

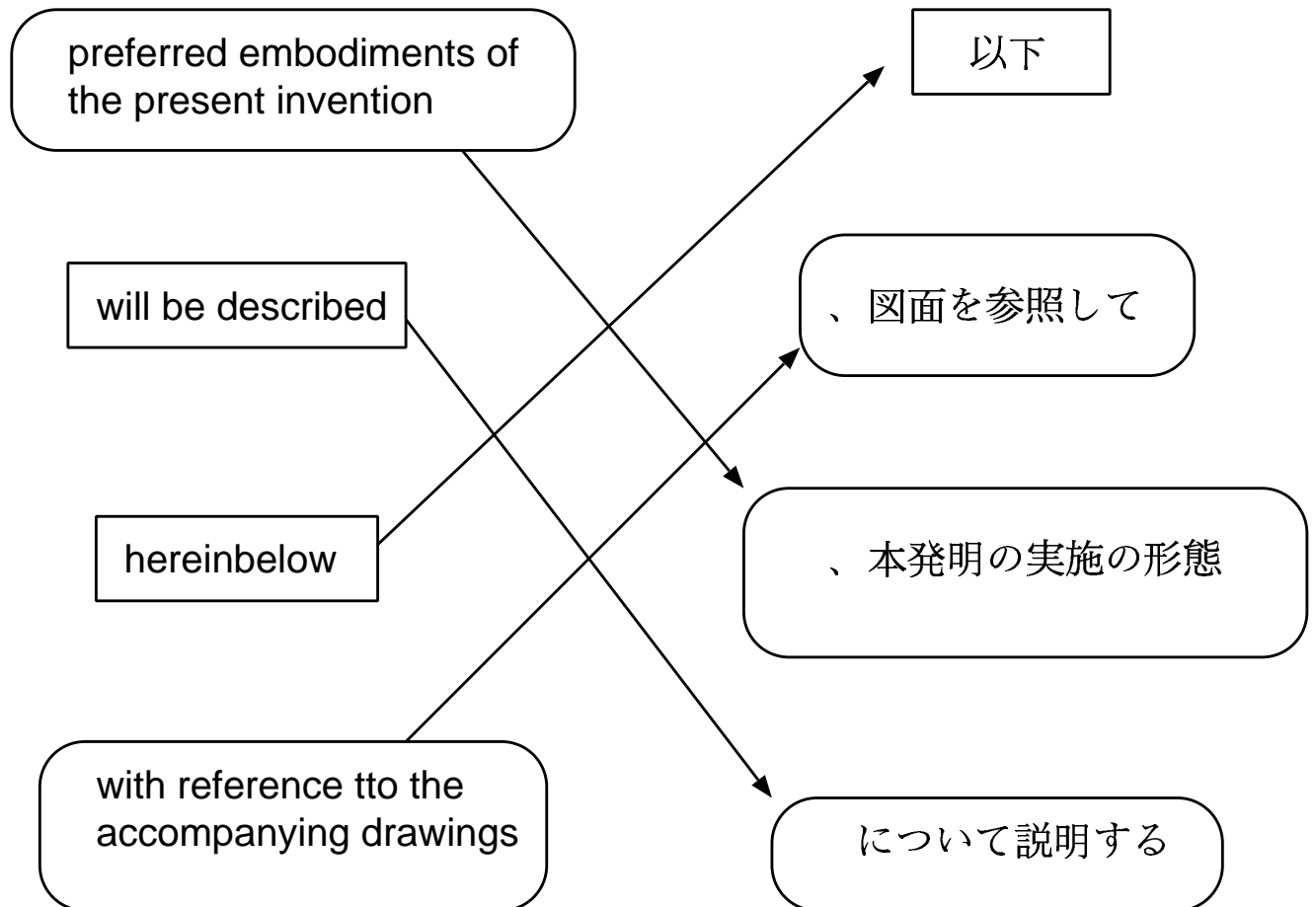
17. 句に基づく統計的機械翻訳 (SMT) における素性

内山将夫@NICT
mutiyama@nict.go.jp

句に基づく統計的機械翻訳 (SMT) における素性

- 句単位の翻訳の例
- 句単位の翻訳を構成する基本要素
- 句単位の翻訳の素性

句単位の英日翻訳の例



- 英語の文を句にわけ
(e.g., will be described)
- 各句を翻訳する
(e.g., について説明する)
- 句を並べ替える

句単位の翻訳の基本要素

- 英語の文 e を句にわけ

$$e = \bar{e}_1 \bar{e}_2 \dots \bar{e}_l$$

- 各句 \bar{e}_i を日本語の句 \bar{n}_i に翻訳する

$$\mathbf{n} = \bar{n}_1 \bar{n}_2 \dots \bar{n}_l$$

- \bar{n}_i を並べかえる .

句単位の翻訳のための素性

原言語を f , 対象言語を e としたとき ,

$$\hat{e} = \arg \max_e \sum_i \lambda_i h_i(e, f)$$

なる \hat{e} を探したい . そのために , e と f の組をスコア付けするために , 素性 $h_i(e, f)$ を利用する .

素性の例

- 言語モデル
- 句単位の翻訳モデル
- 句を構成する単語に基づく翻訳モデル
- 単語ペナルティ
- 句ペナルティ
- 語順ペナルティ

これらの素性の値は , e, f より $h_i(e, f)$ として計算される . 各素性の重みは , 自動推定する (後述) .

言語モデル

$$\begin{aligned}h(\mathbf{e}, \mathbf{f}) &= \log P(\mathbf{e}) \\&= \log P(e_1, e_2, \dots, e_l) \\&= \log \prod_i P(e_i | e_{i-2}, e_{i-1})\end{aligned}$$

3-gram 言語モデルにより対象言語 \mathbf{e} の生成確率の対数を素性とする．これが大きい翻訳文 \mathbf{e} は，対象言語としてのつくりが良いと考えられる．

ところで，

3-gram 言語モデルは，前の2単語しか見ていないので，とても貧弱なモデルに見えるが，これを越えるモデルはあまりない．良い言語モデルを作ることができれば，機械翻訳に与えるインパクトは大きい．

単語ペナルティと句ペナルティ

単語ペナルティ

$$h(e, f) = e \text{ 中の単語数}$$

句ペナルティ

$$h(e, f) = e \text{ 中の句数}$$

これらの重みが正のときには，単語や句がたくさんある e が優先される．負のときには，単語や句は少ない方がよい．

単語ペナルティや句ペナルティは，適切な長さの訳文を出力するために有用である．

句単位の翻訳モデル

$$\begin{aligned}h(\mathbf{e}, \mathbf{f}) &= \log P(\mathbf{e}|\mathbf{f}) \\&= \log \Pi P(\bar{e}|\bar{f}) \\&= \sum \log P(\bar{e}|\bar{f})\end{aligned}$$

これは $\mathbf{f} \rightarrow \mathbf{e}$ の翻訳モデルだが , 逆方向も同様に ,

$$h(\mathbf{e}, \mathbf{f}) = \sum \log P(\bar{f}|\bar{e})$$

ただし ,

$$\begin{aligned}P(\bar{f}|\bar{e}) &= \frac{\text{count}(\bar{f}, \bar{e})}{\sum_{\bar{f}'} \text{count}(\bar{f}', \bar{e})} \\ \text{count}(\bar{f}, \bar{e}) &= \text{句}\bar{f}\text{と句}\bar{e}\text{が対応している回数} \quad (1)\end{aligned}$$

$P(\bar{f}|\bar{e})$ の問題点

$\text{count}(\bar{f}, \bar{e})$ が小さいときに値が信頼できない

→ スムージング

単語に基づく句対応の重み (双方向)

$$h(\mathbf{e}, \mathbf{f}) = \sum \log \text{lex}(\bar{f}|\bar{e}) \quad (2)$$

$$\text{lex}(\bar{f}|\bar{e}) = \max_{\mathbf{a}} P_w(\bar{f}|\bar{e}, \mathbf{a})$$

$$P_w(\bar{f}|\bar{e}, \mathbf{a}) = \prod_{i=1}^n E_w(f_i|\bar{e}, \mathbf{a})$$

$$E_w(f_i|\bar{e}, \mathbf{a}) = \frac{1}{|\{j|(i, j) \in \mathbf{a}\}|} \sum_{(i,j) \in \mathbf{a}} w(f_i|e_j)$$

$$w(f_i|e_j) = \text{単語対応の確率} = \frac{\text{count}(f_i, e_j)}{\sum_{f'} \text{count}(f', e_j)}$$

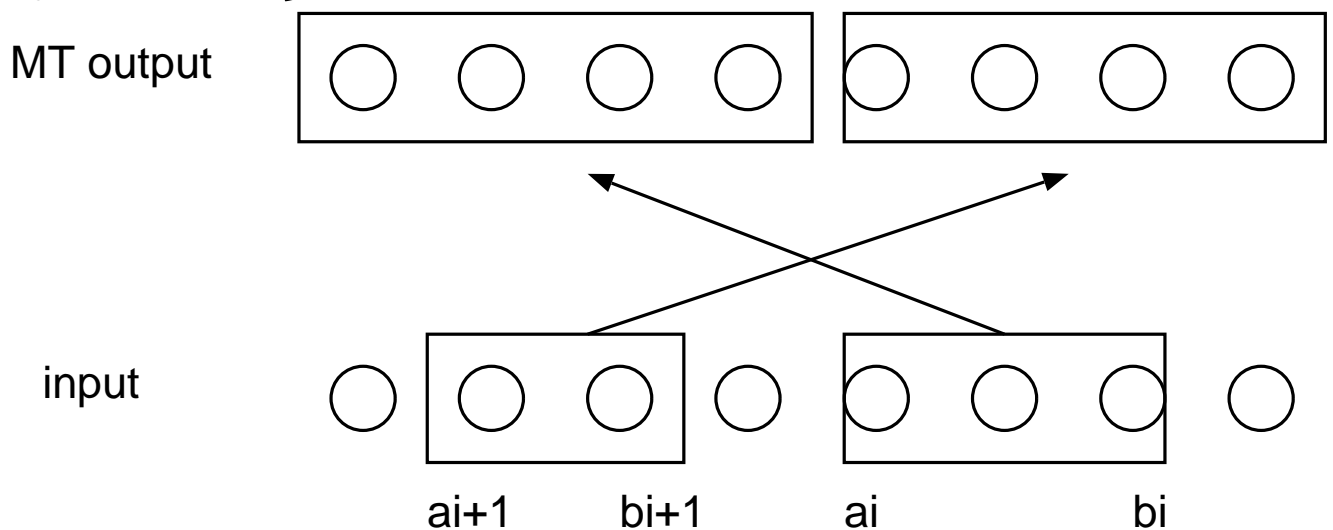
$$E_w(f_i|\bar{e}, \mathbf{a}) = w(f_i|e_j) \text{ の } e_j \text{ に関する平均値}$$

$$P_w(\bar{f}|\bar{e}, \mathbf{a}) = E_w(f_i|\bar{e}, \mathbf{a}) \text{ の積 } f_i \text{ の条件付き独立を仮定}$$

$$\text{lex}(\bar{f}|\bar{e}) = P_w(\bar{f}|\bar{e}, \mathbf{a}) \text{ が最大の } \mathbf{a} \text{ に関する確率を採用}$$

語順

語順については，今，盛んに研究されている．単純なものとしては，原言語の語順と異なるものにペナルティを与えるというものがある．



のとき，

$$|b_i - a_{i+1} + 1| = |\overline{f_i} \text{ に対応する } \overline{e_i} \text{ の最後の単語位置} \\ - \overline{f_{i+1}} \text{ に対応する } \overline{e_{i+1}} \text{ の最初の単語位置} + 1|$$

として，連続するフレーズが連続するときのペナルティを 0 とし，そうでないときに，離れ具合に応じてペナルティをかける．

まとめ

- 対数線形モデルを利用することにより，様々な素性をスコア付けに利用できる
- 各素性は，それぞれ別個に改良できる

特に，

- 言語モデルの改良
- 翻訳モデルの改良
 - － カバー率の高い句表をつくる
 - － 良いスコア付けをする
- 語順の制約

が重要である．