

概要

翻訳システムの一手法として“パターン翻訳”がある。パターン翻訳は大量の対訳文パターンと単語辞書を用いて、翻訳文を出力する手法である。パターン翻訳は、入力文が適切な対訳文パターンに適合した場合に、翻訳精度の高い翻訳文が得られやすいという特徴がある。しかし、パターン翻訳に用いる単語辞書と対訳文パターンは人手で作成するため、開発コストが高くなる。

一方翻訳システムとして“単語に基づく統計翻訳”がある。単語に基づく統計翻訳は、学習データとして対訳文を与えるだけで翻訳ができる。このため、翻訳にかかるコストが低い。さらに、対訳文から単語辞書と単語翻訳確率を自動的に得ることが可能である。

江木らは、パターンに基づく統計翻訳を考案した。パターンに基づく統計翻訳は、統計的手法を用いて、対訳句と句レベル文パターンを自動作成して翻訳を行う。しかし、この手法は翻訳精度が低い。原因の一つとして、不適切な対訳句が翻訳時に選択されていることが挙げられる。また、翻訳時の対訳句の選択にはフレーズ確率を使用している。

本研究では、翻訳精度の向上を目的とし、翻訳時の対訳句の選択における二つの手法を提案する。一つ目は、フレーズ確率の総積を使用する手法である。二つ目は、Dice 係数と類似度の積を使用する手法である。この二つの手法とフレーズ確率を対訳句の選択に使用する翻訳(以下、従来手法)で比較実験を行った。

出力文の精度を比較評価した結果、フレーズ確率の総積を使用する手法の出力文が最も精度が良いことがわかった。しかし従来手法と大きな差がないことも分かった。

目次

第1章	はじめに	1
第2章	翻訳システム	2
2.1	概要	2
2.2	パターン翻訳 ¹	2
2.2.1	概要	2
2.2.2	英日パターン翻訳の手順	3
2.2.3	文パターン辞書	5
2.3	単語に基づく統計翻訳 ²	6
2.3.1	概要	6
2.3.2	IBM 翻訳モデル	6
2.3.3	GIZA++	11
2.3.4	言語モデル	11
2.3.5	デコーダ	11
2.4	パターンに基づく統計翻訳 ³	13
2.4.1	概要	13
2.4.2	対訳単語の作成	14
2.4.3	単語レベル文パターンの作成	15
2.4.4	対訳句の作成	16
2.4.5	句レベル文パターンの作成	18
2.4.6	翻訳文の作成	20
2.4.7	従来手法の問題点	21

¹江木孝史：“句に基づく文パターンを用いた英日翻訳”

²江木孝史：“句に基づく文パターンを用いた英日翻訳”

³村上仁一：“パターンに基づく統計機械翻訳の概要と問題点について”

第3章	提案手法	22
3.1	概要	22
3.2	提案手法1(フレーズ確率の総積(P_2)を使用する手法)	22
3.3	提案手法2(Dice係数と類似度の積(P_3))を使用する手法	23
3.3.1	Dice係数	23
3.3.2	類似度	24
第4章	実験	25
4.1	実験目的	25
4.2	実験データ	25
4.3	評価方法	25
第5章	実験結果	26
5.1	人手評価結果	26
5.2	自動評価結果	33
5.3	実験結果のまとめ	33
第6章	考察	34
6.1	フレーズ確率の総積の有効性	34
6.2	差なしの原因	34
6.3	誤り分析	34
6.4	他の対訳句確率の計算方法	35
第7章	おわりに	36

目次

2.1	英日パターン翻訳の手順	4
2.2	対訳単語作成の例	14
2.3	単語レベル文パターン作成の例	15
2.4	対訳句作成の例	16
2.5	フレーズ確率付与の例 (日英)	17
2.6	句レベル文パターン作成の例	18
2.7	文パターン確率付与の例 (日英)	19
2.8	翻訳文作成の例	20

表 目 次

2.1	英日パターン翻訳の例	5
2.2	言語モデルの例	11
2.3	一方のフレーズ確率が低い例	21
3.1	対訳学習文テストデータ	22
3.2	手順2と手順3の具体例	23
3.3	類似度を求めるときのデータの例	24
4.1	対訳学習文および翻訳実験に使用する文の例	25
4.2	実験データ	25
5.1	テスト文100文の人手評価の結果	26
5.2	従来 の例	27
5.3	提案1 の例	28
5.4	提案2 の例	29
5.5	提案1×の例	30
5.6	提案2×の例	31
5.7	全ての文で意味を読み取れない, ほぼ同一である, または同一の出力例	32
5.8	テスト文1000文の自動評価の結果	33
6.1	提案1×の残りの2文	35
7.1	付録一覧	37

第1章 はじめに

翻訳システムの一手法として“パターン翻訳”がある。パターン翻訳は大量の対訳文パターンと単語辞書を用いて、翻訳文を出力する手法である。パターン翻訳は、入力文が適切な対訳文パターンに適合した場合に、翻訳精度の高い翻訳文が得られやすいという特徴がある。しかし、パターン翻訳に用いる単語辞書と対訳文パターンは人手で作成するため、開発コストが高くなる。

一方、翻訳システムとして“単語に基づく統計翻訳”がある。単語に基づく統計翻訳は、学習データとして対訳文を与えるだけで翻訳ができる。このため、翻訳にかかるコストが低い。さらに、対訳文から単語辞書と単語翻訳確率を自動的に得ることが可能である。

江木らは、パターンに基づく統計翻訳を考案した。パターンに基づく統計翻訳は、統計的手法を用いて、対訳句と句レベル文パターンを自動作成して翻訳を行う。しかし、この手法は翻訳精度が低い。原因の一つとして、対応する原言語と目的言語が不自然な対訳句が翻訳時に選択されていることが挙げられる。そして翻訳時の対訳句の選択には対訳フレーズ確率を使用している。

本研究では、翻訳時の対訳句の選択における二つの手法を提案する。一つ目は、フレーズ確率の総積である。二つ目は、Dice 係数と類似度の積である。

以上の手法により翻訳精度向上を目指す。

本論文の構成を以下に示す。第2章で翻訳システムについて説明する。第3章で提案手法について説明する。そして、第4章では実験環境を、第5章で実験結果を示し、第6章で本研究の考察を述べる。

第2章 翻訳システム

2.1 概要

本章の2.2節および2.3節は、江木の論文 [1] を引用しており、パターンに基づく統計翻訳の節は村上の論文 [2] を参照して記述している。

翻訳システムの一手法として“パターン翻訳”がある。パターン翻訳は大量の対訳文パターンと単語辞書を用いて、翻訳文を出力する手法である。パターン翻訳は、入力文が適切な対訳文パターンに適合した場合に、翻訳精度の高い翻訳文が得られやすいという特徴がある。しかし、パターン翻訳に用いる単語辞書と対訳文パターンは人手で作成するため、開発コストが高くなる。

一方、翻訳システムとして“単語に基づく統計翻訳”がある。単語に基づく統計翻訳は、学習データとして対訳文を与えるだけで翻訳ができる。このため、翻訳にかかるコストが低い。さらに、対訳文から単語辞書と単語翻訳確率を自動的に得ることが可能である。

江木らは、単語辞書と対訳文パターンを統計的手法で自動的に作成し翻訳するパターンに基づく統計翻訳を提案した。パターンに基づく統計翻訳は、句に基づく統計翻訳の特徴である対訳文から単語辞書と単語翻訳確率を自動的に取得できる点に着目し、翻訳に用いる単語辞書と対訳文パターンを統計的手法を用いて自動的に作成する。

2.2 パターン翻訳¹

2.2.1 概要

パターン翻訳とは、機械翻訳手法の一種である。パターン翻訳は、原言語文と目的言語文の対訳文に対して、任意の単語やフレーズを変数化した“文パターン”と“単語辞書”が必要である。原言語入力文と原言語文パターンを照合し、適合する原言語文パターン

¹江木孝史：“句に基づく文パターンを用いた英日翻訳”

に対応する目的言語文パターンを得る．そして，文パターンの変数部に対応する単語やフレーズを，単語辞書を用いて翻訳し，目的言語翻訳文を出力する．

パターン翻訳は適切な文パターンが適合した場合，文全体の構造を保持した翻訳精度の高い翻訳文を得ることができる．しかし，一般的なパターン翻訳は文パターンを人手で作成するため開発に時間がかかる．また，文パターンに辞書に適合しない場合は翻訳ができないため，問題点として，入力文に対するカバー率が低い．

2.2.2 英日パターン翻訳の手順

一般的な英日パターン翻訳の手順を以下に示す．

手順 1 文パターン辞書と単語辞書を用意する．

手順 2 英語入力文と英語側文パターンを照合する．

手順 3 変数部に対応する英単語を単語辞書を用いて日本語単語に翻訳する．

手順 4 英語側文パターンに対応する日本語側文パターンの変数部を，翻訳した日本語単語に置き換える．

手順 5 手順 4 で得た日本語翻訳文を出力する．

英日パターン翻訳の手順を図 2.1 に示す．

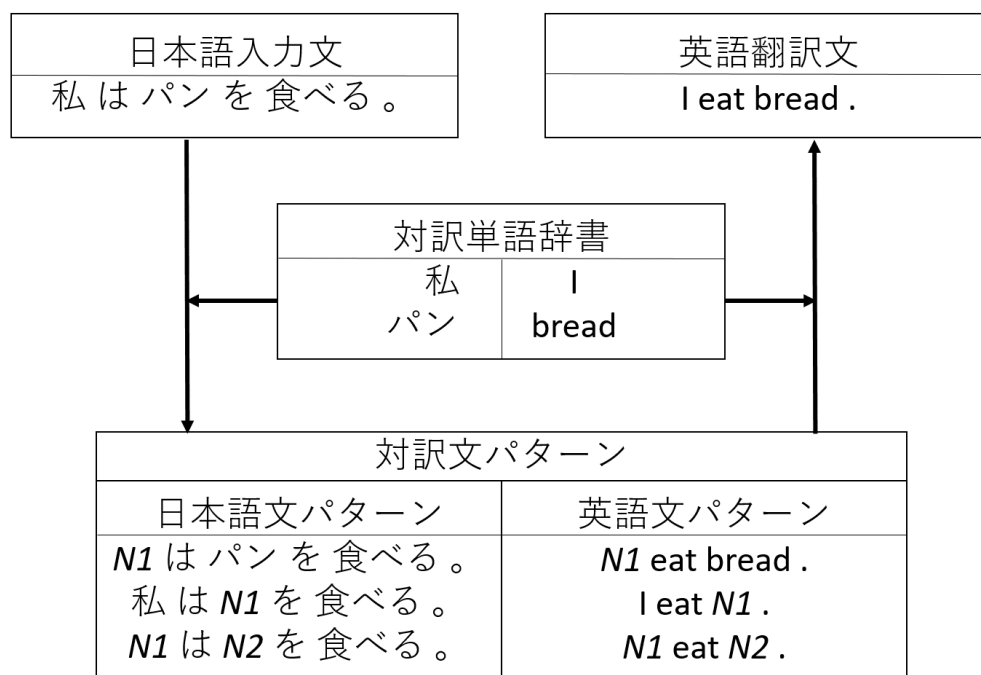


図 2.1: 英日パターン翻訳の手順

2.2.3 文パターン辞書

文パターン辞書とは，大量の対訳文から任意の単語やフレーズを変数化して得られる文パターンの集合である．表 2.1 に例を示す．

表 2.1: 英日パターン翻訳の例

英語入力文	I go to the sea .
英語文パターン	I go to X_1 .
日本語文パターン	私は X_1 に行く。
日本語翻訳文	私は海に行く。

2.3 単語に基づく統計翻訳²

2.3.1 概要

単語に基づく統計翻訳は単語対応の翻訳モデルを用いている．例として，ある英語文を日本語文に翻訳する場合を考える．英単語を日本語に翻訳し，英単語の語順と同じ並びで日本語単語を並べて翻訳する．単語に基づく統計翻訳は単語対応の確率を得る IBM 翻訳モデルが用いられている．

2.3.2 IBM 翻訳モデル

統計翻訳の代表的なモデルとして，IBM の Brown らによる仏英翻訳モデル [3] がある．IBM 翻訳モデルは，単語に基づく統計翻訳を想定して作成された，単語対応の確率モデルである．この翻訳モデルは順に複雑な計算を行うモデル 1 から 5 の 5 つのモデルで構成される．IBM 翻訳モデルでは，フランス語から英語への翻訳を想定しているため，以下の説明では仏英翻訳を前提とする．本章では，原言語であるフランス語文を F ，目的言語である英語文を E として定義する．

IBM モデルでは，フランス語文 E ，英語文 F の翻訳モデル $P(F|E)$ を計算するために，アライメント a を用いる．以下に IBM モデルの基本式を示す．

$$P(F|E) = \sum_a P(F, a|E) \quad (2.1)$$

アライメントとは仏単語と英単語の対応を意味している．IBM モデルのアライメントでは，各仏単語 f に対応する英単語 e は 1 つあり，各英単語 e に対応する仏単語は 0 から

2.3.2.1 モデル 1

(2.1) 式は以下の式に分解することができる． m はフランス語文の長さ， a_1^{j-1} はフランス語文における，1 番目から $j-1$ 番目までのアライメント， f_1^{j-1} はフランス語文における，1 番目から $j-1$ 番目まで単語を表している．

$$P(F, a|E) = P(m|E) \prod_{j=1}^m P(a_j|a_1^{j-1}, f_1^{j-1}, m, E) P(f_j|a_1^j, f_1^{j-1}, m, E) \quad (2.2)$$

(2.2) 式ではとても複雑であるので計算が困難である．そこで，モデル 1 では以下の仮定により，パラメータの簡略化を行う．

²江木孝史：“句に基づく文パターンを用いた英日翻訳”

- フランス語文の長さの確率 ϵ は m, E に依存しない

$$P(m|E) = \epsilon$$

- アライメントの確率は英語文の長さ l に依存する

$$P(a_j|a_1^{j-1}, f_1^{j-1}, m, E) = (l+1)^{-1}$$

- フランス語の翻訳確率 $t(f_j|e_{a_j})$ は、仏単語 f_j に対応する英単語 e_{a_j} に依存する

$$P(f_j|a_1^j, f_1^{j-1}, m, e) = t(f_j|e_{a_j})$$

パラメータの簡略化を行うことで、 $P(F, a|E)$ と $P(F, E)$ は以下の式で表される。

$$P(F, a|E) = \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m t(f_j|e_{a_j}) \quad (2.3)$$

$$P(F|E) = \frac{\epsilon}{(l+1)^m} \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m t(f_j|e_{a_j}) \quad (2.4)$$

$$= \frac{\epsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l t(f_j|e_{a_j}) \quad (2.5)$$

モデル1では翻訳確率 $t(f|e)$ の初期値が0以外の場合、Expectation-Maximization(EM) アルゴリズムを繰り返し行うことで得られる期待値を用いて最適解を推定する。EM アルゴリズムの手順を以下に示す。

手順1 翻訳確率 $t(f|e)$ の初期値を設定する。

手順2 仏英対訳対 $(F^{(s)}, E^{(s)})$ (但し、 $1 \leq s \leq S$) において、仏単語 f と英単語 e が対応する回数の期待値を以下の式により計算する。

$$c(f|e; F, E) = \frac{t(f|e)}{t(f|e_0) + \cdots + t(f|e_l)} \sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^l \delta(e, e_i) \quad (2.6)$$

$\delta(f, f_j)$ はフランス語文 F 中で仏単語 f が出現する回数、 $\delta(e, e_i)$ は英語文 E 中で英単語 e が出現する回数を表している。

手順3 英語文 $E^{(s)}$ の中で1回以上出現する英単語 e に対して、翻訳確率 $t(f|e)$ を計算する。

1. 定数 λ_e を以下の式により計算する .

$$\lambda_e = \sum_f \sum_{s=1}^S c(f|e; F^{(s)}, E^{(s)}) \quad (2.7)$$

2. (2.7) 式より求めた λ_e を用いて , 翻訳確率 $t(f|e)$ を再計算する .

$$\begin{aligned} t(f|e) &= \lambda_e^{-1} \sum_{s=1}^S c(f|e; F^{(s)}, E^{(s)}) \\ &= \frac{\sum_{s=1}^S c(f|e; F^{(s)}, E^{(s)})}{\sum_f \sum_{s=1}^S c(f|e; F^{(s)}, E^{(s)})} \end{aligned} \quad (2.8)$$

手順 4 翻訳確率 $t(f|e)$ が収束するまで手順 2 と手順 3 を繰り返す .

2.3.2.2 モデル 2

モデル 1 では , 全ての単語の対応に対して , 英語文の長さ l にのみ依存し , 単語対応の確率を一定としている . そこで , モデル 2 では , j 番目の仏単語 f_j と対応する英単語の位置 a_j は英語文の長さ l に加えて , j と , フランス語文の長さ m に依存し , 以下のような関係とする .

$$a(a_j|j, m, l) \equiv P(a_j|a_1^{j-1}, f_1^{j-1}, m, l) \quad (2.9)$$

この関係からモデル 1 における (2.4) 式は , 以下の式に変換できる .

$$P(F|E) = \epsilon \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m t(f_j|e_{a_j}) a(a_j|j, m, l) \quad (2.10)$$

$$= \epsilon \prod_{j=1}^m \sum_{i=0}^l t(f_j|e_{a_j}) a(a_j|j, m, l) \quad (2.11)$$

モデル 2 では , 期待値は $c(f|e; F, e)$ と $c(i|j, m, l; F, E)$ の 2 つが存在する . 以下の式から求められる .

$$c(f|e; F, E) = \frac{t(f|e)}{t(f|e_0) + \cdots + t(f|e_l)} \sum_{j=1}^m \delta(f, f_j) \sum_{i=1}^l \delta(e, e_i) \quad (2.12)$$

$$= \sum_{j=1}^m \sum_{i=0}^l \frac{t(f|e) a(i|j, m, l) \delta(f, f_j) \delta(e, e_i)}{t(f|e_0) a(0|j, m, l) + \cdots + t(f|e_l) a(l|j, m, l)} \quad (2.13)$$

$$c(i|j, m, l; F, E) = \sum_a P(a|E, F) \delta(i, a_j) \quad (2.14)$$

$$= \frac{t(f_j|e_i) a(i|j, m, l)}{t(f_j|e_0) a(0|j, m, l) + \cdots + t(f_j|e_l) a(l|j, m, l)} \quad (2.15)$$

モデル2では、EMアルゴリズムで計算すると複数の極大値が算出され、最適解が得られない可能性がある。モデル1では $a(i|j, m, l) = (l+1)^{-1}$ となるモデル2の特殊な場合であると考えられる。したがって、モデル1を用いることで最適解を得ることができる。

2.3.2.3 モデル3

モデル3は、モデル1とモデル2とは異なり、1つの単語が複数対応する単語の繁殖数や単語の翻訳位置の歪みについて考慮する。またモデル3では単語の位置を絶対位置として考える。モデル3では以下のパラメータを用いる。

- 翻訳確率 $P(f|e)$
英単語 e が仏単語 f に翻訳される確率
- 繁殖確率 $n(\phi|e)$
英単語 e が ϕ 個の仏単語と対応する確率
- 歪み確率 $d(j|i, m, l)$
英語文の長さ l 、フランス語文の長さ m のとき、 i 番目の英単語 e_i が j 番目の仏単語 f_j に翻訳される確率

さらに、英単語が仏単語に翻訳されない個数を ϕ_0 とし、その確率 p_0 を以下の式で求める。このとき、歪み確率は $\frac{1}{\phi_0!}$ で、 $p_0 + p_1 = 1$ で p_0, p_1 は0より大きいとする。

$$P(\phi_0|\phi_1^l, E) = \binom{\phi_1 + \dots + \phi_l}{\phi_0} p_0^{\phi_1 + \dots + \phi_l - \phi_0} p_1^{\phi_0} \quad (2.16)$$

したがって、モデル3は以下の式で求められる。

$$P(F|E) = \sum_{a_1=0}^l \dots \sum_{a_m=0}^l P(F, a|E) \quad (2.17)$$

$$= \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \binom{m - \phi_0}{\phi_0} p_0^{m-2\phi_0} p_1^{\phi_0} \prod_{i=1}^l \phi_i! n(\phi_i|e_i) \\ \times \prod_{j=1}^m t(f_j|e_{a_j}) d(j|a_j, m, l) \quad (2.18)$$

モデル3では、全てのアライメントを計算するため、計算量が膨大となるので期待値を近似により求める。

2.3.2.4 モデル4

モデル4では、モデル3と異なり、単語の位置を絶対位置ではなく、相対位置で考える。またモデル3では考慮されていない各単語の位置、例えば形容詞と名詞の関係を考慮する。モデル4では歪み確率 $d(j|i.m, l)$ を2つの場合で考える。

- 繁殖数が1以上である英単語に対応する仏単語の中で、最も文頭に近い場合

$$P(\Pi_{[i]1} = j | \pi_1^{[i]-1}, \tau_0^l, \phi_0^l, E) = d_1(j - \odot_{i-1} | \mathcal{A}(e_{[i-1]}), \mathcal{B}(f_j)) \quad (2.19)$$

\odot_{i-1} は $i-1$ 番目の英単語に対応する仏単語の位置を表している。

- それ以外の場合

$$P(\Pi_{[i]k} = j | \pi_{[i]1}^{k-1}, \pi_1^{[i]-1}, \tau_0^l, \phi_0^l, E) = d_{>1}(j - \pi_{[i]k-1} | \mathcal{B}(f_j)) \quad (2.20)$$

$\pi_{[i]k-1}$ は同じ英単語に対応している直前の仏単語を表している。

2.3.2.5 モデル5

モデル4では、単語の位置に関して直前の単語以外は考慮されていない。したがって、複数の単語が同じ位置に生じたり、単語の存在しない位置が生成される。モデル5では、この問題を避けるために、単語を空白部分に配置するよう改善が施されている。

- 繁殖数が1以上である英単語に対応する仏単語の中で、最も文頭に近い場合

$$\begin{aligned} P(\Pi_{[i]1} = j | \pi_1^{[i]-1}, \tau_0^l, \phi_0^l, E) \\ = d_1(v_j | \mathcal{B}(f_j), v_{\odot_{i-1}}, v_m - \phi_{[i]} + 1)(1 - \delta(v_j, v_{j-1})) \end{aligned}$$

v_j は j 番目までの空白数、 \mathcal{A} は英語の単語クラス \mathcal{B} はフランス語の単語クラスを表している。

- それ以外の場合

$$\begin{aligned} P(\Pi_{[i]k} = j | \pi_{[i]1}^{k-1}, \pi_1^{[i]-1}, \tau_0^l, \phi_0^l, E) \\ = d_{>1}(v_j - v_{\pi_{[i]k-1}} | \mathcal{B}(f_j), v_m - v_{\pi_{[i]k-1}} - \phi_{[i]} + k)(1 - \delta(v_j, v_{j-1})) \end{aligned}$$

2.3.3 GIZA++

GIZA++ [4] とは，統計翻訳で用いることを前提に作られたツールである．IBM 翻訳モデルを用いて，対訳文 (原言語文と目的言語文の対) から対訳単語と単語翻訳確率を自動的に得る．

2.3.4 言語モデル

言語モデルとは，膨大な量の単言語データを用いて単語の列や文字の列が起こる確率を付与するモデルである．統計翻訳では主に N -gram を用いる．以下に言語モデル (tri-gram) の例を示す．表中の w_1, w_2, w_3 はそれぞれ tri-gram の単語列における 1 番目，2 番目，3 番目の単語を示している．

表 2.2: 言語モデルの例

tri-gram の単語列	$\log_{10}(P(w_3 w_1 w_2))$ (スムージングなし)	$\log_{10}(P(w_3 w_1 w_2))$ (バックオフスムージング)
$w_1 w_2 w_3$		
痛み が 伴う	-1.382585	-0.3105274
堤防 が 決壊	-1.061585	-0.1920604
天気 が よかつ	-1.768149	-0.1920604
納得 が いか	-0.6635545	-0.1101559
梅雨 が 明け	-0.7214168	-0.1029072
風 が 吹く	-2.222238	-0.1920604

表の 1 行目の例では，左側の数値が，“痛み” と “が” という文字列が連続した後に，“伴う” が出現する確率を常用対数で表した値 “ $\log_{10}(P(\text{伴う} | \text{痛み が})) = -1.382585$ ” を，中央が tri-gram で表された単語列である “痛み が 伴う” を，右側の数値はバックオフスムージングにより得られる，“痛み が” の後に “伴う” が出現する確率を常用対数で表した値 “ $\log_{10}(P(\text{伴う} | \text{痛み が})) = -0.3105274$ ” を示している．

また，バックオフスムージングとは，高次の N -gram が存在しない場合において，低次の N -gram の値を用いて高次の N -gram の値を推定する方法である．

2.3.5 デコーダ

デコーダは，翻訳モデルと言語モデルを用いて，確率が最大となる翻訳候補を探索し，出力を行う変換器のことである．代表的なデコーダとして，“Moses” [5] がある．

英日統計翻訳において、 $\operatorname{argmax}_e P(e|j)P(j)$ の確率が最大となる日本語文を出力するために、適切な順序で英語と日本語の単語対応を得る必要がある。しかし、適切な英語文を決定するためには、計算量が膨大となり、かつ莫大な時間が必要となる。そこで計算量を削減するために、ビームサーチ法を用いる。

ビームサーチ法とは、翻訳候補の探索において、翻訳確率の低い翻訳候補を枝刈りし、探索範囲を減退する方法である。探索領域の中で一定の確率以上の翻訳候補のみを残し、それ以外の翻訳候補は除外する。

ただし、ビームサーチ法は、切り捨てられた翻訳候補が文章全体で見たときに、最大の確率を持つ翻訳候補であったという可能性がある。そのため選択した翻訳文が最適解であるとは限らないという問題がある。

2.4 パターンに基づく統計翻訳³

2.4.1 概要

パターンに基づく統計翻訳は、大きく5つの手順で翻訳を行う。パターンに基づく統計翻訳の概要を以下に示す。

手順1 対訳単語の作成

GIZA++を用いて、対訳単語を作成

手順2 単語に基づく対訳文パターンの作成

対訳単語を用いて、単語に基づく対訳文パターン(以下、単語レベル文パターン)を作成。

手順3 対訳句の作成

単語レベル文パターンを用いて、対訳句を作成。

手順4 句に基づく対訳文パターンの作成

対訳句を用いて、句に基づく対訳文パターン(以下、句レベル文パターン)を作成。

手順5 翻訳文作成

対訳句と句に基づく対訳文パターンを用いて、翻訳文生成を行う。

以下にそれぞれの手順の詳細を記述する。

³村上仁一：“パターンに基づく統計機械翻訳の概要と問題点について”

2.4.2 対訳単語の作成

GIZA++を用いて、対訳文の単語対応を取り、対訳単語と単語翻訳確率を得る。図 2.2 に対訳単語作成例を示す。

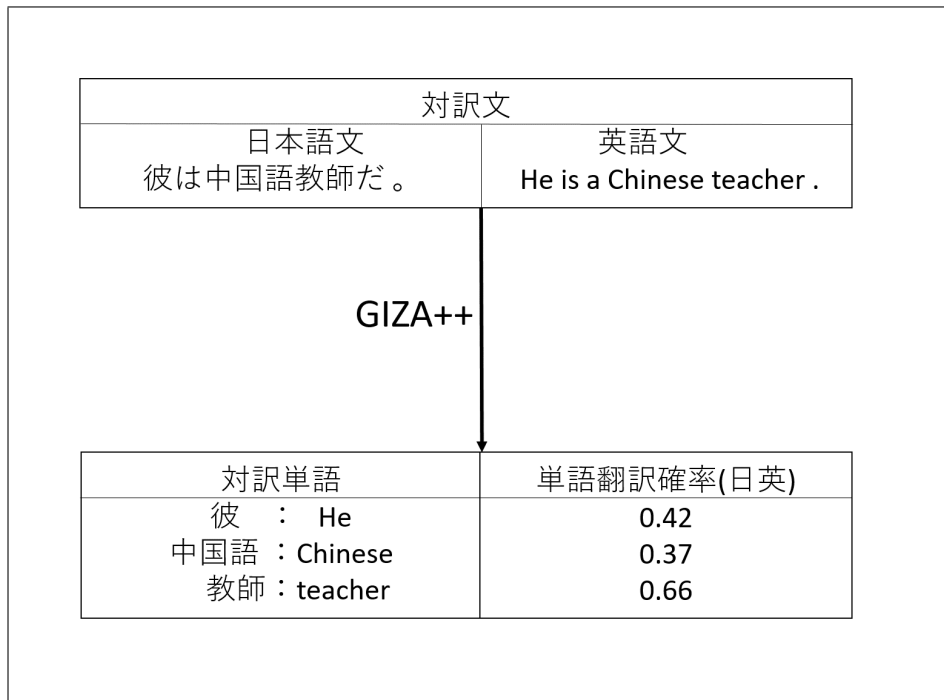


図 2.2: 対訳単語作成の例

2.4.3 単語レベル文パターンの作成

対訳単語と対訳文を用いて，単語レベル文パターンを作成する．まず，対訳単語と対訳文を照合する．そして，対訳文において，適合した対訳単語を変数化する．図 2.3 に単語レベル文パターン作成の例を示す．

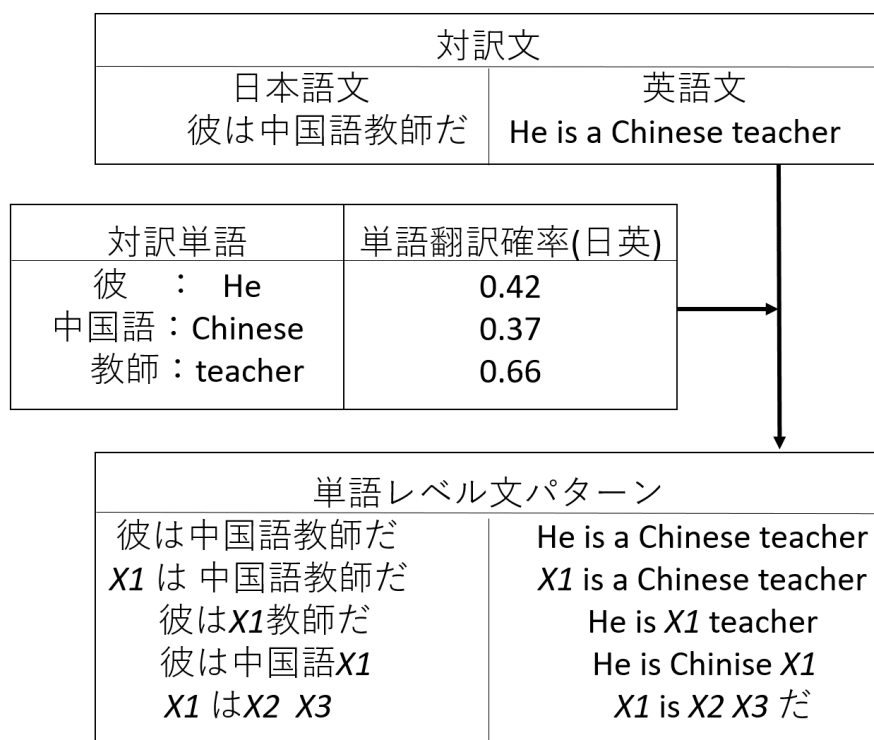


図 2.3: 単語レベル文パターン作成の例

2.4.4 対訳句の作成

1) 対訳句の抽出

単語レベル文パターンと対訳文を照合する．適合した場合，単語レベル文パターンの変数部に対応する単語を，対訳句として対訳文より抽出する．図 2.4 に対訳句抽出の流れを示す．

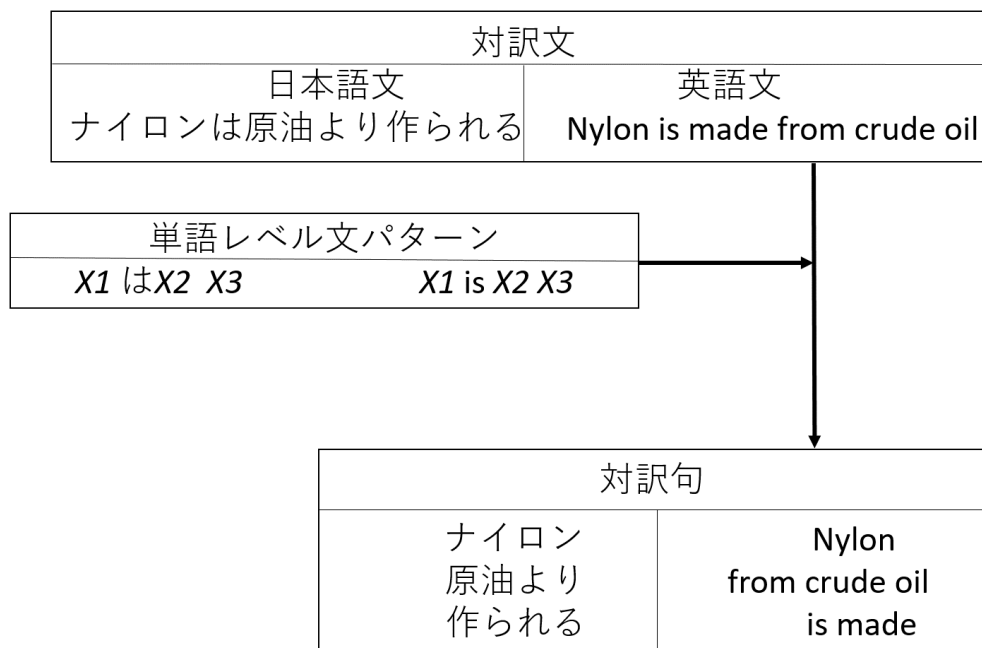


図 2.4: 対訳句作成の例

2) フレーズ確率の付与

対訳単語と単語翻訳確率を用いて、対訳句に確率を付与する。まず、対訳句において日本語句の単語と英語句の単語の全ての組み合わせを得る。次に、日本語句の単語に対応する英語句の単語の中で、単語翻訳確率の最大値を得る。これを各日本語単語に対して行い、得られた値について対数の総和を求める。(以下、フレーズ確率)。同様に対訳句において、英単語に対応する日本語単語の中で単語翻訳確率の最大値を取得し、英日方向のフレーズ確率も求める。日英方向のフレーズ確率付与の例を図 2.5 に示す。

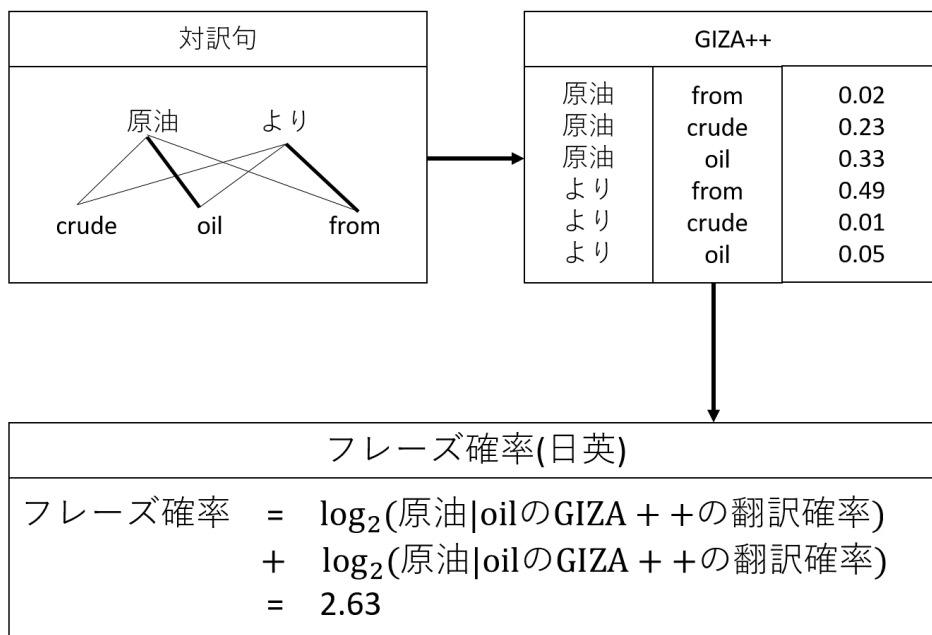


図 2.5: フレーズ確率付与の例 (日英)

2.4.5 句レベル文パターンの作成

1) 対訳文パターンの作成

対訳句と対訳文を用いて，句レベル文パターンを作成する．作成方法は，単語レベル文パターンの作成と同様に変数の組み合わせを考慮して，句レベル文パターンを可能な限り多く作成する．句レベル文パターン作成の例を図 2.6 に示す．

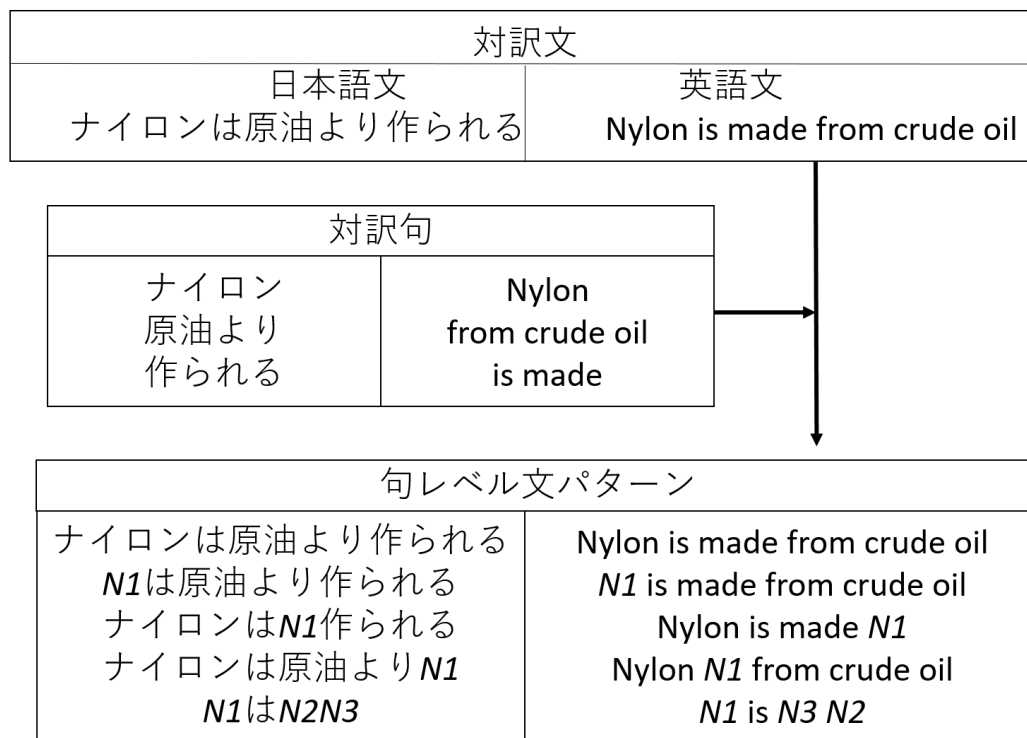


図 2.6: 句レベル文パターン作成の例

2) 文パターン確率の付与

対訳単語と単語翻訳確率を用いて，句レベル文パターンに確率を付与する．句レベル文パターンにおいて字面を用いて，フレーズ確率の付与と同様の計算手法で確率を求める．本研究では，この値を文パターン確率と呼ぶ．日英方向の文パターン確率付与の例を図 2.7 に示す．

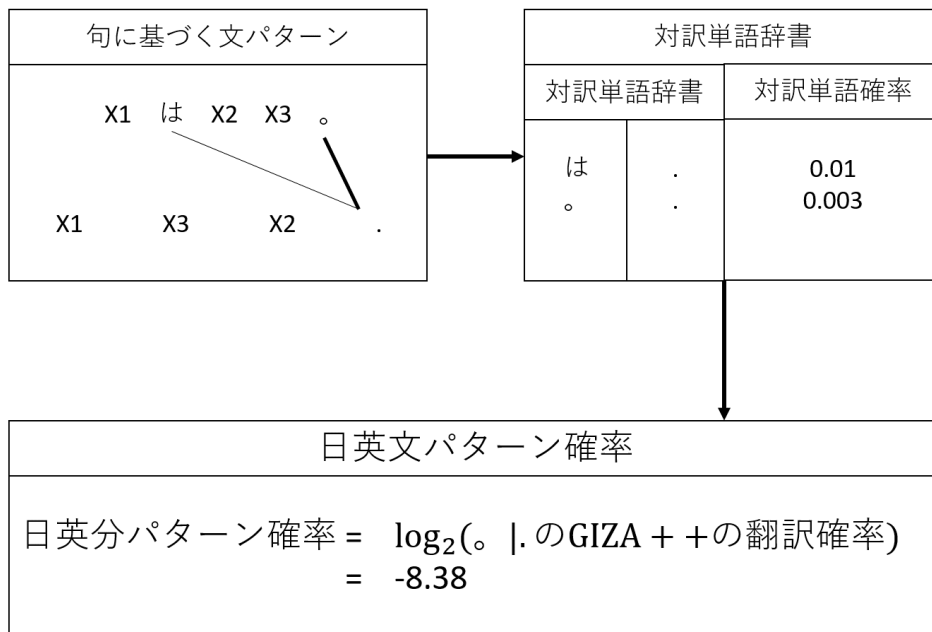


図 2.7: 文パターン確率付与の例 (日英)

2.4.6 翻訳文の作成

句レベル文パターンと対訳句を用いて、翻訳文を生成する。まず、日本語文パターンと入力文を照合し、入力文に適合する日本語文パターンを選択する。なお、文パターンの選択には、入力文と日本語文パターンの字面を比較し、字面が多く一致する文パターンを選択する。そして、選択した文パターンにおいて、英語文パターンの変数部に対訳句を用いて英語句を置換し、翻訳候補文を生成する。この処理を各適合する文パターンに対して同様に行う。最後に、各翻訳候補文から翻訳文を選択するために、句レベル文パターンの文パターン確率()と対訳句のフレーズ確率()、言語翻訳確率(trigram =)の総和を用いる。各翻訳候補文の句レベル文パターンの文パターン確率と対訳句のフレーズ確率、言語翻訳確率(trigram)の総和を求め、翻訳候補文の中で総和が最大となる文を翻訳文として出力する。日英翻訳における翻訳文の生成例を図 2.8 に示す。

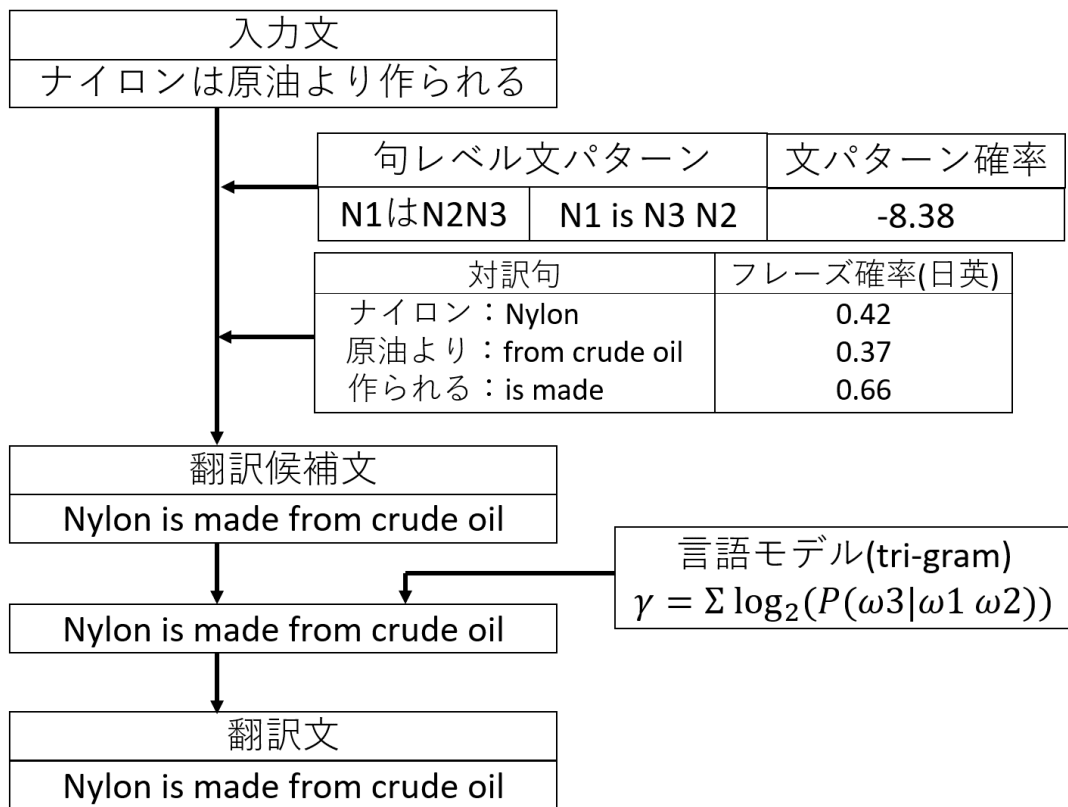


図 2.8: 翻訳文作成の例

2.4.7 従来手法の問題点

従来手法の問題点としては翻訳精度が低い。対訳句の精度が低いことが原因の一つとして挙げられる。そして対訳句はフレーズ確率で選択していた。フレーズ確率はEMアルゴリズムで計算した対訳単語確率を使用するが、1つの日本語に対して複数の英語が相当するとき対訳単語確率が不安定な値を示す場合がある。

表 2.3: 一方のフレーズ確率が低い例

対訳句	対訳句	フレーズ確率
その	The	-0.476996
その	its	-10.440154

第3章 提案手法

3.1 概要

従来手法では，翻訳時に対訳句を選択する際 IBM Model 1 から計算する対訳単語確率を使用したフレーズ確率を使用し，対訳句を選択していた．しかし，この手法における翻訳精度は十分ではない．そこで，本研究では対訳句を選択する際の確率としてフレーズ確率を掛けあわせると値の不安定さが改善されると仮定し，フレーズ確率の総積を使用する手法を提案する．さらに，フレーズ確率を使用しない Dice 係数と類似度の積を使用する手法を提案する．それぞれに対訳句の選択に利用して翻訳精度の向上を試みる．

3.2 提案手法 1(フレーズ確率の総積 (P_2) を使用する手法)

提案手法 1 はフレーズ確率の総積 (P_2) を用いる手法である．フレーズ確率の総積 P_2 は同一パターンの変数部のフレーズ確率 P_1 の総積である．計算式を式 3.1 に示す．

$$P_2(X_i) = \sum_{i=0}^W P_1(x_i) \quad (3.1)$$

W ; 変数の数

以下に例を示す．表 3.1 に対訳句作成時に使用するデータの具体例を示す．

表 3.1: 対訳学習文テストデータ

対訳学習文 (日)	この箱は木より作られる
対訳学習文 (英)	This box is made from wood
対訳文パターン (日)	X_1 は X_2 X_3
対訳文パターン (英)	X_1 X_3 X_2

表 3.1 の対訳学習文は対訳句作成に用いる文である．対訳文パターンは対訳句作成に用いるパターンである．対訳句作成の手順を以下に示す．

手順 1 対訳学習文とパターンを照合

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 3.2 に抽出した対訳句の具体例を示す .

表 3.2: 手順 2 と手順 3 の具体例

変数	対訳句		フレーズ確率 (対数)
X1	この箱	This box	-0.4
X2	木より	from wood	-0.1
X3	作られる	is made	-0.05

手順3 式を用いて対訳句にフレーズ確率を付与

表 3.2 に示しているフレーズ確率は対数値である . フレーズ確率の総積はパターンにおける変数部のフレーズ確率の総積である . ここで , フレーズ確率の総積は 1 つの変数の組から作成された対訳句において同じ値となる . 表 3.2 におけるフレーズ確率の総積の計算式を式 3.2 に示す .

$$\begin{aligned}
 & P_2\left(\frac{\text{この箱}}{\text{this box}}\right) \\
 = & P_1\left(\frac{\text{この箱}}{\text{this box}}\right) * P_1\left(\frac{\text{木より}}{\text{from wood}}\right) * P_1\left(\frac{\text{作られる}}{\text{is made}}\right) \\
 = & 2^{(-0.4)} * 2^{(-0.1)} * 2^{(-0.05)} = 2^{(-0.55)} = 0.683 \\
 = & P_2\left(\frac{\text{木より}}{\text{from wood}}\right) = P_2\left(\frac{\text{作られる}}{\text{is made}}\right) \tag{3.2}
 \end{aligned}$$

最終的なフレーズ確率の総積は全てのパターンと対訳文におけるフレーズ確率の総積の再良値となる .

3.3 提案手法 2(Dice 係数と類似度の積 (P_3)) を使用する手法

提案手法 2 は Dice 係数と類似度の積 (P_3) を用いる手法である . Dice 係数と類似度の積 P_3 は Dice 係数 $Dice(j, e)$ と類似度 P_s の積である . 計算式を式 3.3 に示す .

$$P_3 = Dice(j, e) * P_s(r) \tag{3.3}$$

3.3.1 Dice 係数

Dice 係数は頻度から計算する . Dice 係数を式 3.4 に示す .

$$Dice(j, e) = \frac{2 * count(j, e)}{count(j) + count(e)} \tag{3.4}$$

$count(j, e)$; 日本語句 j , 英語句 e が同じ対訳学習文において共起する頻度

$count(j)$; 対訳学習文の日本語文に日本語句 j が出現する頻度

$count(e)$; 対訳学習文の英語文に英語句 e が出現する頻度

日英の「この箱」と「this box」の共起頻度が4, 「この箱」の出現する頻度が6, 「this box」の出現する頻度が8の時の例を式 3.5 に示す.

$$\begin{aligned} Dice(\text{この箱}, \text{this box}) &= \frac{2count(\text{この箱}, \text{this box})}{count(\text{この箱}) + count(\text{this box})} \\ &= \frac{4}{6 + 8} = 0.28 \end{aligned} \quad (3.5)$$

3.3.2 類似度

類似度は対訳学習文とパターン原文の同一の単語の出現率である. 類似度の計算方法を式 3.6 に示す.

$$P_s(r) = \frac{N_{j1}}{M_{j1}} * \frac{N_{j2}}{M_{j2}} * \frac{N_{e1}}{M_{e1}} * \frac{N_{e2}}{M_{e2}} \quad (3.6)$$

M_{j1} ; 対訳学習文中の日本語単語数 M_{j2} ; パターン原文の日本語単語数

M_{e1} ; 対訳学習文中の英語単語数 M_{e2} ; パターン原文の英語単語数

N_{j1} ; 対訳学習文中の単語とパターン原文の単語が一致している日本語単語数

N_{j2} ; パターン原文の単語と対訳学習文の単語が一致している日本語単語数

N_{e1} ; 対訳学習文中の単語とパターン原文の単語が一致している英語単語数

N_{e2} ; パターン原文の単語と対訳学習文の単語が一致している英語単語数

表 3.3 を用いた場合の類似度の例を, 式 3.7 に示す.

表 3.3: 類似度を求めるときのデータの例

対訳学習文 (日)	この箱は木から作られる
対訳学習文 (英)	This box is made from wood
パターン (日)	X1 は X2 X3
パターン (英)	X1 X3 X2
パターン原文 (日)	この箱は鉄から作られる
パターン原文 (英)	This box is made from iron

$$P_s\left(\frac{\text{この箱}}{\text{this box}}\right) = (5/6) * (5/6) * (6/7) * (6/7) = 0.51 \quad (3.7)$$

したがって上記の例で P_3 は式 3.8 となる.

$$\begin{aligned} P_3 &= Dice(\text{この箱}, \text{this box}) * P_s\left(\frac{\text{この箱}}{\text{this box}}\right) \\ &= 0.28 * 0.51 = 0.142 \end{aligned} \quad (3.8)$$

第4章 実験

4.1 実験目的

対訳句を選択する確率の計算方法はフレーズ確率以外にも考えられる．そこで本研究では対訳句を選択する確率を”フレーズ確率の総積”と”Dice 係数と類似度の積”に変更して日英翻訳を行う．そして”フレーズ確率”を加えた3種の出力文を比較評価する．

4.2 実験データ

実験には，電子辞書などの例文より抽出した単文コーパス [6] を使用する．単文コーパスの日本語文は単文である．しかし英語文は重文と複文を含んでいる．対訳学習文および翻訳実験に使用する文の例を表 4.1 に，使用するデータの内訳を表 4.2 に示す．

表 4.1: 対訳学習文および翻訳実験に使用する文の例

日本語文	運が私に向いてきた。
英語文	Luck turned in my favor .
日本語文	お盆は先祖の霊を祭る行事だ。
英語文	The Bon Festival is a festival of worshipping our ancestors .
日本語文	ポップミュージックがカセットプレーヤーから鳴り響いた。
英語文	Pop music blared from a cassette player .

表 4.2: 実験データ

対訳学習文	160,000 文
入力文	1,000 文

4.3 評価方法

本研究では，出力文の翻訳精度の評価として人手評価と自動評価を行う．人手評価として従来手法と提案手法1と提案手法2の出力文を比較する．また，自動評価にはBLEU[7]，METEOR[8]，TER[9]，RIBES[10]を用いる．

第5章 実験結果

5.1 人手評価結果

3種の確率を用いた出力文1000文ずつから，それぞれランダムに抽出した100文を用いて，人手による対比較評価を行う．評価の基準を以下に示す．

- 従来 : 従来手法の出力文が最も優れている
- 提案1 : 提案手法1の出力文が最も優れている
- 提案2 : 提案手法2の出力文が最も優れている
- 従来× : 従来手法の出力文だけが劣っている
- 提案1× : 提案手法1の出力文だけが劣っている
- 提案2× : 提案手法2の出力文だけが劣っている
- : 全ての文で意味を読み取れない，ほぼ同一である，または同一の出力

出力文の人手評価結果を表5.1に示す．

表 5.1: テスト文 100 文の人手評価の結果

従来	提案 1	提案 2	従来 ×	提案 1 ×	提案 2 ×	
6	14	3	0	6	9	62

出力文の各例を以下の表 5.2 から表 5.7 に示す．なお各表のフレーズ確率，フレーズ確率の総積，Dice 係数と類似度の積は翻訳時に使用する句レベル文パターン中の変数部の確率の和である．

表 5.2: 従来 の例

入力文	日記 を つけ はじめた。
参照文	I've started a journal .
従来出力文	I began to write a diary .
パターン (日)	N00 を つけ N01 た。
パターン (英)	I N 01 write a N 00 .
言語翻訳確率	-29.881766
文パターン確率	-13.009582
フレーズ確率	-3.599263
提案 1 出力文	The Army began the diary .
パターン (日)	N00 を つけ はじめた。
パターン (英)	The Army began the N 00 .
言語翻訳確率	-34.287694
文パターン確率	-23.354406
フレーズ確率の総積	-15.255424
提案 2 出力文	The Army began the diary .
パターン (日)	N00 を N01 。
パターン (英)	N 01 the N 00 .
言語翻訳確率	-34.287694
文パターン確率	-2.476900
Dice 係数と類似度の積	-0.148458
入力文	その 火事 は 10 日 前 に 起こった。
参照文	The fire took place ten days ago .
従来出力文	The fire broke out in 10 days before .
提案 1 出力文	The fire broke out before the tenth .
提案 2 出力文	The fire broke out in front of the tenth .
入力文	彼は 過去 の 経験 に 頼った。
参照文	He drew on his past experience .
従来出力文	He relied on his past experience .
提案 1 出力文	He was past dependent on his experience .
提案 2 出力文	He is economically dependent on his past experience .

表 5.3: 提案1 の例

入力文	神社の境内で落ち葉を燃やしている。
参照文	They are burning fallen leaves in the grounds of the shrine .
従来出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 フレーズ確率	The shrine is Inside burned fallen leaves . N05 N02 N03 N04 を N00 N01 。 N 02 N 05 N 01 N 03 N 00 N 04 . -66.728267 -6.025362 -20.106336
提案1 出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 フレーズ確率の総積	I have burned the fallen leaves in the shrine Inside the premises . N03 の N00 N02 N04 N01 。 I N 01 N 04 N 02 the N 03 N 00 . -72.737649 -7.34012 -98.927199
提案2 出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 Dice 係数と類似度の積	The shrine is Inside burned fallen leaves . N05 N02 N03 N04 を N00 N01 。 ” N 02 N 05 N 01 N 03 N 00 N 04 . -66.728267 -6.025362 -8.828174
入力文	彼は法律の適用を誤った。
参照文	He made a mistake in the application of the law .
従来出力文	He misallocated the laws of minors .
提案1 出力文	He took a wrong to invoke the law .
提案2 出力文	He misallocated the Disaster Relief Law in the law .
入力文	マイケル・ジャクソンが新曲を吹き込んだ。
参照文	Michael Jackson has recorded a new song .
従来出力文	The Michael Jackson inspired a new song recorded .
提案1 出力文	The Michael Jackson recorded a new song .
提案2 出力文	He helped usher in a new Michael Jackson .

表 5.4: 提案 2 の例

入力文	わたしは悔いのない生涯を送りたい。
参照文	I'd like to live my life without regrets .
従来出力文 パターン(日) パターン(英) 言語翻訳確率 文パターン確率 フレーズ確率	I want to see life is not of regret . N01 は N04 の N06 N03 N05 N02 N00 。 N 01 N 00 N 05 N 02 N 03 N 06 of N 04 . -72.387453 -2.649630 -22.870618
提案 1 出力文 パターン(日) パターン(英) 言語翻訳確率 文パターン確率 フレーズ確率の総積	I I wish regretted life . N03 は N02 の ない N01 を N00 。 N 03 N 00 N 02 N 01 . -83.460649 -22.288554 -65.575554
提案 2 出力文 パターン(日) パターン(英) 言語翻訳確率 文パターン確率 Dice 係数と類似度の積	I want to lead a life no my regret . N04 は N02 の N01 N00 を N03 。 N 03 a N 00 N 01 N 04 N 02 . -71.683916 -6.206232 -11.858005
入力文	彼女は足をくじいた。
参照文	She sprained her ankle .
従来出力文	She sprained my ankle .
提案 1 出力文	She sprained my ankle .
提案 2 出力文	She twisted her leg .
入力文	彼はオートバイのエンジンをふかした。
参照文	He gunned the motorcycle engine into life .
従来出力文	He puffed at his motorcycle engine .
提案 1 出力文	He took a pull at his motorcycle engine .
提案 2 出力文	He took his motorcycle engine .

表 5.5: 提案 1 × の例

入力文	主人は客に椅子を勧めた。
参照文	The host waved his guests to a seat .
従来出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 フレーズ確率	The host offered the guest the seat . N03 は N01 に N00 N02 N04 た 。 N 03 N 04 the N 01 N 02 N 00 . -26.036633 -6.831383 -20.436261
提案 1 出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 フレーズ確率の総積	The host offered the guest a glass of the seat . 主人は客に N00 を勧めた。 The host offered the guest a glass of N 00 . The host offered the guest a glass of N 00 . -17.908530 -21.484719
提案 2 出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 Dice 係数と類似度の積	The host offered the guest the seat . N01 は N00 に N02 N03 た 。 N 01 N 03 the N 00 N 02 . -26.036633 -8.155711 -12.526808
入力文	この汽車に犯人が乗り込んだらしい。
参照文	The culprit seems to have got on this train .
従来出力文	The criminal seems to be loaded in the this train .
提案 1 出力文	The criminal seems got is the this train .
提案 2 出力文	The criminal seems to be loaded in the this train .
入力文	その馬には彼が乗る。
参照文	The horse is jockeyed by him .
従来出力文	He holds the horse .
提案 1 出力文	The horse holds him .
提案 2 出力文	He holds the horse .

表 5.6: 提案 2 × の例

入力文	ハンカチに自分の名前を刺しゅうした。
参照文	I embroidered my name on the handkerchief .
従来出力文	She embroidered her name on a handkerchief .
パターン (日)	ハンカチ N00 N02 を N01 た。
パターン (英)	She N 01 her N 02 N 00 a handkerchief .
言語翻訳確率	-51.192848
文パターン確率	-7.890944
フレーズ確率	-10.521583
提案 1 出力文	She embroidered her name on her handkerchief .
パターン (日)	N01 に N00 を刺しゅうした。
パターン (英)	She embroidered her N 00 on her N 01 .
言語翻訳確率	-59.809190
文パターン確率	-11.290342
フレーズ確率の総積	-28.022058
提案 2 出力文	She embroidered her initials on a rock with a handkerchief .
パターン (日)	N02 に N01 を N00 。
パターン (英)	N 00 a N 01 N 02 .
言語翻訳確率	-38.325355
文パターン確率	-8.208456
Dice 係数と類似度の積	-6.702979
入力文	健康にはいつも注意している。
参照文	I am always careful about my health .
従来出力文	There is always pays attention to the health .
提案 1 出力文	Health is always attention .
提案 2 出力文	I always take the warning .
入力文	彼は客から注文をとった。
参照文	He took the customer's order .
従来出力文	He took the order from the customer .
提案 1 出力文	He took the order from the customer .
提案 2 出力文	He took my order for the guest .

表 5.7: 全ての文で意味を読み取れない, ほぼ同一である, または同一の出力例

入力文	契約を守らない方が悪いのだ。
参照文	The fault is with the one has not observed the contact .
従来出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 フレーズ確率	not signed preserve of a bad person . N05 N04 N00 N02 N01 N03 だ。 N 00 N 05 N 04 N 03 N 01 N 02 . -79.228048 -7.824789 -28.636751
提案 1 出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 フレーズ確率の総積	is toword the bad keep their contact . N04 N01 N03 N02 N00 N05 だ。 N 03 N 02 N 05 N 00 N 01 N 04 . -97.000000 -8.485753 -101.084818
提案 2 出力文 パターン (日) パターン (英) 言語翻訳確率 文パターン確率 Dice 係数と類似度の積	The riddance jealously lease of the contact . N01 N04 N02 N00 が N03 。 The N 00 N 02 N 01 N 04 N 03 . -78.458952 -3.281977 -11.793332
入力文	ここから港がよく見える。
参照文	We can see the whole harbor from here .
従来出力文	I see harbor well from here .
提案 1 出力文	I see harbor well from here .
提案 2 出力文	Here you can see the harbor .
入力文	彼らは逃げ腰になった。
参照文	They were ready to flee .
従来出力文	They became 逃げ腰 .
提案 1 出力文	They became 逃げ腰 .
提案 2 出力文	They became 逃げ腰 .

5.2 自動評価結果

テスト文 1,000 文で翻訳実験を行い，出力文に対して自動評価を行う．表 5.8 に，それぞれの手法における自動評価の結果を示す．

表 5.8: テスト文 1000 文の自動評価の結果

パラメータ	BLEU	METEOR	TER	RIBES
フレーズ確率 (P_1)	0.1495	0.4175	0.6467	0.7309
フレーズ確率の総積 (P_2)	0.1552	0.4308	0.6388	0.7354
Dice 係数と類似度の積 (P_3)	0.1459	0.3955	0.6762	0.7139

5.3 実験結果のまとめ

表 5.1 と表 5.8 より”フレーズ確率の総積”を使用したときの結果が最も良く，”Dice 係数と類似度の積”を使用した結果が最も悪いことが分かった．しかし大きな差はないことが分かった．なお未知語が出現する文では 3 種どの確率を使用した文であっても未知語として出力されていた．

第6章 考察

6.1 フレーズ確率の総積の有効性

表 5.1 より，翻訳精度としてフレーズ確率の総積を使用する手法が最も優れていることが分かった．理由は二つ考えられる．一つ目はフレーズ確率の総積は変数と対訳句の組が全て適当であるときのみ確率が高くなる．そして，翻訳時最も確率が高い対訳句を使用するため適切な翻訳が行われたと考えられる．二つ目はフレーズ確率よりフレーズ確率の総積の値が大きくなるので言語翻訳確率の影響を受けにくいことが考えられる．

6.2 差なしの原因

フレーズ確率の総積 P_2 を使用する手法と従来手法の”フレーズ確率 P_1 ”を使用した手法には大きな差がなかった．これは2つの原因が挙げられる．まずフレーズ確率の総積を使用した翻訳の翻訳確率においても言語翻訳確率 (trigram) が大きな割合を占めていること，次に確率の変更をした場所がデコーダー部分であったことである．

なお，現在対訳句から句に基づく文パターンを作成する際に確率を使用して枝刈りを行っているので，句に基づく文パターンを作成する際の確率を変更すると本研究よりも差が出るのではないかと考える．

6.3 誤り分析

表 5.5 の提案 1 × の 5 文について解析すると，1 つ目の文が単語が過剰な文，二つ目の文が動詞連続文，三つ目の文が主語目的語が逆転している文となった．表 6.1 で示していない 2 文の結果を以下に示す．

6.1 の上の文は単語不足の文，下の文は同音異義語を選択している文となった．単語不足の文と単語が過剰な文と同音異義語を選択している文の合計 3 文が出力されたのは誤った対訳句を選択してしまったことが原因である．動詞連続文と主語目的語が逆転している文の合計 2 文が出力された原因は誤ったパターンを選択してしまったことが原因である．翻訳精度の向上において対訳句選択とパターン選択の問題対策が今後の課題と言える．

表 6.1: 提案 1 × の残りの 2 文

入力文	子供達は公園に遊びに行った。
参照文	The children went to play in the park .
従来出力文	The children went to play in the park .
提案 1 出力文	The children went to the park .
提案 2 出力文	The children went to play in the park .
入力文	私は表の通りに出た。
参照文	I went out on to the main street .
従来出力文	I went to the top of the street .
提案 1 出力文	I went to the list .
提案 2 出力文	I went to the top of the street .

6.4 他の対訳句を選択する確率の計算方法

対訳句を選択する確率は他に式 6.1 や式 6.2 で計算可能である．今後比較実験を行っていきたい．

$$P\left(\frac{J_0 \cdots J_{N-1}}{E_0 \cdots E_{M-1}}\right) = \prod_{m=0}^{M-1} \arg \max_{0 \leq n \leq N-1} p(E_m | J_n) * \prod_{n=0}^{N-1} \arg \max_{0 \leq m \leq M-1} p(J_n | E_m) \quad (6.1)$$

$$P_1\left(\frac{J_0 \cdots J_{N-1}}{E_0 \cdots E_{M-1}}\right) = \prod_{m=0}^{M-1} \prod_{n=0}^{N-1} (p(E_m | J_n) * p(J_n | E_m)) \quad (6.2)$$

第7章 おわりに

本研究ではパターンに基づく統計翻訳の問題点を翻訳文の作成時に使用する対訳句の精度の低さと定めた．この問題点を解決するために手法を2つ提案した．提案手法1はフレーズ確率の総積である．また，提案手法2はDice係数と類似度の積である．提案手法1,2と従来手法の翻訳を比較した結果フレーズ確率の総積を使用する手法の翻訳精度が最も優れていた．しかし従来手法と大きな差はなかった．このため翻訳精度の向上には句に基づく文パターン作成時の使用する確率の変更やフレーズ確率の別の計算方法で求めることに着目してパターンに基づく統計翻訳の改善を行なっていきたい．

付録

付録は表 7.1 に示す内容である .

表 7.1: 付録一覧

内容	ページ
入力文ランダム 100 文	1 ~ 3
参照文ランダム 100 文	4 ~ 6
人手評価文:ランダム 100 文	7 ~ 15
人手評価結果:	16 ~ 18

謝辞

最後に，1年間に渡りご指導いただきました鳥取大学工学部知能情報工学科自然言語処理研究室の村上仁一准教授，村田真樹教授，卒業論文発表をご指導してくださいました応用計算知能研究室の徳久雅人講師をはじめ，自然言語処理研究室の方々に厚く御礼申し上げます．

また，参考にさせていただいた論文の著者の方々に対して深く感謝申し上げます．

参考文献

- [1] 江木孝史 “句に基づく文パターンを用いた英日翻訳 ” , 鳥取大学修士論文、 2013.
- [2] 村上仁一 “パターンに基づく統計機械翻訳の概要と問題点について” , 電子情報通信学会技術研究報告 , 言語理解とコミュニケーション , NLC2017-3 , pp.13-18 , 2017.
- [3] Peter F.Brown, Stephen A.Della Pietra, Vincent J.Della Pietra, Robert L.Mercer, “The mathematics of statistical machine translation:Parameter Estimation” , Computational Linguistics, 1993.
- [4] GIZA++ : <http://www.fjoch.com/GIZA++>
- [5] Moses: Philipp Koehn, Marcello Federico, Brooke Cowan, Richard Zens, Chris Dyer, Ondrej Bojar, Alexandra Constantin, Evan Herbst, “Moses: Open Source Toolkit for Statistical Machine Translation” , Proceedings of the ACL 2007 Demo and Poster Sessions, pp.177-180, 2007.
- [6] 村上仁一, 藤波進 “日本語と英語の対訳文対の収集と著作権の考察 ” , 第一回コーパス日本語学ワークショップ, pp.119-130. 2012.
- [7] BLEU: “a Method for Automatic Evaluation of Machine Translation” , Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL), pp.311-318. 2002.
- [8] Meteor: Lavie Alon, and Denkowski Michael “An Automatic Metric for MT Evaluation with High Levels of Correlation with Human Judgments” , Proceedings of the Second Workshop on Statistical Machine Translation, pp.228-231. 2007.

- [9] Richard Schwartz, Linnea Micciulla, John Makhoul: “A Study of Translation Edit Rate with Targeted Human Annotation”, AMTA, 2006.
- [10] Hideki Isozaki, “Automatic Evaluation of Translation Quality for Distant Language Pairs”, Proceedings of the 2010 Conference on Empirical Methods in Natural Language Processing, pp.944-952. 2010.