

帰納的学習を用いた電子メールにおける Subject の自動生成手法の提案

康 錦紅[†] 荒木 健治[†] 栃内 香次[‡]

[†]北海道大学大学院工学研究科電子情報工学専攻
〒060-8628 札幌市北区北13条西8丁目

[‡]北海学院大学大学院経営学研究科
〒062-8605 札幌市豊平区旭町4-1-4

1. はじめに

最近の電子メールの普及は目覚めしく、非常に多くのメールが日々送られてくるという環境にいる人が急激に増大している。電子メールにおいて Subject は相手にまず自分が何を伝えたいのかを表現する部分なので非常に重要ではあるが、時間がなくなってくると十分に考えずに Subject を書いてしまうことも多く、そのことがメールを受け取った人に誤解を招いたり、後ほどメールを保存しておいてメールを見ながら内容を確認する際に問題になる。このような状況を考え、本稿では過去の Subject とメール本文の組より帰納的学習[1]により Subject の生成ルールを獲得し、それらのルールを用いて Subject をメール本文より自動的に生成する手法を提案する。

2. 予備実験

まず、メーリングリストから Subject が入っているメール 100 通を出現順に選択して予備実験を行った。予備実験の結果を表 1 に示す。ここで、Subject を含む文を重要文と定義し、重要文中 5 回以上出現した語をキーワードと定義する。予備実験の結果から重要文は文の位置、日時表現と重要文中で高い確率で使用された語との関係が深いことが分かったので、重要文の決定に使用される要素を表 2 に示すように決定した。

表 1 重要文決定の予備実験結果

段落	第一段落	第二段落	第三段落	他の段落
頻度	8	84	6	2
段落中文の位置	一番目	二番目	三番目	その他
頻度	94	6	0	0
日時表現	ある	なし		
頻度	34	66		
キーワード	ある	なし		
頻度	63	37		

3. 提案手法

本手法の流れを図 1 に示す。

本手法ではあるユーザの過去のメール本文と Subject の組より学習を行うことにより Subject を生成することができる。また、システムがユーザに動的に適應することができる。

A Study on Generation Method of Subject of E-mail Using Inductive Learning

[†]Jinhong Kang, Kenji Araki (Hokkaido University)

[‡]Koji Tochinnai (Hokkai-Gakuen University)

具体的なアルゴリズムとして、以下の 2 段階で Subject を生成する。

(1) Subject を含むと推測される文 (重要文) を抽出する。

(2) 重要文より品詞列パターンを用いて Subject を生成する。

アルゴリズムとしては、まず入力文に対して形態素解析ツール「茶筌」[2]を用いて形態素解析を行う。次に、予備実験の結果より、第何段落目か、段落中文の位置、日時表現の有無、重要文中で高頻度で出現する語の有無などの情報を用いて重要文を決定する。これらの要素の値は評価式を用いて算出し、評価値が最大のものを重要文と決定する。ここで利用される重みは予備実験中の各要素の出現頻度により決定したものである。このようにして決定された重要文より Subject を決定する。Subject の決定は品詞列のパターンを用いて決定する。パターンは以下のようなもので構成されている。

[名詞列] + [助動詞]

[名詞列] + [助詞]

[名詞列] + [助詞] + [名詞-サ変接続]

[名詞列] + [助詞] + [動詞-自立]

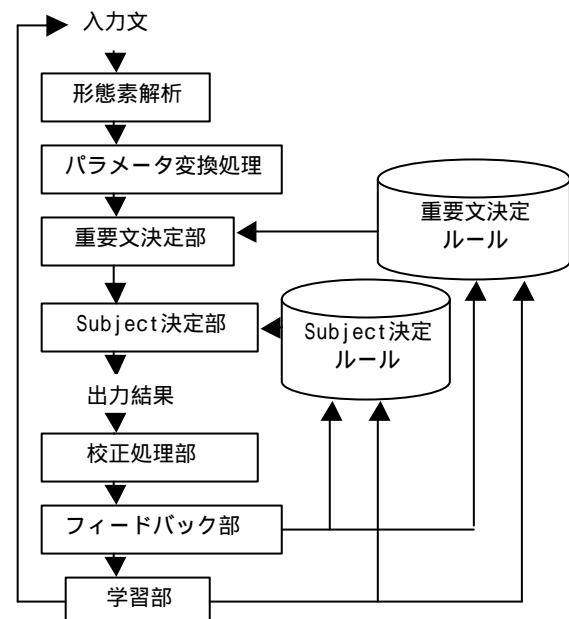


図 1

生成した Subject は校正処理部で正誤の判定後、

誤った結果はユーザにより校正され、学習部で新たなルールとして追加される。また、生成結果の正誤により、フィードバック部で尤度が変化する。尤度は適用したルール正誤の変換度数により変化する。

4. 概要

4.1 パラメータ変換処理

まず、入力文に対して形態素解析ツール茶筌[2]を用いて形態素解析を行う。次に、表2に示す要素で入力文の一文一文を4桁の数字に変換する。例えば、第一段落の一番目の文で、日時表現があってキーワードがある文のパラメータは1111のようにになっている。

表2 パラメータ決定要素

1. 段落
1: 第一段落 2: 第二段落 3: 第三段落 4: その他 他の段落
2. 段落中文の位置
1: 一番目の文 2: 二番目の文 3: 三番目の文 4: その他
3. 日時表現
1: ある 2: なし
4. キーワード
1: ある 2: なし

4.2 重要文の抽出

重要文は、変換されたパラメータを用いて決定する。重要度は式(1)により計算する。

$$\text{重要度} = \text{ルールの重み} \times \text{要素の重み} \dots (1)$$

ここで、ルールの重みとは、ルールと一致した文が重要文になった確率である。入力メールの一文ごとの重要度で重要度が一番大きい文を重要文として抽出する。

4.3 Subject 決定部

Subject はパターンマッチにより決定する。パターンは以下のようなもので構成される。

[名詞列] + [助動詞]

[名詞列] + [助詞]

[名詞列] + [助詞] + [名詞-サ変接続]

[名詞列] + [助詞] + [動詞-自立]

ここで名詞列とは、名詞または接頭詞、名詞、未知語の品詞が連続に並んだ単語列あるいは接頭詞、名詞、未知語が「と」、「の」、「や」で連結されたものである。マッチした文字列は、尤度の高さにより、一番高い尤度を用いる文字列からSubject が生成される。Subject は文字列の名詞列である。

例を以下に示す。

重要文:

[記号-空白]第[接頭詞-名詞接続]3[未知語]回[名詞-一般]国際[名詞-一般]バラ[名詞-一般]と[助詞-格助詞-一般]ガーデンニングショウ[未知語]が[助詞-格助詞-一般]5/18[未知語]([記号-括弧開]金[名詞-一般])[記号-括弧閉]~[記号-一般]23[未

知語]([記号-括弧開]水[名詞-一般])[記号-括弧閉]まで[助詞-副助詞]西武[名詞-固有名詞-組織]ドーム[名詞-一般]で[助詞-格助詞-一般]開催[名詞-サ変接続]する[動詞-自立]れる[動詞-接尾]ます[助動詞]。

マッチした文:

(1) 第[接頭詞-名詞接続]3[未知語]回[名詞-一般]国際[名詞-一般]バラ[名詞-一般]と[助詞-格助詞-一般]ガーデンニングショウ[未知語]が[助詞-格助詞-一般]

(2) 「西武[名詞-固有名詞-組織]ドーム[名詞-一般]で[助詞-格助詞-一般]開催[名詞-サ変接続]する[動詞-自立]

ここで、マッチした文が2つなので、尤度を計算し、計算結果により尤度が高い(1)からSubject が決定される。Subject は、「第3回国際バラとガーデンニングショウ」となる。

5. 評価実験と考察

5.1 実験結果

今回はメールリスト[3]のメールを使って実験を行った。実験結果を表3に示す。

表3 実験結果

メール数	正しい結果数	誤った結果数
100	52	48

5.2 考察

今回、結果の正誤の判定は実例と比較して評価した。システムが出力した結果が実例のSubjectと同じ場合正しいと評価、違う場合誤った結果として評価した。実例のSubjectの中に必ずしも正しいとはいえないものが8通あったので、実験結果に影響を与えた。また、重要文の抽出が正しくないものが19通あったので、Subjectの決定に影響を与えた。

6. おわりに

本稿では、帰納的学習により電子メールの本文からSubjectを自動的に決定する手法を提案し、本手法による評価実験を行った。今回は実験で用いたメールの数が少なかったため、今後、より多いのメールを用いて実験を行う予定である。また、重要文の決定が誤った場合Subjectの出力も必ず間違っているので、重要文の決定ルールの有効性についても検討する必要がある。

参考文献

- (1) 荒木健治, “帰納的学習を用いた自然言語処理の有効性について”, 信学技報, TL99-41, pp.33-40, 2000.
- (2) 松本裕治, 北内啓, 山下達雄, 平野善隆, “日本語形態素解析システム『茶筌』version 2.0 使用説明書”, NAIST Technical Report, NAIST-IS-TR99008, April 1999
- (3) <http://www.egroups.co.jp/messages/eco-school>