

パターンを用いた対訳句の自動作成と翻訳精度の調査

栗下 尚樹 *1 村上 仁一 *2

*1 鳥取大学 工学部 知能情報工学科

{s132024,murakami}@ike.tottori-u.ac.jp

1 はじめに

対訳文から自動的に対訳句を作成する方法として、Barkeley などが挙げられる。これは基本的に Och のヒューリスティック [1] を利用している。しかし従来の方法で作成された対訳句 (特に日英の対訳句) の精度は低い。

そこで本研究では、新たな手法を提案する。具体的には対訳文と文パターンを利用して対訳句を作成する。また翻訳実験を行い翻訳の精度向上を試みる。

2 関連研究

対訳句の作成や精度向上のためにいくつかの先行手法が挙げられる。Och ら [1] は、アライメントテンプレートを用いる方法の中で、2つの方向の対応関係を用いてそれぞれの方向の対応関係をヒューリスティックルールを適用して拡張している。また山田ら [2] は、同じ対訳コーパスから非対称な2つの学習によって構築された2つの翻訳モデルを利用することで、それぞれの翻訳モデルの精度を向上させた。また Chris ら [3] らは、文アライメントだけでなく単語アライメントを学習時に用いることで AER を下げる研究がなされている。

3 提案手法

提案手法としてパターンから対訳句を作成する方法を提案する。

3.1 基本手法

基本的な手法は対訳文とパターンから対訳句を作成する。例として表 1 に具体例を示す。

表 1 テストデータ

対訳文 (日)	ナイロンは原油から作られる
対訳文 (英)	Nylon is made from crude oil
パターン (日)	X1 は X2 X3
パターン (英)	X1 X3 X2
パターン原文 (日)	ワインはブドウから作られる
パターン原文 (英)	Wine is made from grape

表 1 の対訳文は対訳句作成に用いる文であり、パターンは対訳句作成に用いるパターンであり、パターン原文はパターンを作る際に用いた対訳文である。(パターン原文からパターンを作成する手法は文献 [4] を参照) 対訳句作成の手順を以下に示す。

手順 1 対訳文とパターンを照合

日本語文では「ナイロンは原油から作られる」と「X1 は X2 X3」である。

手順 2 単語レベル文パターンの変数部に対応する組み合わ

せの対訳句をすべて抽出

表 2 に具体例を示す。

表 2 手順 2 と手順 3 の具体例

変数	対訳句		対数フレーズ確率
X1	ナイロン	Nylon	-0.5
X2	原油より	from crude oil	-0.1
X3	作られる	is made	-0.05

⋮

X1	ナイロン	Nylon is made from	-0.8
X2	原油より	oil	-0.2
X3	作られる	crude	-0.3

手順 3 GIZA++ を用いて対数フレーズ確率を付与
対数フレーズ確率は以下の式を用いて計算する。

$$\log_2 P\left(\frac{J_0 \dots J_{N-1}}{E_0 \dots E_{M-1}}\right) = \sum_{n=0}^{N-1} (\log_2(p(J_n|E_m)) + \log_2(p(E_m|J_n)))$$

J_n : 日本語の単語 N ; 日本語の単語数

E_m : 英語の単語 M ; 英語の単語数

$p(J_n|E_m)$: 英単語 E_m が日本単語 J_n に翻訳される確率 (GIZA++ の値)

手順 4 全ての対訳句の対を抽出

この手法では、パターンに基づいて対訳文から対訳句を可能な限り抽出していくため、大量の対訳句を出力する。表 2 のパターンでは 40 組 (変数の組み合わせ) × 3 (変数) で 120 句が抽出される。

3.2 パターン制約

基本手法から対訳句の数を減らすために、1つのパターンごとに対数フレーズ確率の総和が最大の対訳句の組み合わせを出力する。その方法をパターン制約 [5] と呼ぶ。以下に具体例を示す。

手順 5 各パターンごとに、手順 3 で付与した各組みの対数フレーズ確率の総和が最大値をとる対訳句を 1 つずつ選出

表 3 に具体例を示す。

表 3 パターン制約の具体例

変数	対訳句		対数フレーズ確率
X1	ナイロン	Nylon	-0.5
X2	原油より	from crude oil	-0.1
X3	作られる	is made	-0.05
総和			-0.65

これにより基本手法の例では 120 句が 3 句に削減される。

3.3 類似制約

パターン制約では1つのパターンに対して対訳句が1組ずつ選ばれるが、実際には1つの対訳句(日)に対し複数のパターンが照合するため複数の対訳句(英)が作成される。またパターン制約を用いても誤った対訳句が存在する。そこで対訳句を作る際、対数フレーズ確率だけでなく、対訳文とパターン原文との類似度を利用する。これを類似制約と呼ぶ。この類似制約はパターン制約を適用しさらに類似度を利用して対訳句を選択する。以下に具体的手順を示す。手順はパターン制約以降の続きである。

手順6 対訳句作成時の対訳文(対訳文)と適合したパターン作成時の対訳文(パターン原文)を抽出

表1 参照

手順7 対訳文から見たパターン原文との類似度 A を作成
類似度は対訳文とパターン原文の同一の単語の出現率を表している。以下の式を用いて計算する。

$$P(r) = \left(\frac{N}{M}\right)$$

M: 対訳文中の単語数

N: 対訳文中の単語がパターン原文中の単語と一致している単語数

表4 から類似度 A の値は 3/5。

表4 類似度 A の具体例

比較する単語	評価
ナイロン	不一致
は	一致
原油	不一致
から	一致
作られる	一致
類似度 A=一致/総語数=3/5	

手順8 パターン原文から見た対訳文との類似度 B を作成

手順7と同様に計算し、類似度 B の値は 3/5

手順9 両者を掛け合わせ、類似度とする

類似度 = 類似度 A * 類似度 B = 3/5 * 3/5 = 9/25

手順10 類似度が最大の対訳句を抽出

表5は「原油より」に対応する対訳句を表す。

表5 類似制約の例

対訳句		類似度
原油より	from crude oil	9/25
	from oil	2/21
	than crude oil	3/14

表5の例に対して類似制約を適用すると表6が出力される。

表6 類似制約の出力

対訳句		類似度
原油より	crude oil	9/25

なお実験においては日英両方向での類似度を計算する。

4 実験環境

4.1 実験データ

本実験では表7のデータを用いて行う。

表7 実験データ

日本語学習文	100,000 文
英語学習文	100,000 文
テスト文	100 文

4.2 評価方法

対訳句の評価において、提案手法の対訳句と先行手法の対訳句から、それぞれランダムに100句抽出し、人手評価を行う。また表8に評価基準を示す。

表8 対訳句の評価方法

評価	評価内容
	適切な対応をとる
	意味が欠落している、もしくは不必要な意味が付与されている
x	不適切な対応をとる

5 実験結果

5.1 対訳句の精度結果

提案手法と先行手法において、対訳句の精度の評価を行う。ただし、先行手法は Moses の phrase-table とし、アライメントは intersection とする。また、人手対訳句[6]も評価する。表9に評価結果を示す。

表9 対訳句の評価結果(100句中)

手法	総数			x
先行手法	30,887,824	5	53	42
パターン制約	94,028	52	14	34
類似制約	77,880	61	27	12
人手対訳句	335,726	77	17	6

表9より、先行手法より2つの提案手法の対訳句の精度の方が高く、提案手法が有効であるといえる。また類似制約はパターン制約に比べて”不適切な対応をとる対訳句”が減少し、”意味が欠落している、もしくは不必要な意味が付与されている対訳句”が増加した。また類似制約の精度は人手対訳句に接近しているがまだ及ばない。

以下に提案手法で作成された対訳句の具体例を示す。表10に先行手法の例を、表11にパターン制約の例を、表12に類似制約の例を示す。

表10 phrase-table の具体例

対訳句(日)	対訳句(日)	評価
途方もなく大きな恩恵 世界を独占している	benefited tremendously have a world monopoly	
煙突から煙が上って 上衣	from the chimney off my coat	
もやっと山 警察署には一種	all downhill station has an atmosphere	x x

先行手法は長い対訳句が多く、対訳句(日)に完全に対応している対訳句(英)は少ない。

表 11 パターン制約の具体例

対訳句 (日)	対訳句 (英)	評価
母からの 向けられ	from my motehr turned on	
捨てた スペースシャトル	He threw away The space	
はこの 円	He has one million	× ×

パターン制約は先行手法に比べ の数が多くなっているがまだ多い。対訳句 (日) に全く対応していない対訳句 (英) が多い。

表 12 類似制約の具体例

対訳句 (日)	対訳句 (英)	評価
全速力 よく聞く	full speed listens well	
かばんのチャック 手に包帯	bag open with a bandage	
は勇敢にもその 32 段	expose the without stoping once	× ×

類似制約は × がパターン制約と比べさらに減っている。また”意味が欠落している, もしくは不必要な意味が付与されている対訳句”が増加し, ”不適切な対応をとる対訳句”がかなり減少した。

6 翻訳精度の調査

対訳句の精度向上による翻訳精度の影響を調べるために翻訳実験を行う。翻訳実験は提案手法を適用した対訳句を用いる。翻訳方法はパターンに基づく統計機械翻訳 (PBSMT) を用いて行う。PBSMT の概要を以下に示す [4]。

- 手順 A 対訳文と GIZA++ から対訳単語を作成。
- 手順 B 対訳文と対訳単語から単語レベル文パターンを作成。
- 手順 C 対訳文と単語レベル文パターンから対訳句を作成。
- 手順 D 対訳文と対訳句から句レベル文パターンを作成。
- 手順 E 対訳句と句レベル文パターンを用いて翻訳を行う。

本研究で提案した 2 つの手法は手順 C で用いる。

6.1 翻訳実験データ

翻訳実験は表 13 のデータを用いる。

表 13 翻訳実験データ

日本語学習文	100,000 文
英語学習文	100,000 文
テスト文	200 文

ただし類似制約とパターン制約の実験において、枝刈りのパラメータが少し異なる。

6.2 翻訳の評価方法

翻訳文の評価方法において、パターン制約の翻訳文と類似制約の翻訳文から、ランダムに 200 文抽出し、人手

による対比較評価を行う。また表 14 に評価基準を示す。

表 14 翻訳の評価方法

評価	評価内容
類似制約	類似制約の方がパターン制約よりも良い
パターン制約	パターン制約の方が類似制約よりも良い
差なし	どちらもあまり差がない
同一文	類似制約とパターン制約が同一

6.3 翻訳結果

翻訳を行った結果、200 文中 76 文が未知語を含み、10 文が翻訳できなかった。よって 114 文で評価を行った。結果を表 15 に示す。

表 15 翻訳の人手評価

類似制約	パターン制約	差なし	同一
36	31	35	12

表 15 より、どちらの提案手法を用いても翻訳精度にあまり差がないことがわかった。また表 16 に類似制約を、表 17 にパターン制約 の出力例を示す。太字は正しいと判断した部分、下線部は間違っていると判断した部分である。

表 16 類似制約 の具体例 1

入力文 1	その店は、新しい経営陣の下で再開した
参照文 1	The store has reopened under new management.
パターン制約	The store resumed the under <u>championship</u> in the, new its management
パターン (日)	X03 は X02 X00 X01 の X04 で X05 した。
パターン (英)	X03 X05 the X04 championship in the X02 X00 X01 .
類似制約	The store has reopen under a new management.
パターン (日)	X03 X01 新しい X04 X02 X00 した。
パターン (英)	X01 X03 has X00 X02 new X04 .
入力文 2	この方法の是非は簡単には決められない。
参照文	We cannot tell at once whether this method is right or wrong .
パターン制約	This method whether not arranged that easily.
パターン (日)	X04 の X02 は X03 X01 決められ X00 。
パターン (英)	X04 X02 X01 arranged X00 X03 .
類似制約	This method can't decided on the easily of right or wrong .
パターン (日)	X03 X02 X00 X01 X05 X04 られない。
パターン (英)	X03 can't X04 X01 X05 X02 X00 .

表 16 において、入力文 1 の結果は類似制約の方が参照文 1 に似ているため、入力文 2 の結果は「是非」が読み取れるため、類似制約 とした。

表 17 において、入力文 3 の結果は文の形として成り立っているため、入力文 4 の結果は文と主語が合っているためパターン制約 とした。

表 17 パターン制約 の具体例

入力文 3 参照文	私の 疑惑 は 大き くな った My suspicion grew .
パターン制約	My Suspicion became louder .
パターン (日)	私 の X00 は X01 X02 。
パターン (英)	My X00 X02 X01 .
パターン原文 (日)	私 の 時計 は 少 し 進 む 。
パターン原文 (英)	My watch gains a little .
類似制約	The grown my suspicion .
パターン (日)	私 の X02 X00 X01 た 。
パターン (英)	X00 X01 my X02 .
パターン原文 (日)	私 の 服 の 趣 味 を 当 て こ す っ た 。
パターン原文 (英)	He made a dig at my taste in clothes .
入力文 4 参照文	冬 は 太陽 が 早 く 沈 む The sun sets early in winter .
パターン制約	The sun sink early winter .
パターン (日)	X03 は X01 が X02 X00 。
パターン (英)	The X01 X00 X02 X03 .
パターン原文 (日)	明 日 は 気 温 が 高 く な る で し ょ う 。
パターン原文 (英)	The temperature will be high tomorrow .
類似制約	Winter sink early in the sun.
パターン (日)	X04 は X02 X00 X03 X01 。
パターン (英)	X04 X01 X03 in X00 X02 .
パターン原文 (日)	彼 は 川 に 釣 り に 行 っ た 。
パターン原文 (英)	He went fishing in the river .

7 考察

7.1 先行手法と類似制約の比較

提案手法の類似制約と先行手法として Moses の翻訳の評価を行う。それぞれの翻訳文からランダムに 200 文抽出する。人手による対比較評価を行う際、200 文中 63 文が未知語を含み、9 文が翻訳できなかったので 128 文で評価を行った。結果を表 18 に示す。

表 18 類似制約 は類似制約の方が Moses よりも良いことを表す。

表 18 Moses の人手評価

類似制約	Moses	差なし	同一
49	14	41	4

表 18 より類似制約は Moses と比較して非常に精度が良いことがわかった。

7.2 対訳句作成の考察

表 9 より、類似制約において対訳句の精度は向上したが、不適切な対応をとる対訳句が 12 文ある。表 12 にある評価 x の「は勇敢にもその」-「expose the」の対応する際に用いられたデータを調べた。表 19 にデータを示す。

類似度の値は表 12 にある評価 の「全速力」-「full speed」よりも低い。「は勇敢にもその」-「expose the」の対応はパターン制約を適用しているためパターンにおいて最大の対数フレーズ確率値であるといえる。しかし、「スキャンダル」-「courage」や「彼」-「scandal」という誤った対応の対訳句を含んでいる。これは対数フレーズ

表 19 類似制約における対訳句の考察

対訳句 (日)	は 勇 敢 に も そ の
対訳句 (英)	expose the
対訳文 (日)	彼 は 勇 敢 に も そ の ス キ ャ ン ダ ル を 暴 露 し た 。
対訳文 (英)	He had the courage to expose the scandal .
パターン (日)	X02 X00 X03 X01 X04 た 。
パターン (英)	He X04 X01 X03 to X00 X02 .
パターン原文 (日)	シ ョ ー ト に ラ イ ナ ー を 打 っ た 。
パターン原文 (英)	He hit a liner to the shortstop .
X00	「は勇敢にもその」-「expose the」
X01	「を」-「the」
X02	「彼」-「scandal」
X03	「スキャンダル」-「courage」
X04	「暴露し」-「had」

確率を付与する際に利用する GIZA++ の推定値に問題がある。解決方法としては対数フレーズ確率の付与の方法を見直す、もしくは学習文のデータ量を増やすことであると考える。

7.3 0 型代名詞を含む文からのパターン作成の問題

入力文に対して形の似ていないパターン原文から作られたパターンを用いて翻訳を行った結果、不適切な文となっている例がある。表 17 の入力文 3 の主語は「私の疑惑」に対してパターン原文は主語がない文である。そのため入力文に対して不適切なパターンを照合したことで不適切な文が出力されている。

8 おわりに

本研究では、対訳文と文パターンを利用して対訳句を作成する手法として、1つのパターンに対し1組ずつの対訳句を出力するパターン制約と、対訳文とパターン原文の類似度を用いて対訳句を選択する類似制約を用いた手法を提案した。その結果、対訳句の精度は大幅に向上した。しかし、これらを適用した対訳句を翻訳に用いても精度はあまり向上しなかった。

参考文献

- [1] Och, F.J., Tillmann, C., and Ney, H. "Improved Alignment Models for Statistical Machine Translation." In Proc. of EMNLP/WVLC-99
- [2] 山田節夫 永田晶明 山田賢治 "再学習により翻訳モデルを用いた単語アライメントの向上"
- [3] Chris Callison-Burch, David Talbot, and Miles Osborne. "Statistical Machine Translation with Word- and Sentence-Aligned Parallel Corpora"
- [4] 江木孝史 "句に基づく対訳文パターンの自動作成と統計的手法を用いた英日パターン翻訳", 言語処理学会 第 20 回年次大会
- [5] 興相 玲架 "パターンに基づく統計翻訳において変数部の確率の総和を使った対訳句の抽出", 鳥取大学卒業論文
- [6] 鳥バンク <http://unicorn.ike.tottori-u.ac.jp/toribank/>
- [7] Moses : <http://www.statmt.org/moses/>.