

日英新聞記事の自動記事対応付け

高橋 大和 白井 諭 藤波 進
NTT コミュニケーション科学研究所

池原 悟
鳥取大学

上田 洋美 松島 英之

NTT アドバンスドテクノロジー(株)

〒238-03 神奈川県横須賀市武1-2356 620C

NTT コミュニケーション科学研究所 翻訳処理研究グループ

TEL: 0468 59 8238 / E-mail: {yamato, shirai, fujinami}@nttkb.ntt.jp

あらまし

機械翻訳などの自然言語処理技術を研究する上で、大量の対訳コーパスは非常に有用である。しかし、大量の一般的なデータの収集は困難である。本稿では、継続的なデータ収集を目的として、日英新聞記事を対象に、数値と名詞をキーワードとして利用し、その対応項目数により自動的な日英記事対応を行う手法を提案する。

この手法を用いた実験により、一日平均90記事中38記事の正しい記事対応を自動的に得られることを確認した。

キーワード 記事対応、対訳コーパス、データベース、新聞記事

Automatically Aligning Japanese & English Newspaper Articles

Takahashi, Y., Shirai, S., Fujinami, S., Ikehara, S., Ueda, H. and Matsusima, H.
NTT Communication Science Lab., Tottori Univ., NTT Advanced Technology Co.
620C 1-2356 Take Yokosuka-Shi Kanagawa 238-03 JAPAN
TEL: +81 468 59 8238 / E-mail: {yamato, shirai, fujinami}@nttkb.ntt.jp

Abstract

Bilingual Corpora are very useful in natural language processing. Unfortunately they are difficult to compile.

We have developed a method in which numerical values and proper nouns in the articles are used as keywords to align the Japanese and English newspaper articles automatically in order to develop a corpus.

In addition, we have developed a way to evaluate the results. We correctly align automatically an average of 38 out of 90 pairs of articles daily.

key words Aligning articles, Bilingual corpora, Database, Newspaper articles

1. はじめに

機械翻訳などの自然言語処理技術を研究する上で、大量の対訳コーパスは非常に有用であるが、大量の一般的なデータの収集は困難である、という問題点がある。

これに対して、新聞記事を対象とすれば継続的にデータを収集することができる利点があるが、従来、対応付けが容易でないと考えられてきた。しかし、日英のように言語のギャップが大きい場合でも、記事内容に着目すれば、基本的に一致すると考えられる。

例えば、日本経済新聞社が提供しているテレコンDBから取得した日英記事を比較検討した例では、部分対応を含めると、ほとんどの英文に内容的には対応する日本語があり、そのうち5割は格要素などの対応もとることができることが報告されている[1]。

対訳コーパスを構築するためには、1. 日英記事対応、2. 日英文対応を行う必要がある。1. に関して、数値と名詞をキーワードとして統計的に対応を行う方法が提案されている[2]が、この手法では、対応の採否の判定基準の決定が困難であるという問題があり、最終的には人手による確認が必要であると考えられる。

そこで、本稿では自動的な記事対応付けの手法の確立をめざし、数値と名詞をキーワードとして利用し、その対応項目数により採否の判定を行う評価アルゴリズムを提案する。

2. 記事の特徴

本稿で対象とした記事は、日本経済新聞社が有料情報サービスとして提供しているテレコンDBから、電話回線経由のパソコン通信により取り寄せることができる。日本文記事は、日経テレコンB I Zに収録されている日経四紙（日本経済新聞、日経産業新聞、日経流通新聞、日経金融新聞）を対象とし、英文記事は、Nikkei Telecom Japan News & Retrievalより、日経四紙の速報訳として提供されている記事を対象とした。

1994年11月1日から10日までの記事数の例を表1に示す。日本文記事は英文記事に比べ7倍近い記事数が収録されていることがわかる。したがって、英文記事を基準として対応する日本

文記事を得るほうが効率がよい。

対訳テキストにおける文の対応付けには、統計的情報や対訳辞書の利用が有効であることが知られている[3][4][5]が、数百記事の中から目的とする記事を見つけるために応用することは困難であると考えられる。一般的には、記事を特定する情報として、DB検索ではキーワードを利用する。この方法を応用した場合、対応する記事の検索に英日の対訳キーワード辞書を利用して、日本文記事の検索を行う。しかし、この場合は対応付けの品質は対訳辞書の品質に大きく左右される。

そこで、通常はキーワードとして利用されていない数値情報に着目する。数値情報は、日英記事のどちらにおいても、記事からの抽出が容易であり、対応も取り易い。よって、記事を特定する情報として有用であると考えられる。例として、数値情報を含む記事数を表1に示す。数値を含む記事は日本語で約80%、英文で約90%あることがわかる。

表1 日本文記事と英文記事の量(1994年11月分)

日付	曜日	日本文記事			英文記事		
		記事数	平均 字数	含数値	記事数	平均 Bytes	含数値
1	火	1001	410.2	843	137	792.9	125
2	水	842	408.5	703	120	775.4	111
3	木	485	396.4	360	39	830.6	34
4	金	738	449.4	647	97	803.0	85
5	土	504	367.7	358	29	905.4	26
6	日	167	588.4	127	12	829.2	8
7	月	629	541.5	519	128	758.8	116
8	火	929	451.9	671	152	754.5	132
9	水	819	410.8	641	146	782.3	131
10	木	941	471.7	756	151	777.3	142
合計		7055	437.7	5625	1011	783.0	910

数値情報と名詞をキーワードとする方法としては、稀少度による評価により、数値情報だけで約50%、名詞キーワードを含めると80%の対応付けを行えることが知られている[2]。しかし、対応記事の採否の判定基準となる稀少度は記事の数やキーワードの出現回数に左右されるため、その基準を決めるのは困難である。

そこで、キーワードの出現回数に代えて、対応項目数に着目する。数値の組は、記事の個性を表す特徴的な情報として考えられるためである。

また、数値による記事対応に名詞キーワードを併用することにより、得られる対応記事を増加させることができる。本稿では、対応項目数により対応記事の採否を判定する日英記事対応アルゴリズムを提案する。

3. 記事対応付けアルゴリズム

対応する記事を発見する際の問題として、検定すべき候補記事はかなり多数にのぼることが表1からわかる。そこで、本稿では字面処理程度の浅い解析による方法を中心に考察した。具体的には、英文記事の本文に含まれる数値と名詞をキーワードとして対応付けを行う。以下に、キーワードの抽出アルゴリズムを述べる。

3. 1 英文記事

記事は1日分を対象として、その範囲内に含まれる数値およびキーワードをすべて抽出する。

3. 1. 1 数値キーワードの切り出し

1. 本文の数値を数値リストとして切り出す。小数点や次の単位が続いている場合はそれらを含めて切り出す。ここで扱う単位は、出現頻度が高く、日本文においても切り出しやすい単位のみを扱う。
例: dollar, yen, %, trillion, billion, million
2. “trillion, billion, million” を含む数値を数値列に正規化する。
例: 4.3 trillion dollar →
4300000000000 dollar
3. 数値リストは記事単位で項目の重複がないように、重複した数値項目を削除する。
4. 数値リストをソートし、項目の出現回数を数え、数値項目重みデータとする。
5. 数値リストに、数値項目重みデータの重みを付加する。これを、英文記事重み付き数値リストとする。

3. 1. 2 名詞キーワードの切り出し

1. 本文と見出しから名詞と推定される単語を名詞リストとして切り出す。以下の条件を満たす単語を一単語とみなす。
 - ・大文字を含む単語列
例1 :NTT Communication Science Lab.
例2 :SL-enhanced Intel i486SX
 - ・大文字を含む単語列の所有格に大文字

を含む単語がある時

例1 : Japan Federation of Employers' Associations

例2 : International Standardization Organization's ISO9001

- ・“of”, “&” を大文字を含む単語間に挟んでいる時

例1 : Social Democratic Party of Japan

例2 : Mitsubishi Trust & Banking Corp.

- ・大文字を含む単語列の所有格に大文字を含む単語列が接続していない時、また、“of” の後ろに大文字を含む単語が接続しない時は、そこまでで切り出す。

例1 : NTT's line → NTT

例2 : Bank of city → Bank

- ・“The” は単語列に含まない。これは、文の途中では小文字になり、切り出し単語が増えてしまうためである。

例: The U.S. → U.S.

2. 名詞リストの項目をキーとして、対訳辞書を検索する。日本語訳があった場合、名詞リストに日英の対として追加する。訳がなかった場合は項目を削除する。
3. 名詞リストは、記事単位で日本語訳の重複がないように、記事単位で単語の一部分を含んでいる日本語単語があれば、長い単語のみ残り短い単語項目を削除する。
例: Tokyo 東京
Bank of Tokyo 東京銀行
この二つの対応項目があった時は、“Bank of Tokyo 東京銀行”のみ残す。
4. これを、英文記事キーワードリストとする。

3. 2 日本文記事

英文1日分に対して、日本文は英文記事の日付前後1日を含めた3日分を対象として対応付けを行う。これは、白井らによって報告されているように[2]、対応する英文記事と日本文記事の日付が同じ記事は63.5%、英文記事が1日遅いものは30.6%のほか、1日早いものが5.9%含まれているからである。

この3日分の日本文記事から数値をすべて抽出する。なお、日本文記事から固有名詞を直接抽出することは容易でないため、対応付けの対象範囲として、見出し文およびリード文（第一段落）を

抽出する。

3. 2. 1 数値キーワードの切り出し

1. 数値を切り出す条件は、数字、漢数字の連鎖、また、小数点、「ドル」、「円」、「銭」、「%」が続いている場合、これらを含めて一単語として扱う。
例：三十五円四十銭、五〇、八%
2. 数値を数字列に正規化する。
例：三十五円四十銭 → 35.40 yen
3. 数値リストは記事単位で項目の重複がないように、重複した数値項目を削除する。
4. 数値リストをソートし、項目の出現回数を数え、数値項目重みデータとする。
5. 数値リストに、数値項目重みデータの重みを付加する。これを、**日本語記事重み付き数値リスト**とする。

3. 2. 2 リード文の抽出

1. 日本語記事三日分のタイトルと第一段落をリード文として切り出す。これを、**日本語記事リード文リスト**とする。

3. 3 対応付け

次のように、数値のみによる対応付け、英文から得たキーワードによる対応付けの2つの方法を試みた。

3. 3. 1 数値の記事対応

英文記事重み付き数値リストと日本語記事重み付き数値リストの対応付けを行う。対応付けは、正規化された数値のマッチングによって行われる。



この結果を、数値対応付けリストとする。

3. 3. 2 キーワードの記事対応

英文記事キーワードリストと日本語記事リード文リストの対応付けを行う。対応付けは、英文記事名詞リストの名詞が日本語記事リード文リストのリード文に含まれているかどうかにより行われる。



この結果を名詞キーワード対応付けリストとする。

4. 対応付けの評価と考察

4. 1 数値による記事対応

前節で述べたアルゴリズムで得られた対応付けリストの評価実験を行った。実験1として、稀少度が最も高くなる日本語記事：英文記事での数値の出現回数が1対1の対応を調べた。結果を表2に示す。ただし、項目の個数が同数の場合は、対応候補を決定しない。この場合は、稀少度を用いて、相対的に判定するためである。ただし、この実験では、稀少度による評価は行わない。

[実験1]

数値対応付けリストにおいて、日英の出現頻度が1対1の数値を含む記事対応で個数が多いもの。

結果を表2に示す。

表2 1対1対応 (1994年11月2日~9日分)

日付	2日	3日	4日	5日	6日	7日	8日	9日
1個	15	4	14	5	0	18	9	23
2個	9	7	22	5	2	21	24	17
以上	9	7	22	5	2	21	24	17

正解数/対応が得られた記事数

結果として、1対1の出現回数による評価では、対応項目の個数が1個の場合は、平均で約70%程の正解率であり、2個以上の組み合わせがある場合は、100%の正解率が得られる。

これより、出現回数による評価より、対応項目個数に着目して評価を行った方がよいと考えられる。そこで、以下の実験2を行った。

[実験2]

数値対応付けリストにおいて、対応する数値項目の個数が多いもの。ただし、項目の個数が同数の場合は、対応候補を決定しない。

結果を表3に示す。

この条件では、全体で332記事の対応が得られ、対応が正しい記事は317記事で、約95.5%の正解率だった。対応個数が少ないと記事の正解率が落ちることが表3から分かる。また、対応が5個あっても間違った記事対応が得られている。

表3 対応項目数と正解数 (1994年11月2日~9日分)

個数	2個	3個	4個	5個	6個以上	合計
2日	7	6	11	6	18	48
3日	2	6	1	2	8	19
4日	2	7	12	2	23	46
5日	2	4	2	2	6	16
6日	0	1	1	0	3	5
7日	6	7	7	8	29	57
8日	7	12	15	6	23	63
9日	1	11	12	7	32	63
	3	13	12	7	32	67

正解数/対応が得られた記事数

記事に含まれている数値で一番多いものは、年度であり(参考資料1を参照)、これらがノイズとして正解率を下げていると考えられる。一日を単位として頻度をみた場合、一日の記事総数の影響を受ける。そこで、対応候補記事間の対応数の差が大きいほど、正しい候補であろうと考え、以下の条件で実験を行った。

[実験3]

数値対応付けリストにおいて、対応記事第一候補と第二候補の対応項目数の差が2個以上のもの。ただし、項目の個数が同数の場合は、対応候補を決定しない。

結果を表4に示す。

表4 対応記事正解数 (1994年11月2日~9日分)

日付	2日	3日	4日	5日	6日	7日	8日	9日	合計
*1	30	14	39	8	2	39	44	46	222
*2	73	25	62	22	6	91	90	98	467
*3	120	39	97	29	12	128	152	146	733

*1:対応記事正解数

*2:数値項目が3個以上含む英文記事数

*3:英文記事数

この条件により得られた日本文記事は全て正しい対応であった。8日間で733記事中222記事(約30.3%)に対して、正しい対応を得ることができた。

4.2 固有名詞による記事対応

次に、3.1.2により切り出した固有名詞のみを用いた記事対応の実験を行った。実験に先立ち、固有名詞の辞書が必要になる。高橋⁶では、機械翻訳システムの持つ代表的な地名と企業名の辞書(1345項目)を用いたが、製品名や人名も記事対応を行うために有用な情報になること、また、企業名に関しても、記事内では略称が使われることが多いことから、ここでは、数値により得られた対応記事222記事から人手により対訳辞書を作成した辞書を使った。内容を表5に示す。

表5 英日対訳辞書

内容	項目数
企業名	588
製品名	107
人名	23
地名・その他	136
合計	1344

この辞書を用いて、以下の条件で記事対応実験を行った。

[実験1]

名詞キーワード対応付けリストにおいて、対応する項目数が一番多いもの。ただし、項目の個数が同数の場合は、対応候補を決定しない。

結果を表6に示す。

実験により、175記事の記事対応が得られ、正しい対応は152記事(約86.9%)であった。まだ対訳辞書が不十分なためか、現状では数値による対応に比べ得られる対応の数も、正解率も低い。特に、対応項目数が2個の場合、正解率が約83.8%と低く、全体の正解率を落としている。内容としては、企業名二つ、もしくは企業名と地名による対応である。名詞キーワードにおいては、地名(参考資料2を参照)が出現頻度が高いため、その影響を受けていると考えられる。

表6 対応記事数 (1994年11月2日～9日分)

個数	2個	3個	4個	5個	6個以上	合計
2日	13 13	11 11	4 4	3 3	0 0	31
3日	2 2	4 4	0 0	1 1	0 0	7
4日	6 7	9 11	1 1	0 0	2 2	18 21
5日	3 3	2 2	0 0	0 0	1 1	6 6
6日	1 1	2 2	0 0	0 0	0 0	3 3
7日	13 17	9 9	2 2	1 2	0 0	25 30
8日	17 19	5 7	8 8	2 2	0 0	32 36
9日	12 21	11 13	3 3	4 4	0 0	30 41

正解記事数/対応が得られた記事数

[実験 2]

名詞キーワード対応付けリストにおいて、対応記事第一候補と第二候補の対応する項目数の差が2個以上のもの。ただし、項目の個数が同数の場合は、対応候補を決定しない。

結果を表7に示す。

表7 対応記事数 (1994年11月2日～9日分)

日付	2日	3日	4日	5日	6日	7日	8日	9日	合計
対応数	14	3	9	1	2	8	12	11	60
	14	3	10	1	2	8	12	11	61

正解記事数/対応が得られた記事数

実験の結果、61記事の対応が得られ、60記事が正解(約98.4%)だった。この条件でも一記事のみ正しい対応は得られていない。また、対応記事の量も約39.5%と、かなり少なくなる。名詞キーワードのみの対応づけでは、現状の対訳辞書ではそれほどよい結果は得られないと考えられる。

4. 3 数値と名詞による記事対応

最後に、対応記事第一候補と第二候補の対応項目数の差が2個以上という条件で、数値対応付け

リストと名詞キーワード対応付けリストをマージしたデータに対して評価した結果を表8に示す。

表8 数値とキーワードによる対応付け (1994年11月2日～9日分)

日付	2日	3日	4日	5日	6日	7日	8日	9日	合計
*1	49	21	46	13	3	53	63	58	306
*2	30	14	39	8	2	39	44	46	222

*1:数値とキーワードによる正解対応記事数

*2:数値のみの正解対応記事数

この結果より、数値のみによる対応づけに比べ、数値とキーワードを併用した場合、84記事(約37.8%増加)の新しい対応を得ることができることがわかる。

5. 問題点と今後の改良

今後は、名詞キーワードの拡張とその効果の確認を行う。しかし、名詞キーワードの対訳の作成を手で行うには、対訳を効率的に見つける方法が難しい。現在、英文記事から対訳がない項目リストを得ることはできるため、その一般的な訳をつけていく方法を行う。

また、対応する個数が同数の場合、日本文記事の候補を決定するには稀少度による相対的評価を行うのがよいと考えられるが、稀少度の計算において、数値の単位による評価値の加減が必要と考えられる。

本手法では、名詞キーワードを数値項目と区別せず、マッチングの個数だけで評価している。この点も、実験により決定していきたい。

6. おわりに

本稿では、並立するデータベースから収集した日英新聞記事を、浅い解析で得られる数値情報と名詞キーワードに着目して対応項目数を評価し、対応記事を自動的にかつ高い精度で得る方法を提案した。日英新聞記事8日分に対する対応付け実験の結果、対応記事候補の第一候補と第二候補の差が2個以上の第一候補を対応記事候補とする、という条件で、1日平均約41.7%の記事対応を正解率100%で得ることができた。今後は固有名詞辞書を拡張していき、その効果の確認と対応

項目が同数の候補記事の評価方法と重みづけを実験により検討する。

本手法により、大量の日英の対訳記事を収集することが可能になり、新語や専門用語の対訳の収集、対訳表現の抽出など、辞書の整備や翻訳表現調査の効率化が図れると考えられる。また、白井[7]に提案されている文対応の方法を実験し、対応情報を SGML タグとして構造化し[8]、継続的な大量の対訳コーパスの構築を目指す予定である。

参考文献

- 1 白井,藤波,池原,上田,井上:新聞記事日英対訳コーパスの構築(1)ー基本構想と検討課題ー,電気関係学会九州支部第48回連合大会(1995)
- 2 白井,上田,阿部,藤波,池原:新聞記事日英対訳コーパスの構築(2)ー並立DBから取得した記事の対応付けー,電気関係学会九州支部第48回連合大会(1995)
- 3 P.F.Brown, J.C.Lai & R.L.Merce: Aligning sentences in parallel corpra, 29th ACL(1991)
- 4 S.F.Chen: Aligning sentences in parallel corpra, 31st ACL(1993)
- 5 T.Utsuro, H.Ikeda, M.Yamane, Y.Matsumoto & M.Nagano: Bilingual text matching using bilingual dictionary and statistics, 15th COLING(1994)
- 6 高橋,白井,藤波,池原,上田,松島: DBから抽出した日英新聞記事の自動対応づけ,言語処理学会第2回年次大会(1996)
- 7 白井,松尾,瀬下,藤波,池原:新聞記事日英対訳コーパスの構築(3)ー記事の特徴分析と文の対応関係の検討ー,電気関係学会九州支部第48回連合大会(1995)
- 8 F.Bond, Y.Takahashi, S.Yamada, M.Nisigaki: Still tagging an aligned Japanese/English corpus, 言語処理学会第2回年次大会(1996)

参考資料1：出現頻度の多い数値項目リスト

日付	1日	2日	3日	4日	5日	6日	7日	8日	9日	10日
日本語	94 (95)	94 (87)	95 (49)	20 (72)	12 (53)	95 (14)	12 (76)	11 (104)	11 (90)	11 (114)
	31 (191)	95 (91)	12 (52)	95 (80)	95 (59)	11 (15)	95 (82)	12 (105)	94 (92)	12 (124)
	11 (240)	11 (201)	11 (70)	11 (140)	11 (65)	12 (26)	11 (84)	95 (115)	95 (108)	95 (150)
英文	93 (17)	93 (11)	100 (4)	20% (8)	94 (4)	300 (2)	93 (10)	93 (10)	93 (12)	96 (14)
	95 (18)	94 (21)	94 (6)	95 (13)	95 (4)	94 (2)	95 (28)	94 (21)	95 (27)	95 (30)
	94 (28)	95 (30)	95 (9)	94 (23)	96 (7)	95 (5)	94 (29)	95 (24)	94 (33)	94 (33)

数値項目(出現数)

参考資料2：出現頻度の多い名詞項目リスト

日付	1日	2日	3日	4日
英文から 切り出した 単語	Bank of Japan 日銀 (4)	Bank of Japan 日銀 (5)	Sakura Bank さくら銀行, さくら銀 (2)	Hitachi 日立 (4)
	Asian アジア (5)	Osaka 大阪府, 大阪市 (5)	Sega Enterprises Ltd. セガ・エンタープライゼス (2)	Ministry of Finance 大蔵省 (5)
	Tokyo 東京, 東京都 (14)	Tokyo 東京, 東京都 (10)	Tokyo 東京, 東京都 (2)	Tokyo 東京, 東京都 (13)

日付	5日	6日	7日
英文から 切り出した 単語	Matsushita Electric Industrial Co. 松下電器産業, 松下電工 (2)	Shoko Chukin Bank 商工中金 (1)	Toshiba Corp. 東芝 (3)
	Tokai Bank 東海銀行, 東海銀 (2)	Small Business Finance Corp. 中小企業金融公庫 (1)	NEC Corp. NEC (5)
	APEC APEC (3)	Tokyo 東京, 東京都 (1)	Tokyo 東京, 東京都 (6)

日付	8日	9日	10日
英文から 切り出した 単語	Kanagawa Prefecture 神奈川県 (6)	Tokyo Stock Exchange 東京証券取引所, 東証 (4)	Mitsubishi 三菱電機 (6)
	Osaka 大阪府, 大阪市 (6)	Hitachi Ltd. 日立製作所, 日立 (5)	Osaka 大阪府, 大阪市 (6)
	Tokyo 東京, 東京都 (10)	Tokyo 東京, 東京都 (13)	Tokyo 東京, 東京都 (14)

英単語列

日本語対訳(出現数)