

# 外国特許文献への分類付与に関する 機械学習活用可能性調査について

Possibility of utilizing machine learning for classification of foreign patent documents

特許庁 企画調査官

富永 泰規

1998年4月特許庁入庁。特許審査、審判に従事のほか、情報技術統括室を経て、2017年4月より現職、審査第一部調整課審査企画室に併任。

## 1 はじめに

特許庁は、「世界最速・最高品質の特許審査」の実現を目標として掲げており、その実現のためには先行技術文献調査を的確、かつ、効率的に実施するための環境整備が重要である。

現在、先行技術文献調査を行うための分類検索キーとして、内国特許文献には、国内分類であるFI及びFターム（以下「Fターム等」という。）が付与されている一方、ほとんどの外国特許文献には、Fターム等が付与されおらず、審査官が内国特許文献及び外国特許文献を検索する際には、Fターム等を用いた国内特許文献の検索とCPC等を用いた外国特許文献の検索とを、異なる分類・検索キーにより行う必要がある。先行技術調査の対象となる外国特許文献が増大し続けている中、外国特許文献を国内特許文献と統一的な付与基準で付与された検索インデックスを用いて横断的に検索できれば、審査官が的確かつ効率的に先行技術を調査できるようになるが、人手により外国特許文献にFターム等を付与する作業は、多大な時間と多大なコストを必要とし、現実的であるとはいえない。

一方、近年、人工知能に関する技術分野の一つである機械学習の分野は、技術進展が著しく、特許や論文、WEBページ等の文章を対象とした、機械学習による文書自動分類の研究も盛んに行われている<sup>[1]</sup>。一般に、機械学習による文書自動分類においては、学習用のデータとして学習用の文献データが必要となるが、外国特許文献についてみると、パテントファミリーという形で、F

ターム等が付与された国内特許文献と紐付いており、これらを活用することにより学習用のデータを準備することができるため、機械学習を用いて外国特許文献にFターム等を付与することは、現実的である可能性がある。

そこで、特許庁では、昨年度の記事<sup>[2]</sup>において紹介させて頂いたように、外国特許文献への機械学習を活用したFターム等付与に関して調査研究を行った。本稿においては、本調査事業の一端を紹介させて頂くこととしたい。

なお、本稿は、昨年度実施した調査事業の報告書を踏まえ、著者の私見に基づいて記載したものであり、特許庁としての意見・見解を表明するものではない。

## 2 分類付与モデルごとの精度比較調査

本調査では、分類付与モデルによる付与精度を検証すべく、特定の技術分野を選定し、英語、独語、中国語、日本語の4つの言語の特許文献を対象として、SVM（サポートベクターマシン）、MLP（多層パーセプトロン）、RNN（リカレントニューラルネットワーク）、NAM（ニューラルアテンションモデル）の4種類の分類付与モデルによる検証を行った。機械学習で使用する学習用データ（教師文献）としては、日本語ファミリー特許文献を持つ外国語特許文献を使用し、外国語特許文献の正解分類として、日本語ファミリー特許文献に付与されているFI・Fタームを流用した。

## 2.1 検証した分類付与モデルの概要

検証を行った4種類の分類付与モデルについて、概要を紹介する。

### (1) SVM (サポートベクターマシン)

SVM (サポートベクターマシン) は、基本的には、2クラスを識別する機械学習モデル (2値分類モデル) であり、2クラスを識別する学習モデルを教師文書から生成し、新規文書に対してどちらのクラスに属するかを識別するものである。

特許分類付与は、一つの特許文書に一つ以上の分類が付与されるマルチラベル分類であり、SVMをマルチラベル分類に適用する方式としては、One-vs-rest、One-vs-one という2種類の方式があるが、本調査では、学習モデルの作成時間を考慮し、One-vs-rest方式 (図1参照) を採用した。

### (2) MLP (多層パーセプトロン)

MLP (多層パーセプトロン) は、最もシンプルでベーシックなニューラルネットワークであり、入力層、出力層と、1層以上の中間層 (隠れ層) から構成される。MLPでは、入力層のノードと中間層のノード、中間層のノードと出力層のノードが全結合した構成となっており、学習時には、出力データと教師データの差分 (誤差) に基づいて、中間層のノードの重み値を更新するバックプロパゲーション (誤差逆伝播学習法) が採用される。

MLPでは、出力層に複数のノードを分類の数だけ設けて、複数の2値分類タスクを同時に解く方式 (以下、「一括付与方式」と呼ぶ) と、MLPを分類の数だけ設けて、分類毎にその分類を付与すべきかを解く方式 (以下、「個別付与方式」と呼ぶ) があるが、本調査では、学習モデルの作成時間を考慮し、一括付与方式 (図2参照) を採用した。

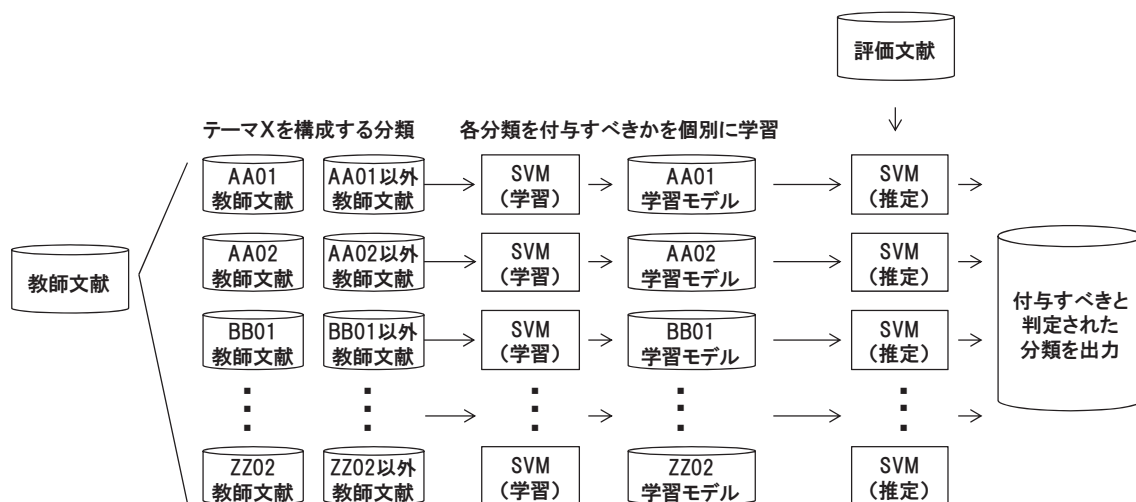


図1 SVM (One-vs-rest 方式) による分類付与方式

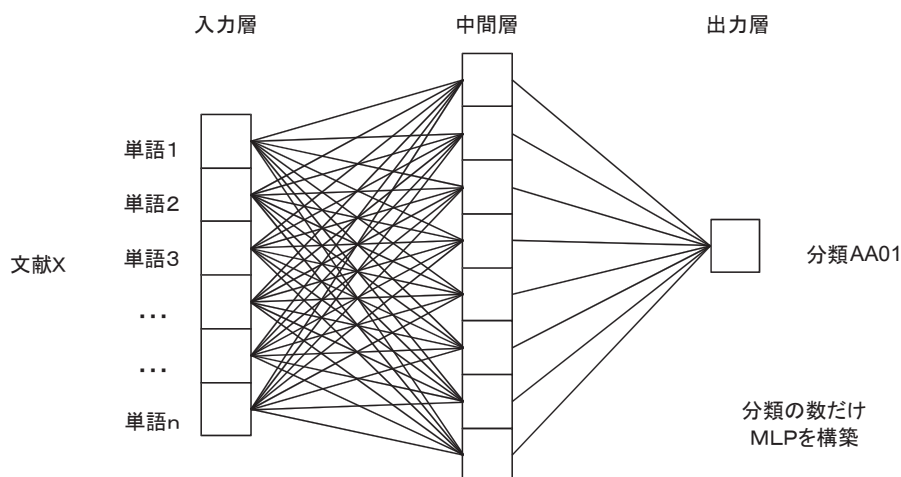


図2 MLP (個別付与方式) による分類付与方式

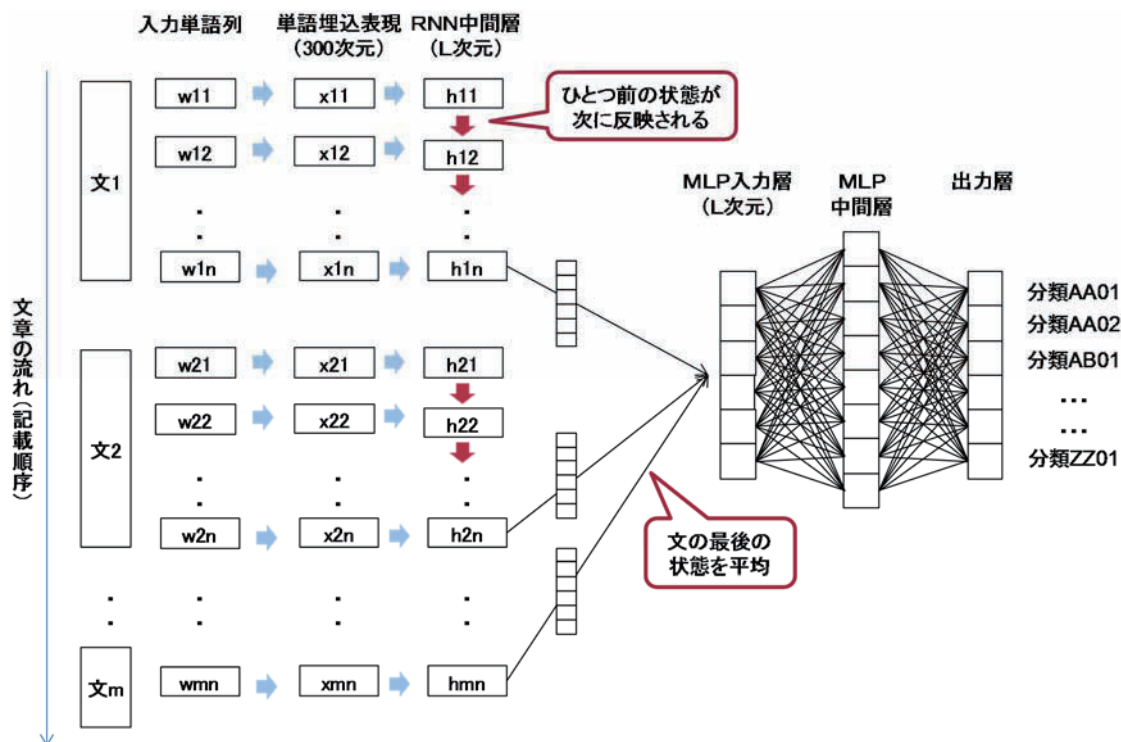


図3 RNNによる特許分類付与方式

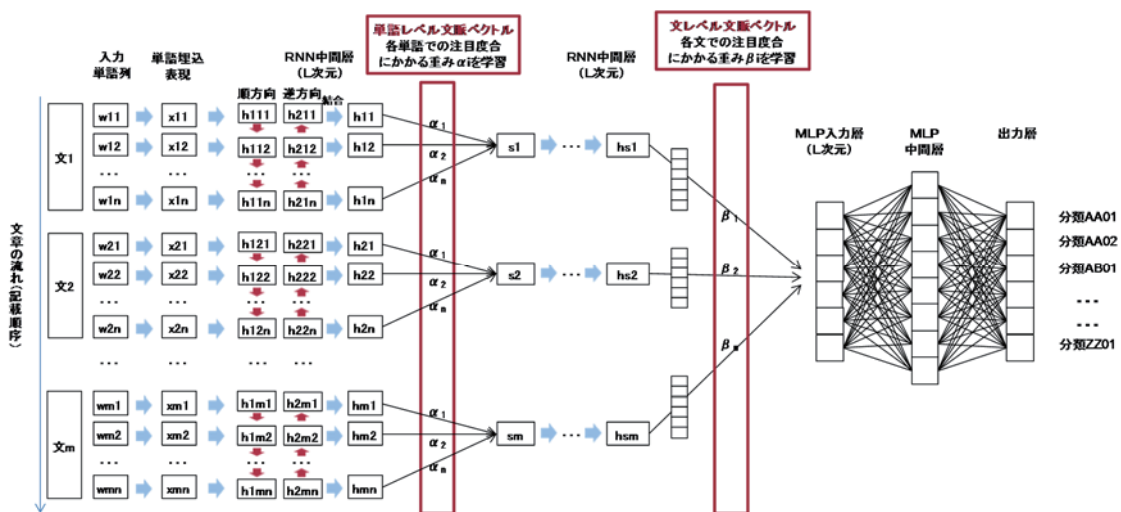


図4 NAMによる特許分類付与方式

### (3) RNN (リカレントニューラルネットワーク)

RNN (リカレントニューラルネットワーク) は、文章の記載順序を保持して学習することにより、文章の記載順序を考慮した学習を可能としたニューラルネットワークであり、中間層を構成するノードにおいて、ひとつ前までのノードの情報 (の一部) を、次のノードに反映させることにより、文章の記載順序を反映した学習ができるものである (図3 参照)。

### (4) NAM (ニューラルアテンションモデル)

NAM (ニューラルアテンションモデル) は、RNN

において、文献に現れる各単語の注目度合 (Attention) を学習し、分類付与に貢献する単語とそうでない単語を識別した上で分類を付与するニューラルネットワークであり、分類を付与する根拠となった記載箇所を学習して分類を付与できるという特長を有する (図4 参照)。

## 2.2 分類付与モデル毎の分類付与結果

外国語特許文献への機械分類付与は、対象言語として原文 (英語、中国語、ドイツ語) を用いる方法と、機械翻訳文 (日本語) を用いる方法とが想定される。また、解析対象として、全文を対象とする方法と要約を対象と

する方法が想定される。SVM、MLP、RNNについては、対象言語として原文と機械翻訳文の双方について、また、解析対象として全文及び要約の双方について、付与精度の検証を行った。NAMについては、学習時間を要したことから、対象言語として英語、解析対象として要約について、付与精度の検証を行った。

その結果、3種類の機械学習モデル（SVM、MLP、RNN）の付与精度は、いずれの対象言語の場合、及び、解析対象の場合においても、SVM、MLP、RNNの順であるという結果が得られた。

また、4種類の機械学習モデル（SVM、MLP、RNN、NAM）の付与精度は、対象言語として英語、解析対象として要約である場合において、概ねSVM、NAM、MLP、RNNの順であるという結果が得られた。

なお、本調査においては、SVMの付与精度がもっとも高いという結果が得られたが、学習データの品質や、学習データの量により、付与精度が変化する可能性があり、また、パラメータのチューニングにより、複雑なモデルであるMLP、RNN、NAMの付与精度が、SVMの付与精度を大幅に超える可能性はあると考えられる。

## 2.3 分類付与向上に向けたエラー分析

外国特許文献に対するFI・Fタームの分類付与技術の実用化のために、その付与精度を向上すべく、機械学習を利用した分類付与モデルによる付与精度の低下原因を分析した。その結果得られた、計算機による分類付与が苦手な分類の特性について以下紹介する。

### （特性1）分類付与根拠が、特許文献のごく一部の記載箇所には現れない分類

分類付与根拠が特許文献全体のごく一部であるため、他の記載箇所に含まれる単語が、ノイズ語となって、誤付与を引き起こす。ただし、分類付与根拠の記載箇所を特定できさえすれば、ノイズ語がなくなるため、付与精度は格段に向上すると考えられる。

解決策としては、教師文献による学習、および、評価文献に対する分類付与の際に、当該分類に対応する分類付与根拠を特定し、そこに現れる単語のみを活用することが考えられる。

### （特性2）分類付与根拠が、句・文・段落レベルで記載されている分類

解決策としては、単語単位で学習するのではなく、フレーズ・文・段落といった単位で、記載内容を表現することが考えられ、例えば、単語の埋込表現ベクトルをフレーズ・文・段落単位で平均したベクトルを素性として使用することが考えられる。

### （特性3）分類付与根拠が、高頻度で使用されている単語の組合せによって記載されている分類

その技術分野においてDF（Document Frequency）の高い単語は、どの文献にも出現しているため、特定の分類を特徴付ける単語となりにくい。

解決策としては、特性2と同様、単語単位で学習するのではなく、フレーズ・文・段落といった単位で、記載内容を表現することが考えられる。

### （特性4）分類付与根拠が、特殊な語彙または表記によって記載されている分類

特許文献は、不特定多数の人によって執筆されているため、他の人が使用しない特殊な語彙あるいは表記によって分類付与根拠が記載されている場合がある。そのような語彙は、教師文献にも含まれていない、あるいは、含まれているとしてもDFが非常に低いため、学習結果に反映されにくい。

一般的な解決策としては、類義語辞書の活用が挙げられるが、類義語辞書にも登録されていないくらい特殊な語彙を吸収することは技術的には難しい。

### （特性5）教師文献における付与文献数が少ない分類

教師文献における付与文献数（正例）が少ないと、その分類の特徴を学習できない。

解決策としては、教師文献を増やしてその分類が付与されている文献数を増やすことが挙げられる。また、正例と負例のバランスを調整することも必要である。

### （特性6）分類の粒度が細かい分類

FI・Fタームにおける分類階層が低い（ドット数が多い）分類になるほど、分類定義の粒度が細くなるため、分類付与が難しくなる。

解決策としては、分類階層に沿った、段階的・局所的

な分類付与が有効である可能性がある。すなわち、まず、分類階層が最も高い分類のみを対象として付与すべき分類を特定し、次に、特定された分類の下位分類のみを対象として、付与すべき分類を特定することが考えられる。

#### (特性 7) 「その他」に相当する分類

「その他」に相当する分類が付与される文献は、他の分類のどれにも該当しない文献の寄せ集めであり、その特徴を学習することは難しい。

解決策としては、「その他」と兄弟関係にある分類のどれも付与されなかった文献に付与する、といった付与方式が考えられる。

#### (特性 8) 発明の傾向が頻繁に変化する分類

発明の傾向が時系列的に頻繁に変化する分類の場合、教師文献の収集が追い付かないため、付与精度が向上しないと考えられる。

## 3 おわりに

本稿では、昨年度実施した機械学習等を利用した分類付与についての調査研究について、その一端を紹介させて頂いた。

検索インデックスは、特許審査の品質の維持・向上を図る観点から、漏れなく均一に付与されていることが望ましいが、現時点において、機械学習等を利用した分類付与の付与精度は、人手付与の付与精度と同等のレベルにあるとは言いがたく、引き続き、付与精度を向上するための付与技術についての調査研究を進めることが必要である。また、上記エラー分析の結果等を踏まえて、機械学習等を利用した分類付与の付与精度を考慮した分類体系の構築についても、検討を進めていくことが望ましい。

また、現時点において、機械学習等を利用した分類付与の付与精度は、人手付与の付与精度と同等のレベルにあるとは言いがたいが、統一した付与基準によって行われるため、均質な付与を行うことができるという利点を有し、その利用方法を工夫することにより活用できる可能性を内在するものと考えられる。そこで、付与精度が一定レベルであるという特長を生かした利用方法についても、併せて検討を進めていくことが望ましい。

最後に、本稿で紹介させて頂いた、機械学習等を利用した分類付与についての調査研究を実施頂いた日立製作所の皆様、及び、特許庁審査第一部調整課審査企画室次期検索システム検討 WG メンバーの皆様に対し、謝意を表します。

#### 参考文献

- [1] 小林英司、“特許分類の自動推定に向けた取り組み—機械学習による自動分類推定の課題と今後の展開—”、Japio YEAR BOOK 2015、2015年、pp.272-275
- [2] 殿川雅也、“内外国特許文献一括検索に向けて”、Japio YEAR BOOK 2016、2016年、pp.44-47

