。

## 8.2 PCSA という手法

PCSA は PLSA を応用したトピック抽出手法となる。PLSA によるトピック抽出のアプローチをおさらいすると、まずテキストデータにテキストマイニングを実行して単語を抽出し、その単語の共起頻度を集計した共起行列を作成する。そしてこの共起行列をインプットとして PLSA を適用することで、使われ方の似ている単語でまとめられたトピックを抽出する。この共起行列は全体のテキストデータから作成するため、当然そこから抽出されるトピックは全体を表現するような代表的なトピックが構成されることになる。

これに対して PCSA では、まず特徴を見たいターゲット変数を設定する。たとえば、特許文書のデータでは、最近トレンドを形成している技術を見なければターゲット変数は出願年が候補となる。知財業務以外の課題で考えると、アンケートデータでは顧客満足度、ユーザレビューデータではレビュー評点、コールセンターの問い合わせデータでは解約希望の有無といった変数が考えられる。PCSA では、このターゲット変数を該当有無 (1, 0) で示される二値変数に変換したものを使用する。

PLSA では共起行列をすべてのデータから 1 つ構築していたが、PCSA ではターゲット変数が該当するデータと該当しないデータ、それぞれから同じ行列構成の共起行列を作成する。そして、この 2 つの共起行列の差分の絶対値を計算した共起行列に対して PLSA を適用するというものである。PCSA で作成する共起行列のイメージを図 35 に示す。

全データをターゲット変数に該当するデータと該当しないデータの 2 つのグループに分割した際、それぞれのデータ件数の規模には違いがあるため、PCSA で適用する差分の共起行列はそのデータ件数の規模の違いを考慮した調整を施して差分を取る。具体的には、件数の少ないグループのデータ件数を基準に、件数の多い方のグループの共起行列の頻度をデータ件数の比率で調整する。データ件数が多い方のグループを調整対象とする理由は、その逆で調整した場合、たまたま共起したような小さな頻度が、調整によって実際には起こりえないような大きな値になってしまう可能性があるためである。つまり、件数の多い方のグループの頻度に合わせると、件数の少ない方のグループの頻度が調整され値が大きくなるが、特に 1 件などの頻度の少ない共起ペアは、全体のデータ件数が多くなればその割合に応じて増加するとは言い切れない。たまたま 1 件出現したというケースは大いに考えられる。これが調整によって大きな値となると現実と乖離した共起頻度となる懸念がある。共起行列では頻度 1 件の共起ペアは非常に多く、これがすべて同じ割合で増加することの影響は大きいと考えられる。一方、件数の少ない方のグループの頻度に合わせると、件数の多い方のグループの頻度が調整され値が小さくなるが、調整後の頻度が 0 件 (存在しない) となることはなく、件数の少ない方のグループの頻度を大きくすることよりも現実的と考えられる。

このようにしてターゲット変数の該当有無によって2つのデータに分割し、それぞれ同じ行列構成で作成された共起行列は、ターゲット変数の該当有無に影響を受ける共起ペアでは頻度の差が大きくなり、その影響を受けないような共起ペアでは頻度の差は小さくなる。つまり、2つの共起行列の差分を取った共起行列では、ターゲット変数の該当有無に影響を受ける共起ペアは頻度が大きくなり、そうでない共起ペアでは頻度が小さくなるということである。そのため、この共起行列にPLSAを適用することでターゲット変数の該当有無に影響を与えるような要因トピックが優先的に抽出されることが期待できる。

なお本手法PCSAは2023年2月6日に特許登録されている(特許第7221526号)<sup>35)</sup>。

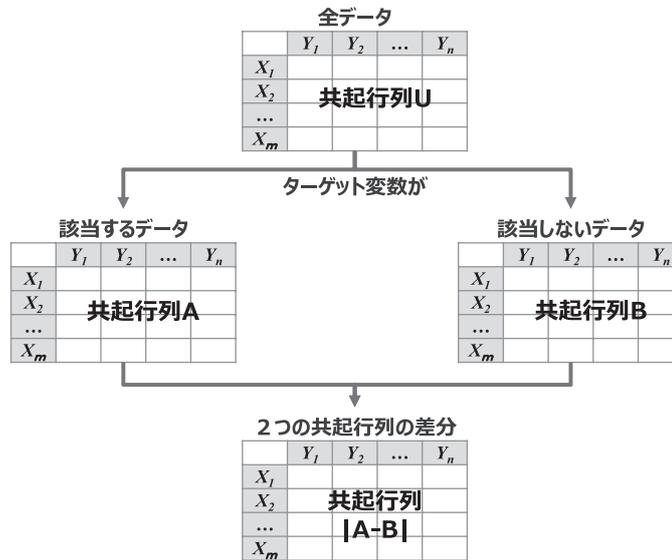


図 35 PCSA で作成する共起行列

### 8.3 PCSA を適用した分析事例

ここでは先述した Nomolytics の二つ目の分析事例で使用した、26,419 件の電気自動車に関連する特許文書データに PCSA を適用した分析事例を紹介し、通常の PLSA で抽出されたトピックの結果と比較する。

#### 8.3.1 ターゲットの設定と共起行列の作成

PCSA ではまず特徴を見たいターゲット変数を設定するが、ここではそれを「出願年」とする。分析データは出願年が 2007 年から 2016 年までとなる 10 年分の特許データであるが、ターゲット変数は出願年が「2013 年以前」か「2014 年以後」という二値変数に設定する。これによって 2014 年前後で変化するような要因トピックを優先的に抽出し、トレンドが 2014 年以降で上昇傾向あるいは下降傾向にあるような技術を把握することを狙いとする。

「2013 年以前」と「2014 年以後」にデータを分割すると、それぞれのデータ件数は 20,014 件と 6,405 件となった。PCSA で作成するそれぞれの共起行列の行列構成は、全データから作成される共起行列の構成を共通して採用する。つまり、Nomolytics の分析事例で作成した「名詞単語 (3,020 語) × 係り受け (2,128 表現)」の共起行列の行列構成をとる。共起頻度は文章単位にカウントした頻度であるが、この共起行列の各共起ペアが 1 つでも存在する文章数は、全データでは 43,836 件であり、そのうち「2013 年以前」のデータでは 33,113 件 (75.5%)、「2014 年以後」のデータでは 10,723 件 (24.5%) であった。「2013 年以前」のデータと「2014 年以後」のデータで、それぞれ「名詞単語 (3,020 語) × 係り受け (2,128 表現)」の共起行列を作成し、件数の多い「2013 年以前」のデータで作成された共起行列の頻度を各データの件数の比率でもって調整した。具体的には、「2013 年以前」のデータで作成された共起行列の全共起頻度に、文章数の比率である (10,723 / 33,113) をかけて調整済みの共起行列を作成した。このようにして作成された 2 つの共起行列の差の絶対値を計算した共起行列を作成し、これをインプットに PLSA を適用してトピックを抽出した。本分析で作成した共起行列のイメージを図 36 に示す。

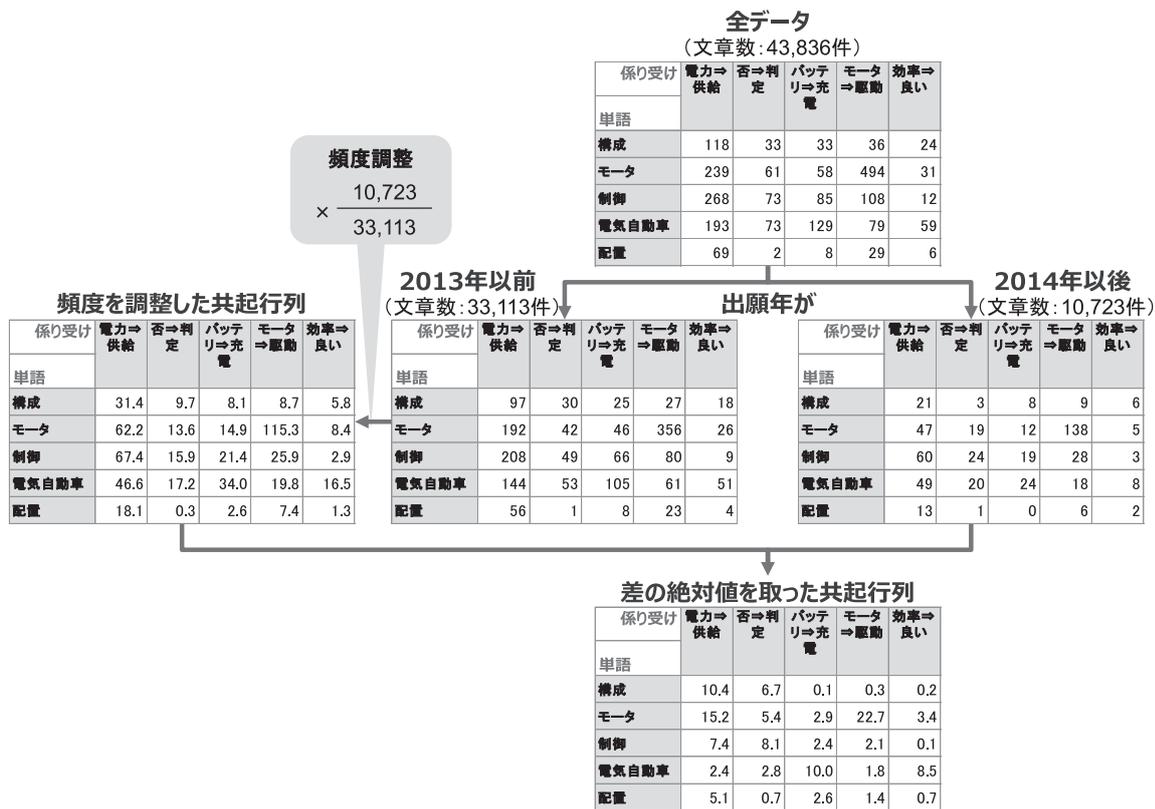


図 36 電気自動車関連の特許データに PCSA を適用するための共起行列作成

### 8.3.2 トピックの抽出と解釈

PLSA の実行は Nomolytics の分析事例と同様にトピック数を 1 刻みで変化させ、それぞれのトピック数に対して初期値をランダムに変えて PLSA を 5 回ずつ実行し、それぞれの解を情報量基準 AIC で評価して最も評価の良い解を採用した。PLSA の実行解の AIC の評価の結果を図 37 に示す。図 37 では、トピック数を 5 から 20 まで 1 刻みで変化させて実行した結果となっている。各トピック数に対して 5 つの黒いプロット一つひとつが初期値を変えて PLSA を実行した結果であり、灰色の折れ線グラフはその 5 つの実行解における AIC の平均値を示している。灰色の折れ線が示す AIC の平均値を比較したとき、それが最小となるトピック数は 11 個となり、このトピック数 11 個における 5 つの実行解の中で AIC が最小となる結果を採用した。

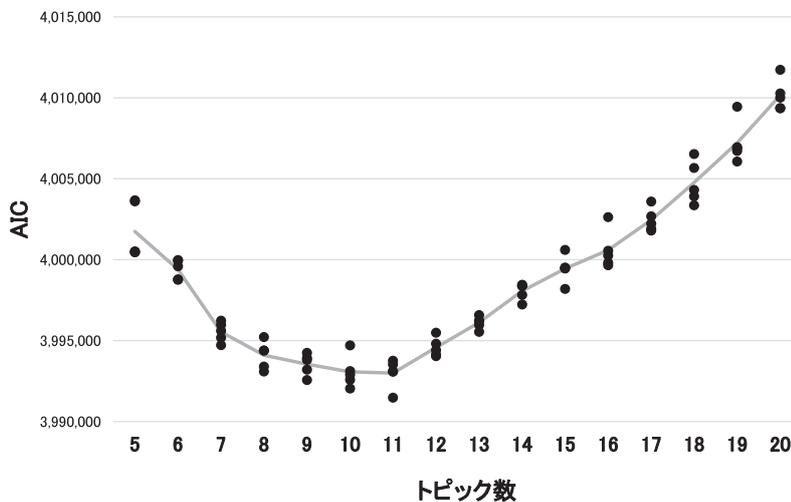


図 37 PCSA の分析における各トピック数での PLSA の実行解の AIC 評価

抽出されたトピックの内容の例を表 13 に示す。単語と係り受けは所属確率の高い順に並べている。表 13 (左) のトピック Zy04 では、単語は、充電、電気自動車、電力、バッテリー、蓄電装置、給電などで所属確率が高く、係り受けは、バッテリー⇒充電、充電⇒行う、電力⇒供給、電気自動車⇒充電、供給⇒電力、蓄電池⇒充電といった表現で所属確率が高い。つまり、この結果は電気自動車への充電、給電装置に関するトピックであると解釈できる。表 13(右) のトピック Zy05 では、単語は、製造、二次電池、負極、正極、製造方法、エネルギー、リチウムイオン電池などで所属確率が高く、係り受けは、水素⇒製造、発電⇒電気、製造方法⇒提供、製造⇒方法、含む⇒リチウムイオン電池といった表現で所属確率が高い。つまり、この結果は二次電池の製造方法に関するトピックであると解釈できる。このように解釈をつけた全 11 個のトピックの一覧を表 14 に示す。

表 13 PCSA の適用で抽出されたトピックの例

トピックZy04		トピックZy05	
確率	単語	確率	係り受け
3.1%	充電	1.9%	バッテリー⇒充電
2.5%	電気自動車	1.9%	充電⇒行う
1.9%	電力	1.8%	電力⇒供給
1.9%	バッテリー	1.5%	電気自動車⇒充電
1.8%	蓄電装置	0.9%	供給⇒電力
1.3%	給電	0.7%	蓄電池⇒充電
1.2%	供給	0.7%	搭載⇒バッテリー
1.0%	蓄電池	0.7%	充電⇒開始
0.8%	制御	0.6%	蓄電装置⇒充電
0.8%	充電スタンド	0.6%	電気自動車⇒接続
0.7%	充電システム	0.6%	給電⇒行う
0.7%	情報	0.6%	用いる⇒充電
0.7%	充電+できる	0.5%	バッテリー⇒供給
0.7%	外部電源	0.5%	電力⇒充電
0.6%	電源	0.5%	充電⇒蓄電装置
...	...	...	...

確率	単語	確率	係り受け
3.1%	製造	4.5%	水素⇒製造
2.2%	二次電池	3.0%	発電⇒電気
1.5%	負極	2.6%	製造方法⇒提供
1.5%	正極	2.6%	製造⇒方法
1.5%	製造方法	1.6%	含む⇒リチウムイオン電池
1.2%	エネルギー	1.3%	成形品⇒提供
1.1%	リチウムイオン電池	1.3%	強度⇒有する
1.1%	水素	1.3%	表面⇒形成
1.1%	セパレータ	1.2%	含む⇒組成物
1.0%	電解液	1.2%	リレー⇒スイッチ
0.9%	形成	1.2%	二次電池⇒備える
0.9%	含有	1.2%	スイッチ⇒適する
0.9%	正極活物質	1.1%	電気部品⇒提供
0.9%	方法	1.1%	段階⇒含む
0.8%	発電	1.1%	二次電池⇒提供
...	...	...	...

表 14 PCSA の適用で抽出されたトピックの一覧

No.	トピック名
Zy01	エンジン駆動・動力伝達の制御
Zy02	モータの回転構成
Zy03	交流・直流の変換
Zy04	電気自動車への充電、給電装置
Zy05	二次電池の製造方法
Zy06	空調などの冷却・加熱
Zy07	情報通信、検出判定システム
Zy08	形成・配置
Zy09	小型化・低コスト化・簡素化・操作性向上
Zy10	重力発電の活用による地球温暖化防止
Zy11	既存エンジンへの警鐘、タービン発電・重力発電

### 8.3.3 トピックのスコアリング

Nomolytics の分析事例と同様に、分析対象とした 26,419 件の特許データに対して、今回抽出された 11 個のトピックのスコア (該当度) を計算した。スコアの計算方法は Nomolytics の分析事例と同様であるため割愛するが、連続値であるスコアから該当有無を示す 1.0 のフラグに変換する閾値は、スコアの分布や実際の文章の内容も確認しながら 5 に設定した。トピックのスコアリングの計算処理により、表 15 に示すようなデータが作成された。26,419 件の特許データに対して、Nomolytics の分析事例で全体のデータから抽出された 34 個のトピックに加え、今回の PCSA の適用によって抽出された 11 個のトピックの 1.0 のフラグ情報が付加されたデータとなる。このデータセットを用いることでこのトピックを軸にしたさまざまな分析を実行することができる。

表 15 PCSA で抽出したトピックのフラグ情報を紐づけた特許データ

特許ID	出願番号	出願年	出願年 グループ	出願人	要約	トピック Z01	トピック Z02	...	トピック Z34	トピック Zy01	トピック Zy02	...	トピック Zy11
1	特願2007-XX	2007	①2013年以前	A社	【課題】電気式変速	1	1		0	0	1		1
2	特願2009-XX	2009	①2013年以前	B社	【課題】従来の電気	0	1		1	0	0		1
3	特願2012-XX	2012	①2013年以前	C社	エンジンのための	0	1		1	0	1		0
4	特願2013-XX	2013	②2014年以後	D社	【課題】駐車場に設	1	0		0	1	0		0
...	...	...	...	...	...	...	...		...	...	...		...
26,419	特願2016-XX	2016	②2014年以後	X社	充電ステーションが	1	0		1	1	0		1

### 8.3.4 PLSA と PCSA の結果の比較

本分析では、ターゲット変数を「出願年」に設定し、全体で 26,419 件の特許データを「2013 年以前」と「2014 年以後」の 2 つのグループに分割して PCSA を適用した。分析の狙いは、出願年が 2014 年前後で変化するような要因トピック優先的に抽出し、トレンドが 2014 年以降で上昇傾向あるいは下降傾向にあるような技術を把握することだった。実際に抽出された 11 個のトピックにそうした傾向があるのか、先の Nomolytics の分析事例で全体のデータから抽出された 34 個のトピックと比較して確認した。

確認の仕方は、表 15 のデータを用いて、先の Nomolytics の PLSA で全体のデータから抽出された 34 個のトピックと、今回の PCSA で抽出された 11 個のトピックについて、それぞれのトピックが該当する特許のうち、出願年が「2014 年以後」となっている割合を比較した。これによって、2014 年前後において PCSA ではどれくらい偏ったトピックを抽出できているか確認した。その結果を図 38 に示す。図 38 のグラフは、横軸をトピック、縦軸を出願年が 2014 年以後となっている割合とし、その割合が高い順に横軸のトピックを並べている。図 38 の上のグラフは先の Nomolytics の PLSA で全体のデータから抽出された 34 個のトピックの結果であり、下のグラフは PCSA で抽出された 11 個のトピックの結果である。この上下のグラフを比較すると、割合の平均値はどちらもほぼ同じだが、そのばらつきには大きな違いが生まれている。PLSA で抽出された全体トピック 34 個の割合の大きさはおおむね 25% 前後となっているが、PCSA で抽出された要因トピック 11 個は、割合の大きさが高いものと低いものに偏っている。つまり、PCSA では出願年が 2014 年前後で特徴を示すトピックが抽出されていることが確認できる。

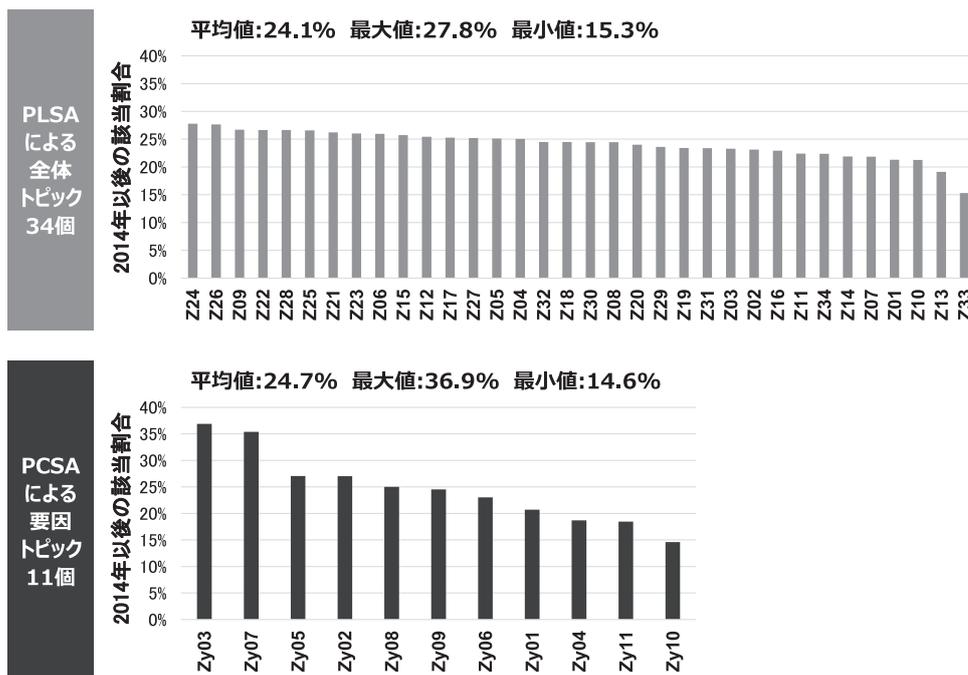


図 38 全体のトピック 34 個と PCSA のトピック 11 個の比較

### 8.3.5 トレンドが上昇傾向にあるトピックの確認

PCSA で抽出されたトピックのうち、図 38（下）で 2014 年以後の割合が高かった上位 2 つのトピック Zy03 と Zy07 について、実際のトレンドを確認した。Nomolytics の分析事例で紹介したトレンド分析と同様に、出願年とトピックの関係の強さを示すリフト値を使ってそれぞれの経年変化を可視化した。その結果を図 39 に示す。図 39 より、確かにどちらもきれいな上昇トレンドを形成していることがわかる。なお、Zy03 は「交流・直流の変換」に関するトピック、Zy07 は「情報通信、検出判定システム」に関するトピックであるが、それぞれのトピックの構成内容と、該当する特許の要約の一例を図 40, 41 に示す。たとえば、Zy03 は、電気自動車が省電力動作をとるときには、コンバータが直流電力を通常よりも低い電圧に変換するが、冷却ポンプへは高い効率で電力を供給することで、十分な冷却性能を得るといった技術がある。Zy07 は、電気自動車の圧力チャンバに設置されている圧力センサで算出される圧力値と、外周面の導体板から形成されるコンデンサで算出される電流値により、車両への衝突を検知する技術がある。



図 39 2014 年以後の割合が高いトピック Zy03 と Zy07 のトレンド

トピックの構成				該当する特許の一例	
トピックZy03				発明の名称	出願年
確率	単語	確率	係り受け	電気自動車用の電源システム	2016
1.9%	供給	1.3%	交流電力-変換	<b>要約文</b> 通常動作と省電力動作の何れにおいても、十分な冷却性能が得るとともに、ポンプを高効率で動作させる電源システムを提供する。コンバータが、変換後の直流電力の電圧を第1電圧に制御する通常動作と、第1電圧よりも低い第2電圧に制御する省電力動作を実行可能である。電力制御ユニットを冷却する冷媒を循環させるポンプの制御ユニットが、省電力動作時に通常動作時よりも変換後の直流電力がポンプに供給される期間の比率を高くする。	
1.5%	変換	1.2%	直流電力-変換		
1.5%	制御	1.0%	電力-供給		
1.3%	電圧	0.9%	変換-供給		
1.2%	インバータ	0.8%	直列-接続		
1.2%	電力	0.8%	並列-接続		
1.1%	制御部	0.6%	供給-電力		
1.1%	検出	0.6%	バッテリー-接続		
1.1%	バッテリー	0.6%	電力-変換		
1.1%	直流電力	0.5%	モーター-供給		
1.1%	電気自動車	0.5%	電圧-検出		
1.0%	モーター	0.5%	交流-変換		
1.0%	交流電力	0.5%	直流-変換		
0.9%	スイッチ	0.5%	負荷-供給		
0.8%	直流	0.5%	直流電圧-変換		
...	...	...	...		

図 40 トピック Zy03「交流・直流の変換」の構成内容と特許の一例

トピックの構成				該当する特許の一例	
トピックZy07				発明の名称	出願年
確率	単語	確率	係り受け	車両用衝突検知装置	2015
1.8%	検出	1.3%	否-判定	<b>要約文</b> 車両への衝突が発生したことを精度よく検出することができる車両用衝突検知装置の提供。バンパーアブソーバーには、圧力チャンバが配置されている。圧力チャンバには圧力センサが接続されている。圧力チャンバの外周面に貼付されている前方導体板および後方導体板により、複数のコンデンサが形成されている。圧力検出値から算出された有効質量が閾値以上であって、かつ、コンデンサの容量の増大により、電流検出値が閾値以上である時に、車両への衝突を検出する。	
1.4%	方法	0.7%	システム-備える		
1.3%	判定	0.6%	基づく-検出		
1.3%	制御	0.6%	基づく-生成		
1.1%	構成	0.6%	方法-備える		
1.1%	受信	0.6%	電気信号-変換		
1.0%	演算	0.5%	信号-受信		
1.0%	制御部	0.5%	制御部-有する		
0.9%	生成	0.5%	車両-走行		
0.9%	信号	0.4%	電気信号-出力		
0.8%	システム	0.4%	情報-基づく		
0.8%	センサ	0.4%	制御-構成		
0.7%	変化	0.4%	有無-判定		
0.7%	測定	0.4%	制御部-含む		
0.7%	制御装置	0.4%	モーター-制御		
...	...	...	...		

図 41 トピック Zy07「情報通信、検出判定システム」の構成内容と特許の一例

### 8.3.6 PCSA のまとめ

PCSA は、特徴を見たいターゲットを定め、そのターゲットに影響を与える要因となり得るトピックを優先的に抽出する手法である。以下に PCSA の手法の特徴を Nomolytics と比較してまとめる。

Nomolytics では、まずテキストデータ全体を表す代表的なトピックを PLSA で抽出し、そのトピックの特徴を、たとえば出願年や出願人といったさまざまな分析軸で探索することができる。一方、PCSA は、たとえば最近トレンドを形成している技術を把握したいというように、特徴を探索したいターゲットが定まっているときに適用する手法となる。そのターゲットに特化したトピックを優先的に抽出し、より顕著な傾向を示す要因を深く分析することでインサイトの獲得を狙う手法である。

どちらが優れているというものではなくて、目的に応じて使い分けることが重要である。PCSA は分析軸を最初から一つに絞り、それに特化したトピックを抽出できるが、Nomolytics のようにさまざまな分析軸で広く特徴を探索したいというときには不向きとなる。ビジネスの課題解決において、特定された課題の要因を探り、有効な施策アイデアを効率的に検討したい場合は、PCSA を適用することで有用な知識が得られるものと期待できる。

## 9. 個性的なトピックの抽出手法 :differential PLSA

differential PLSA<sup>36)</sup> は、より個性的なトピックを抽出することでインサイトの獲得を狙う手法となる。以下に手法の内容とその適用事例について説明する。

### 9.1 典型的でない個性的なトピックの抽出

第7項で述べたとおり、ビジネスにおけるデータ活用場面では、新たな気づきとなるインサイトの獲得が期待されるものの、実際にデータ分析から得られる結果は、経験的によく知っているものも多く、目新しさがないと評価されてしまうことがある。テキストデータから PLSA でトピックを抽出する分析でも、そのトピックはデータ全体を代表するような典型的なトピックが抽出される傾向がある。次元圧縮法である PLSA は、元の高次元のデータに対して、再現性の高い低次元のトピックに変換する手法であるため、元のデータに対して表現力の高いトピックが抽出される。PLSA のインプットとなる共起行列は頻度を値とするデータであるが、元のデータに対して表現力の高いトピックを抽出するには、当然頻度の高い要素を中心に代表的なトピックが形成される傾向がある。結果として高い頻度の要素で構成され、経験的によく知っているような、典型的とも感じられるトピックが抽出されるのである。しかし、これはデータ全体の特徴を把握する上では有用な分析であり、一つひとつのトピックはよく知っている典型的なものであっても、そうしたトピックに類型化することは人間では困難な処理のはずである。実際に結果に対して目新しさが無いと感じる人は、結果を見たからそう感じているだけで、結果を見ないで自分の経験や感覚に基づいて代表的なトピックを網羅的に列挙することはほぼできない。

ここでは、典型的であってもデータ全体を表現するトピックを抽出することで、全体像を把握する分析の重要性を前提とした上で、とは言うものの、より個性的なトピックも抽出できる手法として開発した differential PLSA について説明する。

### 9.2 differential PLSA という手法

differential PLSA も PCSA のように、共起行列を加工して PLSA を適用する手法となる。differential PLSA では、行要素  $X_i$  と列要素  $Y_j$  が共起する実測頻度  $n(X_i, Y_j)$  を値として持つ通常の共起行列  $M$  に加え、行要素  $X_i$  と列要素  $Y_j$  が共起する期待頻度  $n'(X_i, Y_j)$  を値として持つ共起行列  $M'$  を作成する。期待頻度とは、行要素  $X_i$  と列要素  $Y_j$  の総頻度（出現文章数）と全体における出現割合に基づいて計算した理論値となる。具体的には、 $X_i$  の総頻度（出現文章数）を  $n(X_i)$ 、 $Y_j$  の総頻度（出現文章数）を  $n(Y_j)$ 、総文章数を  $N$  とすると、行要素  $X_i$  の総頻度  $n(X_i)$  に対して、列要素  $Y_j$  の出現割合  $n(Y_j)/N$  をかけた値となり、以下の式 (15) のように計算される。あるいは、列要素  $Y_j$  の総頻度  $n(Y_j)$  に対して、行要素  $X_i$  の出現割合  $n(X_i)/N$  をかけた値でも同じである。この実測頻度と期待頻度の2つの共起行

列は、カイ二乗検定をするときに作成する、観測度数と期待度数の2つのクロス集計表に近いイメージのものである。differential PLSA は、共起行列の各共起ペアにおいて、期待頻度に対する実測頻度の比率の対数を取った  $\log(n(X_i, Y_j)/n'(X_i, Y_j))$  を値として持つ共起行列（differential 共起行列）を構築し、これに PLSA を適用する。なお、対数の計算において値が負数となるものは0に置換する。differential PLSA で作成する共起行列のイメージを図 42 に示す。

$$n'(X_i, Y_j) = n(X_i) \cdot n(Y_j) / N \quad (15)$$

実測頻度の共起行列に適用する通常の PLSA では、その解を求める最適化計算において、どうしても頻度が高い要素に高い確率が割り当てられ、結果として抽出されるトピックは典型的なものになる傾向があり、目新しさに欠けてしまう。一方、differential PLSA の共起行列では、実測頻度を期待頻度で除した値を持つが、実測頻度が高い共起ペアでも、元々全体の頻度が高い要素が含まれるときには期待頻度も高い値となるため、実測頻度を期待頻度で除することで値の大きさが制限される。逆に実測頻度が低い共起ペアでも、期待頻度がそれよりも十分小さければ値は大きくなり、これに PLSA を適用した解ではこうした要素にも高い確率が割り当てられる可能性がある。つまり、通常の PLSA では頻度が低い要素には高い確率が割り当てられにくい傾向があるが、differential PLSA ではそうした要素にも高い確率が割り当てられる可能性があり、より個性的なトピックが抽出されることが期待できる。

また、期待頻度に対する実測頻度の比率に対数を取る理由は、極端に高くなる値を制限するためである。特に期待頻度は1未満となるケースも多く、比率のみでは値が高くなりすぎるものもある。この状態では共起行列全体の値の分布は大きくばらつき、極端な値の開きが生まれてしまう。このまま PLSA を適用した場合、この極端に大きな値に引っ張られる結果となるため、必要以上にデフォルメされた歪んだトピックとなることが考えられる。そこでこの比率の値の対数を取ることで値の分布をならし、この問題を緩和する。

なお本手法 differential PLSA は 2023 年 2 月 6 日に特許登録されている（特許第 7221527 号）<sup>37)</sup>。

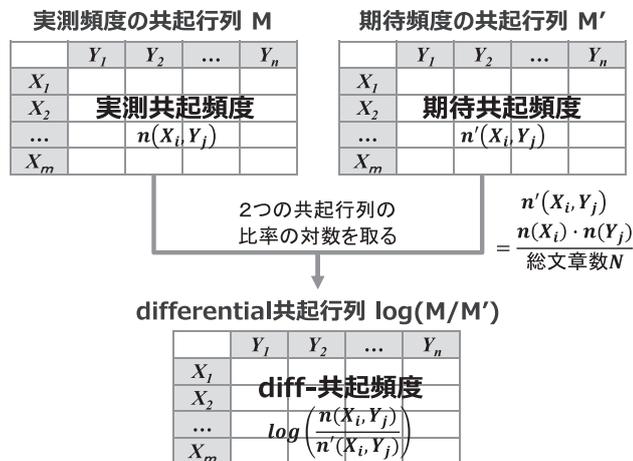


図 42 differential PLSA で作成する共起行列

### 9.3 differential PLSA を適用した分析事例

ここでは先述した Nomolytics の二つ目の分析事例で使用した、26,419 件の電気自動車に関連する特許文書データに differential PLSA を適用した分析事例を紹介し、通常の PLSA で抽出されたトピックの結果と比較する。

#### 9.3.1 共起行列の作成

differential PLSA では実測頻度の共起行列と期待頻度の共起行列を作成するが、実測頻度の共起行列は Nomolytics の分析事例で作成した「名詞単語（3,020 語）×係り受け（2,128 表現）」の共起行列となる。期待頻度の共起行列は、これと同じ行列構成をとり、各共起ペアの期待頻度を計算する。なお、今回のデータにおいて、総文章数 N は

229,598 件であった。この総文章数とは、テキストマイニングの実行の際に、文章の区切り文字に指定した句点（。）や隅付括弧（【】）によって分割された文章の数であり、共起行列を構成する単語や係り受けが一切登場しない文章も含まれる。共起行列の各共起ペアにおいて、期待頻度に対する実測頻度の比率の対数を取った共起行列を作成し、これをインプットに PLSA を適用することでトピックを抽出した。本分析で作成した共起行列のイメージを図 43 に示す。

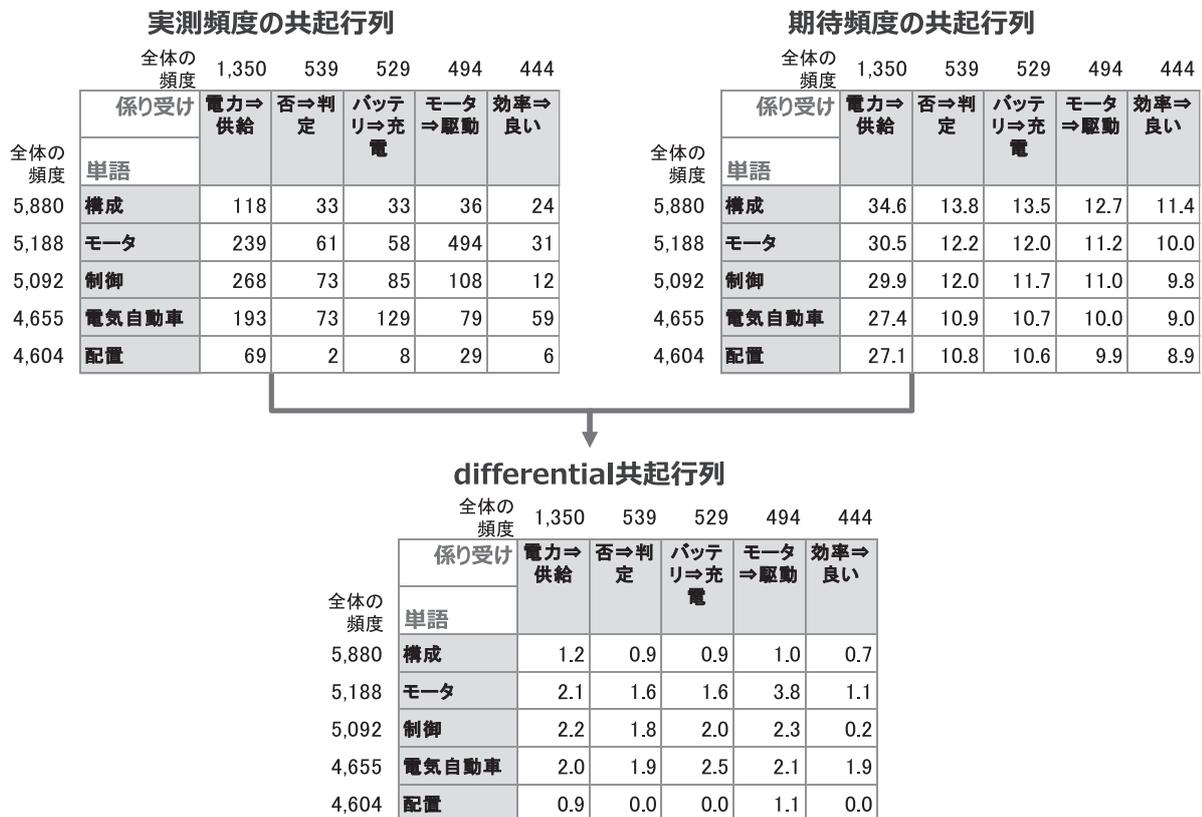


図 43 電気自動車関連の特許データに differential PLSA を適用するための共起行列作成

### 9.3.2 トピックの抽出と解釈

PLSA の実行は Nomolytics の分析事例と同様にトピック数を 1 刻みで変化させ、それぞれのトピック数に対して初期値をランダムに変えて PLSA を 5 回ずつ実行し、それぞれの解を情報量基準 AIC で評価して最も評価の良い解を採用した。PLSA の実行解の AIC の評価の結果を図 44 に示す。図 44 では、トピック数を 35 から 60 まで 1 刻みで変化させて実行した結果となっている。各トピック数に対して 5 つの黒いプロット一つひとつが初期値を変えて PLSA を実行した結果であり、灰色の折れ線グラフはその 5 つの実行解における AIC の平均値を示している。灰色の折れ線が示す AIC の平均値を比較したとき、それが最小となるトピック数は 50 個となり、このトピック数 50 個における 5 つの実行解の中で AIC が最小となる結果を採用した。解釈を付けた 50 個のトピックの一覧を表 16 に示す。

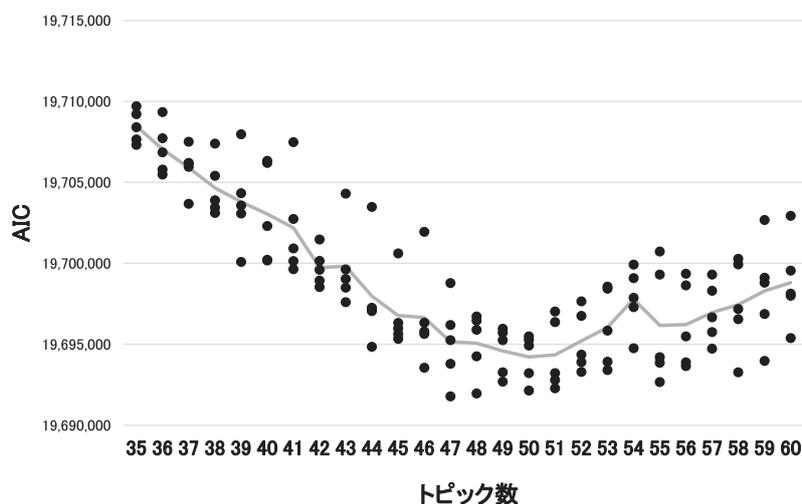


図 44 differential PLSA の分析における各トピック数での PLSA の実行解の AIC 評価

表 16 differential PLSA の適用で抽出されたトピックの一覧

No.	トピック名	No.	トピック名
Z'01	エンジン制御	Z'26	電流・電圧の検出
Z'02	動力伝達	Z'27	温度、電流、充電量などの検出と制御
Z'03	差動機構などを備えた動力伝達の制御	Z'28	演算や推定、測定などのステップを含む方法
Z'04	回転運動	Z'29	情報の取得・提供(位置情報やバッテリー残量等)
Z'05	ロータ・ステータなどモータの構成	Z'30	スイッチなど操作装置
Z'06	モータ制御(トルク制御や回転数制御等)	Z'31	車両用灯具
Z'07	油圧ポンプなどを利用したモータ駆動	Z'32	掃除機
Z'08	ブレーキ	Z'33	基板の構成
Z'09	状態に応じた制御、運転者の操作補助	Z'34	回路の接続(電力変換回路等)
Z'10	コンバータとバッテリー昇降圧	Z'35	端子接続
Z'11	直流と交流の電力変換	Z'36	部品・装置の収容ケース・筐体
Z'12	回転力などの電気エネルギー変換	Z'37	部品・装置の配置
Z'13	エネルギー効率の向上	Z'38	パーツなどの移動、位置
Z'14	発電と蓄電	Z'39	構造の形成・方位
Z'15	電池モジュールの提供	Z'40	支持構造
Z'16	燃料電池	Z'41	装置やユニットの構成
Z'17	二次電池の構成	Z'42	システム・方法の構成
Z'18	バッテリーの充放電	Z'43	その他方法
Z'19	充電システム	Z'44	組成物の製造方法(樹脂や電解液等)
Z'20	充電の接続	Z'45	機能性組成物・成形品(耐熱性や耐衝撃性等)
Z'21	非接触など受給電装置	Z'46	製造の効率化(小型化や低コスト化等)
Z'22	車両用空調など熱交換	Z'47	不具合の防止(損傷、感電、盗難等)
Z'23	冷却装置と放熱	Z'48	その他
Z'24	信号の入出力と検出	Z'49	タービン発電と船舶・飛行機への応用
Z'25	電気信号の取得と変換(センサ検出等)	Z'50	重力発電の活用による地球温暖化防止

### 9.3.3 トピックのスコアリング

Nomolytics の分析事例と同様に、分析対象とした 26,419 件の特許データに対して、今回抽出された 50 個のトピックのスコア(該当度)を計算した。スコアの計算方法は Nomolytics の分析事例と同様であるため割愛する。連続値であるスコアから該当有無を示す 1,0 のフラグに変換する閾値は、スコアの分布や実際の文章の内容も確認しながら 10 に設定した。トピックのスコアリングの計算処理により、表 17 に示すようなデータが作成された。26,419 件の特許データに対して、Nomolytics の分析事例で全体のデータから抽出された 34 個のトピックに加え、今回の differential PLSA の適用によって抽出された 50 個のトピックの 1,0 のフラグ情報が付加されたデータとなる。このデータセットを用いることでこのトピックを軸にしたさまざまな分析を実行することができる。

表 17 differential PLSA で抽出したトピックのフラグ情報を紐づけた特許データ

特許ID	出願番号	出願年	出願人	要約	トピック Z01	トピック Z02	...	トピック Z34	トピック Z'01	トピック Z'02	...	トピック Z'50
1	特願2007-XXX	2007	A社	【課題】電気式変速操	1	1		0	0	0		1
2	特願2009-XXX	2009	B社	【課題】従来の電気自	0	1		1	1	0		0
3	特願2012-XXX	2012	C社	エンジンのための方法	0	1		1	0	1		1
4	特願2013-XXX	2013	D社	【課題】駐車場に設置	1	0		0	1	0		0
...	...	...	...		...	...		...	...	...		...
26,419	特願2016-XXX	2016	X社	充電ステーションが電	1	0		1	0	1		0

### 9.3.4 PLSA と differential PLSA の結果の比較

Nomolytics で適用する通常の PLSA では、データ全体を表現する典型的なトピックが抽出される傾向があることに対して、differential PLSA は、より個性的なトピックを抽出することを狙って開発したものである。実際に differential PLSA で抽出された 50 個のトピックはそうした傾向があるのか、先の Nomolytics の分析事例において通常の PLSA で抽出された 34 個のトピックと比較し、その違いを確認した。注目すべき違いとして、(1) differential PLSA は頻度の低い要素でも高い確率が割り当てられること、(2) differential PLSA でのみ抽出されるトピックがあること、(3) differential PLSA の方がトピックの解釈が難しいことが挙げられる。以下にそれぞれについて結果を交えながら説明する。

#### (1) differential PLSA は頻度の低い要素でも高い確率が割り当てられる

この特徴に関しては、通常の PLSA と differential PLSA で同様の解釈ができるトピックの内容を並べて比較すると分かりやすい。

たとえば、「ブレーキ」に関するトピックでは、通常の PLSA では「Z05. ブレーキ装置」として、differential PLSA では「Z'08. ブレーキ」として抽出されているが、それぞれのトピックの構成内容を比較したものを表 18 に示す。表 18 では、それぞれのトピックにおいて単語と係り受けを所属確率の高い順に並べているが、確率に加えてその単語・係り受けの総頻度（出現文章数）も掲載している。表 18 の比較より、表 18（左）の通常の PLSA によるトピック Z05 は、所属確率の高い単語は全体の出現頻度も高い傾向にあるが、表 18（右）の differential PLSA によるトピック Z'08 は、全体の出現頻度が低い単語でも高い所属確率が割り当てられている。そのため、通常の PLSA によるトピック Z05 では、「ブレーキ」や「制動力」といった、ブレーキの直接的な表現が上位に位置しており、まさにブレーキを代表するトピックとなっている。一方、differential PLSA によるトピック Z'08 では、「マスタシリンダ」や「液圧」といった、頻度は低いがより具体的なブレーキ装置の中身に関する表現が上位にある。

もう一つの比較の例として、電気自動車の「非接触などの給電」に関するトピックでは、通常の PLSA では「Z14. 非接触受電などの給電装置」として、differential PLSA では「Z'21. 非接触など受給電装置」として抽出されている。それぞれのトピックの構成内容を比較したものを表 19 に示す。表 19 の比較より、表 19（左）の通常の PLSA によるトピック Z14 は、やはり所属確率の高い単語は全体の出現頻度も高い傾向にあるが、表 19（右）の differential PLSA によるトピック Z'21 は、全体の出現頻度が低い単語でも高い所属確率が割り当てられている。そのため、通常の PLSA によるトピック Z14 では、「給電」「電力」「電源」といった、給電の直接的な表現ばかりが上位に位置しており、まさに給電を代表するトピックとなっている。一方、differential PLSA によるトピック Z'21 では、「非接触」という単語が最も確率が高く、また「送電コイル」や「受電コイル」といった、頻度は低いがより具体的な非接触給電の装置に関する表現が上位にある。電気自動車において、非接触の充電システムは今後重要な技術と考えられるが、その仕組みは、駐車場の路面に設置した送電コイルと車両に搭載された受電コイルが重なり合うことで電磁誘導が発生し、電力が供給されるという仕組みが代表的である。differential PLSA ではそうした技術のより具体的なキーワードがトピックの上位語となっている。

各トピックにおいて所属確率の高い上位語の頻度が、通常のPLSAよりも differential PLSAの方が少ないことを確認するため統計的検定も行った。その方法は、各トピックにおいて所属確率の高い順に単語および係り受けを並べたときに、その累積確率が50%になるまでの単語および係り受けの平均頻度を検定用データとし、通常のPLSAで抽出された34個のトピックと differential PLSAで抽出されたトピック50個のトピックの平均値の差をWelchのt検定で検定した。通常のPLSAでは単語の頻度は平均1288.2件 (SD = 686.5)、係り受けの頻度は平均85.6件 (SD = 26.0)、differential PLSAでは単語の頻度は平均440.9件 (SD = 137.5)、係り受けの頻度は平均58.8件 (SD = 15.9) となり、単語および係り受けの両方で1%有意の違いがみられた。

表18 「ブレーキ」に関するトピックにおける通常のPLSAと differential PLSAの比較

通常のPLSAで抽出されたトピックZ05			differential PLSAで抽出されたトピックZ'08		
確率	頻度	単語	確率	頻度	係り受け
7.3%	779	ブレーキ	3.2%	92	車両⇒ブレーキ
5.0%	1,384	作動	2.4%	33	ブレーキ液圧⇒発生
3.4%	5,188	モータ	2.2%	53	制動力⇒発生
2.9%	279	制動力	1.8%	49	ブレーキ⇒備える
2.9%	867	運転者	1.7%	32	操作量⇒応ずる
2.7%	1,117	車輪	1.6%	82	ブレーキ⇒提供
2.6%	1,005	操作	1.6%	59	電気信号⇒基づく
1.6%	111	操作量	1.6%	41	運転者⇒操作
1.6%	162	ブレーキペダル	1.6%	55	操作⇒応ずる
1.5%	5,092	制御	1.5%	235	モータ⇒制御
1.4%	485	電気信号	1.4%	32	基づく⇒発生
1.2%	73	ブレーキ液圧	1.3%	494	モータ⇒駆動
1.2%	299	付与	1.3%	34	制動トルク⇒発生
1.2%	3,960	検出	1.3%	21	液圧⇒発生
1.1%	117	液圧	1.3%	22	ブレーキペダル⇒操作
...	...	...	...	...	...

表19 「非接触などの給電」に関するトピックにおける通常のPLSAと differential PLSAの比較

通常のPLSAで抽出されたトピックZ14			differential PLSAで抽出されたトピックZ'21		
確率	頻度	単語	確率	頻度	係り受け
9.7%	1,162	給電	4.3%	1,350	電力⇒供給
5.1%	3,629	電力	2.7%	115	給電⇒行う
3.6%	1,140	電源	2.5%	52	電力⇒受電
3.4%	329	給電装置	2.0%	28	給電⇒電力
2.4%	4,655	電気自動車	1.6%	124	電源⇒接続
2.4%	180	非接触	1.5%	20	駐車装置⇒変化
2.4%	1,052	外部	1.5%	34	非接触⇒受電
2.3%	160	受電	1.3%	85	給電⇒停止
1.3%	158	駐車	1.3%	35	受電⇒電力
1.2%	398	電源装置	1.3%	70	給電⇒制御
1.2%	5,092	制御	1.2%	40	給電装置⇒備える
1.2%	124	受電部	1.2%	57	給電⇒受ける
1.2%	1,064	変化	1.1%	25	電力⇒給電
1.2%	63	駐車装置	1.1%	27	電力⇒送電
1.1%	3,978	供給	1.1%	41	車両⇒給電
...	...	...	...	...	...

(2) differential PLSA でのみ抽出されるトピックがある

differential PLSAでは通常のPLSAでは抽出されないトピックが抽出されていた。これは、differential PLSAの方が抽出トピック数が多かったということもあるが、通常のPLSAのトピック数を50個に設定して計算した解も確認したうえで、differential PLSAのみに現れたトピックが存在していた。その例を図45, 46に示す。図45, 46はそれぞれ differential PLSAでのみに抽出されたトピックZ'09とZ'29であり、そのトピックの構成内容と該当する特許の要約文の一例を掲載したものとなる。図45で示すトピックZ'09は、単語では、シフトレンジ、パーキングレンジ、検出結果、停止などで所属確率が高く、係り受けでは、自動的⇒行う、操作⇒行う、駆動⇒停止といった表現で所属確率が高い。これは運転者の操作を補助したり、自動停止などの運転アシストに関する技術と解釈できる。図46で示すトピックZ'29は、単語では、ナビゲーション装置、情報、目的地、位置情報などで所属確率が高く、係り受

けでは、情報⇒送信、情報⇒含む、情報⇒取得、情報⇒受信といった表現で所属確率が高い。これは位置情報を取得してドライバーにナビ情報として提供するという、情報の獲得と提供に関する技術と解釈できる。

こうした技術は、頻度の高い要素を中心に代表的なトピックを抽出する通常のPLSAでは抽出しきれなかった情報であるが、どちらも自動車の付加価値を高める重要な技術といえる。

トピック'09			該当する特許の要約文の一例		
確率	頻度	単語	確率	頻度	係り受け
1.4%	50	シフトレンジ	1.9%	23	自動的⇒行う
1.0%	38	パーキングレンジ	1.7%	38	操作⇒行う
0.9%	147	検出結果	1.6%	22	駆動⇒停止
0.7%	1,200	停止	1.6%	26	動作⇒行う
0.7%	365	解除	1.6%	40	要する⇒時間
0.6%	34	キースイッチ	1.4%	46	停止⇒状態
0.6%	3,960	検出	1.2%	26	ブレーキ⇒作動
0.6%	237	禁止	1.1%	96	状態⇒検出
0.5%	204	切替	1.1%	62	状態⇒切り替える
0.5%	67	閉状態	1.1%	35	発生⇒検出
0.5%	1,005	操作	1.1%	57	制御装置⇒行う
0.5%	199	移行	1.1%	23	作動⇒行う
0.5%	1,384	作動	1.1%	22	介する⇒制御
0.5%	62	絶縁抵抗	1.0%	25	モータ⇒停止
0.5%	93	シフトレバー	1.0%	40	ロック⇒解除
...	...	...	...	...	...

【課題】モータによる長期間のクリープトルク発生による電力消費を抑えるよう促すとともに、車両停止時の駆動源の切り替えによる車両の動き出しの防止又は車両の動き出しを運転者に知らせることのできるハイブリッド電気自動車を提供すること。【解決手段】アクセルペダル及びブレーキペダルの踏み込みがなく、エンジンの作動中であり、シフト位置が走行レンジ(Dレンジ又はRレンジ)である車両停止時に、クリープ走行制御の駆動源がモータからエンジンに切り替わった際のトルク変動により車両が動き出した際には、パーキングブレーキによる制動を自動的に行い、且つ警報器及び警告灯による運転者への注意喚起を行う。

※例示のための要約文であり、一部内容は筆者が加工している

図 45 differential PLSA でのみ抽出された「運転者の操作補助」に関するトピック

トピック'29			該当する特許の要約文の一例		
確率	頻度	単語	確率	頻度	係り受け
1.2%	122	ナビゲーション装置	2.1%	48	情報⇒送信
1.1%	809	情報	2.1%	42	情報⇒含む
1.0%	165	目的地	1.8%	68	情報⇒取得
1.0%	111	位置情報	1.5%	31	情報⇒受信
0.9%	786	取得	1.5%	35	示す⇒情報
0.9%	819	送信	1.5%	126	情報⇒基づく
0.8%	558	表示	1.4%	25	情報⇒用いる
0.8%	43	検索	1.3%	48	表示部⇒表示
0.8%	72	バッテリー残量	1.2%	31	基づく⇒特定
0.8%	83	サーバ	1.2%	42	無線通信⇒行う
0.7%	511	記憶	1.2%	29	情報⇒提供
0.7%	877	受信	1.1%	28	記憶⇒記憶手段
0.7%	83	無線通信	1.1%	29	基づく⇒取得
0.7%	61	現在地	1.0%	32	車両⇒位置
0.7%	51	地図情報	1.0%	33	情報⇒記憶
...	...	...	...	...	...

【課題】充電スタンドへの電気自動車の有効なナビゲーションを行う。【解決手段】本発明は、複数の充電スタンドの位置情報を記憶するデータベースと、充電スタンドが保持している電池についてその種別およびエネルギー残量を記憶するデータベースと、電気自動車と通信して、電気自動車が搭載する電池の種別、エネルギー残量、エネルギー効率、現在位置を受信し、現在位置から前期エネルギー残量およびエネルギー効率によって到達可能な充電スタンドを見つけ、その位置情報と、その充電スタンドが保持する電池のエネルギー残量の情報を電気自動車に送信する充電スタンドロケータとを備える。

※例示のための要約文であり、一部内容は筆者が加工している

図 46 differential PLSA でのみ抽出された「情報の取得と提供」に関するトピック

### (3) differential PLSAの方がトピックの解釈が難しい

以上のように、differential PLSAで抽出されるトピックは、頻度が低い単語でも高い所属確率が割り当てられる傾向にある。そのため、トピックの所属確率が頻度の低い多様な表現に分散しやすく、通常のPLSAで抽出されるトピックよりも意味の解釈が難しい印象がある。実際にこれまで例示したトピックの構成内容を見ても、通常のPLSAと比較して、differential PLSAのトピックは上位語の所属確率が小さく、確率が分散している傾向がある。通常のPLSAは頻度の高い代表的な表現に所属確率が集中するため、解釈はしやすい。しかし、その結果として通常のPLSAで抽出されるトピックは典型的なトピックとなる傾向がある。この問題を解消する手法としてdifferential PLSAを開発した経緯を踏まえると、このジレンマはやむを得ないものである。逆にそれだけ個性的なトピックが抽出されていることを意味するものといえる。もしトピックの構成内容だけでは解釈しにくい場合は、図 45, 46でも示したように、そのトピックに該当する特許の具体的な文章を確認することが効果的である。

### 9.3.5 differential PLSA のまとめ

differential PLSA は、通常の PLSA で抽出されるトピックが典型的なものになりがちであるという課題に対して、より個性的なトピックを抽出する手法として開発した。通常の PLSA は、データ全体を表現するような代表的なトピックを抽出するため、頻度の高い表現を中心にトピックが構成される傾向があり、結果として典型的なトピックという印象を持つことがある。これに対して differential PLSA は、期待頻度に対する実測頻度の比率の対数を取った共起行列をインプットにすることで、頻度は低いが高具体度の高い表現にも高い確率を割り当ててトピックを構成することができる。これによって、より個性の強いエッジの立ったトピックを抽出でき、データ全体では埋もれがちな特徴を発見することができる。

したがって、differential PLSA はビジネスの課題解決の場面において、新たな気づきとなるようなインサイトが得られる手法として期待できるが、あくまでもデータ全体を表現するトピックを把握した前提で用いることが望ましい。データ分析の取り組みでは、木を見て森を見ずとならないように、まずは全体像を把握することが第一であり、全体像の現状を把握した上で個性のある特徴を探ることが重要となる。そのため、実際には通常の PLSA と differential PLSA を併用した分析が有効であると考えられる。

## 10. Nomolytics, PCSA, differential PLSA の比較

ここまで紹介した Nomolytics, PCSA, differential PLSA という 3 つの手法について、その違いと共通する考え方をまとめる。

### 10.1 各手法の違い

この 3 つの手法の目的と特徴を比較したものを表 20 に示す。Nomolytics は、テキストデータ全体を表現する代表的なトピックを抽出し、そのトピックをベースに傾向の分析や要因関係の分析をすることで、テキストデータ全体に潜む特徴を俯瞰する手法となる。PCSA と differential PLSA は、Nomolytics の中でもトピックを抽出する機能を拡張させた手法となる。Nomolytics がデータ全体を表現するトピックを抽出する手法であるのに対して、PCSA と differential PLSA はより特徴的で個性的なトピックを優先して抽出する手法となり、より深いインサイトの獲得を狙って開発したものである。PCSA は、特徴を見たいターゲットを定め、そのターゲットの特徴を左右する要因となり得るトピックを抽出する手法となる。differential PLSA は、通常の PLSA が頻度の高い要素を中心とした典型的なトピックを抽出しがちであるのに対して、頻度の大小に依存しないより個性的なトピックを抽出し、データ全体では埋もれがちな特徴を発見する手法となる。どの手法が優れているかということではなく、目的に応じて使い分けることが重要である。たとえば PCSA と differential PLSA は深いインサイトの獲得は期待できるが、トピック抽出のインプットとなる共起行列は加工を施しているため、そこで抽出されるトピックはデータ全体を表現するものとはいえずなくなっている。先述したように、データ分析の取り組みでは、まずは全体像を把握することが第一である。データ全

表 20 Nomolytics, PCSA, differential PLSA の比較

	Nomolytics	PCSA	differential PLSA
目的	テキストデータ全体を表現する代表的なトピックを抽出する。そのトピックをベースに、テキスト情報に潜む傾向や要因関係を可視化し俯瞰する。	特徴を見たいターゲットを定め、そのターゲットに影響を与える要因となり得るトピックを優先的に抽出する。	頻度の大小に依存しない、より個性的なトピックを抽出し、データ全体では埋もれがちな特徴を発見する。
PLSAを適用する共起行列の特徴	「単語 × 係り受け」の行列構成をとり、全データにおいて、それぞれの単語と係り受けの共起頻度（同時に出現する文章数）をカウントした共起行列を採用する。	全データをターゲットが該当するデータと該当しないデータに分割し、それぞれのデータで同じ構成の「単語 × 係り受け」共起行列を作成し、その 2 つの共起行列の差分を取った共起行列を採用する。	全データから集計される実測頻度の共起行列に加え、期待頻度の共起行列を作成し、期待頻度に対する実測頻度の比率の対数を取った「単語 × 係り受け」共起行列を採用する。

体を表現するようなトピックを把握し、全体像を俯瞰したいときには、通常のPLSAを適用することが合理的である。したがって、繰り返しになるが、目的に応じた手法の使い分けが重要ということである。

## 10.2 各手法で共通する分析のコンセプト

一方、この3つの手法で共通することは、テキストマイニングで抽出した大量の単語をいくつかのトピックに集約し、その集約されたトピックをベースに傾向を可視化することで、大量のテキストデータに潜む特徴をシンプルに理解しようという分析の考え方である。その分析の考え方のイメージを図47に示す。テキストデータを分析する際に、以前から使用されているテキストマイニングという手法は、テキストデータの文書情報に含まれる単語を抽出し、その単語を単位とした可視化によってデータの特徴を分析する。しかし、対象がビッグデータとなると、マイニングによって抽出される単語の数は膨大となるため、その単語をベースとした可視化では結果が複雑になり、解釈が難しくなってしまう。そこで、トピックモデルを応用することで、テキストマイニングで抽出された大量の単語をいくつかのトピックに集約する。これにより、従来の単語ベースではなく、トピックをベースとして、テキストデータに潜む特徴をシンプルに分析することができる。つまり、テキストデータからマイニング（採掘）された単語を、トピックというかたまりにリファイニング（精製）し、さまざまな可視化の形にプロセッシング（加工）するというアプローチである。これが本節で紹介したNomolytics、PCSA、differential PLSAに共通する重要なコンセプトとなる。

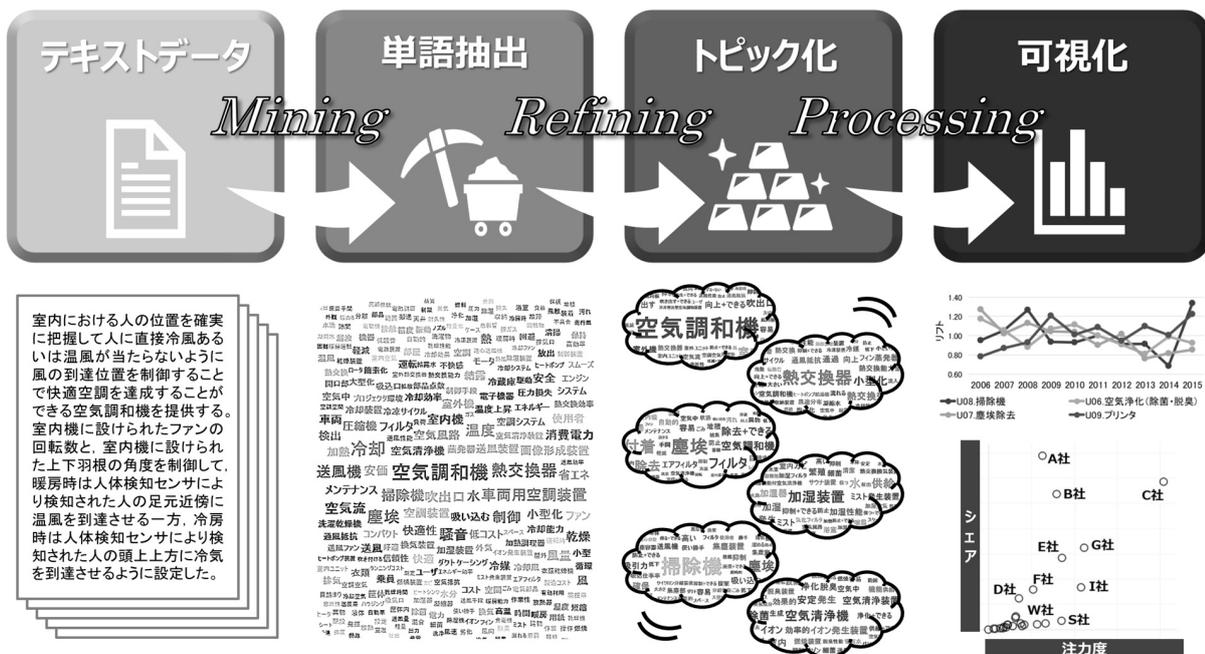


図47 トピックをベースとしたテキストデータの分析

おわりに

本節では、特許の文書情報を用いた分析について、パテントマップと呼ばれるシンプルで分かりやすい従来の分析アプローチの有用性を認識しながら、大量の特許データを対象とする場合の分析の課題を挙げ、それを解決するアプローチとして、テキストマイニングにトピックモデルのPLSAとベイジアンネットワークという2つのAI技術を応用した分析手法Nomolyticsを解説し、それを実際の特許文書データに適用した分析事例を紹介した。

Nomolyticsを適用した特許文書分析のメリットとして、①単語ではなく集約されたトピックを軸にした分析を実行することで、膨大な特許文書に潜む傾向を分かりやすく理解することができること、②用途と技術の統計的な関係を分析することで、各用途を実現する上で重要な要素技術を把握できたり、自社技術を有効活用できるような新規用

途のアイデアを創出できることが挙げられる。このように、従来のテキストマイニングにトピックモデルやベイジアンネットワークという2つのAI技術を組み合わせることで、特許文書データを分析することで、人間では読み切れない膨大な特許文書に潜む傾向や要因関係を把握でき、企業の技術戦略の検討において有益な気づきを得る新たな切り口を提供することができる。

また本節では、より特徴的で個性的なトピックを抽出することでインサイトの獲得を狙って開発した、PCSAとdifferential PLSAというPLSAの応用手法も紹介した。それぞれの手法は目的が異なっている。Nomolyticsはデータ全体を表現するトピックを抽出して全体像を理解する手法となる。そのトピックをベースにさまざまな軸で傾向を可視化したり、そのトピックの周辺に潜む要因関係をベイジアンネットワークで分析する。PCSAは特徴を見たいターゲットの要因となるトピックを抽出する。differential PLSAは頻度の大小に依存することのない個性的なトピックを抽出する。なお、これらの手法で共通する本質的な考え方は、テキストマイニングで抽出した大量の単語を、いくつかのトピックに集約し、そのトピックをベースに傾向を可視化することで、大量のテキストデータに潜む特徴をシンプルに理解するという点である。

昨今の第3次AIブームでは、特に機械学習のアルゴリズムの開発が急速に進み、技術先行型で高度で高性能な分析アプローチが次々に提案されており、特許分析の分野でもそのトレンドは例外ではない。一方で、「データのビジネス活用」という側面では、どれだけ進化しているかは疑問がある。以下に活用的側面からAIを使った特許分析における注意や懸念点をいくつか挙げる。

企業では最先端のAI技術、一番いいAI技術を適用したいという願望が強く、より高度で難しいアプローチを採用しようとする傾向がある。しかしそれに従って分析はどんどん難解なものとなっていく、それを扱える専門性の高い人材は限定的になってしまう。つまり、高度な機械学習アルゴリズムのようにリテラシーのハードルが高すぎると、データ分析に参加できる人間が限られ、ビジネスにおける活用が進みにくくなる。

また、AIによって結果だけが自動で出力されたり、ワンタッチでさまざまな可視化を提示してくれる操作性の高いツールも企業では好まれる。こうしたAIツールは多様な分析結果をすぐに自動で出力してくれるので、分析業務自体が効率化されるというメリットがある。しかし、その結果のプロセスはブラックボックスになってしまい、それがどうして得られたのかというロジックを人間が理解できないことが多く、より期待する結果を得るための探索の仕方も分かりにくい。その結果が得られた妥当性を確認することが難しいので、組織内で議論をしにくく、納得感のある意思決定が難しくなる。

さらに生成AIについては、その使い方は特に注意が必要である。生成AIは、プロンプトとして入力されたあらゆる質問や指示に対して、回答や分析結果を瞬時に返してくれるため、非常に使い勝手の良いアプリである。特許分析においても、特許の分類やラベル付け、要約などを生成でき、さらには指定した特許データの中から何が有益なインサイトとなるのか、生成AIに直接分析させるような使い方もある。まるで生成AIは全知全能で、これさえあれば分析は十分と思ってしまうくらいである。生成AIは大規模言語モデルという、大量のデータから事前学習された汎用的な言語モデルに基づいて推論した結果を生成している。あくまでも推論結果であるので、ハルシネーションと呼ばれる誤った情報を生成してしまうこともある。特に特許文書のように専門性の高いデータでは、汎用的な言語特徴を学習した生成AIにとってハルシネーションを起こしやすい対象といえる。また、生成AIに特許データの中から重要な傾向やインサイトを回答させ、それに該当する特許文献の詳細情報を確認したい場合、実際の特許文献を正しく参照できないこともある。他にも、大量の特許を要約させる場合、生成AIのモデルの汎用性から、重要な気づきとなるような細かい技術要素が省略されてしまう可能性もある。ただし、生成AIの技術は日進月歩で進化しており、こうした課題に対して今後対応可能な生成AIが登場する可能性もあるが、常にその信憑性を確認する意識を持つことは利用者として必須である。

ビジネスにおいてデータを分析する目的は課題解決に活用するためであり、その最終的な課題解決の意思決定をするのは人間である。従来のパテントマップは手動で手間がかかるかもしれないが、ビジネス活用という側面では有用性に優れた手法といえる。古典的な集計分析ではあるものの、Excelだけでも分析と可視化が十分実行できるため、誰もがデータ分析に参加できる。また、そのプロセスは明確なため、結果を信頼でき、チームで議論をしながら分析

の探索を柔軟に設計できる。生成 AI のように推論に基づいた分析ではなく、実データに基づいた分析なので、誤った情報が含まれていたり、実際の特許文献を参照できないというような事態にはならない。

本節で紹介した Nomolytics, PCSA, differential PLSA では、その最も重要なアウトプットは表 7, 12, 15, 17 に示したデータセットである。これは元々の特許データにトピックという新たな分析軸が加わった表形式のデータであり、Excel でも分析を十分実行できるデータである。テキストマイニングと PLSA でトピックを抽出し、各データに対するトピックのスコアを計算するところまでは確かに専門知識が求められる。しかし、そうした専門家がこのような表形式のデータさえ作成すれば、あとはそのデータを使って従来のパテントマップと同様の分析と探索を誰もが実行できる。分析は業務担当者のリテラシーに応じて自身が操作しやすく理解しやすい方法を自由に選択すればよく、例えば Excel で単純に該当件数を集計してグラフ化したり、ピボットテーブルでクロス集計したり、統計解析ソフトで分析したり、ベイジアンネットワークのような要因関係のモデリングを実行することもできる。また抽出したトピックの構成内容は単語や係り受けに対するシンプルな所属確率であり、誰もが十分解釈できる。このようにトピックという新たな分析の切り口を得たデータを用いて、誰もがデータ分析に参加し、そして議論しながら納得感のある技術戦略を検討できる。

AI の技術的側面だけで見ればその発展スピードは目まぐるしく、現在人間はデータを分析する強力な武器がたくさん用意されている状況ではある。一方で、データを活用してビジネスの課題解決を達成するにはその活用的側面も評価しなければならない。今直面しているビジネスの課題の本質を明確に認識した上で、どの武器をどう使えばその課題を効率的にかつ効果的に解決に導くことができるのか、その武器の技術的側面だけでなく活用的側面も考え、戦略的に AI 技術の利用を検討することが望まれる。

## 文 献

- 1) 新井喜美雄：特許情報分析とパテントマップ，情報の科学と技術，Vol. 3, No. 1, pp. 16-21 (2003)
- 2) 安藤俊幸：テキストマイニングと統計解析言語 R による特許情報の可視化，情報管理，Vol. 52, No. 1, pp. 20-31 (2009)
- 3) 小池孝幸，石川徹也：知的財産戦略経営のための特許情報分析手法：パテントマップ作成および読解支援ツールの開発について，Japio YEAR BOOK 2008, pp. 244-249 (2008)
- 4) 松尾豊：『人工知能は人間を超えるか ディープラーニングの先にあるもの』，角川 EPUB 選書 (2015)
- 5) T. Mikolov, K. Chen, G. Corrado, and J Dean : "Efficient Estimation of Word Representations in Vector Space," Proc. of the International Conference on Learning Representations (ICLR) (2013)
- 6) Q. V. Le, and T. Mikolov : "Distributed Representations of Sentences and Documents," Proc. of the 31st International Conference on Machine Learning (ICML-14) (2014)
- 7) T. Mikolov, M. Karafiat, L. Burget, J. Cernocky, and S. Khudanpur: "Recurrent neural network based language model," ISCA Interspeech, pp. 1045-1048 (2010)
- 8) S. Hochreiter, and J. Schmidhuber : "Long Short-Term Memory," Neural Computation, Vol. 9, No. 8, pp. 1735-1780 (1997)
- 9) I. Sutskever, O. Vinyals, and Q. V. Le : "Sequence to sequence learning with neural networks," Proc. of the 27th International Conference on Neural Information Processing Systems (NIPS 2014), pp. 3104-3112 (2014)
- 10) D. Bahdanau, K. Cho, and Y. Bengio : "Neural machine translation by jointly learning to align and translate," arXiv preprint arXiv:1409.0473 (2014)
- 11) A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin : "Attention is All You Need," Proc. of the 30th International Conference on Neural Information Processing Systems (NeurIPS 2017), pp. 6000-6010 (2017)
- 12) J. Devlin, M. W. Chang, K. Lee, and K. Toutanova : "BERT: Pre-training of Deep Bidirectional Transformers for

- Language Understanding," arXiv preprint arXiv:1810.04805 (2018)
- 13) A. Radford, K. Narasimhan, T. Salimans, and I. Sutskeve : "Improving Language Understanding by Generative Pre-Training," Open AI technical report (2018)
  - 14) 安藤俊幸：機械学習を用いた効率的な特許調査方法ーディープラーニングの特許調査への適用に関する基礎検討ー, Japio YEAR BOOK 2018, pp. 238-249 (2018)
  - 15) 坪田匡史, 宮村祐一, 神津友武：深層学習を利用した特許請求項ベースの特許技術俯瞰マップ, 第 34 回人工知能学会全国大会論文集 (2020)
  - 16) 難波英嗣：類似内容の特許請求項の自動対応付け, 情報処理学会 第 142 回情報基礎とアクセス技術・第 120 回ドキュメントコミュニケーション合同研究発表会 (2021)
  - 17) 川上成年：生成 AI の特許データ分析への活用について, 日本マーケティング学会ワーキングペーパー, Vol. 9, No. 19, pp. 1-12 (2023)
  - 18) 伊藤孝佑, 久慈渉, 後藤昌夫：特許庁における AI 技術の活用の現状と最新の取組, 情報の科学と技術, Vol. 73, No. 7, pp. 256-261 (2023)
  - 19) 那須川哲哉：『テキストマイニングを使う技術／作る技術：基礎技術と適用事例から導く本質と活用法』, 東京電機大学出版局 (2006)
  - 20) 金明哲：『テキストアナリティクス』, 共立出版 (2018)
  - 21) S. Deerwester, S. T. Dumais, G. W. Furnas, T. K. Landauer, and R. Harshman : "Indexing by Latent Semantic Analysis," Journal of the American Society for Information Science, Vol. 41, No. 6, pp. 391-407 (1990)
  - 22) D. D. Lee and H. S. Seung : "Learning the parts of objects by non-negative matrix factorization," Nature, vol.401, pp. 788-791 (1999)
  - 23) T. Hofmann : "Probabilistic Latent Semantic Analysis," Proc. of the 15th Conference on Uncertainty in Artificial Intelligence, pp. 289-296 (1999)
  - 24) D. M. Blei, A. Y. Ng, and M. I. Jordan : "Latent Dirichlet Allocation," Journal of Machine Learning Research, Vol. 3, pp. 993-1022 (2003)
  - 25) 柳井啓司：確率トピックモデルによる Web 画像の分類, 人工知能学会全国大会論文集, vol. 22 (2008)
  - 26) 石垣司, 竹中毅, 本村陽一：百貨店 ID 付き POS データからのカテゴリ別状況依存的変数間関係の自動抽出法, オペレーションズ・リサーチ, Vol. 56, No. 2, pp. 77-83 (2011)
  - 27) 繁榊算男, 植野真臣, 本村 陽一：『ペイジアンネットワーク概説』, 培風館 (2006)
  - 28) 野守耕爾, 北村光司, 本村陽一, 西田佳史, 山中龍宏, 小松原明：大規模傷害テキストデータに基づいた製品に対する行動と事故の関係モデルの構築ーエビデンスベースド・リスクアセスメントの実現に向けてー, 人工知能学会論文集, Vol. 25, No. 5, pp. 602-612 (2010)
  - 29) 野守耕爾：テキストマイニングに複数の人工知能技術を応用した特許文書分析と技術戦略の検討, 情報の科学と技術, Vol. 68, No. 8, pp. 32-337 (2018)
  - 30) 特許第 6085888 号：分析方法, 分析装置及び分析プログラム (2017 年 2 月 10 日登録)
  - 31) 野守耕爾, 神津友武：口コミビッグデータに人工知能を応用した地域観光の次世代マーケティング, 人工知能学会全国大会論文集, Vol. 30 (2016)
  - 32) H. Akaike : "Information Theory and an Extension of the Maximum Likelihood Principle," Proc. of the 2nd International Symposium on Information Theory, pp. 267-281 (1973)
  - 33) 若杉徹, 高橋勲男：医薬品調剤履歴に関する確率的構造解析に基づく適応症の推定, 人工知能学会全国大会論文集, Vol. 28 (2014)
  - 34) 野守耕爾：確率的因果意味解析 (PCSA) ーテキストデータを用いたターゲット事象の要因トピックの抽出ー, 人工知能学会全国大会論文集, Vol. 32 (2018)
  - 35) 特許第 7221526 号：分析方法, 分析装置及び分析プログラム (2023 年 2 月 6 日登録)

- 36) 野守耕爾：differential PLSA –テキスト情報の典型的なトピックではないより個性的なトピックの抽出–, 人工知能学会全国大会論文集, Vol.33 (2019)
- 37) 特許第 7221527 号：分析方法, 分析装置及び分析プログラム (2023 年 2 月 6 日登録)