

文章の書き手の同定における分類法の精度比較

三品 光平* 松田 眞一†

E-Mail: matsu@nanzan-u.ac.jp

金・村上 [13] は文章の書き手の同定における分類法の精度比較を行っているがその際、書き手の特徴が顕著に現れないようにするために変数には単語の相対頻度のみを用いている。本論文では変数について、書き手の同定に有効性が示されている複数の変数を追加し比較検証を行う。また、分類法については金・村上 [13] が有効性を示したランダムフォレスト法と同時期に提案された MART 法を新たに追加する。検証の結果、MART 法は正例と負例に差がある不均衡なデータの 2 値判別に対して有効であることがわかり、ランダムフォレスト法が均衡なデータでの多値判別に最も有効であることがわかった。また、変数は単語の長さ (動詞) の分布と助詞の分布においてはあまり良い結果が出なかったが、読点前の文字の分布と品詞の n-gram 分布においては高い判別精度を示した。

1 はじめに

テキストを計量的に分析する研究は 100 年以上前から行われてきており、著者不明の文章の書き手の同定に関する研究の歴史も長い。しかし、日本は欧米に比べて始まるのが遅かったこともあり研究が十分とは言い難い。

日本語は英文と違い、分かち書き (単語と単語の間をスペースで区切る) がされておらず、文字の種類も多いため人の手での計量が難しかった。しかし、コンピュータを用いた形態素解析技術の発展により、文章情報を自動で読み取り処理することができるようになった。コンピュータによる解析は完璧ではないが、高速にテキストデータを処理できるようになったことは言うまでもなく、文章に含まれている様々な情報を客観的に集計できるようになったことは大きな前進である。著者推定の研究は、著作者不明の文献の推定や文献の真贋判定だけでなく、裁判における被告人の上申書と日記の作成者の同一性の検証、さらにはインターネット上のメールや情報分野にまで拡大し、ウェブ文書の推定やスパムメールやウェブスパム (Google のページランクをあげるためのダミーページによる強リンクネットワークの構築) への対策といった実社会の様々な場所で需要がある。本論文では、文章の書き手の同定に焦点を当てて研究を行う。先行研究として金・村上 [13] があり、そこでは文章の書き手の同定における集団学習法の有効性を示している。集団学習法の中でもランダムフォレスト法 (Breiman[3] 参照) が特に高い精度を示している。本論文では金・村上 [13] で検証が行われておらず、ランダムフォレスト法と同時期に提案された MART 法 (Friedman[5, 6]) を比較する分析手法に追加し精度比較を行う。

*南山大学大学院数理情報研究科数理情報専攻

†南山大学情報理工学部情報システム数理工学科

2 分析について

2.1 用いる文章データ

金・村上 [13] で分析に用いられたデータは 10 人各 20 編の合計 200 編の小説, 11 人 10 タイトルの 110 編の作文, 6 人の 10 日間文の日記 60 編である。ただし, 小説データにおいては長い作品の場合分割したものを独立した文章として扱っている。

本論文では小説データは青空文庫から金・村上 [13] とできるだけ同様のものを使用した。

表 1: 分析に用いる小説のリスト

著者名 [†]	作品名
芥川 竜之介 (1892-1927)	或阿呆の一生, 玄鶴山房, 齒車, 芋粥, 煙管, 或日の大石内蔵助, 偷盗, 地獄変, 毛利先生, 路上, お律と子等と, 奇怪な再開, 杜子春, 將軍, 母, おぎん, 保吉の手帳から, 少年, 春, 彼
菊池 寛 (1888-1948)	芥川の事ども, 仇討禁止令, 仇討三態, 青木の出京, 勲章を貰う話, 身投げ救助業, 三浦右衛門の最後, M 侯爵と写真師, 無名作家の日記, 大島が出来る話, 恩を返す話, 恩讐の彼方に, 乱世, 船医の立場, 俊寛, 勝負事, 出世, 忠直卿行状記, 若杉裁判長, ゼラール中尉
夏目 漱石 (1867-1916)	それから, 一夜, 三四郎, 倫敦塔, 吾輩は猫である 1, 吾輩は猫である 2, 吾輩は猫である 3, 坊ちゃん, 幻影の盾, 彼岸過迄 1, 彼岸過迄 2, 琴のそら音, 草枕, 薙露行, 虞美人草 1, 虞美人草 2, 虞美人草 3, 行人 1, 行人 2, 趣味の遺伝, 門,
森 鷗外 (1862-1922)	かのように, じいさんばあさん, カズイヌチカ, 冴た・セクスアリス, 二人の友, 余興, 堺事件, 妄想, 寒山拾得, 山椒大夫, 普請中, 最後の一句, 杯, 百物語, 護持院原の敵討, 阿部の一族, 雁, 青年, 高瀬舟, 鶏
島崎 藤村 (1872-1943)	ある女の生涯, 三人, 並木, 伸び支度, 分配, 刺繍, 千曲川のスケッチ, 家-上巻, 岩石の間, 嵐, 旧主人, 春, 桃の雫, 桜の実の熟する時, 海へ, 熱海土産, 船, 芽生, 藁草履, 食堂
泉 鏡花 (1873-1939)	七宝の柱, 伯爵の釵, 化鳥, 半島の一奇抄, 国貞えがく, 売色鴨南蛮, 女客, 婦系図, 小春の狐, 怨霊借用, 木の子説法, 歌行燈, 眉かくしの霊, 絵本の春, 縁結び, 草迷宮, 葉草取, 遺稿, 高野聖, 鷲狩
岡本 綺堂 (1872-1939)	ゆず湯, 宣室志(唐), 搜神後記(六朝), 搜神記(六朝), 白猿伝・其他, 西陽雑俎(唐), お化け師匠, お文の魂, 勘平の死, 半鐘の怪, 湯屋の二階, 石燈籠, 寄席と芝居と, 影を踏まれた女, 心中浪華の春雨, 異妖編, 穴, 箕輪心中, 青蛙堂奇談, 鳥辺山心中
海野 十三 (1897-1949)	あの世から便りをする話, ある宇宙塵の秘密, 奇賊は支払う, 奇賊悲願, 宇宙の迷子, 宇宙尖兵, 宇宙戦隊, 怪星ガン, 恐ろしき通夜, 暗号の役割, 暗号音盤事件, 海底都市, 火薬船, 生きている腸, 化学者と夜店商人, 英本土上陸作戦前夜, 鍵から抜け出した女, 鞆らしくない鞆, 骸骨館, 鬼仏洞事件
佐々木 味津 (1896-1934)	なぞの八卦見, へび使い小町, 七化け役者, 京人形大尽, 干柿の鐔, 卍のいれずみ, 南蛮幽霊, 明月一夜騒動, 曲芸三人娘, 村正騒動, 毒色のくちびる, 生首の進物, 笛の秘密, 耳のない浪人, 血染めの手形, 袈裟切り大夫, 足のある幽霊, 身代わり花嫁, 達磨を好く遊女, 青眉の女
太宰 治 (1909-1948)	二十世紀旗手, 八十八夜, 愛と美について, 小さいアルバム, 老ハイデルベルヒ, 兄たち, 美少女, 地球図, 千代女, 断崖の錯覚, 男女同権, 誰, 誰も知らぬ, 服装に就いて, 玩具, 逆行, 恥, 花吹雪, 春の盗賊, 皮膚と心, 富嶽百景

現在ダウンロードできない作品もあるためそれらの代わりに同じ作者の他の作品を使用

[†]() 内は生没年を示す

し、10人各20編の合計200編の小説データを揃えた。また、作文と日記に関しては同様のデータが手に入らなかったためインターネット上のブログから5人各10編の合計50編のブログ記事を使用する。

2.2 文章のクリーニング

文章を分析する上で重要になるのが、電子化しデータの形式をそろえたり、不要なものを削除したりする文章のクリーニング作業である。本論文で用いる青空文庫では以下の様な処理を行った。

1. ルビの様な本文以外の内容を削除する。
2. コンピュータ上で正常に表記されない外字を識別可能な字に置き換え、その単語を解析ソフトの辞書に登録する。
3. 地の文以外の会話文などを削除する。

2.3 形態素解析

文章を統計的に分析するためには文章情報を読み取り集計する必要がある。しかし、膨大な量の文章データの処理を人の手で行う事は難しい。そこで、コンピュータを使った自然言語処理技術である形態素解析を用いる。形態素解析とは文を形態素つまり意味の最小の単位に分割することをいう。例えば「本を読んだ」という文は「本」、「を」、「読んだ」と分割できると思うかもしれないが、言語学では「本」、「を」、「読む」、「だ」と分割され、「読んだ」を動詞の「読む」と助動詞の「だ」に分割する。

形態素解析では文を形態素に分割すると同時に、形態素の品詞を特定するところまで行われる。本論文では形態素解析にフリーの形態素解析ソフトである MeCab を用い、データの集計には統計解析ソフト R 上で MeCab を実行し集計することができる RMeCab を用いる。(石田 [8] 参照)

2.4 変数

金・村上 [13] では分類法の精度と小サンプルにおける書き手の同定に関するアルゴリズムの適応性に焦点をあてており、用いたデータに書き手の特徴が顕著に現れないようにしている。そのため、変数にはノイズが多く含まれていると思われる単語の相対頻度を用いている。また各単語すべてを変数として用いるとデータセット内の値の多くが0となってしまう分類手法によっては正常に動かないため頻度がある値以下のものは「その他」の項目にまとめている。

本論文では単語の相対頻度だけではなく、金 [9, 10, 11, 12] で書き手の同定における有効性が示された読点前の文字の分布、単語の長さの分布、助詞の分布、品詞の n-gram 分布を使用し分類の精度検証を行う。

- 読点前の文字
読点前の文字の分布は読点の前の文字の出現頻度をその総数で割ったものである。(金 [9] 参照)
- 単語の長さ
金 [10] は単語の長さの分布は品詞別に分けることでより明確に書き手の特徴が現れることを示している。その理由として、書き手の個性が単語の長さに出にくい助詞や文章の内容に依存する名詞などがノイズになっていることが挙げられている。最も書き手の特徴が現れる品詞として動詞が挙げられている。
- 助詞
助詞の分布は各助詞の出現頻度をその総数で割ったものである。本論文では、金 [11] にならって出現頻度が一定以下のものはまとめてその他とした。
- 品詞の n-gram
n-gram とは隣接している n 個の文字の共起関係をを現すものであり、品詞の n-gram の場合は形態素解析を行い形態素ごとに分けられ品詞のタグを付けられたものを品詞に関して n-gram をとったものである。
例文：「今日 < 名詞 > は < 助詞 > 雪 < 名詞 > の < 助詞 > 降る < 動詞 > 寒い < 形容詞 > 日 < 名詞 > な < 助動詞 > ので < 助詞 > 家 < 名詞 > に < 助詞 > い < 動詞 > ます < 助動詞 > 。 < 記号 > 」において、 $N = 2$ で集計したものを表 2 に示す。

表 2: 品詞の n-gram(N=2) の例

品詞	度数	相対頻度
[形容詞 – 名詞]	1	0.077
[助詞 – 動詞]	2	0.154
[助詞 – 名詞]	2	0.154
[助動詞 – 記号]	1	0.077
[助動詞 – 助詞]	1	0.077
[動詞 – 形容詞]	1	0.077
[動詞 – 助動詞]	1	0.077
[名詞 – 助詞]	3	0.231
[名詞 – 助動詞]	1	0.077

2.5 分析手法

先行研究である金・村上 [13] では分類手法として k 最近傍法，学習ベクトル量子化法，サポートベクターマシン法，Bagging 法，Boosting 法，RandomForest 法[§]が用いられており，RandomForest 法が最もよい結果を示し，その次により結果を示したのが Bagging 法

[§]金・村上 [13] ではカタカタ表記のランダムフォレスト法を用いているが、他の方法と統一感を出すため以降はこの表記を用いる。

と Boosting 法 (AdaBoost) であった。本論文ではそれらの三つの手法に加えて、Boosting 法の一つである勾配 Boosting に学習器として CART 型樹木を用い拡張した MART 法を用いる。

3 分類手法

3.1 CART 法

CART 法は Breiman et al.[1] によって提案された樹木に基づく方法 (樹木構造接近法) であり、分割規則に従いデータを複数の群に分割するものである。CART 法は次の三つの手順に大きく分かれる。(杉本ら [14] 参照)

1. 前進過程：樹木の成長過程
規則に従いデータを 2 群 (ノード) に分割していき停止規則に達するまで分割を行う。
2. 後退過程：樹木の刈り込み過程
成長させた大きな樹木は過剰適合を起こすため、弱い枝を切り落とす。
3. 最適モデル選択過程：最適な樹木の決定
最適な樹木を決定には樹木の刈り込み過程において作られた候補となる部分樹木に、樹木の作成に使われていないテストデータを当てはめたとき、誤差率が少ない部分樹木を選択する。

3.2 Boosting 法

Boosting 法とは、逐次学習データの調整しながら複数の弱分類器を構築しそれらを組み合わせることによって精度の高い強分類器を構築する方法である。

3.2.1 AdaBoost

AdaBoost は Freund and Schapire[4] によって提案された Boosting 法の一つであり、逐次弱分類器の重み付き誤り率から求めた信頼度を更新していくことで強分類器を構築する手法である。

3.2.2 MART 法

MART 法は Friedman[5, 6] によって提案された Boosting 法の一つである勾配 Boosting に CART 型樹木を分類器に用いた手法であり、勾配 Boosting は逐次損失関数 $L(y, f(x))$ の傾きにより重みを更新していき強分類器を構築する手法である。損失関数とその傾きを表 3 に示す。本論文では、 \mathbb{R} 上で実行できる関数がなかったため、作成して使用した。

表 3: 損失関数とその傾き

種類	損失関数	傾き
回帰	$\frac{1}{2}[y - f(\mathbf{x})]^2$	$y - f(\mathbf{x})$
回帰	$ y - f(\mathbf{x}) $	$\text{sign}[y - f(\mathbf{x})]$
2分類 (2項ロジット)	$\log(1 + \exp(-2yf(\mathbf{x})))$	$\frac{2y}{(1 + \exp(2yf(\mathbf{x})))}$
多分類 (S クラス)	$-\sum_{s=1}^S y_s \log p_s(\mathbf{x})$	$y_s - p_s(\mathbf{x})$

本論文では多分類に対応した損失関数を用いた。

3.3 Bagging 法

Bagging 法とはアンサンブル学習法の一つであり, Breiman[2] によって提案された。Bagging 法はブートストラップと呼ばれる復元抽出法で複数の学習データセットを作成し, 各学習データで分類器を作成し, 多数決を取ることで精度の向上を図っている手法である。

3.4 RandomForest 法

RandomForest 法とはアンサンブル学習法の一つであり, Breiman[3] によって提案された。RandomForest 法も Bagging 法と同様に復元抽出によりサンプリングされた複数の学習データセットを作成しそれらから分類器を作成し, 多数決を取る。Bagging 法との違いは変数もサンプリング (非復元抽出) されたものを用いる点である。また, Bagging では分類器を作成する際, 樹木の刈り込み過程があったが, RandomForest 法では刈り込みを行わず最大の樹木を用いる。

4 分析結果

4.1 検証方法について

小説データ作者 10 人各 20 編とブログデータ 5 人各 10 編のデータに関して, それぞれ分類器の学習に用いる学習データと分類器の評価に用いるテストデータにサンプリングで分けることにする。標本サイズの違いによる判別精度をみるために, 著者一人あたりの標本サイズを S としたとき学習用データを $(S - 1, S - 2, \dots, 3)$ 個ずつ各著者からランダムサンプリングし, それ以外をテストデータとした。分類器内で用いられている乱数やサンプリングされる標本データの違いにより判別精度が違ってしまいますので, 評価には実験を 100 回繰り返した評価指標の平均を用いることとした。

4.2 再現率・精度

評価指標として再現率 (recall) と精度 (precision), そして, それらから求められる調和平均 F を用いる。著者 i ($i = 1, 2, \dots, n$) とその他の著者を A, B とラベルをつけたグルー

ブを G_i とし A を正しく分類したい場合を例とすると、再現率 R_i は A と判断されるべきものの内どれだけ正しく A と判断されたかを表し、精度 P_i は A と判断されたものの内どれだけ A と正しく判断されたかを表す。同定結果を表 4 とおいたとき再現率・精度・ F 値は下記の式で表される。一般的に精度は式 (2) の計算で求めるが、本論文の場合、小さい標本サイズで精度を出すことがあり、その際 a_i, b_i とともに 0 となり計算ができないことがでてくるため、精度を式 (3) で求めることとする。

$$\text{再現率} : R_i = \frac{a_i}{a_i + c_i} \quad (1)$$

$$\text{精度 1} : P_i = \frac{a_i}{a_i + b_i} \quad (2)$$

$$\text{精度 2} : P_i = \frac{a_i + d_i}{a_i + b_i + c_i + d_i} \quad (3)$$

多数の分類の場合は再現率と精度の平均を用いる。それは次の式で定義される。

$$\text{再現率} : \hat{R} = \frac{1}{n} \sum_{i=1}^n \frac{a_i}{a_i + c_i}, \quad \text{精度} : \hat{P} = \frac{1}{n} \sum_{i=1}^n \frac{a_i + d_i}{a_i + b_i + c_i + d_i} \quad (4)$$

調和平均 F は次の式で定義される。

$$F = \frac{2 \times \hat{P} \times \hat{R}}{\hat{P} + \hat{R}} \quad (5)$$

表 4: 同定結果のクロス表

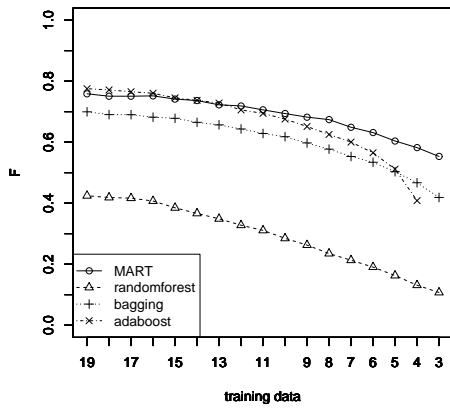
G_i		分類法の結果	
		A	B
データ	A	a	c
	B	b	d

4.3 単語の相対頻度

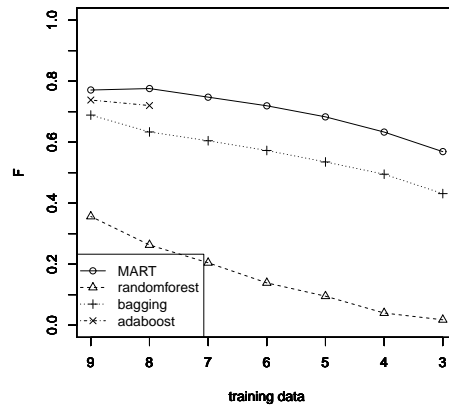
小説データでは出現頻度が 50 以上を基準[¶]とし、それより下の単語はその他の項目にまとめたところ 587 項目となり、ブログデータでは出現頻度が 10 以上を基準とし、それより下の単語はその他の項目にまとめたところ 56 項目となった。これらの変数を用いて各分類法で分析を行う。また、 F 値は各著者とその他の 2 値に分類し判別・同定し求めたものと各著者を同時に (多分類) 判別・同定し求めたものの 2 つの結果を図 1, 2 に示す。まず、図 1 の小説データとブログデータと、図 2 のブログデータでの AdaBoost の分析結果が途中から途切れていることについて述べる。

金・村上 [13] でデータセット内の値の多くが 0 となると正常に動かない分類法があることが示されていたので、内部に 0 を含まないデータセットでも試したが上手くいかず、他

[¶]基準は定量的な方法がなく、試行錯誤によって決めている。



小説

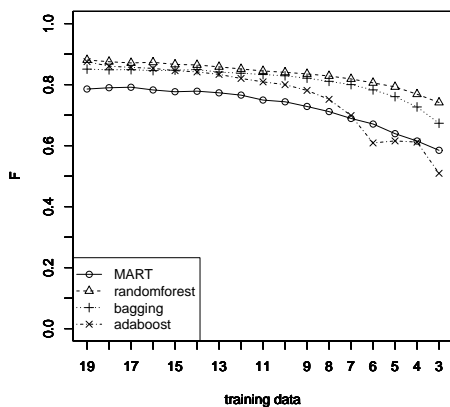


ブログ

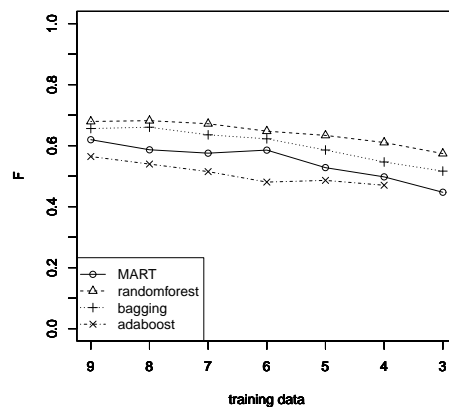
図 1: 単語の相対頻度 (2 値判別・ F 値の平均)

の変数での分析でも同様の標本サイズで正常に動かなくなることからその原因に標本サイズが関係していると思われる。ある程度、標本サイズのある状態では判別が正常に行われているのでその点で比較を行う。

2 値判別での F 値の図 1 から小説・ブログデータともに, MART, AdaBoost がほぼ同じで Bagging がそれに続き,そして RandomForest という順番になっている。特に RandomForest 法が他の分類手法よりも低くなっていることがわかる。これは RandomForest 法の欠点である正例と負例の数が異なる場合判別精度が下がるという性質が原因と考えられる。(sfchaos[7] 参照) 正例と負例の比は小説データ 1:9, ブログデータ 1:4 となる。



小説



ブログ

図 2: 単語の相対頻度 (多値判別・ F 値の平均)

多値判別での F 値の図2から小説データでは RandomForest , Bagging , AdaBoost , MART , プログデータでは RandomForest , Bagging , MART , AdaBoost , の順に良く , ともに RandomForest 法が一番良い結果となった。また , 小説データでの標本サイズの減少に伴う AdaBoost の F 値の変化をみてみると 2 値判別 , 多値判別ともに他の分類法よりも標本サイズの減少による判別精度の低下が大きいことがわかる。以上のことから正例と負例に差のある 2 値判別においては MART 法と AdaBoost が有効であるが , 標本サイズの減少に対して影響を受けにくい MART 法の方が有効だといえる。また均衡データでの多値判別では RandomForest 法が最も有効だといえる。

4.4 読点の前の文字の分布

読点の前の文字の分布において , 小説データでは出現頻度 50 以上を基準とし , それより下の単語はその他の項目にまとめたところ 24 項目となり , プログデータでは出現頻度が 3 以上を基準とし , それより下の単語はその他の項目にまとめたところ 22 項目となった。

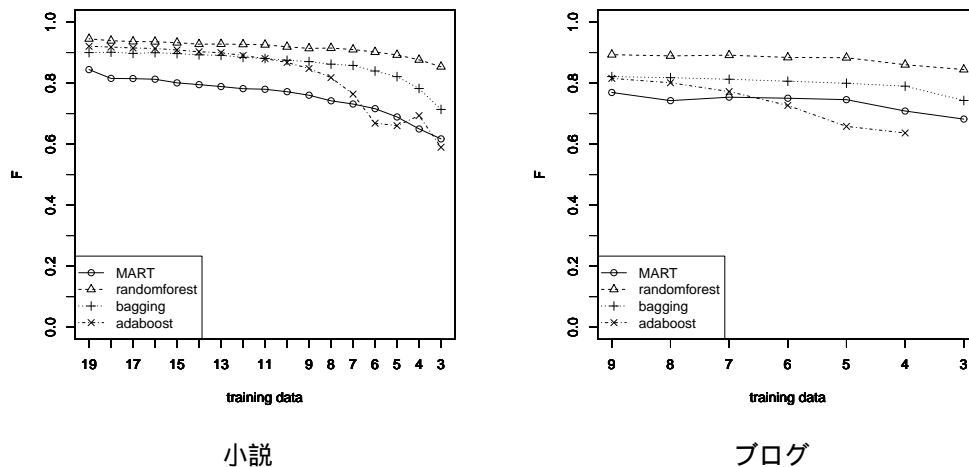


図 3: 読点前の文字の分布 (多値判別・ F 値の平均)

多値判別での分析結果を図3に示す。2 値判別では MART と AdaBoost が同じくらいの F 値となり , それらに続き Bagging , RandomForest の順となった。多値判別においては RandomForest 法が最も良く , Bagging , AdaBoost が続き , それらと少し離れて MART の順となった。AdaBoost は標本サイズが大きいときは Bagging と同等の判別精度であるが , 標本サイズの減少に伴い急激な判別精度の低下がみられる。全体的に高い判別精度となり , 特に多値判別での RandomForest 法での F 値が小説データで最大 0.944 , 各著者の標本サイズが 3 でも 0.853 と高く , プログデータでも最大 0.892 標本サイズ 3 でも 0.844 と小説データに比べ文字数の少ないプログデータでも高い判別精度を示した。

4.5 単語の長さの分布

単語の長さの分布には最も書き手の特徴が現れるとされる動詞を用いた。単語の長さの分布を集計したところ、小説データでは8項目、ログデータでは6項目となった。2値判別を行い求めた F 値で、小説データでは MART, そして, Bagging, AdaBoost と続き, それらから離れて RandomForest となり, ログデータでは MART, Bagging, RandomForest, AdaBoost となった。このことから, 単語の相対頻度と同様に MART 法が不均衡データでの2値判別に有効であるといえる。次に, 多値判別においては小説データ, ログデータともにだいたい RandomForest, Bagging, AdaBoost, MART となっており, RandomForest 法が最も良い結論となった。しかし, 全体的に判別精度が低く, 最も良い多値判別においての小説データでの F 値でもどの方法も 0.65 を超えなかった。そのため, 単語の長さの分布は本論文で用いた分類法ではあまり有効でないという結果になった。金 [10] では小説データ 3 人 21 編と少ない人数で検証しているのに対し, 本論文では小説データ 10 人 200 編, ログデータでも 5 人 50 編としており, 判別対象や標本サイズが大きかったことや, それに対する変数の項目数が少ないことも理由として考えられる。

4.6 助詞の分布

助詞の分布において, 小説データでは出現頻度 50 以上を基準とし, それより下の単語はその他の項目にまとめたところ 38 項目となり, ログデータでは出現頻度が 5 以上を基準とし, それより下の単語はその他の項目にまとめたところ 28 項目となった。2値判別では MART, Bagging, AdaBoost, RandomForest の順に良く, 多値判別においては小説データで RandomForest, Bagging, AdaBoost, MART の順に良く, ログデータでは RandomForest, MART, Bagging, AdaBoost の順に良いという結果となった。単語の相対頻度や品詞の n -gram 分布に比べて項目数が少ないが, 多値判別においての小説データでの F 値は 0.8 以上をとっており, 助詞の分布の書き手の同定における有効性はある程度あるといえる。しかし, ログデータでは最大でも 0.591 と低い値をとなり文字数の少ないデータではあまり有効ではない可能性がある。

4.7 n -gram 分布

品詞の n -gram 分布を集計したところ, 小説データでは 151 項目, ログデータでは 116 項目となった。多値判別での分析結果を図 4 に示す。2値判別では MART と AdaBoost が同じくらいの値を示し, それらに続き Bagging, RandomForest の順となった。多値判別においては, 小説データは RandomForest 法が最も良く, Bagging, AdaBoost が続き, それらと少し離れて MART の順となり, F 値も最大で 0.941, 学習用の標本を減らしていても, 各著者の標本サイズが 6 になるまで 0.9 以上, 標本サイズ 3 でも 0.845 という高い判別精度を示した。また, ログデータでは最も良かった RandomForest が次に良い Bagging と F 値で平均 0.147 も差があり, RandomForest 自体の F 値も最大で 0.916, 標本サイズ 4 まで 0.8 以上と小説データにくらべ文字数が少ないにも関わらず高い判別精度を示した。これらのことから多値判別では品詞の n -gram 分布を用いた書き手の同定には RandomForest 法が非常に有効であり, 文字数が少ないデータに対しても有効であることがわかった。

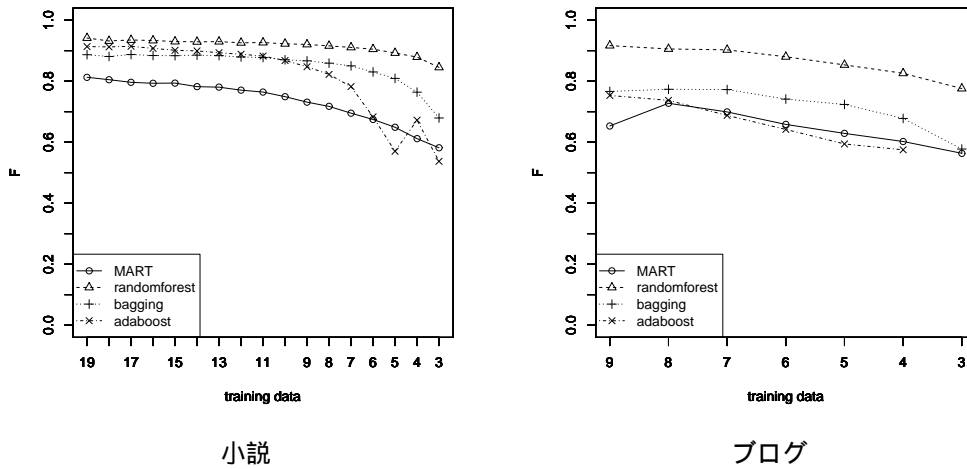


図 4: 品詞の n-gram 分布 (多値判別・F 値の平均)

5 まとめ

各変数の結果から不均衡データでの2値判別ではMART法やAdaBoostといったBoosting法が有効であるといえる。しかし、AdaBoostはMART法に比べて、変数によっては標本サイズの減少に伴い判別精度が下がる傾向が見られるため不均衡データでの2値判別においてはMART法が良いといえる。また、均衡データでの多値判別においてはRandomForest法が最も有効であるといえる。RandomForest法はどの変数においても常に最も良い結果を出した。特に読点前の文字の分布、品詞のn-gram分布におけるブログデータでの判別精度は他の分類法よりも圧倒的に良い結果となった。このことから文字数が少ない文章での書き手の同定において有効であるといえるだろう。

変数に関しては読点前の文字の分布と品詞のn-gram分布がとても高い判別精度をとり、2値判別、多値判別それぞれで良い結果となった。上記の分類手法と組み合わせることで、書き手の同定における有効な手段になるといえるだろう。

6 おわりに

本論文では、正例と負例に差がある場合でもMART法は高い判別精度を持つことと、均衡データでの多値判別におけるRandomForest法の判別精度の高さを示すことができた。しかし、AdaBoostが正常に動かない問題を解決できずに終わってしまった。また、テキストのクリーニング作業は文献やインターネット上の情報を基に独学で行っているため完璧とはいえない。より、正確にテキストクリーニングを行うことで今回あまり良い結果とならなかった変数の判別結果も変わってくる可能性がある。

参考文献

- [1] Breiman, L., Friedman, J. H., Olshen, R. A. & Stone, C. j.(1984): *Classification And Regression Trees*, Wadsworth.
- [2] Breiman, L.(1996): Bagging predictors , *Machine Learning* , **26**(2) , 123-140.
- [3] Breiman, L.(2001): Random Forests , *Machine Learning* , **45**(1) , 5-32.
- [4] Freund, Y. and Schapire, R. E.(1996): Experiments with a new boosting algorithm , *Machine Learning , Proceedings of the Thirteen International Conference* , 148-156.
- [5] Friedman, J. H.(2001): Greedy function approximation: a gradient boosting machine , *The Annals of Statistics* , **29**(5) , 1189-1232.
- [6] Friedman, J. H.(2002): Stochastic gradient boosting: Nonlinear methods and data mining , *Computational Statistics and Data Analysis* , **38** , 367-378.
- [7] sfchaos(2012): 不均衡データのクラス分類, www.slideshare.net/sfchaos/ss-11307051 .
- [8] 石田基広 (2008) : RMeCab の使い方, rmecab.jp/wiki/index.php?plugin=attach&refer=RMeCab&openfile=manual.pdf .
- [9] 金明哲 (1993) : 読点の情報に基づく文献の分類 ,『全国大会講演論文集 第 46 回平成 5 年前期 (3)』 , 131-132.
- [10] 金明哲 (1996) : 日本語における単語の長さの分布と文章の著者, 『社会情報』 , 5(2), 13-21.
- [11] 金明哲 (2002) : 助詞の分布における書き手の特徴に関する計量分析, 『社会情報』 , 11(2), 15-2.
- [12] 金明哲 (2004) : 品詞のマルコフ遷移の情報を用いた書き手の同定, 『日本行動計量学会大会発表論文抄録集』 , 32, 384-385.
- [13] 金明哲・村上 征勝 (2007) : ランダムフォレスト法による文章の書き手の同定 ,『統計数理』 , 55(2), 255-268.
- [14] 杉本知之・下川敏雄・後藤昌司 (2005) : 樹木構造接近法と最近の発展, 『計算機統計学』 , 18(2), 123-164.