

人手作成パターンとコーパスにおける統計情報の融合による固有表現抽出システム

宮本昌幸 松尾衛 森辰則

横浜国立大学工学部電子情報工学科, E-mail: {laforge,mamoru,mori}@forest.dnj.ynu.ac.jp

1 はじめに

膨大な文書情報が氾濫している昨今、文書から特定の情報を抜き出す情報抽出が注目されている [MUC96, MUC98, IREX99]。文書中の組織名、人名、地名、商品名のような固有名や、時間表現、数値表現、割合表現などを同定する課題である、固有表現抽出タスク (Named Entity タスク, NE タスク) はその基礎技術として重要である。昨年開催された IREX においては、NE タスクに参加したシステムの約半数がパターン駆動型であり、好成績をおさめているものが多い。しかし、パターン駆動型システムの問題点はパタンの大部分が人手で書かれており、作成や調整のための労力が大きいことである。

そこで本稿では、この労力を軽減するために、人手や機械的に作成した各々のパターンに対し、正解タグ付きコーパスから得られた固有表現に関する統計情報を付与することにより、パターンが同定する固有表現の推定を行なう枠組を提案する。さらに、この統計情報を用いることにより、パタンの性能に関する評価やそれにもとづく選別を行なえることを示す。

2 パターンと統計量

パターン駆動型 NE システムにおいて、パタンのなす役割とは、形態素解析等により得られた形態素列 (ならびに品詞などの付加情報) において、ある特別な部分列を発見することである。そして同時に、その部分形態素列の全部又は一部が、ある固有表現であるという情報を付加する役目もなす。これは、図1に示すとおり、該当する各形態素の先頭と末尾に対して、ある固有表現の開始 (例えば組織名の開始 ORG-ST)、中間 (例えば組織名の中間 ORG-CN)、終了 (例えば組織名の終了 ORG-ED) というクラスを付与することである。図1は三形態素により一つの組織名が構成されている場合に相当する。このクラスを以下では「NE クラス」あるいは単に「クラス」と呼ぶ。

ORG-ST ORG-CN ORG-CN ORG-CN ORG-ED

図1: NE クラスと形態素の関係

図2に示す人名ならびに組織名を同定するパターンを例にとってみよう。このパターンでは左辺が特定の形態素列を発見するための条件部であり、右辺が照合した形態素の先頭と末尾に付与される NE クラスを表す¹。

ここで、図2の右辺に見られるような NE クラス情報付与について注目する。通常、この情報はパターンに

¹このパターンは説明のためのものであり、実際の記述とは異なる。

0 1 2 3 4 5 6
 (*) (人名+) (・) (名詞+) '社長' (*)
 → 1=PER-ST, 2=PER-ED, 3=ORG-ST, 4=ORG-ED

図2: NE 用パタンの例

より記述された特定の形態素列に対して、人間が割当を行ない、固有表現同定に利用する。これは、人間が NE の正解を与えていることになる。

しかし、もしも、固有表現に関するタグを付与したコーパスがあるならば、それに対して各パターンを照合させた時の結果を蓄積することにより、NE クラスに関する情報が収集できるはずである。さらに、この方法では、パターンに対して人間が NE クラスの情報を付与した箇所以外 (図2においては、区切り 0,5,6) についても情報を収集できるので、その箇所にも NE が有意に出現した場合にはそれを自動的に獲得できる可能性がある。

すなわち、パタンの持つ機能のうち「特定の形態素列を発見する機能」のみを用い、そのパターンに現れる全ての部分パターンについてのコーパスにより NE に関する統計量を求めれば、人間が記述したパターンを包摂する情報が獲得されると期待される。

3 パターンに付与された統計情報に基づく NE 抽出

前節で述べた考え方に従うと、未知の文書に対する NE 同定は以下のステップからなる。

1. NE タグ付コーパスからの統計量の収集

- NE のためのパターンを用意する。NE のクラスに関する情報はあってもなくてもよい。
- 正解タグ付コーパスよりパターンに照合する形態素に現れる NE クラスの頻度を調べる。
- NE クラスの頻度情報から、そのパターンが適用された時の NE の生起確率を求める。

2. 未知の文書における NE の同定

- 未知の文書に対して、各々のパタンの適用を試みる。パターンが適用された場合、そのパターンが持っている NE クラスの生起確率を対応する形態素に割り当てる。
- 付与された NE クラスの生起確率に基づき、スムージングを行ないつつ、NE クラスの接続可能性を満足しつつ文全体の生起確率が最も高くなるような NE クラスの選択をし、NE クラスの推定値とする。

3.1 正解タグ付きコーパスによるボタンへの統計情報の付与

ステップ2においては、ボタンの各部分がどのような確率である特定のNEクラスになるかを推定する。これには、正解タグ付コーパスに対して各ボタンの照合を繰り返すことにより、実際に照合したNEクラスの頻度情報を蓄積する。具体的には、各ボタン中の各々の部分ボタンについて、START,END,CENTER-ST,CENTER-EDという名前の頻度リスト(NE頻度リスト)を持たせ、どのようなNEのクラスが何回現われたかという頻度情報を保存する。NE頻度リストSTARTにはその部分ボタンの先頭に現れるNEのクラスと頻度を、ENDには末尾に現れるNEのクラスと頻度を保存する。ボタンにおいて繰り返し表現が現れる場合には、その繰り返しの途中の形態素に付与されるNEのクラスについても集計する。これにはNE頻度リストCENTER-STとCENTER-EDが用いられ、それぞれ、繰り返しの内側に現れる形態素の先頭、末尾に対応する。

例として、図2のボタンが「PER-ST「山田」PER-CN「太郎」PER-ED「。」ORG-ST「日本」ORG-CN「電子」ORG-CN「機器」ORG-ED「社長」】というNE情報付きの文字列に照合した場合を考える²。区切番号3,4の間にある繰り返しボタンは3形態素(「日本」、「電子」、「機器」)に照合するが、この時にはその繰り返しボタンが以下のように展開されたと考え、文書中のNEのクラスの頻度情報に対応するリストに保存する。

START CENTER-ED	CENTER-ST CENTER-ED	CENTER-ST END
-----------------	---------------------	---------------

この1回の照合の結果、このボタンに付随するNE頻度リストはそれぞれ次のようになる。ただし、各要素は「NEクラス * 頻度」の形式である。

START: (ORG-ST*1), END: (ORG-ED*1), CENTER-ST: (ORG-CN*2), CENTER-ED: (ORG-CN*2)

3.2 ボタンにおけるNEクラスの生起確率の推定

あるボタンが文書に照合した時にNEクラスが生起する確率は、各NE頻度リストにおけるNEクラスの相対頻度で推定される。例えば、ある部分ボタンの持つNE頻度リストが次のようになっていたとしよう³。

START: (ORG-ST*1 NON*1), END: (ORG-ED*1 PER-ED*1), CENTER-ST: (ORG-CN*2 PER-ST*1), CENTER-ED: (ORG-CN*2 NON*1)

これより、各部分の生起確率は以下のように推定される。

START: (ORG-ST=0.50 NON=0.50), END: (ORG-ED=0.50 PER-ED=0.50), CENTER-ST: (ORG-CN=0.66 PER-ST=0.33), CENTER-ED: (ORG-CN=0.66 PER-ST=0.33)

しかし、このようなナイーブな方法では統計量をとるためのコーパスの規模によるデータスパースネスの問題が生じる。つまり、正解タグコーパスにおいて照応できる事例が少なかったボタンについては、ほとんど

²ここでは、見やすさのため、品詞情報を省略している。

³これは、前節の例において注目していた部分ボタンが、さらに、

NON NON	PER-ST PER-ED
---------	---------------

という形態素列に1回照合した状況である。

のNEクラスで推定確率が0になってしまう。ここでは、以下のm推定(m-estimate)を用いて、平滑化を行なう[Mit97]。ただし、nはそのボタンが適用された回数、n(c)はその中のNEクラスcの頻度、p(c)はNEクラスcの事前確率(正解付コーパス全体での出現確率)、Mは定数⁴である。これらを各クラスについて求めている。以下では、M=m*n, m=0.2としている。

$$\frac{n(c) + M * p(c)}{n + M} \quad (1)$$

3.3 未知の文書に対するNEクラスの生起確率の推定

正解タグ付コーパスからNEクラスの生起確率に関する情報を付与されたボタンがあれば、未知の文書に対しても各形態素におけるNEクラスの生起確率を推定可能である。すなわち、文書全体に対して全てのボタンを順次照合可能性を調べる。照合した場合には、対応する形態素に対し、ボタンの持つNEクラスの生起確率情報をそのまま付与する。このとき文書のある部分に複数のボタンが照合することがある。この場合、厳密には条件付確率の重ね合わせを計算すべきであるが、ここでは簡単のために各ボタンのNEクラスの推定性能が等しいと仮定し、各ボタンのもつ生起確率の平均値を採用している。

3.4 NEクラスの生起確率に基づく固有表現の抽出

前記の方法により各形態素にNEクラスの生起確率を求められるが、最も生起確率が高いNEクラスを採用すると、その前後にどのようなNEクラスが採用されたかとは無関係にクラスが推定されてしまう。このため、各形態素に付与された生起確率とNEクラスの接続可能性の両者を考慮しつつ、一文全体としてNEクラスの生起確率が最大となる候補を選択する。探索方法としてはビタビアルゴリズムを用いた。

4 人手で付与したNEクラス情報に基づくNE抽出との比較

本稿で提案する枠組によるNE抽出手法の基本的な特性を調べるために、人間がボタンに付与したNEクラスによるNE抽出と比較する実験を行なった。使用するボタンはIREX NEタスクに参加するために作成したものをを用いた。従来のボタンに基づくNE抽出では、これをそのまま用いる。本手法の枠組においては、同じボタンについて、付与されているNEクラスの情報を使わず、正解付コーパスから得たNEクラスの生起確率を用いる。なお、形態素解析器として両システムともJUMAN3.6を用いた。正解タグ付コーパスとしては郵政省通信総合研究所により作成、公開されたCRL固有表現データ(毎日新聞1174記事)を用いた。また、評価対象としては、IREX本大会のNE

⁴等価サンプル数(equivalent sample size)と呼ばれる。

タスクで使用された文書の自由トピック記事(72記事)により行なった。結果を表1,2に示す。

表 1: 従来のボタン型システムによる抽出結果

	再現率	適合率	F 値
ORGANIZATION	44.99	76.75	56.73
PERSON	49.86	56.91	53.15
LOCATION	61.06	72.36	66.23
ARTIFACT	4.08	15.38	6.45
DATE	88.04	83.79	85.86
TIME	96.55	76.71	85.49
MONEY	86.67	86.67	86.67
PERCENT	80.95	100	89.47
ALL	59.34	72.19	65.14

表 2: 統計情報に基づくシステムによる抽出結果

	再現率	適合率	F 値
ORGANIZATION	41.39	74.19	53.14
PERSON	52.68	57.72	55.08
LOCATION	65.38	71.77	68.43
ARTIFACT	4.08	33.33	7.27
DATE	87.68	78.57	82.88
TIME	87.93	96.23	91.89
MONEY	86.67	86.67	86.67
PERCENT	80.95	100	89.47
ALL	59.85	71.65	65.22

これらの表に示される通り、ボタンのもつ特定の形態素を同定する機能だけに注目し、人間の付与したNEクラス情報を用いずとも、同等以上の精度で抽出できることがわかる。個別の固有表現クラスについて見ると、統計を用いたことによりPERSON, LOCATION, ARTIFACT, TIMEの表現について改善が見られた。一方、ORGANIZATION, DATEについては精度が落ちている。

5 機械的に作成したボタンでの評価

前節ではボタンに対してNEクラスの確率情報を付与することにより、人間がNEクラスを付与した場合と同程度の精度でNEを抽出できることを示した。そのことはまた、正解コーパスから機械的に半網羅的に作りだしたボタンであっても、統計情報を用いることによりある程度の精度でNEを抽出できるものと考えられる。ここでは、人名抽出の場合に注目して、その評価を行なう。

5.1 評価に用いる人名抽出ボタン

ボタンを機械的に作成するにあたって、抽出の戦略として以下のものを採用した。

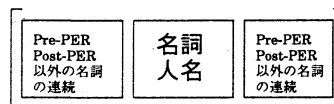
- 人名を構成し得る形態素として、名詞、片仮名列と‘.’, 平仮名列のみを考える。これらの集合をMid-PERとする。
- 人名に接続し、なおかつ、人名にならない名詞を集める。そのうち、前方に現れる表現の集合を

Pre-PER、後方に現れる表現の集合をPost-PERと表す。

- 人名の一部となる形態素を発見する。
- 人名の一部がわかっている時その前(後)の形態素について、Mid-PERでありなおかつ、Pre-PER(Post-PER)に属していないものであれば、そこも人名と認定する。

なおPre-PERならびにPost-PERについては、CRL固有表現データから集めた。

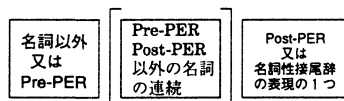
上記戦略に基づき人名抽出ボタンを作成する。一つのボタンを図3に示す。これは、形態素解析の結果、人名と判定された箇所から人名候補を探索するものである。次に考えるべきは、人名であることが形態素解



[]の部分をPERSONとして抽出する

図 3: 人名抽出ボタン 1

析の結果からはわからなかった場合に対処するボタンのグループである。これは、図4に示すように、人名の直後(直前)に現れ、その直前の表現が人名であると特定できるような表現に基づくボタンである⁵。人名の直後(直前)に来る表現としては、Post-PER(Pre-PER)の他にも表現がある。ここでは、その多数を占める名詞性接尾辞(接頭辞相当のもの)を追加して考える。実際のボタンはその表現を埋め込むことにより生成されるので、多数のボタンのグループからなる。また、表層表現をそのまま展開しただけでは正解付コーパスに特化し過ぎているものも数多くみうけられる。そこで、人名の直後の表層表現をその品詞などに置き換えて一般化したボタンもも含めている。これらのパタ



[]の部分をPERSONとして抽出する

図 4: 人名抽出ボタングループ 2

ンを用いて、節4と同じ評価実験を行なった。結果を表3,4に示す。

これらの表をみると、まず、従来のボタン型システムでは、網羅的に書かれたボタンによるため、再現率は高いものの適合率が非常に低くなっていることがわかる。一方、統計情報に基づくシステムでは、パタ

⁵図4は、人名の直後の表現に基づくボタンのグループである。これと同様に人名の直前の表現に基づくボタンのグループもある。

ンごとに付与された NE クラスの統計情報が有効に働き、適合率が大幅に向上している。最終的に F 値では約 10 ポイントの改善が見られた。

表 3: 従来のボタン型システムによる抽出結果

再現率	適合率	F 値
71.27	45.67	55.67

表 4: 統計情報に基づくシステムによる抽出結果

再現率	適合率	F 値
61.41	69.21	65.08

6 ボタンの選別

前節では、統計情報に基づくシステムを用いると、例えば、網羅的に展開されたボタンを用いても適合率の低下を招くことなく、十分な精度が保つことができることを示した。しかし、不要なボタンを数多く残しておく、ボタンの照合に非常に時間がかかり有用なシステムとならない。例えば、図 3 に示すボタンが 1 つのボタンで記述できるのとは対象的に、図 4 のボタングループは表層表現の組合せにより非常に多数のボタンのグループとなる。さらに、人名の直後の表層表現をその品詞などに置き換えて一般化したボタンもあるので、これらも含めるとその数は増大するので問題となる。

そこで、ボタンに付与されている統計情報を用い、確率的に有用でないボタンを排除することを考える。その尺度としては、エントロピーを用いる。この選別ができれば、ボタン数を減少させることが可能であり、統計情報に基づくシステムでは処理時間の短縮が見込まれる。更に、選別されたルールを使えば、ボタン型でも十分な精度が得られることが期待される。

6.1 ボタンのエントロピー

各ボタンの各々の部分ボタンに対して、正解付コーパスから獲得された NE クラスの生起確率が付与されているとすれば、頻度リスト START, END, CENTER-ST, CENTER-ED ごとにエントロピー次式で求めることができる。ただし、 c は NE クラスを表す。

$$-\sum_c P(c) \log P(c) \quad (2)$$

さて、1 つのボタンは複数の部分ボタンを持ち、また、一つの部分ボタンは 4 つの頻度リストを持つためボタン全体のエントロピーはそれらから求める必要がある。ここでは、それらの平均値を採用した。

ここでボタンが持つエントロピーの意味を考える。エントロピーが大きい場合には、NE クラス分布が均一化されているということを表すので、ある形態素列に対してそのボタンが照合したとしても、ある特定の NE クラスを指すことがなく、有用な情報となり得

ない。逆に、エントロピーが小さいボタンは、照合した部分について、ある特定のクラスになる可能性が高い。

6.2 エントロピーによる選別の効果

エントロピーによるボタン削減の効果を示すために、大量に生成された図 4 に示される型のボタンについて、エントロピーを計算し、ある閾値以上のボタンを削除する実験を行なった。結果を表 5 に示す。

表 5: エントロピーによる選別後の従来のボタン型システムによる抽出結果

閾値	0.10	0.15	0.20	0.25	0.30
再現率	64.23	64.51	67.89	67.89	67.89
適合率	65.71	65.43	64.96	60.86	58.50
F 値	64.96	64.97	66.39	64.18	62.85

閾値が 0.20 のとき F 値は 66.39 となり、ボタンを削減する前から約 11 ポイントの改善がみられた。

7 まとめ

本稿では人手で作成された固有表現抽出ボタンに対して、正解タグ付コーパスから取得された統計情報を付与することにより、確率に基づいて固有表現を特定する手法を提案し、その有効性を確認した。特に、正解タグ付コーパスから半網羅的に作成した人名抽出ボタンにおいて、統計情報を用いた手法で固有表現を特定することにより、ボタン単独による抽出と比べて、F 値において約 10 ポイントの改善が見られ、適合率の低下を軽減した。また、ボタンに付与された統計情報から求めたエントロピーを用い閾値をもうけて精度の低いボタンを削除することで、ボタン型システムにおいても約 11 ポイントの F 値の改善がみられた。

今後の課題としては他の学習型システムとの融合などがある。

参考文献

- [IREX99] IREX 実行委員会 (編). IREX ワークショップ予稿集. IREX 実行委員会, 1999.
- [Mit97] Tom M. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- [MUC96] MUC6, editor. *Proceedings of the sixth Message Understanding Conference (MUC-6)*. Morgan Kaufmann Publishers Inc., 1996.
- [MUC98] MUC7, editor. *Message Understanding Conference Proceedings*. <http://www.muc.saic.com/>, 1998.