

スタッキングアルゴリズムを用いた 特許拒絶理由類型の判別

本 橋 永 至 高 橋 省 吾 真 鍋 誠 司
鈴 井 智 史 伊 田 英 紀 松 井 重 明

1. はじめに

近年、企業における知財マネジメントの重要性が高まっている。あらゆる市場がグローバルに展開されるようになり、企業は他社に対する技術的な競争優位を保持するために、発明した知財を適切に守らなければならなくなった。また、自社が保有する知財を部分的に開放したり、反対に、他社が開放した知財を活用し、新しい製品・サービスを生むオープン・イノベーションも注目されている (Chesbrough 2013)。

一方、AIの進歩により、AIを活用した特許業務の範囲が拡大している (野崎 2018)。具体的には、自然言語処理や機械学習手法を用いて特許情報を分析することにより、発明する技術分野を検討したり、特許業務を自動化する業務などが挙げられる。近年、特許情報の分析に関するAIツールが多数開発されており、実際に多くの企業において利用されている。宇野ら (2016) は、今後、すべての特許業務をAIに置き換えることは不可能だが、機械が人間に勝る部分をうまく活用し、人間が更に生産性を上げられるように協調すべきだと主張している。

特許業務において、発明の内容を説明する明細書の作成も重要な業務の1つである。たとえ優れた技術を発明しても、発明の内容が明細書に的確に表現できなければ、出願は拒絶される。出願が拒絶されても明細書の開示範囲の内容で補正を行うことが可能であるが、拒絶理由の一部の記載不備については、補正が困難な場合があるため、可能な限り記載不備のない明細書を作成することが求められる。

出願が拒絶される際、特許庁から出願人に対して、拒絶理由が記載された拒絶理由通知書という文書が送られる。過去の拒絶理由通知書を分析することにより、明細書作成業務の課題を発見し、適切な対策を講じることができる。具体的には、特許の審査基準となる特許法の各条項は特許・実用新案審査基準に規定されたいくつかの類型に分類されるため、どの類型に違反することが多いかを把握することにより、明細書作成業務を改善することが期待できる。しかしながら、拒絶理由通知書には、拒絶理由条項は記載されているが、具体的にどの拒絶理由類型に該当するかは明記されていないため、通知書の内容から判断しなければならない。

企業における明細書作成業務を効率化するためには、知財に関わる部署が拒絶理由の傾向を把握し、全社的な対策をとる必要があるが、現状、そのような取組みはほとんど行われていな

い。出願件数が少ない企業においては、拒絶理由類型を手動で分類しても大した労力を必要としないが、年間数千件の出願を行う企業においては、人手による分類は多大な労力を必要とするため現実的ではない。また、同一の拒絶理由類型であっても拒絶理由通知書に記載される表現方法は多岐に渡るため、拒絶理由通知書のテキストデータから拒絶理由類型を正確に判別することは容易ではない。そのため、多くの出願を行う企業にとって、拒絶理由通知書の拒絶理由類型を自動的に判別することができれば、人的リソースを必要とせず、明細書作成業務を改善することが可能となる。

本研究の目的は、特許出願における拒絶理由通知書から、拒絶理由条項ごとの拒絶理由類型を判別する手法を提案し、その有用性を示すことである。本研究では、複数のモデルから予測値を求め、その予測値を特徴量とするメタモデルから最終的な予測値を求めるスタッキングアルゴリズムを用いる。実験では、通知書の実データを用いて、提案手法の有用性を検証する。

次節以降の本論文の構成は、以下の通りである。2節では、本論に関連する研究を整理する。3節では、提案手法を紹介する。4節では、実験方法と使用するデータを説明した後、実験結果について考察する。最後に、5節で本論のまとめと今後の課題を整理する。

2. 関連研究

2.1 特許情報とテキストマイニング

これまでに特許情報を活用する研究は多数行われてきた (Abbas et al. 2014)。特許情報の分析においては、特許公報などのテキストデータが主たる分析対象となるため、自然言語処理によるテキストマイニングに関する研究が多い (那須川 2018)。例えば、西山ら (2009) は、用言ベースの構文パターンを用いて、技術的優位点の特徴の記述表現を特許文書から自動抽出する手法を提案した。福田ら (2013) は、論文と特許から技術動向に関する情報を抽出し、マップを自動作成して可視化する方法を提案した。野守 (2018) は、特許文書にテキストマイニングを応用した研究の目的として、1) 全体像の把握、2) トレンドの把握、3) 競合他社の把握、4) 用途と技術の関係の把握、などが挙げられると整理している。さらに、確率的潜在意味解析 (Probabilistic Latent Semantic Analysis: PLSA) とベイジアンネットワークに基づくテキスト分析技術とそれを特許文書データに適用した分析事例を紹介した。これら研究のように、テキストマイニングによる特許業務の効率化は様々な視点から行われているが、過去に本研究で扱うような明細書の品質向上に関する研究は行われていない。

2.2 スタッキング

機械学習の分野において、1990年代からスタッキングと呼ばれる予測手法の研究が行われている。スタッキングとは、複数のモデルから予測値を求め、その予測値を特徴量とするメタモデルから最終的な予測値を求める手法である。スタッキングは、KDD CUPやkaggleなどの予測精度を競うコンペにおいて、近年、高く評価されるようになったが、ビジネスへの応用研究は少ない。

これまで、スタッキングに関する研究として、メタモデルで用いる特徴量の生成やメタモデル自体の構築方法などに関する研究が行われてきた。Wolpert (1992) は、メタモデルで用いる特徴量の型とメタモデルの学習アルゴリズムの型を考慮することが重要であることを示した。

発送番号 ○○○○ 発送日 令和○年○月○日	拒絶理由通知書	特許出願の番号 ○○○○ 起案日 令和○年○月○日 特許庁審査官 ○○○○ 特許出願人代理人 ○○○○ 適用条文 第36条
1. (サポート要件) 拒絶理由条項について	理 由	
2. (明確性) 拒絶理由条項について		
●理由1について 拒絶理由の説明	記	
●理由2について 拒絶理由の説明		
----- <先行技術文献調査結果の記録>		
審査第○部 ○○○○		

図1 拒絶理由通知書のイメージ

Breiman (1996a) は，異なる大きさの回帰木，もしくは，異なる次元の特徴量を用いた線形回帰を第1段階のモデルとし，すべての回帰係数が非負となる制約を課した線形回帰をメタモデルとするスタック回帰を提案した. Ting and Witten (1999) は，メタモデルにおける特徴量として，判別値ではなく，クラス確率を用いることを推奨している. さらに，Breiman (1996a) が提案した非負の制約は，分類においては有効ではないことを明らかにした. 田村ら (2018) は，テクニカル分析のモデルとファンダメンタル分析のモデルの双方の予測値を用いることにより，スタッキングが株主価値推定に有効であることを示した. メタモデルで用いる特徴量を生成する個々のモデルについては，互いに異なるべきであり，多様性を持たすことが性能の向上に繋がることが知られている (Tumer and Ghosh 1996).

3. 提案手法

提案手法は，拒絶理由通知書から拒絶理由に関するテキストの抽出，形態素出現行列の作成，スタッキングアルゴリズムによる判別モデルの構築の3つのステップから構成される. 以降，各ステップを順に説明する.

3.1 拒絶理由に関するテキストの抽出

拒絶理由通知書には，最初に，特許法のどの条項に違反するかが記載されており，その後，条項ごとに具体的な拒絶理由が記載されている. 図1は，拒絶理由通知書で用いられる典型的

表1 実施可能要件・委任省令要件 (36条4項1号)

類型	内容
1	技術的手段の記載が抽象的又は機能的である場合
2	技術的手段相互の関係が不明確である場合
3	製造条件等の数値が記載されていない場合
4	発明の詳細な説明に、請求項に記載された上位概念に含まれる一部の低位概念についての実施の形態のみが実施可能に記載されている場合
5	発明の詳細な説明に、特定の実施の形態のみが実施可能に記載されている場合
6	マーカッシュ形式で記載された請求項の場合
7	達成すべき結果によって物を特定しようとする記載を含む請求項の場合
8	委任省令要件で記載することが求められる事項が記載されていない場合

表2 サポート要件 (36条6項1号)

類型	内容
1	請求項に記載されている事項が、発明の詳細な説明中に記載も示唆もされていない場合
2	請求項及び発明の詳細な説明に記載された用語が不統一であり、その結果、両者の対応関係が不明瞭となる場合
3	出願時の技術常識に照らしても、請求項に係る発明の範囲まで、発明の詳細な説明に開示された内容を拡張ないし一般化できるとはいえない場合
4	請求項において、発明の詳細な説明に記載された、発明の課題を解決するための手段が反映されていないため、発明の詳細な説明に記載した範囲を超えて特許を請求することになる場合

表3 明確性要件 (36条6項2号)

類型	内容
1	請求項の記載自体が不明確である結果、発明が不明確となる場合
2	発明特定事項に技術的な不備がある結果、発明が不明確となる場合
3	請求項に係る発明の属するカテゴリーが不明確であるため、又はいずれのカテゴリーともいえないため、発明が不明確となる場合
4	発明特定事項が選択肢で表現されており、その選択肢同士が類似の性質又は機能を有しないため、発明が不明確となる場合
5	範囲を曖昧にし得る表現がある結果、発明の範囲が不明確となる場合
6	機能、特性等を用いて物を特定しようとする記載がある場合
7	サブコンビネーションの発明を「他のサブコンビネーション」に関する事項を用いて特定しようとする記載がある場合
8	製造方法によって生産物を特定しようとする記載がある場合

なフォーマットのイメージである。記載不備に関する代表的な条項として、36条4項1号の実施可能要件・委任省令要件、36条6項1号のサポート要件、36条6項2号の明確性要件などが挙げられる。さらに、各条項はいくつかの類型に分類される。具体的には、36条4項1号は8類型(表1)、36条6項1号は4類型(表2)、36条6項2号は8類型(表3)に分類される。

各拒絶理由は、必ずいずれかの類型に分類されるが、1つの条項違反に対して、複数の類型が該当する場合がある。例えば、ある通知書において、36条4項1号の違反が存在する場合、類型1の「技術的手段の記載が抽象的又は機能的である場合」と類型4の「発明の詳細な説明に、請求項に記載された上位概念に含まれる一部の低位概念についての実施の形態のみが実施可能に記載されている場合」の双方に該当する場合がある。

各拒絶理由がどの類型に該当するかを判別するためには、まず当該通知書において、どの条項に違反するかを識別し、さらに、各条項において、どの類型に該当するかを識別する必要がある。どの条項に違反するかについては、拒絶理由通知書に明記されているため、比較的容易に識別することが可能である。しかしながら、どの類型に該当するかを識別するためには、該当する条項の説明文を抽出し、そのテキストデータから判断しなければならない。拒絶理由通知書は、ある程度決まったフォーマットで記載されているが、審査官によって書き方が微妙に異なる。そのため、拒絶理由の説明文のみを抽出する自然言語処理アルゴリズムを用いてテキストの抽出を行う。

3.2 形態素出現行列の作成

条項ごとに拒絶理由の説明文を抽出した後、形態素解析により、テキストを形態素単位に分割する。次に、各通知書の条項ごとに、各形態素の出現回数を算出し、形態素出現行列を作成する。本研究では、通知書の拒絶理由類型を判別することが目的であるため、判別に影響を与えない明細書の記載からの引用された発明技術に関するテキストは分析に使用しない。通知書において、通常、そのような内容の文章は鍵括弧で囲われている。

3.3 判別モデルの構築

各拒絶理由は、条項ごとにいくつかの類型に分類されるが、前述の通り、同一の通知書において、1つの条項違反に対して、複数の類型が該当する場合がある。したがって、提案手法では、各通知書の各条項について、各類型に該当するか否かの2値変数を目的変数とする判別モデルを構築する。例えば、ある通知書において、36条4項1号違反に該当するとき、当該条項は8つの類型があるため、最大で8つの判別モデルを構築することになる。

提案モデルでは、複数のモデルの予測値を用いて最終的な予測値を求めるスタッキングアルゴリズムを用いる。スタッキングアルゴリズムでは、まず使用する各モデルから予測値を求め(第1段階)、その予測値を特徴量とするメタモデルにより、最終的な予測値を求める(第2段階)。図2はスタッキングアルゴリズムのパラメータ推定の流れを表す概念図である。

具体的には、モデルの多様性を考慮して、スパースロジスティック回帰 (Friedman et al. 2010)、ニューラルネットワーク (Ripley 1996)、決定木 (Breiman et al. 1984)、バギング (Breiman 1996b)、ランダムフォレスト (Breiman 2001)、XGBoost (Chen and Guestrin 2016) の6つのモデルの予測値をメタモデルの特徴量として用いる。ただし、Ting and Witten (1999)と同様に、個々のモデルの予測値は、2値の判別結果ではなく目的変数が1となる確率とする。

個々のモデルを学習する際、すべての訓練用データを同時に用いてしまうと、予測値は正解ラベルを用いて得られたものであるため、過学習を起こしてしまう。そのため、訓練用データを k 個のデータセットに分割し、 j 番目のデータセットの予測値については、 j 番目以外の

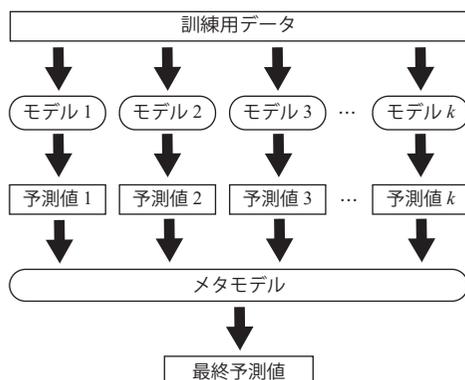


図2 スタッキングアルゴリズムの概念図

データセットを用いて求めるという作業を k 回繰り返し、すべてのデータセットの予測値を求める。この作業を第1段階のすべてのモデルに対して実行し、メタモデルの特徴量を生成する。

次に、それらの予測値を特徴量とするメタモデルをスパースロジスティック回帰により構築する。つまり、個々のモデルから得られた6次元特徴量ベクトルを $\mathbf{x} = (x_1, \dots, x_6)^T$ として、目的変数 Y が1をとる確率 $P(Y = 1 | \mathbf{x}_i)$ を

$$\pi(\mathbf{x}_i) = \frac{\exp(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta})}{1 + \exp(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta})}, \quad (i = 1, \dots, n) \quad (1)$$

と定式化する。ここで、 β_0 は切片、 $\boldsymbol{\beta} = (\beta_1, \dots, \beta_6)^T$ は回帰係数ベクトルである。目的変数 Y は2値の値をとる確率変数であるので、 Y の確率関数は

$$f(Y = y | \mathbf{x}) = \pi(\mathbf{x})^y \{1 - \pi(\mathbf{x})\}^{1-y} \quad (2)$$

である。したがって、パラメータ $\boldsymbol{\beta}$ のいくつかの成分が0になるようスパース推定するためには、正則化対数尤度関数

$$\frac{1}{n} \sum_{i=1}^n [y_i \log \pi(\mathbf{x}_i) + (1 - y_i) \log \{1 - \pi(\mathbf{x}_i)\}] - \lambda \|\boldsymbol{\beta}\|_1 \quad (3)$$

を最大化する $\boldsymbol{\beta}$ を求めればよい (川野ら 2018, Park and Hastie 2007)。ここで、 λ は正の値をとる正則化パラメータである。 λ の値を適切に設定することにより、回帰係数のいくつかの値を正確に0に推定することができる。

提案モデルの利点として、既存のモデルを組み合わせているため実装が容易であり、かつ、

表4 混同行列

		真の値	
		正例	負例
予測値	正例	真陽性 TP	偽陽性 FP
	負例	偽陰性 FN	真陰性 TN

メタモデルにロジスティック回帰を用いているため推定結果の解釈がしやすい点が挙げられる。判別モデルの選択においては、判別精度が最も重要であるが、判別結果について説明を求められることもあるため、モデルの解釈性も考慮すべき要素の1つである。

4. 実験

4.1 実験方法

提案モデルの判別精度を検証するために、提案モデルの第1段階で用いるスパースロジスティック回帰（LR）、ニューラルネットワーク（NN）、決定木（TR）、バギング（BG）、ランダムフォレスト（RF）、XGBoost（XG）を比較モデルとする。判別精度の指標として、正答率、F値、AUC値の3つを用いる。予測値と真の値の組み合わせは、予測値を正例としたか負例としたか、その予測が正しいか誤りかによって4つに分けられる。それらの事例数を表4のように集計すると、正答率（*accuracy*）とF値（ F_1 ）はそれぞれ

$$accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (4)$$

$$F_1 = \frac{2TP}{2TP + FP + FN} \quad (5)$$

と定義される。AUC値とは、ROC曲線を作成した時のグラフの曲線より下の部分の面積である。AUC値は0から1までの値をとり、値が1に近いほど判別能が高いことを示す。判別能力がランダムであるとき、AUC値は0.5となる。

4.2 データ概要

実験には、三菱電機株式会社が2018年から2019年に受領した拒絶理由通知書のテキストデータを用いる。分析に用いる通知書は、36条4項1号、36条6項1号、36条6項2号のいずれかの条項違反が含まれるものである。各拒絶理由がどの類型に分類されるかという正解ラベルについては、特許業務に関して専門的な知識を持つ者が手動により作成した。

テキストデータには、助詞や記号などの判別に寄与しない形態素も含まれる。本実験では、すべての形態素の中から、以下の条件に当てはまる形態素を抽出し、形態素出現行列を作成する。

表5 サンプルサイズ

	36.4.1_1	36.4.1_2	36.6.1_1	36.6.1_3	36.6.1_4	36.6.2_1	36.6.2_2
該当	98	64	357	143	302	495	1,194
非該当	117	151	421	635	476	1,172	473

表6 データ分割後のサンプルサイズ

	36.4.1	36.6.1	36.6.2
訓練用データ	150	544	1,166
検証用データ	65	234	501

表7 特徴量重要度

順位	36.4.1_1	36.4.1_2	36.6.1_1	36.6.1_3	36.6.1_4	36.6.2_1	36.6.2_2
1	拒絶	記載	範囲	拡張	反映	指し示す	指し示す
2	記載	段落	課題	一般	課題	誤記	前
3	経済	詳細	記載	制御	超える	前	誤記
4	もの	本願	超える	常識	いえる	こと	する
5	場合	ため	拡張	事務所 A	内容	記載	よう
6	特定	示す	一般	本願	いる	日本語	れる
7	実施	場合	化	こと	規定	する	の
8	可能	よう	詳細	明確	発明	の	記載
9	こと	開	なる	いる	拡張	項	こと
10	よう	理由	引用	小さい	する	以前	いる

1. 名詞, 形容詞, 動詞のいずれか
2. 全文書を通じて, 3回以上出現
3. ひらがな, もしくは, 漢字

データには, 11の特許事務所から出願した特許に関する通知書が含まれる。そのため, 特徴量として, 各形態素の出現回数に加えて, 各特許事務所のダミー変数も含める。各条項の最終的な特徴量の次元は, 36条4項1号が1,016, 36条6項1号が1,964, 36条6項2号が1,934である。

本研究では, 各条項の類型ごとに, その類型に該当するか否かを判別する。各条項には, 頻出する類型とそうではない類型がある。モデルの学習において十分なサンプルサイズを確保するために, 36条4項1号については類型 [1, 2], 36条6項1号については類型 [1, 3, 4], 36条6項2号については類型 [1, 2] の計7つの類型を分析対象とする。表5は, 各類型の該当と非該当の件数を示す。

提案モデルの判別精度を検証するために, 訓練用データと検証用データの割合がそれぞれ70%, 30%になるようにデータを分割する。訓練用データと検証用データのサンプルサイズは, それぞれ表6の通りである。

表8 正答率

	36.4.1_1	36.4.1_2	36.6.1_1	36.6.1_3	36.6.1_4	36.6.2_1	36.6.2_2	平均
LR	0.4615	0.5846	0.8846	0.9478	<u>0.9316</u>	0.8423	0.8044	0.7797
NN	0.5385	0.6000	0.8846	0.9578	0.9231	0.8443	0.7984	0.7911
TR	<u>0.5692</u>	0.6308	0.8590	0.9658	0.9274	0.8403	0.8124	0.8007
BG	<u>0.5692</u>	0.6000	0.8590	0.9615	0.9188	0.8463	0.8244	0.7970
RF	0.5385	0.6308	0.8974	0.9673	0.9231	0.8423	0.8343	0.8034
XG	0.4769	<u>0.6923</u>	0.8974	<u>0.9701</u>	<u>0.9316</u>	<u>0.8563</u>	0.8144	0.8056
ST	0.5538	0.6308	<u>0.9060</u>	<u>0.9701</u>	0.9231	<u>0.8563</u>	<u>0.8383</u>	<u>0.8112</u>

表9 F値

	36.4.1_1	36.4.1_2	36.6.1_1	36.6.1_3	36.6.1_4	36.6.2_1	36.6.2_2	平均
LR	0.4776	0.7216	0.8898	0.9695	0.9459	0.8879	0.6370	0.7899
NN	0.5946	0.7292	0.8851	0.9701	0.9384	0.8883	0.6667	0.8103
TR	<u>0.6667</u>	0.7647	0.8675	0.9798	0.9431	0.8870	0.6328	0.8202
BG	0.6585	0.7292	0.8653	0.9772	0.9365	0.8917	0.6480	0.8152
RF	0.6154	0.7647	0.8992	0.9750	0.9400	0.8942	0.6612	0.8214
XG	0.5405	<u>0.7959</u>	0.9008	<u>0.9825</u>	<u>0.9467</u>	0.9000	0.6517	0.8169
ST	0.6506	0.7692	<u>0.9091</u>	0.9824	0.9396	<u>0.9003</u>	<u>0.6897</u>	<u>0.8344</u>

4.3 実験結果

4.3.1 特徴量重要度

実験結果に関して、まず特徴量重要度について考察する。表7は、XGBoostによって得られた重要度が高い上位10の特徴量を示す。重要度の高い特徴量は、条項ごとに異なっていることがわかる。例えば、36条4項1号では、技術的手段の記載が適切か否かが問われるので、「記載」という単語が上位に入っている。36条6項1号では、「拡張」、「一般」、「課題」といった単語が上位に入っているが、これは類型3に該当する場合、「拡張ないし一般化」という表現がよく用いられ、また、類型4に該当する場合、「課題を解決するための手段」といった表現がよく用いられるため、妥当な結果といえる。最後に、36条6項2号では、「指し示す」、「誤記」、「前」といった単語が上位に入っているが、これらは類型1に該当する場合、「前記○が何を指しているか不明確」といった表現がよく用いられるためである。事務所のダミー変数については、36条6項1号の類型3のみにおいて上位10に入っていたので、重要度は高くはないといえる。

4.3.2 モデル比較

次に、提案モデルと比較モデルの判別精度について考察する。表8は正答率、表9はF値、

表10 AUC値

	36.4.1_1	36.4.1_2	36.6.1_1	36.6.1_3	36.6.1_4	36.6.2_1	36.6.2_2	平均
LR	0.5140	<u>0.5724</u>	0.9019	0.9216	0.9605	0.8531	0.8031	0.7895
NN	0.5180	0.5169	0.9250	0.9596	0.9744	0.8623	0.8209	0.7967
TR	0.5585	0.5174	0.9072	0.9556	0.9026	0.8201	0.7503	0.7731
BG	0.5890	0.4868	0.9204	0.9557	0.9711	0.8425	0.7905	0.7937
RF	0.6075	0.4831	0.9363	0.9857	0.9718	0.8933	0.8268	0.8149
XG	0.5290	0.5243	<u>0.9465</u>	0.9743	<u>0.9807</u>	0.8884	0.8169	0.8086
ST	<u>0.6250</u>	0.5048	0.9417	<u>0.9868</u>	0.9735	<u>0.8986</u>	<u>0.8290</u>	<u>0.8228</u>

表11 メタモデルの推定値

	36.4.1_1	36.4.1_2	36.6.1_1	36.6.1_3	36.6.1_4	36.6.2_1	36.6.2_2
切片	-2.1391	-1.7457	-3.5092	-4.3043	-4.4581	-3.3881	-4.8228
LR	-1.9029	-	0.3781	0.1313	0.2789	1.4659	0.7763
NN	-	-	-	-	-	-	-0.0160
TR	1.6828	1.7182	1.2066	2.4570	-	0.2852	0.2161
BG	0.5591	0.9363	-	-	2.6322	0.5228	0.1873
RF	3.8465	-	4.0715	3.4233	6.1004	4.4773	6.5261
XG	-	-	1.4250	1.8034	0.3773	1.0695	0.5632

表10はAUC値を示す。条項別の傾向としては、36条6項1号が最も判別精度が高かった。次に36条6項2号の判別精度が高く、36条4項1号が最も低かった。提案モデルと比較モデルの差異に着目してみると、提案モデルは、正答率とAUC値において、7類型中4類型で最も判別精度が高く、F値においては、7類型中3類型で最も判別精度が高かった。各指標における7類型の平均値については、すべての指標において、提案モデルが最も判別精度が高く、正答率、F値、AUC値の平均値はそれぞれ0.8112、0.8344、0.8228だった。総合的に見て、提案モデルは他の比較モデルと比較して最も判別精度が高く、次に、ランダムフォレスト、XGBoostの精度が高かった。

4.3.3 メタモデルの推定値

最後に、メタモデルの推定値について考察する。メタモデルでは、第1段階の各モデルの予測値を特徴量として、最終的な予測値を求める。提案モデルのメタモデルでは、スパースロジスティック回帰を用いたため、最終的な予測値を求める際、各類型の判別に有効なモデルが選択される。

表11は、各類型において、どのモデルが選択されたか、および、それらの推定値を示しており、ハイフンは当該モデルが選択されなかったことを意味している。

類型ごとに、メタモデルにおいて選択されたモデルが異なることがわかる。例えば、36条4項1号の類型2では、決定木とバギングのみが選択されたのに対して、36条6項2号の類型2では、すべてのモデルが選択された。最も多く選択されたモデルは、スパースロジスティック回帰、決定木、ランダムフォレストであり、それらは7類型中6類型で選択された。今回の実験では、7つの類型しか用いなかったため、どのようなデータに対して、どのモデルが選択されるのかについては明言できない。しかしながら、36条4項1号については、他の条項と比べてサンプルサイズが小さかったため、サンプルサイズの大きさが選択されたモデルの数に影響した可能性が考えられる。

5. おわりに

5.1 本研究の貢献

本研究では、拒絶理由通知書の拒絶理由類型を判別する手法を提案した。提案モデルは、スタッキングアルゴリズムを用いており、実験において、既存のモデルと比べて判別精度が優れていることを示した。特許出願は、たとえ技術が優れていても、適切な明細書が作成できていなければ登録されない。すなわち、明細書に記載不備があれば、拒絶理由通知書が送付され、その対処に無用の費用と時間が浪費されることになる。提案手法を用いて、企業や特許事務所ごとに、どの拒絶理由が多いかを把握することにより、明細書作成における課題を把握することができる。その課題に応じた対策を講じることにより、出願が拒絶される確率を減少させ、特許出願費用の削減につながることを期待される。

本研究の貢献として、大きく2点が挙げられる。第1に、本研究では、拒絶理由通知書のテキストデータを用いて、拒絶理由類型を判別する手法を初めて提案した。過去に特許情報のデータを用いた研究は数多く行われてきたが、拒絶理由通知書を分析することによる明細書の品質向上を目的とする研究は、過去に行われていない。しかしながら、特許業務において、明細書の品質向上は、出願費用の削減に貢献できるため、重要な課題である。本研究が特許明細書の品質向上に寄与する研究の契機となり、実務においても、有効に活用されることが期待される。第2に、本研究は、スタッキングアルゴリズムによる判別モデルが既存のモデルを単体で用いるよりも判別精度が高いことを示した。近年、スタッキングの有用性に注目が集まっているが、実データに対してスタッキングモデルを適用し、判別精度を検証した研究は少ない。本研究で提案した判別モデルは実装が容易であり、分析結果の解釈が可能である点も実務上有益であるといえる。

5.2 今後の課題

今後の課題として、3点が挙げられる。第1に、スタッキングアルゴリズムに用いるモデルの選択についてである。スタッキングでは、第1段階で用いる複数のモデルと第2段階で用いるメタモデルを選択する必要がある。第1段階で用いるモデルは、多様性が高いほど判別精度が高いと言われているが、具体的にどのようなモデルの組み合わせが最適であるかは明らかではない。また、第2段階で用いるメタモデルについても、最適なモデルは不明である。一方、

kaggleなどのコンペにおいて、第1段階において数百ものモデルを用いることにより、予測精度が向上したという報告がある。スタッキングアルゴリズムに用いるモデルの選択方法に関する研究が期待される。第2に、正解データの作成についてである。本研究では、手動により正解データを作成したが、拒絶理由がどの類型に分類されるかの判断は専門知識がないと難しい。そのため、正解データを作成できる者が限られる点は、提案モデルを実務に応用する際の課題になるだろう。第3に、本研究では、3つの拒絶理由条項に限定して分析を行ったため、他の条項については、判別精度が不明である。本モデルを実務でより有効に活用するためには、他の条項についても検証する必要があるだろう。

参 考 文 献

- [1] H. W. Chesbrough, *Open Innovation: The New Imperative for Creating and Profiting from Technology*, Harvard Business School Press, 2003.
- [2] 野崎篤志, “特許情報と人工知能 (AI): 総論,” 『情報の科学と技術』, 68, pp. 316-325, 2018.
- [3] 宇野毅明, 野崎篤志, 那須川哲哉, 小川延浩, “人工知能が知財業務に及ぼす影響,” 『パテント』, 69, pp. 10-18, 2016.
- [4] A. Abbas, L. Zhang and S. U. Khan, “A Literature Review on the State-of-the-art in Patent Analysis,” *World Patent Information*, 37, pp. 3-13, 2014.
- [5] 那須川哲哉, “テキストアナリティクスと特許情報分析,” 『情報の科学と技術』, 68, pp. 326-331, 2018.
- [6] 西山莉紗, 竹内広宣, 渡辺日出雄, 那須川哲哉, “新技術が持つ特長に注目した技術調査支援ツール,” 『人工知能学会論文誌』, 24, pp. 541-548, 2009.
- [7] 福田悟志, 難波英嗣, 竹澤寿幸, “論文と特許からの技術動向情報の抽出と可視化,” 『情報処理学会誌』, 6, pp. 16-29, 2013.
- [8] 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』, 68, pp. 332-337, 2018.
- [9] D. H. Wolpert, “Stacked Generalization,” *Neural Networks*, 5, pp. 241-260, 1992.
- [10] L. Breiman, “Stacked Regressions,” *Machine Learning*, 24, pp. 49-64, 1996a.
- [11] K. M. Ting and I. H. Witten, “Issues in Stacked Generalization,” *Journal of Artificial Intelligence Research*, 10, pp. 271-289, 1999.
- [12] 田村浩一郎, 上野山勝也, 飯塚修平, 松尾豊, “深層学習を用いたアンサンブルモデルによる株主推定モデルの提案,” 『人工知能学会論文誌』, 33, pp. 1-11, 2018.
- [13] K. Tumer and J. Ghosh, “Theoretical Foundations of Linear and Order Statistics Combiners for Neural Pattern Classifiers,” Technical Report, IEEE Transactions on Neural Networks, 1996.
- [14] J. Friedman, T. Hastie and R. Tibshirani, “Regularization Paths for Generalized Linear Models via Coordinate Descent,” *Journal of Statistical Software*, 33, pp. 1-22, 2010.
- [15] B. D. Ripley, *Pattern Recognition and Neural Networks*, Cambridge University Press, 1996.
- [16] L. Breiman, J. Friedman, C. J. Stone and R. A. Olshen, *Classification and Regression Trees*, Chapman and Hall/CRC, 1984.
- [17] L. Breiman, “Bagging Predictors,” *Machine Learning*, 24, pp. 123-140, 1996b.
- [18] L. Breiman, “Random Forest,” *Machine Learning*, 45, pp. 5-32, 2001.
- [19] T. Chen and C. Guestrin, “XGBoost: A Scalable Tree Boosting System,” *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785-794, 2016.
- [20] 川野秀一, 松井英俊, 廣瀬慧, 『スパース推定法による統計モデリング』, 共立出版, 2018.
- [21] M. Y. Park and T. Hastie, “L1-regularization Path Algorithm for Generalized Linear Models,” *Journal of the Royal Statistical Society Series B*, 69, pp. 659-677, 2007.

〔もとはし えいじ 横浜国立大学大学院国際社会科学研究院准教授〕

〔たかはし しょうご 鹿児島大学産学・地域共創センター教授〕

〔まなべ せいじ 横浜国立大学大学院国際社会科学研究院教授〕

〔すずい さとし 三菱電機株式会社知的財産センター特許業務管理部業務推進第一グループ専任〕

〔いだ ひでのり 三菱電機株式会社知的財産センター特許業務管理部業務推進第四グループマネージャー〕

〔まつい じゅうめい 三菱電機株式会社知的財産センター副センター長〕

〔2021年11月8日受理〕