

# Семинары по EM-алгоритму

Евгений Соколов  
sokolov.evg@gmail.com

11 марта 2015 г.

## 1 EM-алгоритм

### §1.1 Смеси распределений

Говорят, что распределение  $p(x)$  является *смесью распределений*, если его плотность имеет вид

$$p(x) = \sum_{k=1}^K \pi_k p_k(x), \quad \sum_{k=1}^K \pi_k = 1, \quad \pi_k \geq 0, \quad (1.1)$$

где  $p_k(x)$  — распределения компонент смеси,  $\pi_k$  — априорные вероятности компонент,  $K$  — число компонент. Будем считать, что распределения компонент смеси принадлежат некоторому параметрическому семейству:  $p_k(x) = \varphi(x | \theta_k)$ .

Рассмотрим следующий эксперимент: сначала из дискретного распределения  $\{\pi_1, \dots, \pi_K\}$  выбирается номер  $k$ , а затем из распределения  $\varphi(x | \theta_k)$  выбирается значение  $x$ . Покажем, что распределение переменной  $x$  будет представлять собой смесь вида (1.1).

Введем *скрытую переменную*  $z$ , отвечающую за выбор компоненты смеси. Пусть она представляет собой  $K$ -мерный бинарный случайный вектор, ровно одна компонента которого равна единице:

$$z \in \{0, 1\}^K, \quad \sum_{i=1}^K z_k = 1.$$

Вероятность того, что единице будет равна  $k$ -я компонента, равна  $\pi_k$ :

$$p(z_k = 1) = \pi_k.$$

Запишем распределение сразу всего вектора:

$$p(z) = \prod_{k=1}^K \pi_k^{z_k}.$$

Если номер компоненты смеси известен, то случайная величина  $x$  имеет распределение  $\varphi(x | \theta_k)$ :

$$p(x | z_k = 1) = \varphi(x | \theta_k),$$

или, что то же самое,

$$p(x | z) = \prod_{k=1}^K [\varphi(x | \theta_k)]^{z_k}.$$

Запишем совместное распределение переменных  $x$  и  $z$ :

$$p(x, z) = p(z)p(x | z) = \prod_{k=1}^K [\pi_k \varphi(x | \theta_k)]^{z_k}.$$

Чтобы найти распределение переменной  $x$ , нужно избавиться от скрытой переменной:

$$p(x) = \sum_z p(x, z).$$

Суммирование здесь ведется по всем возможным значениям  $z$ , то есть по всем  $K$ -мерным бинарным векторам с одной единицей:

$$p(x) = \sum_z p(x, z) = \sum_{k=1}^K \pi_k \varphi(x | \theta_k).$$

Мы получили, что распределение переменной  $x$  в описанном эксперименте представляет собой смесь  $K$  компонент.

## §1.2 Модели со скрытыми переменными

Рассмотрим вероятностную модель с наблюдаемыми переменными  $X$  и параметрами  $\Theta$ , для которой задано правдоподобие  $\log p(X | \Theta)$ . Предположим, что в модели также существуют *скрытые переменные*  $Z$ , описывающие ее внутреннее состояние. Тогда правдоподобие  $\log p(X | \Theta)$  называется *неполным*, а правдоподобие  $\log p(X, Z | \Theta)$  — *полным*. Они связаны соотношением

$$\log p(X | \Theta) = \log \left\{ \sum_Z p(X, Z | \Theta) \right\}.$$

Как правило, знание скрытых переменных существенно упрощает правдоподобие и позволяет достаточно просто оценить параметры  $\Theta$ .

Рассмотрим пример со смесями распределений. В качестве наблюдаемых переменных здесь выступает выборка  $X^\ell = \{x_1, \dots, x_\ell\}$ , в качестве скрытых переменных — номера компонент, из которых сгенерированы объекты  $Z = \{z_1, \dots, z_\ell\}$  (здесь каждый из  $z_i$  является  $K$ -мерным вектором), в качестве параметров — априорные вероятности и параметры компонент  $\Theta = (\pi_1, \dots, \pi_K, \theta_1, \dots, \theta_K)$ . Неполное правдоподобие имеет вид

$$\log p(X | \Theta) = \sum_{i=1}^{\ell} \log \left\{ \sum_{k=1}^K \pi_k p(x_i | \theta_k) \right\}.$$

Правдоподобие здесь имеет вид «логарифм суммы». Если приравнять нулю его градиент, то получатся сложные уравнения, не имеющие аналитического решения. Данное правдоподобие сложно вычислять, оно не является вогнутым и имеет много локальных экстремумов, поэтому применение итерационных методов для его непосредственной максимизации приводит к медленной сходимости.

Рассмотрим теперь полное правдоподобие:

$$\log p(X, Z | \Theta) = \sum_{i=1}^{\ell} \sum_{k=1}^K z_{ik} \left\{ \log \pi_k + \log \varphi(x_i | \theta_k) \right\}.$$

Оно имеет вид «сумма логарифмов», и позволяет аналитически найти оценки максимального правдоподобия на параметры  $\Theta$  при известных переменных  $X$  и  $Z$ . Проблема же заключается в том, что нам не известны скрытые переменные  $Z$ , поэтому их необходимо оценивать одновременно с параметрами, что никак не легче максимизации неполного правдоподобия. Решение данной проблемы предлагается в *EM-алгоритме*.

### §1.3 EM-алгоритм

EM-алгоритм решает задачу максимизации полного правдоподобия путем попеременной оптимизации по параметрам и по скрытым переменным.

Опишем сначала «наивный» способ оптимизации. Зафиксируем некоторое начальное приближение для параметров  $\Theta^{\text{old}}$ . При известных наблюдаемых переменных  $X$  и параметрах  $\Theta^{\text{old}}$  мы можем оценить скрытые переменные, найдя их наиболее правдоподобные значения:

$$Z^* = \arg \max_Z p(Z | X, \Theta^{\text{old}}) = \arg \max_Z p(X, Z | \Theta^{\text{old}}).$$

Зная скрытые переменные, мы можем теперь найти следующее приближение для параметров:

$$\Theta^{\text{new}} = \arg \max_{\Theta} p(X, Z^* | \Theta).$$

Повторяя итерации до сходимости, мы получим некоторый итоговый вектор параметров  $\Theta^*$ . Данная процедура, однако, не гарантирует сходимости и может выдать неоптимальное решение.

Гораздо лучшие результаты можно получить, воспользовавшись байесовским подходом. Как и прежде, зафиксируем вектор параметров  $\Theta^{\text{old}}$ , но вместо точечной оценки вычислим апостериорное распределение на скрытых переменных  $p(Z | X, \Theta^{\text{old}})$ . В этом заключается *E-шаг* EM-алгоритма.

Усредним логарифм полного правдоподобия по всем возможным значениям скрытых переменных  $Z$  с весами, равными апостериорным вероятностям этих переменных  $p(Z | X, \Theta^{\text{old}})$ :

$$Q(\Theta, \Theta^{\text{old}}) = \mathbb{E}_{Z \sim p(Z | X, \Theta^{\text{old}})} \log p(X, Z | \Theta) = \sum_Z p(Z | X, \Theta^{\text{old}}) \log p(X, Z | \Theta).$$

Формально говоря, мы нашли матожидание логарифма полного правдоподобия по апостериорному распределению на скрытых переменных. На *M-шаге* новый вектор параметров находится как максимизатор данного матожидания:

$$\Theta^{\text{new}} = \arg \max_{\Theta} Q(\Theta, \Theta^{\text{old}}) = \arg \max_{\Theta} \sum_Z p(Z | X, \Theta^{\text{old}}) \log p(X, Z | \Theta).$$

Далее мы рассмотрим условия сходимости описанного итерационного процесса и увидим, что при достаточно общих условиях любая из его сходящихся подпоследовательностей сойдется к стационарной точке неполного правдоподобия.

## §1.4 Вывод формул для задачи разделения смеси

На лекциях EM-алгоритм для разделения смеси гауссиан был выведен следующим образом: неполное правдоподобие продифференцировано по параметрам и приравнено к нулю, в полученных уравнениях сгруппированы и зафиксированы выражения, имеющие смысл вероятности получения объектов из определенных компонент смеси; после этого из уравнений стало возможно аналитически выразить параметры. В данном разделе мы выведем формулы, исходя из описанной общей схемы EM-алгоритма, и увидим, что они совпадают с формулами из лекций.

### 1.4.1 Разделение смеси нормальных распределения

Запишем полное правдоподобие для смеси нормальных распределений:

$$\log p(X, Z | \Theta) = \sum_{i=1}^{\ell} \sum_{k=1}^K z_{ik} \left\{ \log \pi_k + \log \varphi(x_i | \theta_k) \right\}.$$

На E-шаге вычисляется апостериорное распределение на скрытых переменных:

$$p(Z | X, \Theta^{\text{old}}) = \frac{p(X, Z | \Theta^{\text{old}})}{p(X | \Theta^{\text{old}})} \propto \prod_{i=1}^{\ell} \prod_{k=1}^K \left[ \pi_k^{\text{old}} \mathcal{N}(x_i | \mu_k^{\text{old}}, \Sigma_k^{\text{old}}) \right]^{z_{ik}}.$$

Заметим, что данное распределение распадается в произведение распределений, соответствующих отдельным объектам  $p(z_i | x_i, \Theta)$ :

$$p(Z | X, \Theta^{\text{old}}) = \prod_{i=1}^{\ell} p(z_i | x_i, \Theta^{\text{old}}).$$

Иными словами, величины  $\{z_i\}$  независимы при известных объектах  $\{x_i\}$ . Вектор  $z_i$  имеет  $K$  возможных значений. Запишем их вероятности, воспользовавшись формулой Байеса:

$$\begin{aligned} g_{ik} \equiv p(z_{ik} = 1 | x_i, \Theta^{\text{old}}) &= \frac{p(z_{ik} = 1) p(x_i | z_{ik} = 1, \Theta^{\text{old}})}{\sum_{j=1}^K p(z_{ij} = 1) p(x_i | z_{ij} = 1, \Theta^{\text{old}})} = \\ &= \frac{\pi_k^{\text{old}} \mathcal{N}(x_i | \mu_k^{\text{old}}, \Sigma_k^{\text{old}})}{\sum_{j=1}^K \pi_j^{\text{old}} \mathcal{N}(x_i | \mu_j^{\text{old}}, \Sigma_j^{\text{old}})}. \end{aligned}$$

Вычислим теперь матожидание полного правдоподобия:

$$\begin{aligned}
 Q(\Theta, \Theta^{\text{old}}) &= \mathbb{E}_{Z \sim p(Z|X, \Theta^{\text{old}})} \log p(X, Z | \Theta) = \\
 &= \mathbb{E}_{Z \sim p(Z|X, \Theta^{\text{old}})} \sum_{i=1}^{\ell} \sum_{k=1}^K z_{ik} \left\{ \log \pi_k + \log \mathcal{N}(x_i | \mu_k, \Sigma_k) \right\} = \\
 &= \sum_{i=1}^{\ell} \sum_{k=1}^K \mathbb{E}_{Z \sim p(Z|X, \Theta^{\text{old}})} [z_{ik}] \left\{ \log \pi_k + \log \mathcal{N}(x_i | \mu_k, \Sigma_k) \right\}.
 \end{aligned}$$

Нам понадобится вспомогательная величина:

$$\begin{aligned}
 \mathbb{E}_{Z \sim p(Z|X, \Theta^{\text{old}})} [z_{ik}] &= 1 * p(z_{ik} = 1 | x_i, \Theta) + 0 * p(z_{ik} = 0 | x_i, \Theta) = \\
 &= \frac{\pi_k^{\text{old}} \mathcal{N}(x_i | \mu_k^{\text{old}}, \Sigma_k^{\text{old}})}{\sum_{j=1}^K \pi_j^{\text{old}} \mathcal{N}(x_i | \mu_j^{\text{old}}, \Sigma_j^{\text{old}})} = \\
 &= g_{ik}.
 \end{aligned}$$

Получаем:

$$Q(\Theta, \Theta^{\text{old}}) = \sum_{i=1}^{\ell} \sum_{k=1}^K g_{ik} \left\{ \log \pi_k + \log \mathcal{N}(x_i | \mu_k, \Sigma_k) \right\}.$$

Дифференцируя данный функционал, нетрудно получить формулы М-шага:

$$\begin{aligned}
 \pi_k &= \frac{1}{\ell} \sum_{i=1}^{\ell} g_{ik}; \\
 \mu_k &= \frac{1}{\ell \pi_k} \sum_{i=1}^{\ell} g_{ik} x_i; \\
 \Sigma_k &= \frac{1}{\ell \pi_k} \sum_{i=1}^{\ell} g_{ik} (x_i - \mu_k)(x_i - \mu_k)^T.
 \end{aligned}$$

## §1.5 Дивергенция Кульбака-Лейблера

Дивергенция Кульбака-Лейблера — это мера расстояния между двумя вероятностными распределениями, которая определяется как

$$\text{KL}(q \parallel p) = \int q(x) \log \frac{q(x)}{p(x)} dx.$$

В случае с дискретными распределениями она принимает вид

$$\text{KL}(q \parallel p) = \sum_x q(x) \log \frac{q(x)}{p(x)}.$$

KL-дивергенция определена только в том случае, когда из  $p(x) = 0$  следует  $q(x) = 0$ . При вычислении мы считаем, что  $0 \log 0 = 0$  и  $0 \log \frac{0}{0} = 0$ .

**Задача 1.1.** Покажите, что KL-дивергенция неотрицательна.

**Решение.** Нам понадобится неравенство Йенсена в интегральной форме: для любой вогнутой функции  $f(x)$  выполнено

$$f\left(\int \alpha(x)y(x)dx\right) \geq \int \alpha(x)f(y(x))dx, \quad \int \alpha(x)dx = 1, \quad \alpha(x) \geq 0.$$

Пользуясь данным неравенством и вогнутостью логарифма, получаем:

$$\text{KL}(q \parallel p) = -\int q(x) \log \frac{p(x)}{q(x)} dx \geq -\log\left(\int q(x) \frac{p(x)}{q(x)} dx\right) = -\log\left(\int p(x) dx\right) = 0.$$

Можно доказать данное утверждение и без использования неравенства Йенсена, если в определении используется натуральный логарифм. Заметим, что при  $y > 0$  имеет место неравенство  $\ln y \leq y - 1$ , которое обращается в равенство только при  $y = 1$ . Тогда:

$$\text{KL}(q \parallel p) = -\int q(x) \log \frac{p(x)}{q(x)} dx \geq -\int q(x) \left(\frac{p(x)}{q(x)} - 1\right) dx = \int q(x) dx - \int p(x) dx = 0.$$

■

Неравенство Йенсена обращается в равенство тогда и только тогда, когда  $y(x) = \text{const}$ . В нашем случае это означает, что  $\frac{p(x)}{q(x)} = \text{const}$ . Поскольку  $q$  и  $p$  — вероятностные распределения, это возможно только при их равенстве. Мы получили важное свойство KL-дивергенции: она обращается в нуль тогда и только тогда, когда ее аргументы равны.

**Задача 1.2.** Пусть заданы выборка  $X^\ell$  и распределение на объектах  $p(x | \theta)$ . Эмпирическим распределением называется дискретное распределение на объектах, присваивающее каждому объекту из обучающей выборки вероятность  $1/\ell$ :

$$\hat{L}(x | X^\ell) = \sum_{i=1}^{\ell} \frac{1}{\ell} [x = x_i].$$

Покажите, что максимизация правдоподобия эквивалентна минимизации дивергенции Кульбака-Лейблера между эмпирическим распределением и модельным распределением:  $\text{KL}\left(\hat{L}(x | X^\ell) \parallel L(x | X^\ell, \theta)\right)$ .

**Решение.** Распишем указанную дивергенцию:

$$\begin{aligned} \text{KL}\left(\hat{L}(x | X^\ell) \parallel L(x | X^\ell, \theta)\right) &= \sum_{i=1}^{\ell} \frac{1}{\ell} \log \frac{1/\ell}{p(x_i | \theta)} = \\ &= \sum_{i=1}^{\ell} \frac{1}{\ell} \log \frac{1}{\ell} - \frac{1}{\ell} \sum_{i=1}^{\ell} \log p(x_i | \theta) \rightarrow \min_{\theta}. \end{aligned}$$

Отбросим константные члены:

$$\sum_{i=1}^{\ell} \log p(x_i | \theta) \rightarrow \max_{\theta}.$$

Мы получили задачу максимизации логарифма правдоподобия. ■

Таким образом, метод максимума правдоподобия старается подобрать такие параметры модели, чтобы она давала равномерное распределение на объектах выборки и присваивала нулевую вероятность всем остальным объектам.

## §1.6 Обоснование EM-алгоритма

Представим неполное правдоподобие в виде суммы двух функций:

$$\log p(X | \Theta) = \mathcal{L}(q, \Theta) + \text{KL}(q \| p), \quad (1.2)$$

где

$$\begin{aligned} \mathcal{L}(q, \Theta) &= \sum_Z q(Z) \log \frac{p(X, Z | \Theta)}{q(Z)}, \\ \text{KL}(q \| p) &= - \sum_Z q(Z) \log \frac{p(Z | X, \Theta)}{q(Z)}. \end{aligned}$$

Здесь  $q(Z)$  — это произвольное распределение на скрытых переменных.

**Задача 1.3.** Докажите, что это представление корректно.

**Решение.**

$$\begin{aligned} &\sum_Z q(Z) \log \frac{p(X, Z | \Theta)}{q(Z)} - \sum_Z q(Z) \log \frac{p(Z | X, \Theta)}{q(Z)} = \\ &= \sum_Z q(Z) \log \frac{p(X, Z | \Theta)}{p(Z | X, \Theta)} = \\ &= \sum_Z q(Z) \log p(X | \Theta) = \\ &= \log p(X | \Theta) \sum_Z q(Z) = \\ &= \log p(X | \Theta). \end{aligned}$$

Заметим, что  $\mathcal{L}(q, \Theta)$  — это нижняя оценка на логарифм правдоподобия: ■

$$\log p(X | \Theta) = \mathcal{L}(q, \Theta) + \underbrace{\text{KL}(q \| p)}_{\geq 0} \geq \mathcal{L}(q, \Theta).$$

Чем «правильнее» выбрано распределение  $q(Z)$ , тем точнее эта оценка. Будем по очереди максимизировать нижнюю оценку  $\mathcal{L}(q, \Theta)$  по  $q$  и по  $\Theta$ . Зафиксируем сначала вектор параметров  $\Theta^{\text{old}}$  и найдем максимум по  $q$ . Заметим, что в разложении (1.2) левая часть не зависит от  $q$ , поэтому нижняя оценка будет максимальна тогда, когда KL-дивергенция будет минимальна. Мы знаем, что минимум дивергенции равен

нулю и достигается на равных распределениях. Таким образом, нижняя оценка достигнет своего максимума на  $q = p(Z | X, \Theta^{\text{old}})$ . Мы получили E-шаг EM-алгоритма — вычисление апостериорного распределения на скрытых переменных.

Зафиксируем теперь  $q$  и найдем максимум нижней оценки по  $\Theta$ . Преобразуем задачу:

$$\begin{aligned} \mathcal{L}(q, \Theta) &= \sum_Z q(Z) \log \frac{p(X, Z | \Theta)}{q(Z)} = \\ &= \sum_Z q(Z) \log p(X, Z | \Theta) - \sum_Z q(Z) \log q(Z) = \\ &= \sum_Z p(Z | X, \Theta^{\text{old}}) \log p(X, Z | \Theta) + \text{const}(\Theta) = \\ &= Q(\Theta, \Theta^{\text{old}}) + \text{const}(\Theta) \rightarrow \max_{\Theta}. \end{aligned}$$

Мы получили оптимизационную задачу с M-шага EM-алгоритма.

Описанный способ вывода E- и M-шагов позволяет получить важное свойство EM-алгоритма — на каждой его итерации значение правдоподобия не уменьшается. Действительно, после E-шага значение нижней оценки совпадает со значением правдоподобия, а значит, максимизация оценки на M-шаге приведет и к максимизации правдоподобия:

$$\log p(X | \Theta^{\text{new}}) = \mathcal{L}(q, \Theta^{\text{new}}) + \text{KL}(q \| p) \geq \mathcal{L}(q, \Theta^{\text{new}}) \geq \mathcal{L}(q, \Theta^{\text{old}}) = \log p(X | \Theta^{\text{old}}).$$

Если правдоподобие ограничено сверху, то последовательность значений правдоподобия  $\{p(X | \Theta^i)\}_i$  обязательно сойдется. Здесь мы обозначили последовательность параметров, генерируемую EM-алгоритмом, через  $\{\Theta^i\}_i$ .

Существуют и более сильные утверждения о сходимости.

**Теорема 1.1 ([1]).** Пусть  $Q(\Theta, \Theta^{\text{old}})$  непрерывна по  $\Theta$  и  $\Theta^{\text{old}}$ . Тогда все предельные точки последовательности  $\{\Theta^i\}_i$  являются стационарными точками неполного правдоподобия  $p(X | \Theta)$ , а последовательность  $\{p(X | \Theta^i)\}_i$  монотонно сходится к значению правдоподобия  $L^* = p(X | \Theta^*)$  в одной из стационарных точек  $\Theta^*$ .

Обратим внимание на тот факт, что сходимость последовательности  $\{\Theta^i\}_i$  не гарантируется — у нее может быть несколько подпоследовательностей, каждая из которых будет сходиться к своей стационарной точке. Также отметим, что речь идет только о сходимости к стационарной точке; сходимость к локальному максимуму гарантируется лишь для некоторых семейств распределений (например, для экспоненциальных [1]).

Покажем одно из свойств EM-алгоритма.

**Задача 1.4.** Докажите, что если  $\Theta^i$  не является стационарной точкой логарифма правдоподобия, то следующее приближение  $\Theta^{i+1}$ , выданное EM-алгоритмом, будет отличаться от  $\Theta^i$ .

**Решение.** Пусть  $\Theta^i$  не является стационарной точкой, то есть

$$\nabla_{\Theta} \log p(X | \Theta)|_{\Theta^i} \neq 0.$$



Выполним E-шаг, найдем апостериорное распределение  $q(\Theta^i)$ , и запишем разложение правдоподобия:

$$\log p(X | \Theta^i) = \mathcal{L}(q, \Theta^i) + \underbrace{\text{KL}(q(\Theta^i) \| p)}_{=0}.$$

KL-дивергенция здесь равна нулю в силу выбора распределения  $q(\Theta^i)$ . Поскольку на данном распределении достигается минимум дивергенции, ее градиент равен нулю:

$$\nabla_{\Theta} \text{KL}(q(\Theta) \| p)|_{\Theta^i} = 0.$$

Получаем:

$$\nabla_{\Theta} \mathcal{L}(q, \Theta)|_{\Theta^i} = \nabla_{\Theta} \log p(X | \Theta)|_{\Theta^i} - \underbrace{\nabla_{\Theta} \text{KL}(q(\Theta) \| p)|_{\Theta^i}}_{=0} = \nabla_{\Theta} \log p(X | \Theta)|_{\Theta^i} \neq 0.$$

Таким образом, точка  $\Theta^i$  не является максимумом нижней оценки, и поэтому на M-шаге будет сделан переход к новой точке  $\Theta^{i+1} \neq \Theta^i$ . ■

## Список литературы

- [1] *Wu, C. F. Jeff* (1983). On the Convergence Properties of the EM Algorithm. // *Annals of Statistics*, 11(1), p. 95-103.