

ЕМ - алгоритм для задачи перевода

Обозначим: за R число предложений в корпусе;

n_i - длина i -го предложения

m_i - длина i -го перевода

$S_i = (S_{i1}, \dots, S_{in_i})$ - исходное i -е предложение

$T_i = (T_{i1}, \dots, T_{im_i})$ - его перевод

$A_i = (a_{i1}, \dots, a_{im_i})$ - выравнивание каждого слова в i -м предложении (скрытые переменные)

Е шаг:

$$p(a_i = z | t_i, S_i) = \frac{P(a_i = z, t_i | S_i)}{P(t_i | S_i)} =$$

$$= \frac{P(a_i = z, t_i | S_i)}{\sum_{m=1}^M P(a_i = m, t_i | S_i)} = \frac{\theta(t_i | S_i)}{\sum_{m=1}^M \theta(t_i | S_{im})} = q_i(z)$$

М шаг:

$$L(q_i; z) = \sum_{i=1}^R \sum_{j=1}^m q_i(z) \cdot \log \frac{P(a_i = z, t_i | S)}{q_i(z)} =$$

$$= \sum_{i=1}^R \sum_{j=1}^m q_i(z) \cdot \log \frac{\theta(t_i | S_k)}{n q_i(z)} =$$

$$= \sum_{i=1}^R \sum_{j=1}^m q_i(z) \cdot \log \theta(t_i | S_k) - \sum_{i=1}^R \sum_{j=1}^m q_i(z) \cdot \log n q_i(z) \rightarrow \max_{\theta}$$

$$\Rightarrow \text{задача сводится к } \sum_{i=1}^R \sum_{j=1}^m q_i(z) \log \theta(t_i | S_k) \rightarrow \max_{\theta}$$

Лагранжиан:

$$L = \sum_{i=1}^R \sum_{j=1}^m q_i(z) \cdot \log \theta(t_i | S_k) - \sum_{x \in X} \lambda_x \left(\sum_{y \in Y} \theta(y | x) - 1 \right)$$

$$\frac{\partial L}{\partial z} = \sum_{i=1}^R \sum_{j=1}^m q_i(z) \frac{\partial \log \theta(t_i | S_k)}{\partial \theta(y | x)} - \frac{\partial \sum_{x \in X} \left(\sum_{y \in Y} \theta(y | x) - 1 \right)}{\partial \theta(y | x)} =$$

$$= \sum_{i=1}^R \sum_{j=1}^m q_i(z) \cdot \frac{[t_i = y] [S_k = x]}{\theta(y | x)} - \lambda_x = 0$$

Нужно найти λ_x . Воспользуемся св-м бер-ми:

$$\begin{aligned}
 \sum_{y \in Y} \theta(y|x) &= 1 = \sum_{y \in Y} \frac{1}{\lambda_x} \sum_{z=1}^n \sum_{i=1}^m q_i(z) [t_i = y] [S_k = x] = \\
 &= \frac{1}{\lambda_x} \sum_{z=1}^n \sum_{i=1}^m q_i(z) [S_k = x] \underbrace{\sum_{y \in Y} [t_i = y]}_{=1} \\
 \Rightarrow \theta^*(y|x) &= \frac{\sum_{z=1}^n \sum_{i=1}^m q_i(z) [t_i = y] [S_k = x]}{\sum_{z=1}^n \sum_{i=1}^m q_i(z) [S_k = x]}
 \end{aligned}$$