

Лекция 15

ЕМ-алгоритм

Е. А. Соколов
ФКН ВШЭ

9 февраля 2020 г.

Существуют задачи, в которых помимо параметров и известных данных имеются скрытые переменные — некоторые величины, которые существенно упрощают модель, но неизвестны. Ниже мы рассмотрим смеси распределений, которые являются примером такой модели. Общим методом обучения моделей со скрытыми переменными является ЕМ-алгоритм, о котором и пойдёт речь в этой лекции.

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

Каждую компоненту распределения $p_k(x)$ можно рассматривать как кластер, а значение данной плотности на объекте — как вероятность принадлежности данному кластеру. Таким образом, с помощью смеси распределений можно описывать *мягкую кластеризацию*, в которой каждый объект относится к каждому из кластеров с некоторой вероятностью.

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:

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

Вероятность того, что единице будет равна k -я компонента, равна π_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,

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

Запишем совместное распределение переменных x и z :

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

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.

2 Модели со скрытыми переменными

Рассмотрим вероятностную модель с наблюдаемыми переменными X и параметрами Θ , 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 Θ .

Рассмотрим пример со смесями распределений. В качестве наблюдаемых переменных здесь выступает выборка $X = \{x_1, \dots, x_\ell\}$, 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 Θ 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 *ЕМ-алгоритме*.

3 ЕМ-алгоритм

ЕМ-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 Θ^{old} . При известных наблюдаемых переменных X и параметрах Θ^{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 Θ^* . Данная процедура, однако, далека от идеальной — ниже мы предложим подход, который приводит к более качественным результатам.

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

4 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 4.1. 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 4.2. Пусть заданы выборка X^ℓ и распределение на объектах $p(x \mid \theta)$, параметр которого мы хотим настроить под данную выборку. Эмпирическим распределением называется дискретное распределение на объектах, присваивающее каждому объекту из обучающей выборки вероятность $1/\ell$:

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

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

Решение. Распишем указанную дивергенцию:

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

Отбросим константные члены:

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

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.

5 Обоснование ЕМ-алгоритма

Представим неполное правдоподобие в виде суммы двух функций:

$$\log p(X | \Theta) = \mathcal{L}(q, \Theta) + \text{KL}(q \parallel p), \quad (5.1)$$

где

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

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

Задача 5.1. 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:

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

Чем «правильнее» выбрано распределение $q(Z)$, тем точнее эта оценка. Будем по очереди максимизировать нижнюю оценку $\mathcal{L}(q, \Theta)$ по q и по Θ . Зафиксируем сначала вектор параметров Θ^{old} и найдем максимум по q . Заметим, что в разложении (5.1) 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 Е-шаг ЕМ-алгоритма — вычисление апостериорного распределения на скрытых переменных.

Зафиксируем теперь q и найдем максимум нижней оценки по Θ . 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 М-шага ЕМ-алгоритма.

Описанный способ вывода Е- и М-шагов позволяет получить важное свойство ЕМ-алгоритма — на каждой его итерации значение правдоподобия не уменьшается. Действительно, после Е-шага значение нижней оценки совпадает со значением правдоподобия, а значит, максимизация оценки на М-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 ЕМ-алгоритмом, через $\{\Theta^i\}_i$.

Существуют и более сильные утверждения о сходимости.

Теорема 5.1 ([1]). Пусть $Q(\Theta, \Theta^{\text{old}})$ непрерывна по Θ и Θ^{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 Θ^* .

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 ЕМ-алгоритма.

Задача 5.2. Докажите, что если Θ^i не является стационарной точкой логарифма правдоподобия, то следующее приближение Θ^{i+1} , выданное ЕМ-алгоритмом, будет отличаться от Θ^i .

Решение. Пусть Θ^i не является стационарной точкой, то есть

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

Выполним Е-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:

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

Получаем:

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

Таким образом, точка Θ^i не является максимумом нижней оценки, и поэтому на М-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.