

Оценка параметров смеси распределений

К. В. Павлов

kirill.pavlov@phystech.edu

Московский физико-технический институт

В работе рассматриваются способы построения смеси моделей и экспертов. Предлагается *EM*-алгоритм для совместного нахождения параметров моделей и их весов в смеси, а так же для нахождения параметров смеси обобщенных линейных моделей.

Ключевые слова: смеси моделей, обобщенно-линейные модели, смеси экспертов.

Введение

При решении задачи анализа данных строится модель — отображение известных характеристик объекта в неизвестные. Часто оказывается, что качество алгоритма можно улучшить с помощью комбинирования нескольких моделей [3, р. 653–676]. Например, можно обучить l моделей и в качестве ответа выводить усредненный ответ по всем моделям. Подобные комбинации моделей называются комитетами. Один из наиболее важных случаев комитета является бустинг. Алгоритмы в комитет добавляются последовательно и их параметры зависят от уже созданного на момент добавления комитета. Другим важным частным случаем комитета является смесь экспертов. В этом случае ответы алгоритмов взвешиваются в зависимости от области пространства, в которой находится объект. Рассмотрим способы построения композиций.

Общий подход к оценке параметров моделей

В случае, когда одной модели для описания данных не хватает, используют смеси моделей. Предполагается, что исходная зависимость $p(\mathbf{y} | \mathbf{x})$ выражается как композиция моделей $p(y | \mathbf{x}, \mathbf{w}_k)$ формулой:

$$p(\mathbf{y} | \mathbf{x}) = \sum_{k=1}^l p(\mathbf{w}_k | \mathbf{x}) p(y | \mathbf{x}, \mathbf{w}_k) = \sum_{k=1}^l \pi_k p(y | \mathbf{x}, \mathbf{w}_k), \quad (1)$$

где $\pi_k = p(\mathbf{w}_k | \mathbf{x})$ — вероятность принадлежности к модели k . На π_k накладываются условия нормировки: вероятность каждой модели неотрицательна и сумма вероятностей равна единице.

$$\sum_{k=1}^l \pi_k = 1, \quad \pi_k \geq 0 \quad \forall k. \quad (2)$$

Далее предполагается, что объекты в выборке независимы и плотность совместного распределения преобразуется в произведение плотностей распределения каждого объекта.

$$p(\mathbf{y} | \mathbf{x}) = \sum_{k=1}^l \pi_k \prod_{i=1}^n p(y^i | \mathbf{x}^i, \mathbf{w}_k) = \prod_{i=1}^n \sum_{k=1}^l \pi_k p(y^i | \mathbf{x}^i, \mathbf{w}_k). \quad (3)$$

В формуле 3 произведена смена порядка суммирования перемножения. Используя принцип максимума правдоподобия, будет максимизировать $p(\mathbf{y} | \mathbf{x})$. Проще это делать, введя функцию правдоподобия $Q(\mathbf{w}_1, \dots, \mathbf{w}_l, \boldsymbol{\pi})$ как логарифм плотности вероятности данных.

$$Q(\mathbf{w}^1, \dots, \mathbf{w}^l, \boldsymbol{\pi}) = \ln p(\mathbf{y} | \mathbf{x}) = \sum_{i=1}^m \ln \left[\sum_{k=1}^l \pi_k p(y^i | \mathbf{x}^i, \mathbf{w}_k) \right]. \quad (4)$$

Обозначим через $p(y, \mathbf{w}_k | \mathbf{x})$ вероятность того, что объект (\mathbf{x}, y) был порожден компонентой \mathbf{w}_k , $\gamma_{ik} = p(\mathbf{w}_k | y^i, \mathbf{x}^i)$ — вероятность того, что i -объект порожден j -компонентой. Каждый объект был порожден какой-либо моделью, по формуле полной вероятности

$$\sum_{k=1}^l \gamma_{ik} = 1, \quad \forall i. \quad (5)$$

Для произвольного объекта (\mathbf{x}, y) вероятность его получения моделью w_k по формуле условной вероятности равна:

$$p(y, \mathbf{w}_k | \mathbf{x}) = p(\mathbf{w}_k | \mathbf{x}) p(y | \mathbf{x}, \mathbf{w}_k) \equiv \pi_k p(y | \mathbf{x}, \mathbf{w}_k). \quad (6)$$

Подставим это равенство в формулу Байеса для γ_{ik}

$$\gamma_{ik} = \frac{\pi_k p(y^i | \mathbf{x}^i, \mathbf{w}_k)}{\sum_{s=1}^l \pi_s p(y^i | \mathbf{x}^i, \mathbf{w}_s)}. \quad (7)$$

Для определения параметров смеси необходимо решить задачу максимизации правдоподобия $Q(\mathbf{w}^1, \dots, \mathbf{w}^l, \boldsymbol{\pi}) \rightarrow \max$, это можно сделать с использованием функции Лагранжа [1], которая имеет вид:

$$L = \sum_{i=1}^m \ln \left[\sum_{k=1}^l \pi_k p(y^i | \mathbf{x}^i, \mathbf{w}^k) \right] - \lambda \left(\sum_{k=1}^l \pi_k - 1 \right). \quad (8)$$

Необходимым условием экстремума функции является равенство нулю первых производных. Приравняем производную функции Лагранжа по π_k к нулю:

$$\frac{\partial L}{\partial \pi_k} = \sum_{i=1}^m \frac{p(y^i | \mathbf{x}^i, \mathbf{w}^k)}{\sum_{s=1}^l \pi_s p(y^i | \mathbf{x}^i, \mathbf{w}^s)} - \lambda = 0. \quad (9)$$

Умножим обе части равенства на π_k и просуммируем по $k = 1..l$

$$m = \sum_{k=1}^l \sum_{i=1}^m \frac{\pi_k p(y^i | \mathbf{x}^i, \mathbf{w}^k)}{\sum_{s=1}^l \pi_s p(y^i | \mathbf{x}^i, \mathbf{w}^s)} = \lambda \sum_{s=1}^l \pi_s = \lambda. \quad (10)$$

Получилось, что $\lambda = m$ необходимое условие минимума. В выражении для производной $\frac{\partial L}{\partial \pi_k}$ заменим λ на m и домножим обе части равенства на π_k :

$$\pi_k = \frac{1}{m} \sum_{i=1}^m \frac{\pi_k p(y^i | \mathbf{x}^i, \mathbf{w}^k)}{\sum_{s=1}^l \pi_s p(y^i | \mathbf{x}^i, \mathbf{w}^s)} = \frac{1}{m} \sum_{i=1}^m \gamma_{ik}. \quad (11)$$

Равенство 11 позволяет находить коэффициенты π_k смеси модели при известных γ_{ik} . Вычислим производную функции Лагранжа по параметрам k -й модели:

$$\frac{\partial L}{\partial \mathbf{w}^k} = \sum_{i=1}^m \frac{\pi_k \frac{\partial p(y^i | \mathbf{x}^i, \mathbf{w}^k)}{\partial \mathbf{w}^k}}{\sum_{s=1}^l \pi_s p(y^i | \mathbf{x}^i, \mathbf{w}^s)} = \sum_{i=1}^m \frac{\pi_k p(y^i | \mathbf{x}^i, \mathbf{w}^k)}{\sum_{s=1}^l \pi_s p(y^i | \mathbf{x}^i, \mathbf{w}^s)} \frac{\partial \ln p(y^i | \mathbf{x}^i, \mathbf{w}^k)}{\partial \mathbf{w}^k}. \quad (12)$$

Преобразуем выражение:

$$\frac{\partial L}{\partial \mathbf{w}^k} = \frac{\partial}{\partial \mathbf{w}^k} \sum_{i=1}^m \gamma_{ik} \ln p(y^i | \mathbf{x}^i, \mathbf{w}^k) = 0. \quad (13)$$

Полученное равенство совпадает с необходимым условием максимума в задаче максимизации взвешенного правдоподобия:

$$\sum_{i=1}^m \gamma_{ik} \ln p(y^i | \mathbf{x}^i, \mathbf{w}^k) \rightarrow \max_{\mathbf{w}^k}. \quad (14)$$

В общем случае задача оптимизации $Q(\mathbf{w}^1, \dots, \mathbf{w}^l, \boldsymbol{\pi}) \rightarrow \max$ трудна, для её решения используют EM-алгоритм, заключающийся в итеративном повторении двух шагов. На E -шаге вычисляются ожидаемые значения вектора скрытых переменных γ_{ik} по текущему приближению параметров моделей $(\mathbf{w}_1, \dots, \mathbf{w}_l)$. На M -шаге решается задача максимизации правдоподобия Q при начальном приближении параметров моделей и значений γ_{ik} .

E -шагу соответствует выражение

$$\gamma_{ik} = \frac{\pi_k p(y^i | \mathbf{x}^i, \mathbf{w}_k)}{\sum_{s=1}^l \pi_s p(y^i | \mathbf{x}^i, \mathbf{w}_s)}. \quad (15)$$

M -шаг заключается в оптимизации параметров распределений.

$$Q(\mathbf{w}^1, \dots, \mathbf{w}^l | \boldsymbol{\pi}) \rightarrow \max \quad (16)$$

Формула на M -шаге может упроститься для случая конкретного распределения. Для упрощения дальнейших рассуждений введем обозначения

$$G = (\gamma_1, \dots, \gamma_l) = \begin{pmatrix} \gamma_{11} & \dots & \gamma_{1l} \\ \vdots & \ddots & \vdots \\ \gamma_{m1} & \dots & \gamma_{ml} \end{pmatrix}, \quad G_k = \begin{pmatrix} \gamma_{1k} & & 0 \\ & \ddots & \\ 0 & & \gamma_{mk} \end{pmatrix}. \quad (17)$$

Перейдем к рассмотрению линейных и обобщенных линейных моделей.

Оценка параметров смеси линейных моделей

Линейная модель имеет вид:

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\varepsilon}, \quad (18)$$

где $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, B)$ — вектор нормально распределенных ошибок. В данной постановке вектор \mathbf{y} является нормальным с математическим ожиданием $E(y | \mathbf{x}) = \boldsymbol{\mu} = \mathbf{x}^T \mathbf{w}$, и корреляционной матрицей B . Плотность распределения \mathbf{y} задается формулой:

$$p(\mathbf{y} | X, \mathbf{w}) = \frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{|\det B|}} \exp\left(-\frac{1}{2}(\mathbf{y} - X\mathbf{w})^T B(\mathbf{y} - X\mathbf{w})\right). \quad (19)$$

Применим для задачи описанный EM-алгоритм. Шаг E сводится к применению формулы 15, а шаг M алгоритма принимает следующий вид:

$$G_k \ln \left[\frac{1}{(2\pi)^{\frac{n}{2}} \sqrt{|\det B|}} \right] - \frac{1}{2} (G_k(\mathbf{y} - X\mathbf{w})^T B(\mathbf{y} - X\mathbf{w})) \rightarrow \max_{\mathbf{w}} \quad (20)$$

Первое слагаемое не зависит от \mathbf{w}_k , его можно не учитывать. Преобразование второго слагаемого дает

$$\frac{1}{2} \mathbf{w}^T X^T G_k B X \mathbf{w} - \mathbf{w}^T X^T G_k B \mathbf{y} \rightarrow \min_{\mathbf{w}} \quad (21)$$

Задача квадратична по \mathbf{w} , решение находится аналитически

$$\mathbf{w}^* = (X^T G_k B X)^{-1} G_k B X \mathbf{y}. \quad (22)$$

Оценка параметров смеси обобщенно-линейных моделей

В случае обобщенных линейных моделей функция плотности распределения имеет вид

$$p(\mathbf{y} | \boldsymbol{\theta}) = \exp(\mathbf{T}(\mathbf{y})^T \boldsymbol{\eta}(\boldsymbol{\theta}) - b(\boldsymbol{\theta}) + c(\mathbf{y})). \quad (23)$$

M -шаг алгоритма сводится к максимизации

$$\mathbf{T}(\mathbf{y})^T G_k \boldsymbol{\eta}(\boldsymbol{\theta}) - b(G_k \boldsymbol{\theta}) + c(G_k \mathbf{y}) \rightarrow \max_{\boldsymbol{\theta}}. \quad (24)$$

Последнее слагаемое не зависит от параметров модели $\boldsymbol{\theta}$, что позволяет упростить функционал

$$\mathbf{T}(\mathbf{y})^T G_k \boldsymbol{\eta}(\boldsymbol{\theta}) - b(G_k \boldsymbol{\theta}) \rightarrow \max_{\boldsymbol{\theta}}. \quad (25)$$

Дальнейшая минимизация зависит от конкретного семейства из обобщенного класса распределений.

Оценка параметров смеси экспертов

Понятие смеси экспертов было введено Джорданом и Якобсом в 1991г [2]. Предполагается, что параметры смеси π являются функциями от объекта, т.е.

$$p(\mathbf{y} | \mathbf{x}) = \sum_{k=1}^l \pi_k(\mathbf{x}) p(\mathbf{y} | \mathbf{x}, \mathbf{w}_k). \quad (26)$$

Компоненты $\pi_k(\mathbf{x})$ называются функциями селективности, а $p(\mathbf{y} | \mathbf{x}, \mathbf{w}_k)$ экспертами. Функция селективности отвечает за компетентность эксперта в определенной области.

Оказывается [4], что наличие функции компетенции допускает решение задачи с помощью EM -алгоритма, причем, E -шаг остается прежним:

$$\gamma_{ik} = \frac{\pi_k(\mathbf{x}^i) p(y^i | \mathbf{x}^i, \mathbf{w}_k)}{\sum_{s=1}^l \pi_s(\mathbf{x}^i) p(y^i | \mathbf{x}^i, \mathbf{w}_s)}. \quad (27)$$

M -шаг принимает вид:

$$\pi_k = \frac{1}{m} \sum_{i=1}^m \gamma_{ik}. \quad (28)$$

$$\sum_{i=1}^m \gamma_{ik}(\mathbf{x}^i) \ln p(y^i | \mathbf{x}^i, \mathbf{w}^k) \rightarrow \max_{\mathbf{w}^k}. \quad (29)$$

Уравнение 29 можно решить с помощью метода итеративно перевзвешенных наименьших квадратов (IRLS).

Литература

- [1] Воронцов К. В. Курс лекций: Линейные методы классификации. — 2009. — 01. <http://www.machinelearning.ru/wiki/images/6/68/Voron-ML-Lin.pdf>.

- [2] Adaptive mixtures of local experts / R. A. Jacobs, M. I. Jordan, S. J. Nowlan, G. E. Hinton // *Neural Computation*. — 1991. — no. 3. — Pp. 79–87.
- [3] *Bishop C. M.* Pattern Recognition and Machine Learning. — Springer, Series: Information Science and Statistics, 2006. — 740 pp.
- [4] *Jordan M. I., Jacobs R. A.* Hierarchical mixtures of experts and the EM algorithm // *Neural Computation*. — 1994. — no. 6. — Pp. 181–214. citeseer.ist.psu.edu/article/jordan94hierarchical.html.