

品詞bigram(Bi-POS)モデルを用いたコスト最小法による形態素解析は、以下に示されるビタービアルゴリズムによって効率的に計算される。

記号の定義：

長さ  $L$  の入力文字列： $S_L = c_1, c_2, \dots, c_L$

出力単語列： $W = w_1, w_2, \dots, w_{|W|}$  (すべての  $w_i$  は辞書に入っている単語であること)

出力品詞列： $T = t_1, t_2, \dots, t_{|T|}$  (ちなみに  $|W| = |T|$  である。)

単語  $w_i$  の品詞： $t(w_i)$  単語列  $W$  の品詞列： $t(W) = t(w_1), t(w_2), \dots, t(w_{|W|})$

品詞  $t_{i-1}$  の直後に品詞  $t_i$  が続くコスト(品詞バイグラムのコスト)： $con(t_{i-1}, t_i)$

品詞  $t_i$  の場合に単語  $w_i$  が出現する単語コスト： $word(w_i, t_i)$

解くべき問題(品詞bigramモデルを用いたコスト最小法による形態素解析の形式的定義)：

$$(\hat{W}, \hat{T}) = \arg \min_{W, T} \sum_{i=1}^{|W|} \{con(t_{i-1}, t_i) + word(w_i, t_i)\} \quad (\text{ここで } t_0 \text{ は「文頭」という特殊な品詞})$$

上記問題のビタービアルゴリズムによる計算方法：

(0) 記憶領域 (変数)

$\#w(x)$  : 文字位置  $x$  で終わる単語の数 (簡単のためすべての位置で1以上とする)

$w_x(j)$  : 文字位置  $x$  で終わる単語のうち  $j$  番目の単語

$t_x(j)$  :  $w_x(j)$  の品詞

$\delta_x(j)$  :  $w_x(j)$  を採用した場合の先頭から  $w_x(j)$  までの最小コスト

$B_x(j)$  : バックポインタ ( $\delta_x(j)$  を達成するために選ばれた左隣の単語)

(1) 初期化：

・文頭部分の初期化： $\#w(0) = 1, t_0(1) = \text{文頭}, \delta_0(1) = 0$

・辞書を用いて単語ラティスを作る ( $\#w(x), w_x(i), t_x(i)$  の設定)

(2) 繰り返し： $x = 1: L$

for  $i = 1: \#w(x)$

$y = x - |w_x(i)|$

$$\delta_x(i) = \min_{j=1: \#w(y)} \{ \delta_y(j) + con(t_y(j), t_x(i)) + word(w_x(i), t_x(i)) \}$$

$$B_x(i) = \arg \min_{j=1: \#w(y)} \{ \delta_y(j) + con(t_y(j), t_x(i)) + word(w_x(i), t_x(i)) \}.$$

(3) 終了：

$$bp(0) = \arg \min_{i=1: \#w(L)} \delta_L(i)$$

$$w'_1 = w_L(bp(0)), \text{left}(1) = L - |w'_1|, bp(1) = B_L(bp(0))$$

(4) バックトラック

$$w'_k = w_{\text{left}(k-1)}(bp(k-1))$$

$$\text{left}(k) = \text{left}(k-1) - |w'_k| \quad k = 2, 3, \dots, M \quad (\text{left}(M) = 0)$$

$$bp(k) = B_{\text{left}(k-1)}(bp(k-1))$$

最後に  $w'_k$  を逆順にして出力： $\hat{W} = \hat{w}_1, \hat{w}_2, \dots, \hat{w}_M$  ( $\hat{w}_i = w'_{M-i+1}$ )

(品詞列は  $w \rightarrow t$  と置き換えればよい)

もっと簡潔に書けるが、Tri-POSモデルのときにこの方が考えやすい。