

ЕМ-алгоритм для задачи машинного перевода

Пусть $S = (s_1, \dots, s_n)$ исходное предложение, $T = (t_1, \dots, t_m)$ — его перевод. В роли латентных переменных будут выступать выравнивания $A = (a_1, \dots, a_m)$ каждого слова в целевом предложении, причём $a_i \in \{1, \dots, n\}$ (считаем, что каждое слово в t является переводом какого-то слова из s). Параметрами модели является матрица условных вероятностей перевода: каждый её элемент $\theta(y|x) = p(y|x)$ отражает вероятность того, что переводом слова x с исходного языка на целевой является слово y (нормировка, соответственно, совершается по словарю целевого языка). Правдоподобие латентных переменных и предложения на целевом языке в этой модели записывается так:

$$p(A, T|S) = \prod_{i=1}^m p(a_i) p(t_i|a_i, S) = \prod_{i=1}^m \frac{1}{n} \theta(t_i|s_{a_i}).$$

Обозначим за V^s - словарь исходного языка, а за V^t словарь целевого.

ЕМ-алгоритм

В ЕМ-алгорите мы максимизируем $p(T|S)$. Причём по всем парам предложений из нашего корпуса данных. Пусть таких пар всего R штук и длины в каждой r -ой паре равны n_r и m_r для исходного и целевого языка соответственно. Так как \log - монотонная функция, то исходная задача эквивалентна максимизации лог-правдоподобия. Для начала получим выражение для подсчёта нижней оценки правдоподобия. Для некоторого распределения $q(A)$ на скрытых переменных верно:

$$\log\left(\prod_{r=1}^R p(T^r|S^r)\right) = \sum_{r=1}^R \int q(A^r) \log p(T^r|S^r) dA^r =$$

$$\begin{aligned}
&= \sum_{r=1}^R \int q(A^r) \log \frac{p(A^r, T^r | S^r)}{p(A^r | T^r, S^r)} dA^r = \sum_{r=1}^R \int q(A^r) \log \frac{p(A^r, T^r | S^r) q(A^r)}{p(A^r | T^r, S^r) q(A^r)} dA^r = \\
&= \sum_{r=1}^R \int q(A^r) \log \frac{p(A^r, T^r | S^r)}{q(A^r)} dA^r + \sum_{r=1}^R \int q(A^r) \log \frac{q(A^r)}{p(A^r | T^r, S^r)} dA^r = \\
&= \mathcal{L}(q, \theta) + KL(q || p)
\end{aligned}$$

Будем считать, что $q(A^r) = p(A^r | T^r, S^r)$. Тогда $KL(q || p) = 0$. Также учтём, что в данной задаче $p(A^r | T^r, S^r)$ имеет дискретное распределение, поэтому интегралы перейдут в суммы. Таким образом получим, что:

$$\mathcal{L} = \sum_{r=1}^R \sum_{i=1}^{m_r} \sum_{j=1}^{n_r} p(a_i = j | T^r, S^r) \log \frac{\theta(t_i^r | s_{a_i^r}^r)}{n_r p(a_i = j | T^r, S^r)}$$

Е-шаг

Мы хотим получить апостериорное распределение латентных переменных $p(A, |T, S)$ (для каждого предложения в нашем корпусе данных). Тогда по формуле Байеса получим:

$$p(a_i = j | T, S) = \frac{p(T | a_i = j, S) p(a_i = j | S)}{p(T | S)} = \frac{p(T, a_i = j | S)}{p(T | S)}$$

Таким образом, мы получили, что апостериорная вероятность представляет собой сумму вероятностей всех выравниваний, содержащих связь между t_i и s_j ($a_i = j$), деленную на сумму вероятностей

всех возможных выравниваний:

$$\begin{aligned}
\frac{p(T, a_i = j | S)}{p(T | S)} &= \frac{\sum_{A: a_i = j} p(A, T | S)}{\sum_A p(A, T | S)} = \frac{\sum_{A: a_i = j} \prod_{k=1}^m \frac{1}{n} \theta(t_k | s_{a_k})}{\sum_A \prod_{k=1}^m \frac{1}{n} \theta(t_k | s_{a_k})} = \\
&= \frac{\theta(t_i | s_j) \sum_{a_1=1}^n \dots \sum_{a_{i-1}=1}^n \sum_{a_{i+1}=1}^n \dots \sum_{a_m=1}^n \prod_{k=1, k \neq i}^m \theta(t_k | s_{a_k})}{\sum_{a_1=1}^n \dots \sum_{a_m=1}^n \prod_{k=1}^m \theta(t_k | s_{a_k})} =^* \\
&=^* \frac{\theta(t_i | s_j) \prod_{k=1, k \neq i}^m \sum_{a_1=1}^n \dots \sum_{a_{i-1}=1}^n \sum_{a_{i+1}=1}^n \dots \sum_{a_m=1}^n \theta(t_k | s_{a_k})}{\prod_{k=1}^m \sum_{a_1=1}^n \dots \sum_{a_m=1}^n \theta(t_k | s_{a_k})} = \frac{\theta(t_i | s_j)}{\sum_{a_i=1}^n \theta(t_i | s_{a_i})}
\end{aligned}$$

Таким образом, получим апостериорные вероятности для скрытых переменных для каждой пары предложений:

$$p(a_i = j | T, S) = \frac{\theta(t_i | s_j)}{\sum_{k=1}^n \theta(t_i | s_k)}$$

М-шаг

На М-шаге мы максимизируем правдоподобие $p(A, T | S)$. Причём по всем парам предложений из нашего корпуса данных. Тогда получим следующую оптимизационную задачу:

$$\mathbb{E}_A \log \left(\prod_{r=1}^R p(A^r, T^r | S^r) \right) \rightarrow \max_{\theta}$$

Выведем формулы для обновления наших параметров θ :

$$\begin{aligned} \mathbb{E}_A \log\left(\prod_{r=1}^R p(A^r, T^r | S^r)\right) &= \mathbb{E}_A \log\left(\prod_{r=1}^R \prod_{i=1}^{m_r} \frac{1}{n_r} \theta(t_i^r | s_{a_i^r}^r)\right) = \mathbb{E}_A \sum_{r=1}^R \sum_{i=1}^{m_r} \left[\log \frac{1}{n_r} + \right. \\ &\quad \left. + \log \theta(t_i^r | s_{a_i^r}^r) \right] = - \sum_{r=1}^R m_r \log n_r + \sum_{r=1}^R \sum_{i=1}^{m_r} \mathbb{E}_A \log \theta(t_i^r | s_{a_i^r}^r) \quad (=) \end{aligned}$$

Обозначим $C = - \sum_{r=1}^R m_r \log n_r$. Тогда получим:

$$(\quad) \quad C + \sum_{r=1}^R \sum_{i=1}^{m_r} \sum_{j=1}^{n_r} p(a_i = j | T^r, S^r) \log \theta(t_i^r | s_{a_i^r}^r)$$

Тогда получим следующую оптимизационную задачу:

$$\begin{cases} C + \sum_{x,y} count(x, y) \log \theta(y|x) \rightarrow \max_{\theta} \\ \sum_y \theta(y|x) = 1, \quad \forall x \in V^s \end{cases}$$

Где $count(x, y) = \sum_{r=1}^R \sum_{i=1}^{m_r} \sum_{j=1}^{n_r} \mathbb{1}\{t_i^r = x\} \mathbb{1}\{s_j^r = y\} p(a_i = j | T^r, S^r)$,

$\sum_{x,y}$ - сумма по всевозможным парам (x, y) , где $x \in V^s$ - произволь-

ное слово из словаря исходного языка, а $y \in V^t$ из целевого.

Будем искать максимум через лагранжиан системы. Запишем его:

$$L = C + \sum_{x,y} count(x, y) \log \theta(y|x) - \sum_x \lambda_x \left(\sum_y \theta(y|x) - 1 \right)$$

Продифференцируем и приравняем к 0 частные производные:

$$\frac{\partial L}{\partial \theta(y|x)} = \frac{count(x, y)}{\theta(y|x)} - \lambda_x = 0 \quad \Rightarrow \quad \theta(y|x) = \frac{count(x, y)}{\lambda_x}$$

$$\frac{\partial L}{\partial \lambda_x} = \sum_y \theta(y|x) - 1 = 0 \Rightarrow \sum_y \theta(y|x) = 1$$

Просуммировав первое уравнение по всем $y \in V^t$, учитывая второе, получим:

$$1 = \sum_y \theta(y|x) = \frac{\sum_y count(x, y)}{\lambda_x} \Rightarrow \lambda_x = \sum_y count(x, y)$$

Таким образом, получим аналитический максимум по параметрам нашей модели:

$$\theta(y|x) = \frac{count(x, y)}{\sum_y count(x, y)}$$

References

Word Alignment and the Expectation-Maximization Algorithm. Adam Lopez. Johns Hopkins University.

Appendix

* Докажем, что знак суммы и произведения можно поменять местами в знаменателе. Для числителя всё аналогично.

$$\begin{aligned} \sum_{a_1=1}^n \dots \sum_{a_m=1}^n \prod_{k=1}^m \theta(t_k|s_{a_k}) &= \theta(t_1|s_1) \dots \theta(t_m|s_1) + \theta(t_1|s_1) \dots \theta(t_{m-1}|s_1) \theta(t_m|s_2) + \\ &\dots + \theta(t_1|s_n) \dots \theta(t_m|s_n) = [\theta(t_1|s_1) + \dots + \theta(t_1|s_n)] \left[\sum_{a_2=1}^n \dots \sum_{a_m=1}^n \prod_{k=2}^m \theta(t_k|s_{a_k}) \right] = \\ &\left[\sum_{a_1=1}^n \theta(t_1|s_{a_1}) \right] \cdot \left[\sum_{a_2=1}^n \dots \sum_{a_m=1}^n \prod_{k=2}^m \theta(t_k|s_{a_k}) \right] = \{ \text{Аналогично проделывая} \\ &\text{такие шаги по всем } t_i, \text{ получим} \} = \left[\sum_{a_1=1}^n \theta(t_1|s_{a_1}) \right] \dots \left[\sum_{a_m=1}^n \theta(t_m|s_{a_m}) \right] = \\ &\prod_{k=1}^m \sum_{a_1=1}^n \dots \sum_{a_m=1}^n \theta(t_k|s_{a_k}) \end{aligned}$$