



機械翻訳

東京大学 情報基盤センター
情報理工学系研究科、情報学府 (兼任)
中川裕志

昔の機械翻訳

- ▶ 入力文: 私はりんごを食べた。
- ▶
 - ▶ →形態素解析→構文解析
 - ▶ noun verb noun → subj predicate object
 - ▶ →意味解析
 - ▶ (action=食べる, agent=私, target=りんご, time=past)
 - ▶ 英語の語彙に変換(つまり意味表現のレベルないしはそれに近い深さで変換 ←対訳辞書利用)
 - ▶ (action=eat, agent=I, target=an apple, time=past)
 - ▶ 構文および形態素の生成(語順の変換)して翻訳出力を得る。
←対訳辞書利用
 - ▶ noun=I, verb (past)=ate, noun=an apple
- ▶ 出力文: I ate an apple.

昔の機械翻訳

- 意味のレベルで精密に日英が同一であることが前提だった。
- また、形態素解析、構文解析、意味解析が正確に動作すると想定している。
- しかし、なかなかそうとも言い切れない
 - 意味レベルでの概念が一致しない例
 - 湯 → hot water、
 - もったいない → ?、
 - check という習慣が日本にない！

対訳辞書

▶ 日本語 → 意味

▶ りんご → APPLE

▶ 意味 → 英語

▶ APPLE → if bear noun or singular: apple
if plural: apples

◆ 単数の場合には an apple, 複数なら apples
を選ぶのは、構文および形態素のレベル

少し前の機械翻訳: example based translation

- ▶ 翻訳対の例文が類似検索可能な形でデータベース化
 - ▶ 例: 私はみかんを食べた。 ↔ I ate an orange.
- ▶ 入力文: 私はりんごを食べた。
 - ▶ 翻訳対データベースから類似した日本語例文を検索
 - ▶ 私はみかんを食べた。
 - ▶ 違っている部分みかんをりんごに置き換え
 - ▶ さらに日英辞書でりんごをan appleに置き換え
- ▶ 結果出力: I ate an apple.
 - ▶ 当然ながら、冠詞の選択などは文法規則によって行う。つまり、相当程度に従来の構文規則や、形態素解析技術と共同することになる。

少し前の機械翻訳: example based translation

- ▶ 類似検索の部分が重要。ここで構文解析を使うことも可能だが、だんだん古典的な機械翻訳に近づく。
- ▶ 翻訳対を集めれば集めるほどが翻訳の質があがる。
 - ▶ この収集作業は機械的にできる。

統計的機械翻訳

Statistic Machine Translation (SMT)

- 言語的知識を全く使わずに対訳を得る。アンチ言語学理論
- 2言語並行コーパスが蓄積
- 文どうしの対応付けされた aligned corpus
- これを使って単語や句どうしの対応付け、すなわち対訳を自動的に抽出
 - 文同士の対応はあるが、単語列同士の対応は不明
 - 探索空間が膨大
- IBMの Peter Brown, S. Della Pietra, V. Della Pietra, Robert Mercerらの1993年のCLの論文“The Mathematics of Statistical Machine Translation:Parameter Estimation”を中心に解説

Bayesの定理

- Canadian Hansard : French-English Bilingual corpus
- フランス語の単語列: f に対して妥当な英語の単語列: e を求める
- Given French string: f , find $e^{\wedge} = \arg \max_e \Pr(e/f)$
 - 種々の f に対応しそうな e はやたらと多い！！

■ then
$$\Pr(e/f) = \frac{\Pr(e)\Pr(f/e)}{\Pr(f)}$$

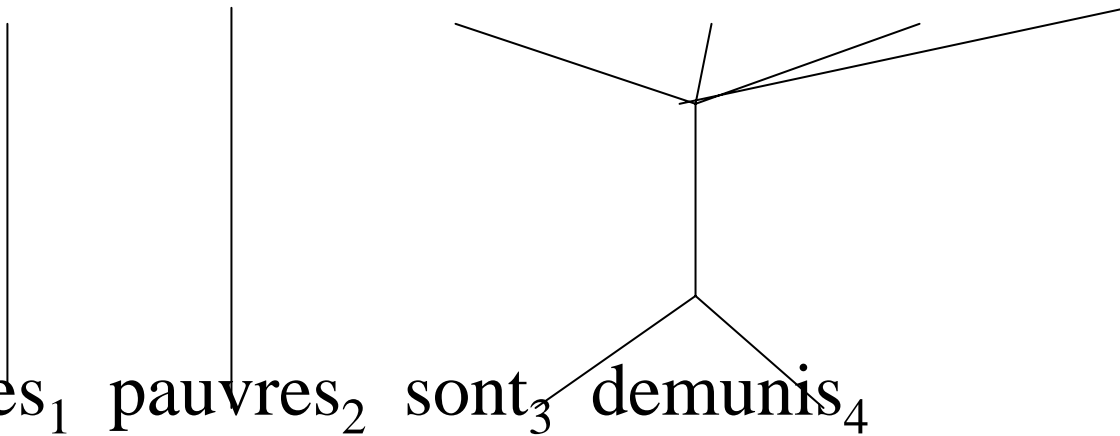
$$\Rightarrow e^{\wedge} = \underset{e}{\operatorname{argmax}} \Pr(e)\Pr(f/e)$$

なぜ $\Pr(e/f)$ ではなく、 $\Pr(f/e) \times \Pr(e)$ か？

- 種々の f に対応しそうな e はやたらと多い！！
 - 対訳コーパスの対訳文はやはり少数
- 無尽蔵に多くあるフランス語の単語列 f に対して、対応すべき正しい英語を求めるのが目的
- $\Pr(e/f)$ 直接では、正しい英単語列 e に高い確率が割り当てられることが保証されない。
- 正しい英文という要因を直接考慮するために $\Pr(e)$ を別個の情報源から得て利用する。

Alignment: 対応

- The₁ poor₂ don't₃ have₄ any₅ money₆



- Les₁ pauvres₂ sont₃ demunis₄

(Les pauvres sont demunis |

The(1) poor(2) don't(3,4) have(3,4) any(3,4)
money(3,4))

=A(e,f)=a

→ e,fはここでは文

記法

- Alignmentも考慮した $\Pr(f, a | e)$

$$\Pr(f | e) = \sum_a \Pr(f, a | e)$$

$e = e_1^l \equiv e_1 e_2 \dots e_l$ where l English words

$f = f_1^m = f_1 f_2 \dots f_m$ where m french words

$a = a_1^m = a_1 a_2 \dots a_m$, where $a_j = i$

f, e は単語列 a はalignment f_i, e_{a_j} は単語

$$\Pr(f, a | e) = \Pr(m | e) \prod_1^m \Pr(a_j | a_1^{j-1}, f_1^{j-1}, m, e) \Pr(f_j | a_1^j, f_1^{j-1}, m, e)$$

- 以後は $\Pr(f, a, | e)$ を評価する方法

IBM Model 1

- ◆ このモデルでは、英、仏文の単語の出現順序には関係がないとしている。－(1)
- ◆ また対訳は個々の単語にだけ依存する－(2)

$$\varepsilon \equiv \Pr(m / e)$$

$$\Pr(a_j | a_1^{j-1}, f_1^{j-1}, m, e) = (l + 1)^{-1} \quad - (1)$$

$$\Pr(f_j | a_1^{j-1}, f_1^{j-1}, m, e) = t(f_j | e_{a_j}) \quad - (2)$$

$$\Pr(f, a | e) = \frac{\varepsilon}{(l + 1)^m} \prod_1^m t(f_j | e_{a_j}) \quad - (3)$$

f, e は単語列 a は *alignment* f_i, e_{a_j} は単語

Model 1

- このモデルでは、Alignment a_j は0から m の任意の値をとるから $1/(l+1)$ 。ラグランジュ未定乗数法によって $\Pr(f|e)$ を最大化する。

$$\Pr(f | e) = \frac{\varepsilon}{(l+1)^m} \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m t(f_j | e_{a_j}) \quad -(4)$$

$$\text{constraint: } \sum_f t(f | e) = 1 \quad f, e \text{ は単語 } f_j, e_{a_j} \text{ のいずれかを表す。}$$

$$h(t, \lambda) \equiv \frac{\varepsilon}{(l+1)^m} \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m t(f_j | e_{a_j}) - \sum_e \lambda_e (\sum_f t(f | e) - 1) \quad -(5)$$

$$0 = \frac{\partial h}{\partial t(f | e)} = \frac{\varepsilon}{(l+1)^m} \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \sum_{j=1}^m \delta(f, f_j) \delta(e, e_{a_j}) t(f | e)^{-1} \prod_{k=1}^m t(f_k | e_{a_k}) - \lambda_e \quad -(6)$$

$$\begin{aligned}
t(f | e) &= \lambda_e^{-1} \frac{\varepsilon}{(l+1)^m} \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \sum_{j=1}^m \delta(f, f_j) \delta(e, e_{a_j}) \prod_{k=1}^m t(f_k | e_{a_k}) \\
&= \lambda_e^{-1} \sum_a \Pr(f, a | e) \sum_{j=1}^m \delta(f, f_j) \delta(e, e_{a_j}) \quad - (7)
\end{aligned}$$

$$c(f | e; f, e) = \sum_a \Pr(a | f, e) \sum_{j=1}^m \delta(f, f_j) \delta(e, e_{a_j}) \quad - (8)$$

- $c(\dots)$ とは翻訳($f|e$)において、英単語 e がフランス語単語 f に翻訳される回数。2番目の Σ はあるalignment a において f, e の対訳された回数で総和。

$$\Pr(a | f, e) = \Pr(f, a | e) / \Pr(f | e) \quad f, a, e \text{ は単語列}$$

$$\rightarrow c(f | e; f, e) = \frac{\sum_a \Pr(f, a | e) \sum_{j=1}^m \delta(f, f_j) \delta(e, e_{a_j})}{\Pr(f | e)} \quad - (9)$$

$$\rightarrow t(f | e) = \lambda_e^{-1} \Pr(f | e) c(f | e; f, e) \quad - (10)$$

- (9) 式の $\sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^l \delta(e, e_i)$ の部分は f と e の対訳回数になることが分かる。(alignment a がないのでこの式) 下図参照。

- $f =$ f1, f2(=f), f3, ..., f7(=f), fm
- e1(=e) * *
- e2
- $e =$:
- e8(=e) * *
- :
- e1

- 教師データとして S 個の翻訳 $(f^{(s)} | e^{(s)})$ $s=1, \dots, S$ がコーパスから知られているので、 S 個の翻訳の総和を用いて式(10)を書き換えた以下の式を使うので覚えておいてください。ただし、(10)の $\lambda_e Pr(f/e)$ を λ_e と置き換えた。

$$t(f | e) = \lambda_e^{-1} \sum_{s=1}^S c(f | e; f^{(s)} e^{(s)})$$

$t(f|e)$ を求めるまではもう一工夫

◆ $t(f_j | e_{a_i})$ は、単項式だから

$$\sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m t(f_j | e_{a_j}) = \prod_{j=1}^m \sum_{i=0}^l t(f_j | e_i) \quad - (11)$$

◆ 例 $t_{10}t_{20} + t_{10}t_{21} + t_{11}t_{20} + t_{11}t_{21} = (t_{10} + t_{11})(t_{20} + t_{21})$

◆ これによると

$$\Pr(f | e) = \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l t(f_j | e_i) - (12)$$

◆そこで、またラグランジュ未定乗数法で

$$\text{constraint: } \sum_f t(f | e) = 1$$

$$h(t, \lambda) \equiv \frac{\varepsilon}{(l+1)^m} \prod_{j=1}^m \sum_{i=0}^l t(f_j | e_i) - \sum_e \lambda_e (\sum_f t(f | e) - 1) \quad - (13)$$

$$0 = \frac{\partial h}{\partial t(f | e)} = \frac{\varepsilon}{(l+1)^m \sum_{i=0}^l t(f_j | e_i)} \prod_{j=1}^m \sum_{i=0}^l t(f_j | e_i) \sum_{j=1}^m \delta(f, f_j) \sum_{i=1}^l \delta(e, e_i) - \lambda_e$$

-(14)

$$\rightarrow \lambda_e^{-1} \Pr(f | e) = \frac{\sum_{i=0}^l t(f_j | e_i)}{\sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^l \delta(e, e_i)} \quad - (15)$$

$$\text{by(10)(15)} \rightarrow c(f | e; f, e) = \frac{t(f | e)}{t(f | e_0) + \dots + t(f | e_l)} \sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^l \delta(e, e_i) - (16)$$

いよいよEMで $t(f/e)$ を推定 - 1

1. $t(f/e)$ の初期値を適当に決める
2. 各 $(f^{(s)}, e^{(s)})$, $1 \leq s \leq S$ に対して、

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

を計算する。

$$\sum_{j=1}^m \delta(f, f_j) \sum_{i=0}^l \delta(e, e_i) \text{ の値は } f, e \text{ が } f^{(s)}, e^{(s)} \text{ の要素}$$

のときだけ0でない。

いよいよEMで $t(f/e)$ を推定 - 2

3. $t(f | e) = \lambda_e^{-1} \sum_{s=1}^S c(f | e; f^{(s)}, e^{(s)})$ を \sum_f すると
左辺が1になるので、 $\lambda_e = \sum_f \sum_{s=1}^S c(f | e; f^{(s)}, e^{(s)})$

この λ_e の値を用いて $t(f/e)$ の新たな値を推定する。

$$t(f | e) = \frac{\sum_{s=1}^S c(f | e; f^{(s)}, e^{(s)})}{\sum_f \sum_{s=1}^S c(f | e; f^{(s)}, e^{(s)})}$$

4. $t(f/e)$ が収束するまで2,3を繰り返す。

Model 2

- Alignmentが位置に依存する。つまり、

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

が j a_j m l に依存

$$\sum_{i=0}^l a(i | j, m, l) = 1$$

$$\text{then } \Pr(f | e) = \varepsilon \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \prod_{j=1}^m t(f_j | e_{a_j}) a(a_j | j, m, l)$$

ラグランジュ

$$h(t, a, \lambda, \mu)$$

$$= \Pr(f | e) - \sum_e \lambda_e (\sum_f t(f | e) - 1) - \sum_j \mu_{jml} (\sum_i a(i | j, m, l) - 1)$$

ラグランジュ未定乗数法で h を微分し計算すると

$$t(f | e) = \lambda_e^{-1} \sum_a \Pr(f, a | e) \sum_{j=1}^m \delta(f, f_j) \delta(e, e_j)$$

$$c(f | e; f, e) = \sum_a \Pr(a | f, e) \sum_{j=1}^m \delta(f, f_j) \delta(e, e_{a_j})$$

$$t(f | e) = \lambda_e^{-1} \Pr(f | e) c(f | e; f, e)$$

$$t(f | e) = \lambda_e^{-1} \Pr(f | e) \sum_{s=1}^S c(f | e; f^{(s)}, e^{(s)})$$

ここまでは同じだが、さらに意味的に考えて

$$c(i | j, m, l; f, e) = \sum_a \Pr(a | e, f) \delta(i, a_j)$$

$$a(i | j, m, l) = \mu_{jml}^{-1} \sum_{s=1}^S c(i | j, m, l; f^{(s)}, e^{(s)})$$

注： $\Pr(a | f, e) \Pr(f | e) = \Pr(a, f | e)$

$$\mu_{jml} \Pr(f | e) \rightarrow \mu_{jml}$$

Model 1と同じように計算し

- Model 1 では $(1+1)^{-1}$ だった $a(i|j,m,l)$ を Model 2 では変数と見ているので、

$$\Pr(f | e) = \varepsilon \prod_{j=1}^m \sum_{i=0}^l t(f_j | e_i) a(i | j, m, l)$$

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

$$c(i | j, m, l; f, e) = \frac{t(f_j | e_i) a(i | j, m, l)}{t(f_j | e_0) a(0 | j, m, l) + \dots + t(f_j | e_l) a(l | j, m, l)}$$

- 後は同じくEMアルゴリズムで $t(f/e)$ を求める
- 初期値には Model 1 の結果を用いる。



Model 3

- 1 単語が n 単語に翻訳 not \Rightarrow ne ... pas
- $n=0$ (翻訳されない) の場合も考慮。冠詞は日本語にはない。
- 対応する単語の出現場所がねじれる
 - 日英での語順の差
- こういった現象に対応するモデル

- 繁殖確率 $n(\phi/e)$: 英語単語 e が ϕ 個のフランス語単語に接続される確率
- 翻訳確率 $t(f|e)$: 英語単語 e がフランス語単語 f に翻訳される確率
- 歪確率 $d(j|i, m, l)$: 英文長さ l , フランス文長さ m , 英文の単語位置 i がフランス文の単語位置 j に接続される確率
- 空の英単語の繁殖数 = ϕ_0

- 空でない英単語から生成されたフランス語単語の後に空から生成された ϕ_0 個の単語が確率 p_1 で挿入されるとすると、

$$\Pr(\phi_0 \mid \phi_1^l, e) = \binom{\phi_1 + \dots + \phi_l}{\phi_0} p_0^{\phi_1 + \dots + \phi_l - \phi_0} p_1^{\phi_0}$$

$$p_0 + p_1 = 1$$

$$\sum_{i=0}^m \phi_i = m$$

以上の準備の下

$$\begin{aligned}\Pr(\mathbf{f} | \mathbf{e}) &= \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \Pr(\mathbf{f}, \mathbf{a} | \mathbf{e}) \\ &= \sum_{a_1=0}^l \dots \sum_{a_m=0}^l \binom{m - \phi_0}{\phi_0} p_0^{m-2\phi_0} p_1^{\phi_0} \prod_{i=1}^l \phi_i! n(\phi_i | e_i) \times \\ &\quad \prod_{j=1}^m t(f_j | e_{a_j}) d(j | a_j, m, l) \quad - (32)\end{aligned}$$

$\phi_i!$ は e_i から生成された ϕ_i 個の単語列における順番を入れ替えた単語列の
数え上げ

- (32)式を用いて、 $n, t, d, p_{0,1}$ に関する各々の総和=1の確率による条件をつけてラグランジュ未定乗数法で $\Pr(e|f)$ を最大化すればよい。
- しかし、model1,2と異なり和積の交換ができないので、直接計算する。
- 組み合わせの数が多いので、ビタビ法で近似計算する。

二言語コーパスからの対訳抽出 -- Aligned corpus の場合 --

- Parallel Corpus (平行、対訳コーパス)
 - Aligned Corpus: 種々の研究あり。要はどのようにして2つの言語のコーパスにおける文、単語、複合語、句を対応付ける(align) するか集中。
 - 90年代前半にきれいな2言語対訳コーパスを対象にした多数の研究があり。
 - 90年代後半に、Noisy Parallel Corpus への展開が試みられた (Fung94,Fung98)

対訳コーパスからの対応文のペアを求める

- ◆ Gale and Church 1993
- ◆ 2言語の文書 S, T から対応付け (Alignment) A を求める。
- ◆ S と T の対応する文のペアを bead という。
 - ◆ 例 $B = (\text{言語 language},$
 $B = (\text{les eaux mineral, mineral water})$
- ◆ $\text{Alignment} = \operatorname{argmax}_A P(A/S, T) = \operatorname{argmax}_A P(A, S, T)$
- ◆ B_k は 文書先頭から k 番目の bead

$$P(A, S, T) \approx \prod_{k=1}^K P(B_k)$$

対訳コーパスからの対応文のペアを求める

◆例 $B=(\text{言語 language})$,

$B=(\text{les eaux mineral, mineral water})$

$$\text{Align} = \operatorname{argmax}_A P(A, S, T) \approx \operatorname{argmax}_A \prod_{k=1}^K P(B_k)$$

◆これらから分かるように、よい対応付けは、単語やcollocationレベルでの2言語での対応付けから得られるbead: B_k の確からしさ $P(B_k)$ が大きいものを集めて作る。

◆具体的なアルゴリズムはDPによる

対訳コーパスからの対応文のペアを求める

- ◆このようにして、文内部の対応がうまく付く文同士が高い対応付け確率を持つとする。
- ◆文書同士で文同士が高い対応付けを持つペアを連ねるように選ぶと対応文ペアが求まる。
 - ◆文書を構成する各文同士の対応付けもまた、文の対応付け確率を用いたDPによる。
- ◆Beadを作る単語同士の対訳の他に、文の長さが類似しているほど、よく対応する文のペアと考える方法も加味できる。

対訳コーパスからの対訳文の概要

- ◆ 動的計画法による探索 Gale&Chru93,
Nissen et al 98
- ◆ A*探索 Wang & Waibel 97
- ◆ 確実な対訳と見られるアンカー一点を決めつつ
の反復法 Kay & Roscheisen 93

Aligned corpus からの 単語やcollocation 対訳抽出の概要

◆ 対訳判定の尺度：

- ◆ 各言語の対訳文毎に対訳かどうか判定したい表現 w_1, w_2 (=各言語から抽出された表現) が出現する文としない文を数える。
- ◆ w_1, w_2 が対訳文の出現する確率による。ただし、この確率 = 0 ならEMアルゴリズムで推定。Kupiec93
- ◆ w_1, w_2 の相互情報量 Haruno93, Smadja96
- ◆ w_1, w_2 のDice係数 Smadja96
- ◆ ヒューリスティック
- ◆ Likelihood ratio Melamed97

語の2言語コーパスでの共起と その有意さの尺度

- あるテキストの単位(文、段落など)へのW1(言語A),w2(言語B)のcontingency matrix (頻度) $a+b+c+d=N$

	W2出現	W2非出現
W1出現	a	b
W1非出現	c	d

- 相互情報量 $MI(w1,w2)=\log\frac{P(w1,w2)}{P(w1)P(w2)}=\log\frac{aN}{(a+b)(a+c)}$

- Dice係数 $Dice(w1,w2)=\frac{2a}{(a+b)+(a+c)}$

Champollion (Smdja et al 96)

- Translating collocations: based on sentence aligned bilingual corpus
- 1つの言語においてXtractで collocation を抽出
- 相手側の言語の対応文である統計的フィルタをパスするような分布で出現している collocation を対訳とする。
- フィルタとしては、相互情報量、Dice係数を比較した結果、

Champollion

- Dice係数は、 $X=0, Y=0$ (双方の言語の文に collocationの要素がない場合)に影響されないなので、精度が高いので、これを使う。

$$Dice(X, Y) = \frac{2 \text{freq}(X=1, Y=1)}{\text{freq}(X=1) + \text{freq}(Y=1)}$$

- 大雑把には collocation (の要素たち)の Dice係数が閾値以上のもののみを残す。
- なお、極端に出現頻度の低い要素語は捨てる。

Champollion

- Canadian Hansards (50MB order)
 - 3000, 5000 collocations extracted by Xtract
 - 中頻度の300 collocations で評価
 - Xtractのerror rate = 11%
 - Incorrect translations = 24%
 - Correct translations = 65%
 - Champollion's precision = 73%

Likelihood ratio: Melamed 97

- Melamed 97 の方法は、対訳はone-to-oneという仮定をおき、
 - 2言語の単語 u, v が対訳であるかどうかの尤度比を計算し、
 - 最大の尤度比のものを対訳とする。
 - ただし、 u, v の個別の出現例から
 - $\lambda_+ = P(\text{対訳} | \text{共起かつ対訳})$ 、 $\lambda_- = P(\text{対訳} | \text{共起かつ非対訳})$ 、をコーパスで評価しておく
 - $\lambda_+ \lambda_-$ を $\prod_{u,v} P(k(u,v) | n(u,v), \lambda_{\pm})$ を最大化するようにして求める
- 次の尤度比を用いる。
$$L(u,v) = B(k(u,v) | n(u,v), \lambda_+) / B(k(u,v) | n(u,v), \lambda_-)$$
 - λ_{\pm} 機能語、内容語で分類する
 - recall=90% -- precision=87%

二言語コーパスからの対訳抽出 -- non aligned corpus の場合 --

- Non-aligned Corpus の場合
 - Alignさせながら対訳を探す方法 Fung95ACL
 - Alignment を使わず別の情報源による方法
 - 複合語の構成 Fung95WVLC
 - 文脈に出現する語の分布
Rapp95, Tanaka96, Fung98,
 - 特定の分野のコーパスであること

二言語コーパスからの対訳抽出 -- non aligned corpus の場合 --

- 基本的な方法は、まず双方の言語で用語あるいは collocation を抽出し、次に何らかの方法で対応を見つけようとする。

Noisy Parallel Corpus の場合

- Fung94 ACL では English-French parallel corpus (Hansards) から次のような方法で対訳を求めた。
 1. 同じ内容をもつ英仏テキスト対をK等分する
 2. 英単語 W_{en} と仏単語 W_{fr} のK個のセグメントでの現れ方($\langle 1, 0, 1, 0, 0 \rangle$ などと表現)の相関を相互情報量MIで測る
 3. MIの高いものが対訳とみなす。
 4. 低頻度単語でMIが高くでる傾向の是正策(t-scoreによる検定)も行う。
- K分割というところが味噌。Fung95につながる。

高頻度の確実な対訳語対を探してAlignment させながら低頻度語の対訳も探す方法 ——Noisy Parallel Corpus の場合——

- Fung95 ACL95
- Alignment と 対訳抽出を同時に行う。
- Step1. 両言語のコーパスで頻度の高い語についての対応付けをする。うまくいけば、この対応付けが同時に対訳になる。
 - 各語の相対出現位置ベクトルを求める
 - このベクトルをDP風にマッチさせ対応する対を選ぶ。
なお、ベクトルの平均と分散のによるユークリッド距離の離れている対は除去
 - この対応付けをアンカー点と呼ぶ
 - 高頻度の対応付けだけを残すとかなりきれいな alignment になっている

Fung 95 ACL の続き

- Step2: 低頻度の語の対応付け
 - アンカー一点で区切ったセグメントを s_1, s_2, \dots とする (K等分ではないところが進歩)
 - 低頻度の語の相対位置ベクトルをセグメントの順番 (i of s_i) の列に変換する。これによって、ラフな alignment でも誤差がセグメント単位に切り上げられる。
 - その上で相互情報量の小さい対を対訳として選ぶ。

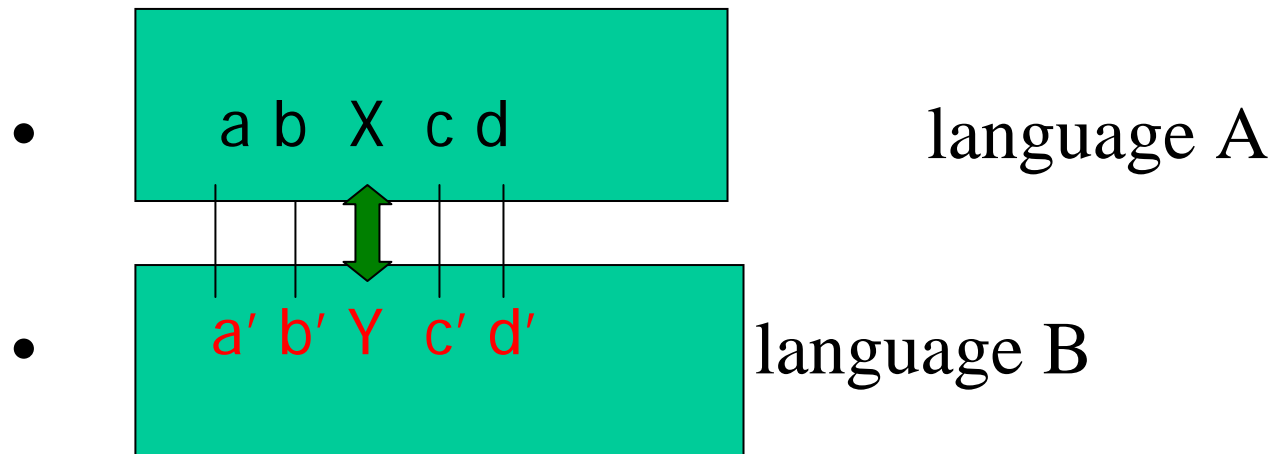
Fung 95 ACL の続き

□ 結果

- 6000語ほどの小さな英中コーパス
 - 高頻度の対訳 128語 80%弱の精度
 - 低頻度の対訳 533語 70%程度の精度
 - 平均して 73%の精度
- 必ずしも実用レベルの実験ではない

非平行な同一分野の2言語コーパスからの訳語抽出

- Non-parallel comparable corpora
- Similarity of context is cue.



- Calculation of contexts similarity is heavy

Context Heterogeneity (Fung95WVLC)

- 扱う分野が同じなら parallel でなくてよい。
- 単語 trigram により文脈を作る。つまり
- 単語L 単語0 単語R
 - ここで、単語0 の context heterogeneity は
 - 単語Lの異なり数 = a, 単語Rの異なり数 = b
 - 単語0 の出現頻度 = c のとき
 - Left-heterogeneity = a/c, right-heterogeneity = b/c
 - このアイデアは語基の接続数(中川)に似ている
- 二つの言語の単語 w_1, w_2 の context heterogeneity x_1, y_1 (for w_1), x_2, y_2 (for w_2) の距離は
$$\varepsilon = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Context heterogeneity の評価

- ε は統計量なのでとにかくコーパスが大きくなると信頼性の高いデータが得られない。
 - 大規模なデータ量は、non-aligned corpus の場合、必然かも。
- Context heterogeneity のような局所的な文脈だけで対訳を絞りこめるのか疑問
 - 評価結果が論文に書いてない
- 既存の辞書を使っていない。ゼロから対訳を得ようというのはかなり無理では。

より本格的に文脈を利用する方法 Rapp

- Rapp95,99 では、共起する語の分布の近さで訳語の曖昧性を解消しようとする。
- 小規模な英独対訳辞書は利用する。
- ドイツ語コーパスである単語wdの前後2単語づつでの共起語を使う。
- 上で述べた位置関係において二つの単語 w_a, w_b の共起を頻度や確率、MIでなく次にしめす対数尤度比(これが最良と報告)によってあらわす。

より本格的に文脈を利用する方法 Rapp

□ 対数尤度比

□ あるテキストの単位（文、段落など）への w_a, w_b （言語B）の contingency matrix（頻度） $a+b+c+d=N$,

□ w_a, w_b が共起 = a , w_a だけ = b , w_b だけ = c , 両方ともない = d

$$\begin{aligned} -2 \log \lambda = & a \log \frac{aN}{(a+b)(a+c)} + b \log \frac{bN}{(a+b)(b+d)} \\ & + c \log \frac{cN}{(a+c)(c+d)} + d \log \frac{dN}{(b+d)(c+d)} \end{aligned}$$

□ この対数尤度比をベクトルの要素にして、 w_a に共起する $w(\text{ger})$ とその対訳 $w(\text{eng})$ のベクトルの距離を次式で測る

$$s(w(\text{ger}), w(\text{eng})) = \frac{(\lambda_{-2} \lambda_{-1} \lambda_{+1} \lambda_{+2})_{w(\text{ger})} - (\lambda_{-2} \lambda_{-1} \lambda_{+1} \lambda_{+2})_{w(\text{eng})}}{}$$

より本格的に文脈を利用する方法 Rapp

- 前式で $w(\text{ger})$ の対訳 $w(\text{eng})$ は既存の辞書で求める。
- 距離 $s(w(\text{ger}), w(\text{eng}))$ を w_a に近接する語基全てについて加算し、これを S とする。
- S の最小のものが w_a の訳語であるとする。
- ドイツ語の語基から英語への対訳を100語のサンプルで評価した結果、
- 第1順位の対訳の精度は72%、
- 上位10位以内に正しい訳語がある場合も加えると89%の精度

より本格的に文脈を利用する方法 Fung

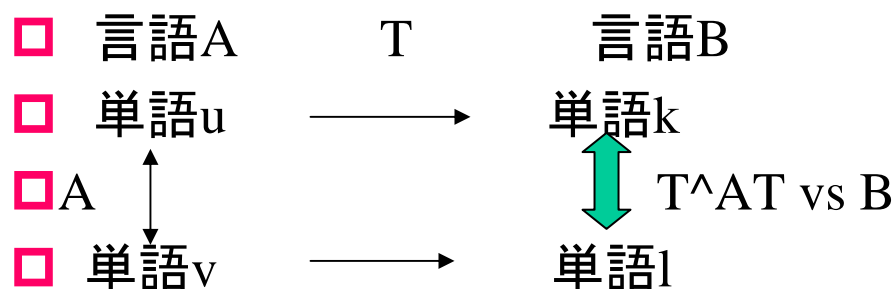
- Fung98 : Rappと同様のアイデア
- 既存の中英対訳辞書を利用して中国語未知語 W の英訳語を求める。
- W の出現した文と同じ文における共起する既存の辞書の単語 W_d をひとつの次元とみなすベクトル空間法
- W と同じ文における W_d の出現数を tf
- W と同じ文を document と見なした場合の idf (W と同じ文に出現する総単語数で正規化

より本格的に文脈を利用する方法 Fung

- $tf \cdot idf$ を直接 W_d の重みにせずに、 W_d の英訳語の既存の辞書での重要度順位 k で割る。
- cosine, Dice 係数などで類似度の高い英単語を W_d の対訳とする。
- 第1位の訳語の精度は30%
- 第20位までの訳語だと76%の正解をカバー

より本格的に文脈を利用する方法 Tanaka K

□ Tanaka96 類似性はMI、辞書はEDICT



□ $T_{ij} = p(\text{言語Bの単語}j \mid \text{言語Aの単語}i)$

□ $|T^AT - B|$ を最急下降法で最小化し、求めたTが対訳マトリクス。

□ 378語の対訳実験で、曖昧さ解消成功率 = 82%。一方、ふるい落とした対訳の85%は誤訳語

Context Heterogeneity (Fung95WVLC)

- 扱う分野が同じなら parallel でなくてよい。
- 単語 trigram により文脈を作る。つまり
- 単語L 単語0 単語R
 - ここで、単語0 の context heterogeneity は
 - 単語Lの異なり数 = a, 単語Rの異なり数 = b
 - 単語0 の出現頻度 = c のとき
 - Left-heterogeneity = a/c, right-heterogeneity = b/c
 - このアイデアは語基の接続数(中川)に似ている
- 二つの言語の単語 w_1, w_2 の context heterogeneity x_1, y_1 (for w_1), x_2, y_2 (for w_2) の距離は
$$\varepsilon = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

Context heterogeneity の評価

- ε は統計量なのでとにかくコーパスが大きくなると信頼性の高いデータが得られない。
 - 大規模なデータ量は、non-aligned corpus の場合、必然かも。
- Context heterogeneity のような局所的な文脈だけで対訳を絞りこめるのか疑問
 - 評価結果が論文に書いてない
- 既存の辞書を使っていない。ゼロから対訳を得ようというのはかなり無理では。

既存の対訳辞書と語基の接続数との併用 nakagawa2000 LREC WTRC, 2001 NLPRS

- 各言語のコーパスから用語候補を抽出
- 言語毎に、抽出された語基を接続数の大きさの順に並べる
- 日本語の語基 w_j の英語の対訳候補 w_{e1}, w_{e2}, \dots を既存の対訳辞書(EDICT)から求める。当然、曖昧
- w_j の対訳として、 w_j の接続数の順位に近い接続数の順位を持つ w_{ei} ($i=1, 2, \dots$)を選ぶ

抽出、順位つけされた
用語集合(英語)

抽出、順位つけされた
用語集合(日本語)

1位

1位

N位: memory system

N-2位: メモリシステム

N+3位: 記憶システム

○

×

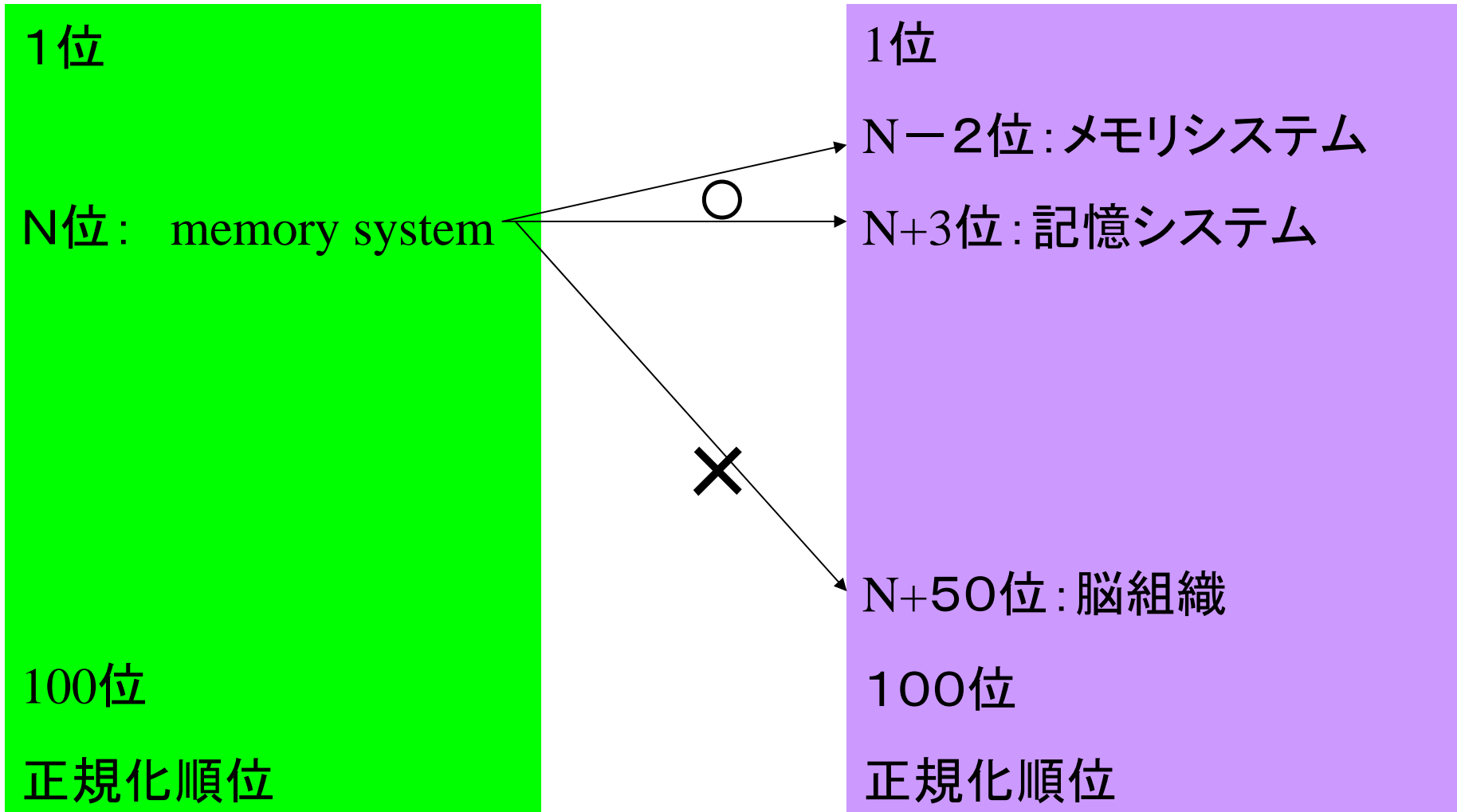
N+50位: 脳組織

100位

100位

正規化順位

正規化順位



Distance

- $\text{distance}(X_e, X_j) = | \text{rank}(X_e) - \text{rank}(X_j) |$
- If $\text{distance}(X_e, X_j)$ is small, X_e is the translation of X_j .
- $\text{distance}(X_{e1}, X_j) < \text{distance}(X_{e2}, X_j) < \dots$
then X_{e1} is most likely translation of X_j

Example of distance

memory system	メモリシステム	0.051493
	記憶システム	0.956459
	メモリ方式	1.234347
	記憶方式	3.809609
	脳組織	63.498688

日英, 英日の情報通信関連専門用語の対訳
で60%から80%の精度

より本格的に文脈を利用する方法 まとめ

- Rapp、Fung, Tanaka の方法は
 1. タネになる小さな対訳辞書によって
 2. 対訳を求めたい語基の文脈(に現れる既知語)をベクトル表現し
 3. 相手側テキストで、このベクトルによく一致するものを探し、その未知語を訳語にする
- この方法は、行き着くところまで行っている (local minimum) ので、新規なアイデアが必要??