

構造を持った定型表現の自動獲得と機械翻訳での利用

京都大学大学院情報学研究所 中澤 敏明

PROFILE

2010年京都大学大学院情報学研究所知能情報学専攻博士課程修了。博士（情報学）。機械翻訳の研究に従事。

✉ nakazawa@nlp.kuee.kyoto-u.ac.jp

☎ 075-753-5346

京都大学大学院情報学研究所教授 黒橋 禎夫

PROFILE

1994年京都大学大学院工学研究科電気工学第二専攻博士課程修了。博士（工学）。2006年4月より京都大学大学院情報学研究所教授。自然言語処理、知識情報処理の研究に従事。

1 はじめに

現在の機械翻訳では単語よりも大きなフレーズを単位として翻訳を行うことが一般的である。しかし、ここでのフレーズは対訳コーパスから自動的に推定された単語対応をヒューリスティックなルールを用いて拡張することで得られたものであり、言語的に意味のある句とは必ずしも一致しない。そのため、複数語で一つの意味を持つ定型表現を扱う際に問題が起こる。定型表現の多くは“方が良い”や“in order to”などのように一つ以上の機能語を伴っているが、これらの機能語は単体では意味を成さないか、他の語と結びついて異なる意味を表すので、対訳文中で出現した場合、相手言語文に一对一で対応する単語がないことがほとんどである。そのため、対応の推定を誤ってしまい、正しい翻訳知識が獲得されず、翻訳誤りを引き起こしてしまう。

対応を誤ってしまう原因は各単語を個別に扱っていることであり、これを解決するためには意味を持つ単位として正しい表現を考慮する方法が考えられ、既にいくつかの手法が提案されている [3, 1]。しかしこれらの研究では単語列上で連続した表現しか考慮しておらず、中国語の“在～中”（“～において”）のような単語列上は連続しない定型表現を扱うことができない。そこで本研究では依存構造木から定型表現を自動的に獲得し、

知識として機械翻訳で利用する手法を提案する。定型表現の獲得はその表現の出現頻度や周辺語の異なり数に基づいたスコアを用いて行う。依存構造木を用いることで単語列では不連続であっても、直接の依存関係が存在していれば定型表現として獲得することができる。また、自動獲得された定型表現を対訳文内の単語・句アライメントで利用することにより、アライメント精度の向上を目指す。

2 依存構造木からの定型表現の獲得

本手法では任意の表現に対してスコア付けを行い、その値が閾値以上の表現を定型表現として獲得する。定型表現はまとまりで頻繁に出現し、接続する単語の種類が多いと考えられる。この特徴を取り入れた指標として C-value [2] を依存構造木に拡張したものをを用いる。

2.1 C-value

C-value とは単語列を対象としたコーパス中のコロケーションを判定するためのスコアであり、以下の式で定義される。

$$C\text{-value}(a) = \begin{cases} \log_2(|a|) \cdot f(a) & (T_a = f) \\ \log_2(|a|) \cdot \left(f(a) - \frac{\sum_{b \in T_a} f(b)}{\#(T_a)} \right) & (\text{otherwise}) \end{cases}$$

a は対象とする表現、 $f(a)$ と $|a|$ はそれぞれ a の頻度と a を構成する単語数である。 T_a は a を内部に含むより大きな表現の集合であり、その異なり数を $\#(T_a)$ とする。式の形から頻度 $f(a)$ が高い表現であっても周辺の単語の種類 $\#(T_a)$ が少なければ、値が小さくなることがわかる。

C-value では大きさに関係なく a を含む全ての表現を T_a として扱っている。例えば、“in spite” の C-value は “in spite of”、“increased in spite”、“in spite of the” などを T_a として計算する。しかし、本研究では【図 1】のように a よりも一単語大きい表現だけを T_a とし、文頭側と文末側で別々に C-value の計算を行い、両方が閾値以上の表現を獲得する。なお、図の A,B,C,X,Y は単語を表す。つまり、“in spite” の計算を行う際は “in spite of” と “increased in spite” のように文頭又は文末側に一単語が接続した表現だけを T_a とする。“in spite” はほとんどの場合 “in spite of” に含まれて出現するため、文末側の値が低くなり、定型表現として獲得されず、“in spite of” のみが獲得される。

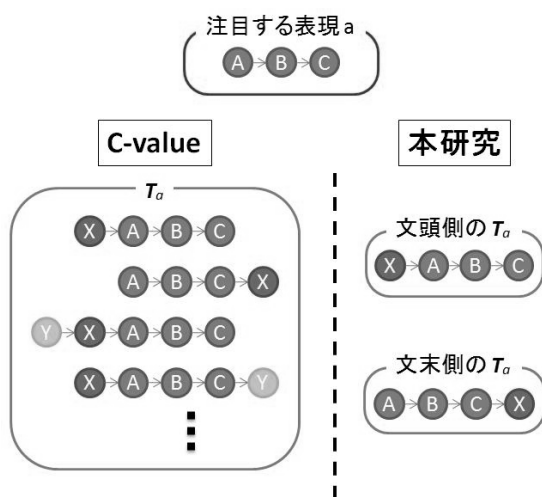


図 1：単語列の T_a

2.2 依存構造木からの獲得

依存構造木では文頭側と文末側の代わりに root 側と leaf 側で T_a を区別し、計算を行う。root 側に関しては係り先は一つしかないので単語列と同様に計算できるが、leaf 側については単語列の場合と異なるため、 T_a

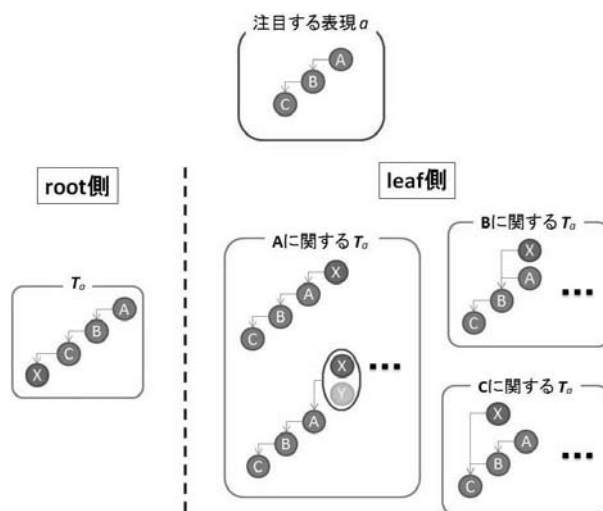


図 2：依存構造木の T_a

をさらに詳細に区別する必要がある。そこで本研究では以下の変更を加え、【図 2】のように T_a を区別した。なお、今後は依存構造木上の単語はノードと呼ぶ。

単語列の場合、単語は常に注目する表現の端のノードに接続していた。しかし依存構造木では全てのノードに単語が接続する可能性があり、同じ単語でも異なるノードに接続することが考えられる。このような異なりを全て T_a の一つとしてしまうと、あるノードに接続しやすい単語があっても他のノードに接続する単語の種類が多いため、 $\#(T_a)$ が大きくなり誤って獲得されてしまう。そこで、【図 2】に示すように a のノードごとに T_a を区別しそれぞれについて計算を行い、その最小値を leaf 側のスコアとする。これにより、あるノードに接続しやすい単語が存在すれば、そのノードの C-value は小さくなるので、定型表現としては獲得されず、全てのノードについて接続する単語の種類が多い表現のみを獲得することができる。

依存構造木では一つのノードに複数の単語が接続する場合があります。そのような表現をどの表現の T_a として扱うかが問題になる。本研究では接続している単語をまとめて扱い、【図 2】に示すように a の任意のノードに複数の単語が接続する表現も T_a として計算を行う。また、同じ単語でも接続するノードに対し単語列上で前から接続している場合と後ろから接続している場合で区別して扱っている。



本研究のように依存構造木を利用した表現獲得の研究には葛原ら [7] や Martens ら [4] の研究がある。葛原らは本研究と同じようにある表現のノードごとにスコアを計算することで獲得を行っている。しかし、英文作成を支援する表現の獲得が目的であり、節や句なども含んだ大きな表現も獲得されている。機械翻訳ではできるだけ小さい単位を扱う方が望ましいため、本研究とは獲得したい表現の大きさが異なる。また、Martens らはいくつかのスコアを用いて表現を獲得しているが、獲得した表現の具体的なアプリケーションへの応用は行っていない。

3 定型表現の獲得実験

提案手法による定型表現の獲得実験を行った。利用したコーパスは小規模論文コーパスと大規模 Web コーパスである。論文コーパスは内山・井佐原らの方法により作成した JST 日英抄録 (約 100 万文対) [6] と日中科学論文 (約 70 万文対) であり、各言語を単言語のコーパスとみなして獲得を行った。ただし、Web コーパスからの獲得を行ったのは日本語と英語のみである。次に単語の区別であるが、単語は以下の情報が全て異なるものを一種類として扱った。

日本語	代表表記、品詞、活用形
英語	原形、品詞 (“it” 以外の代名詞、三単現、複数形は区別しない)
中国語	表層語、品詞

今回は獲得の対象を 6 単語以下の表現に限定した。また、獲得されたものの中には機械翻訳で扱うには不適切な表現があったので、以下に示す簡単なルールでフィルタリングを行った。

日本語：

root が名詞または格助詞

例：“システムが”、“を解析” など

leaf が “する” または名詞性接尾辞

例：“されている”、“性について” など

英語：

be 動詞が動詞または形容詞と接続していない

例：“the system is”、“and ~ is” など

中国語：

“进进行” が動詞に接続していない

例：“对 ~ 进行” など

論文コーパスから獲得された定型表現の例をそれぞれ【表 1】に示す。論文コーパスでは “ことができる” や “in order to” など翻訳で有用な定型表現が獲得できていることが分かる。また、“在 ~ 中” や “as ~ as” など単語列上では不連続な定型表現も本手法により獲得できる。また、Web コーパスを用いた場合も “ことができる” や “due to” などの定型表現が獲得されたが、論文コーパスの結果に比べると、“という” や “at least” などの一般的な定型表現が多く獲得されていた。

表 1：獲得された定型表現の例

日本語	英語	中国語
について	this paper	在 ~ 中
では	be carried out	是 ~ 的
本稿では	based on	的方法
を提案する	in order to	一个
ことができた	this paper described	不能
について述べる	due to	就 ~ 是
示唆した	as ~ as	高的
...

4 定型表現の利用

本研究ではベースラインシステムとして中澤らの用例ベース機械翻訳システム [5] を利用し、定型表現をアライメント時の制約として用いた。

4.1 ベースラインシステム

中澤らのモデルでは依存構造木上で統計的句アライメントを行っている。依存構造木を用いることで、言語構造が大きく異なる言語対でも柔軟に対応することができ。アライメント手法を簡単に説明すると、まず既存の

統計的単語アライメントモデルにより単語レベルでの対応を推定し、これをヒューリスティックなルールで依存構造木上での句対応にマッピングする。これを初期状態とし、句対応確率と句の依存関係の確率を考慮して EM アルゴリズムにより繰り返しモデル推定を行う。EM アルゴリズムの途中により大きな句を獲得するステップがあることが特徴である。

4.2 定型表現による制限

機械翻訳では定型表現を構成する単語を個別に扱うと翻訳誤りの原因になってしまうため、定型表現はまとめて扱う必要がある。しかし、定型表現を構成する単語のアライメントが誤っているために依存構造木上で対応先が不連続になり、まとめて扱えない場合がある。その例を【図 3】に示す。ここでは黒い四角(■)が対応関係を表し、薄いマスと濃いマスの部分が正解の対応である。定型表現である“方が良い”が“方⇔me”と“良い⇔should”という対応を持っているため、“方が良い”の対応先は依存構造木上で不連続になってしまう。そこで、本研究では定型表現はまとまりで一つの意味を持つので対応先でもまとまっていると考え、定型表現を構成する単語の対応先が依存構造木上で不連続になる対応を禁止するという制約を用いた。具体的には、定型表現を構成する単語の対応先が不連続になってしまう場合、対

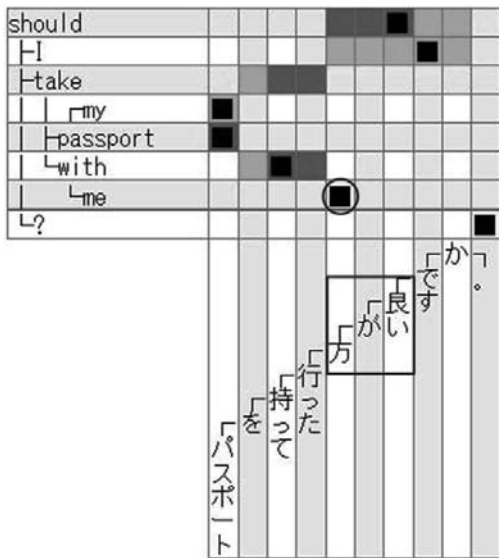


図 3：定型表現による制約の例

応確率が低い方を利用しない。こうすることで、【図 3】の“方 ⇔ me”がなくなり、“方が 良い”をまとめて扱うことができる。

また、文内のどのまとまりを定型表現とするかも問題となる。本研究では前章で獲得した定型表現のうち文内に存在するものは基本的に全てを定型表現として採用する。ただし、候補がオーバーラップした場合は長い表現を優先する。もし、同じ長さであった場合は C-value の高い方を採用する。この時の C-value は root 側と leaf 側の平均を取った値である。例えば、獲得された定型表現の C-value が【表 2】であった場合、“ことができますか”という文では“ことができます”と“ますか”の二つの定型表現が採用される。

表 2：定型表現のスコア

定型表現	スコア
ことができます	22830
ができます	22605
ますか	20714

5 アライメント実験

5.1 実験設定

定型表現を利用したアライメントを行い、精度への影響を調べた。定型表現は以下の 4 種類から獲得したものを利用し、ベースラインの結果と比較した。

- ・論文コーパスの単語列
- ・論文コーパスの依存構造木
- ・Web コーパスの単語列 (日英のみ)
- ・Web コーパスの依存構造木 (日英のみ)

実験には定型表現の獲得で利用した対訳コーパスを用いた。アライメントの評価には人手で正解を与えた日英 480 対訳文と日中 500 対訳文を利用し、以下の式で示される Precision、Recall、Alignment Error Rate(AER) を用いた。AER はアライメントの総合的な精度を示す指標であり、その値が低い程精度が良い。

$$\text{Precision} = \frac{|A \cap P|}{|A|} \quad \text{Recall} = \frac{|A \cap S|}{|S|} \quad \text{AER} = 1 - \frac{|A \cap S| + |A \cap P|}{|A| + |S|}$$

A がシステムの出力、P と S が正解である。S(Sure) は必ず必要な正解であり、【図 3】の濃いマスである。P(Possible) は英語の冠詞や日本語の助詞などのようにあっても誤りではない正解であり、【図 3】の薄いマスである。

5.2 実験結果

日英と日中の定型表現を利用したアライメントの精度をそれぞれ【表 3】と【表 4】に載せる。

結果を比較すると日英では Web コーパスの単語列と依存構造木から獲得した定型表現を用いた場合、日中に関しては単語列から獲得したものをを用いた場合の精度が最も良かった。この結果から定型表現がアライメントに

有効であることが分かる。実際に改善した例を【図 4】に載せる。この例では、定型表現“在 ~ 中”があるので、誤った“在 ⇔ は”の対応が禁止され、アライメントが改善している。

また、依存構造木を用いた場合、全体的な精度である AER は単語列より低下しているが、Precision が上昇していることがわかる。翻訳においては Precision が高い方が正確な翻訳が行えるので、依存構造木から獲得した定型表現は翻訳に有効であると考えられる。しかし、言語構造の違いなどによって定型表現の対応が正しくても対応先が不連続になってしまう場合があり、正しい対応が禁止され、Recall 低下の原因となっている。定型表現の情報をどのようにアライメントの制約として利用すべきかについては今後さらに検討する必要がある。

表 3：日英アライメント精度

定型表現	Pre.	Rec.	AER
なし	84.53	65.53	25.78
単語列 (論文)	85.10	65.18	25.80
依存構造木 (論文)	85.22	65.08	25.81
単語列 (Web)	84.90	65.59	25.63
依存構造木 (Web)	85.13	86.46	25.63

表 4：日中アライメント精度

定型表現	Pre.	Rec.	AER
なし	86.49	77.31	18.09
単語列 (論文)	86.74	77.25	18.00
依存構造木 (論文)	86.78	77.17	18.02

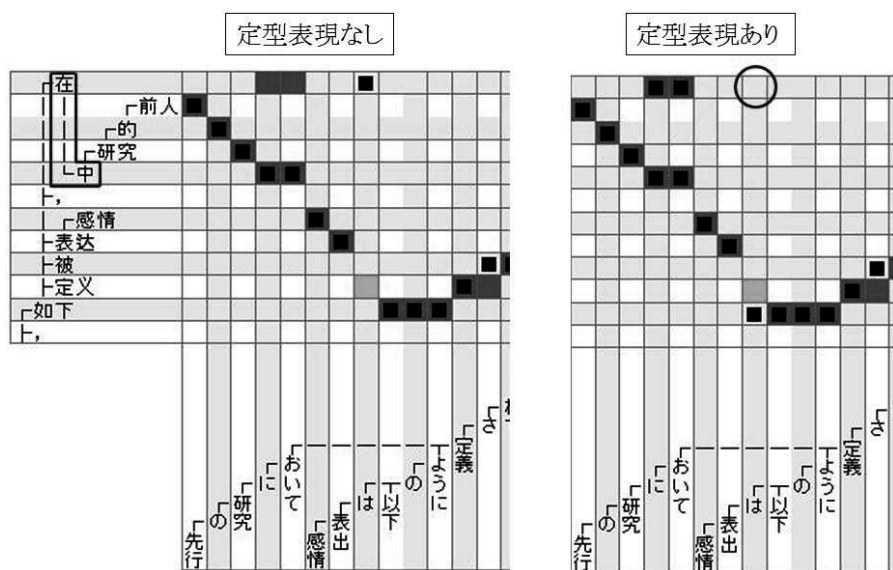


図 4：定型表現による制約の例

6 おわりに

本研究では定型表現をコーパスから自動的に獲得し、アライメント時に制約として利用する手法を提案した。定型表現はコロケーション獲得の指標である C-value を拡張することで依存構造木から獲得しており、単語列上では連続しない表現も獲得できる。また、アライメントでは定型表現を構成する単語の対応先が依存構造木上で不連続になる対応を禁止するという制約を用いた。実験は日英、日中間で行い、そのどちらでも定型表現を利用するとアライメントの精度が向上することを確認した。今後の課題は、制約により誤って禁止される問題を解決するために定型表現の利用方法を検討することと定型表現を利用した翻訳を行いその有効性を検証することである。また、定型表現の獲得の精度を向上させることも検討する必要があり、その方法には現在獲得されている定型表現を候補として対訳コーパスの情報をを用いることなどが考えられる。また今回は論文コーパスと Web コーパスを用いたが、特許コーパスからの定型表現抽出および機械翻訳での利用を行なう。

参考文献

- [1] Xiangyu Duan, Min Zhang, and Haizhou Li. Pseudoword for phrase-based machine translation. In Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, pp. 148-156, 2010.
- [2] Katerina T. Frantzi and Sophia Ananiadou. Extracting nested collocations. In Proceedings of the 16th conference on Computational linguistics, pp. 41-46, 1996.
- [3] Zhanti Liu, Haifeng Wang, Hua Wu, and Sheng Li. Improving statistical machine translation with monolingual collocation. In Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, pp. 825-833, 2010.
- [4] Scott Martens and Vincent Vandeghinste. An efficient, generic approach to extracting multiword expressions from dependency trees. In Proceedings of the Multiword Expressions: From Theory to Applications(MWE 2010), pp. 85-88, 2010.
- [5] Toshiaki Nakazawa and Sadao Kurohashi. Fully syntactic ebmt system of kyoto team in ntcir-8. In Proceedings of the 8th NTCIR Workshop Meeting on Evaluation of Information Access Technologies(NTCIR-8), pp. 403-410, 2010.
- [6] Masao Utiyama and Hitoshi Isahara. Reliable measures for aligning japanese-english news articles and sentences. In Proceedings of the 43rd Annual Meeting of the Association for Computational Linguistics, pp.72-79, 2003.
- [7] 葛原和也, 加藤芳秀, 松原茂樹. 構文構造を利用した英語論文からの表現の自動獲得. 研究報告自然言語処理 (NL), pp. 1-7, 2010.