

Смесь экспертов

Грабовой Андрей Валериевич

Московский физико-технический институт
Факультет управления и прикладной математики
Кафедра интеллектуальных систем

Москва
2019 г

① ЕМ–алгоритм:

- Классический,
- Вариационный,

② Смесь моделей:

- Постановка задачи,
- Итерационные формулы полученные при помощи вариационного ЕМ–алгоритма,
- Иллюстрация сходимости,

③ Смесь экспертов

- Постановка задачи,
- Итерационные формулы полученные при помощи вариационного ЕМ–алгоритма,
- Иллюстрация сходимости,

Максимизация обоснованности:

$$\Theta = \arg \max_{\Theta} p(\mathbf{y}|\mathbf{X}, \Theta) \quad (1)$$

ELBO:

$$\begin{aligned} \mathcal{L}(q(\mathbf{Z}), \Theta) &= \int q(\mathbf{Z}) \log \frac{p(\mathbf{y}, \mathbf{Z}|\Theta, \mathbf{X})}{q(\mathbf{Z})} d\mathbf{Z} \\ &= p(\mathbf{y}|\mathbf{X}, \Theta) - D_{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{y}, \mathbf{X}, \Theta)) \end{aligned} \quad (2)$$

EM-алгоритм:

❶ E-шаг:

$$q^s(\mathbf{Z}) = \arg \max_{q(\mathbf{Z}) \in Q} \mathcal{L}(q(\mathbf{Z}), \Theta^{s-1}) \quad (3)$$

❷ M-шаг:

$$\Theta^s = \arg \max_{\Theta} \mathcal{L}(q^s(\mathbf{Z}), \Theta) \quad (4)$$

Вариационный EM-алгоритм (Mean Field Approximation¹):

❶ E-шаг:

$$\log q(\mathbf{Z}_k^s) \propto \mathbb{E}_{q/k} \log p(\mathbf{y}, \mathbf{Z}|\mathbf{X}, \Theta^{s-1}) \quad (5)$$

❷ M-шаг:

$$\Theta^s = \arg \max_{\Theta} \mathbb{E}_{q^s} \log p(\mathbf{y}, \mathbf{Z}|\mathbf{X}, \Theta) \quad (6)$$

¹<https://github.com/andriygav/EMprior/blob/master/Lecture/Grabovoy2019MeanField.pdf>

Definition

Смесь моделей — мультимодель, ответы которой представляют собой взвешенную сумму ответов всех задействованных моделей независимо от объекта.

$$\hat{\mathbf{f}} = \sum_{k=1}^K \pi_k \mathbf{f}_k, \quad \pi_k = \text{const}, \quad \sum_{k=1}^K \pi_k = 1, \quad (7)$$

где \mathbf{f} — мультимодель, а \mathbf{f}_k — локальная модель.

Пример 1:

- ① Веса моделей в смеси $\boldsymbol{\pi}$ получены из априорного распределения

$$p(\boldsymbol{\pi}|\mu); \quad (8)$$

- ② Вектора параметров \mathbf{w}_k получены из нормального распределения

$$p(\mathbf{w}_k|\mathbf{A}_k) = \mathcal{N}(\mathbf{w}_k|\mathbf{0}, \mathbf{A}_k), \quad k = 1, \dots, K; \quad (9)$$

- ③ Для каждого объекта \mathbf{x}_i существует модель \mathbf{f}_{k_i} , которой он описывается, причем $p(k_i = k) = \pi_k$;

- ④ Для каждого объекта \mathbf{x}_i класс y_i определен в соответствии с моделью

$$\mathbf{f}_{k_i} : y_i \sim \mathcal{N}(\mathbf{w}_{k_i}^\top \mathbf{x}_i + b_k, \beta^{-1}) \quad (10)$$

Правдоподобие модели:

$$p(\mathbf{y}, \mathbf{W}, \boldsymbol{\pi} | \mathbf{X}, \mathbf{A}, \beta, \boldsymbol{\mu}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\mu}) \prod_{k=1}^K N(\mathbf{w}_k | \mathbf{0}, \mathbf{A}_k) \prod_{i=1}^N \left(\sum_{k=1}^K \pi_k \mathcal{N}(y_i | \mathbf{w}_k^T \mathbf{x}_i, \beta^{-1}) \right) \quad (11)$$

Введем скрытые переменные $\mathbf{Z} = ||z_{ik}||$, где $z_{ik} = 1 \Leftrightarrow k_i = k$:

$$p(\mathbf{y}, \mathbf{W}, \boldsymbol{\pi}, \mathbf{Z} | \mathbf{X}, \mathbf{A}, \beta, \boldsymbol{\mu}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\mu}) \prod_{k=1}^K N(\mathbf{w}_k | \mathbf{0}, \mathbf{A}_k) \prod_{i=1}^N \prod_{k=1}^K \left(\pi_k \mathcal{N}(y_i | \mathbf{w}_k^T \mathbf{x}_i, \beta^{-1}) \right)^{z_{ik}} \quad (12)$$

Вариационный EM-алгоритм $q(\mathbf{Z}, \mathbf{W}, \boldsymbol{\pi}) = q(\mathbf{Z}) q(\mathbf{W}) q(\boldsymbol{\pi})$:

① E-шаг:

$$\begin{aligned} \log q(\mathbf{Z}^s) &\propto \mathbb{E}_{q/\mathbf{Z}} \log p(\mathbf{y}, \mathbf{W}, \boldsymbol{\pi}, \mathbf{Z} | \mathbf{X}, \mathbf{A}^{s-1}, \beta^{s-1}, \boldsymbol{\mu}) \\ \log q(\mathbf{W}^s) &\propto \mathbb{E}_{q/\mathbf{W}} \log p(\mathbf{y}, \mathbf{W}, \boldsymbol{\pi}, \mathbf{Z} | \mathbf{X}, \mathbf{A}^{s-1}, \beta^{s-1}, \boldsymbol{\mu}) \\ \log q(\boldsymbol{\pi}^s) &\propto \mathbb{E}_{q/\boldsymbol{\pi}} \log p(\mathbf{y}, \mathbf{W}, \boldsymbol{\pi}, \mathbf{Z} | \mathbf{X}, \mathbf{A}^{s-1}, \beta^{s-1}, \boldsymbol{\mu}) \end{aligned} \quad (13)$$

② M-шаг:

$$\mathbf{A}^s, \beta^s = \arg \max_{\mathbf{A}, \beta} \mathbb{E}_{q^s} \log p(\mathbf{y}, \mathbf{W}, \boldsymbol{\pi}, \mathbf{Z} | \mathbf{X}, \mathbf{A}, \beta, \boldsymbol{\mu}) \quad (14)$$

Итерационные формулы EM-алгоритма¹:

① E-шаг:

$$p(z_{ik} = 1) = \frac{\exp \left(\mathbb{E} \log \pi_k - \frac{\beta}{2} \left[y_i^2 - 2y_i \mathbf{x}_i^\top \mathbb{E} \mathbf{w}_k + \mathbf{x}_i^\top (\mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top) \mathbf{x}_i \right] \right)}{\sum_k p(z_{ik} = 1)},$$

$$q(\boldsymbol{\pi}) = \text{Dir}(\boldsymbol{\pi} | \boldsymbol{\mu} + \boldsymbol{\gamma}), \quad q(\mathbf{w}_k) = \mathcal{N}(\mathbf{w}_k | \mathbf{m}_k, \mathbf{B}_k),$$

$$\gamma_k = \sum_{i=1}^N \mathbb{E} z_{ik}, \quad \mathbf{m}_k = \beta \mathbf{B}_k \left(\sum_{i=1}^N \mathbf{x}_i y_i \mathbb{E} z_{ik} \right), \quad \mathbf{B}_k = \left(\mathbf{A}_k^{-1} + \beta \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \mathbb{E} z_{ik} \right)^{-1}. \quad (15)$$

② M-шаг:

$$\begin{aligned} \mathbf{A}_k &= \mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top \\ \frac{1}{\beta} &= \frac{\sum \sum [y_i^2 - 2y_i \mathbf{x}_i^\top \mathbb{E} \mathbf{w}_k + \mathbf{x}_i^\top \mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i] \mathbb{E} z_{ik}}{\sum \sum \mathbb{E} z_{ik}} \end{aligned} \quad (16)$$

Некоторые математические ожидания:

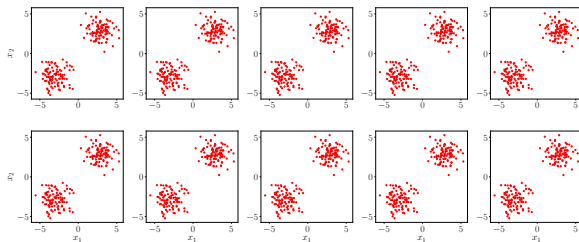
① $\mathbb{E} z_{ik} = p(z_{ik} = 1),$

② $\mathbb{E} \log \pi_k = \psi^0(\mu_k + \gamma_k) - \psi^0(K\mu_k + N),$

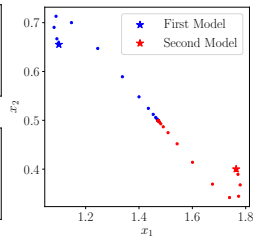
③ $\mathbb{E} \mathbf{w}_k \mathbf{w}_k^\top = \mathbf{B}_k + \mathbf{m}_k \mathbf{m}_k^\top.$

¹<https://github.com/andriygav/EMprior/blob/master/paper/Grabovoy2019Draft.pdf>

Пример 1



(a)



(b)

На рис. 1a показано, что в случае смеси моделей, предсказать то, какую модель использовать в каждой точке нельзя.

На рис. 1b показана зависимость векторов \mathbf{w}_k — параметры локальных моделей в процессе обучения.

¹<https://github.com/andriygav/MixtureLib>

Definition

Смесь экспертов — мультимодель, определяющая правдоподобие веса π_k каждой локальной модели \mathbf{f}_k на признаковом описании объекта \mathbf{x} .

$$\hat{\mathbf{f}} = \sum_{k=1}^K \pi_k \mathbf{f}_k, \quad \pi_k(\mathbf{x}, \mathbf{V}) : \mathbb{R}^{n \times |\mathbf{V}|} \rightarrow [0, 1], \quad \sum_{k=1}^K \pi_k(\mathbf{x}, \mathbf{V}) = 1 \quad (17)$$

где $\hat{\mathbf{f}}$ — мультимодель, а \mathbf{f}_k является некоторой моделью,
 π_k — параметрическая модель, \mathbf{w}_k — параметры k -й локальной модели,
 \mathbf{V} — параметры шлюзовой функции.

Пример 2:

- ❶ Правдоподобие k -й локальной модели $p_k(y_i | \mathbf{w}_k, \mathbf{x}_i) = \mathcal{N}(y_i | \mathbf{w}_k^T \mathbf{x}_i, \beta^{-1})$,
- ❷ Априорное распределение параметров k -й локальной модели $p^k(\mathbf{w}_k) = \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k)$,
- ❸ Шлюзовая функция $\pi(\mathbf{x}, \mathbf{V}) = \text{softmax}(\mathbf{V}_1^T \boldsymbol{\sigma}(\mathbf{V}_2^T \mathbf{x}))$.

Правдоподобие модели:

$$p(\mathbf{y}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \beta) = \prod_{k=1}^K \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k) \prod_{i=1}^N \left(\sum_{k=1}^K \pi_k \mathcal{N}(y_i | \mathbf{w}_k^\top \mathbf{x}_i, \beta^{-1}) \right). \quad (18)$$

Введем скрытые переменные $\mathbf{Z} = \{z_{ik}\}$, где $z_{ik} = 1 \Leftrightarrow k_i = k$:

$$p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \beta) = \prod_{k=1}^K \mathcal{N}(\mathbf{w}_k | \mathbf{w}_k^0, \mathbf{A}_k) \prod_{i=1}^N \prod_{k=1}^K \left(\pi_k \mathcal{N}(y_i | \mathbf{w}_k^\top \mathbf{x}_i, \beta^{-1}) \right)^{z_{ik}}. \quad (19)$$

Вариационный EM-алгоритм $q(\mathbf{Z}, \mathbf{W}) = q(\mathbf{Z}) q(\mathbf{W})$:

① E-шаг:

$$\begin{aligned} \log q(\mathbf{Z}^s) &\propto \mathbb{E}_{q/\mathbf{Z}} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}^{s-1}, \mathbf{A}^{s-1}, \mathbf{W}^{0,s-1}, \beta^{s-1}) \\ \log q(\mathbf{W}^s) &\propto \mathbb{E}_{q/\mathbf{W}} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}^{s-1}, \mathbf{A}^{s-1}, \mathbf{W}^{0,s-1}, \beta^{s-1}) \end{aligned} \quad (20)$$

② M-шаг:

$$\mathbf{W}^{0,s}, \mathbf{A}^s, \beta^s = \arg \max_{\mathbf{W}^0, \mathbf{A}, \beta} \mathbb{E}_{q^s} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \beta) \quad (21)$$

Итерационные формулы EM-алгоритма¹:

① E-шаг:

$$p(z_{ik} = 1) = \frac{\exp(\log \pi_k(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (\mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k))}{\sum_{k'=1}^K \exp(\log \pi_{k'}(\mathbf{x}_i, \mathbf{V}) - \frac{\beta}{2} (\mathbf{x}_i^\top \mathbf{E} \mathbf{w}_{k'} \mathbf{w}_{k'}^\top \mathbf{x}_i - \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_{k'}))},$$

$$q(\mathbf{w}_k) = \mathcal{N}(\mathbf{w}_k | \mathbf{m}_k, \mathbf{B}_k),$$

$$\mathbf{m}_k = \mathbf{B}_k \left(\mathbf{A}_k^{-1} \mathbf{w}_k^0 + \beta \sum_{i=1}^N \mathbf{x}_i y_i \mathbf{E} z_{ik} \right), \quad \mathbf{B}_k = \left(\mathbf{A}_k^{-1} + \beta \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^\top \mathbf{E} z_{ik} \right)^{-1}.$$
(22)

② M-шаг:

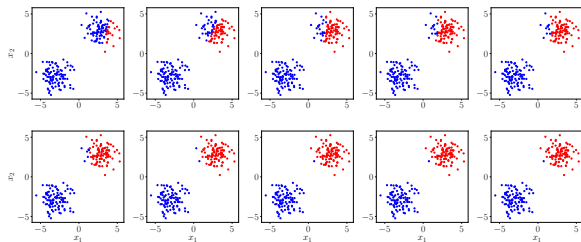
$$\mathbf{A}_k = \mathbf{E} \mathbf{w}_k \mathbf{w}_k^\top - \mathbf{w}_k^0 \mathbf{E} \mathbf{w}_k^\top - \mathbf{E} \mathbf{w}_k \mathbf{w}_k^{0\top} + \mathbf{w}_k^0 \mathbf{w}_k^{0\top},$$

$$\frac{1}{\beta} = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K \left[y_i^2 - 2 y_i \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k + \mathbf{x}_i^\top \mathbf{E} \mathbf{w}_k \mathbf{w}_k^\top \mathbf{x}_i \right] \mathbf{E} z_{ik},$$
(23)

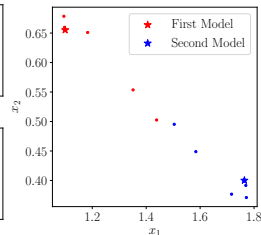
$$\mathbf{w}_k^0 = \mathbf{E} \mathbf{w}_k,$$

$$\mathbf{V} = \arg \max_{\mathbf{V}} \mathbf{E}_{q^s} \log p(\mathbf{y}, \mathbf{Z}, \mathbf{W} | \mathbf{X}, \mathbf{V}, \mathbf{A}, \mathbf{W}^0, \beta).$$

¹<https://github.com/andriygav/EMprior/blob/master/paper/Grabovoy2019MixtureOfExpert.pdf>



(a)

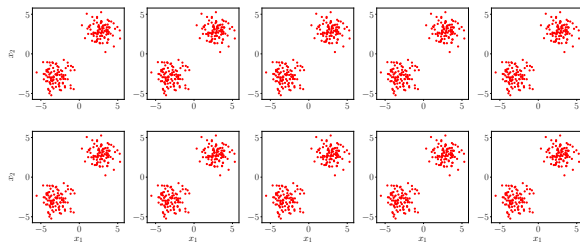


(b)

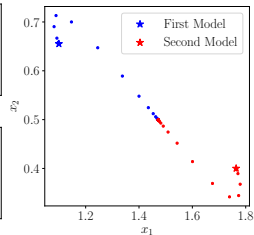
На рис. 1a показано, что в случае смеси экспертов гипермодель предсказывает к какому классу относится каждая точка в пространстве объектов.

На рис. 1b показана зависимость векторов \mathbf{w}_k — параметры локальных моделей в процессе обучения.

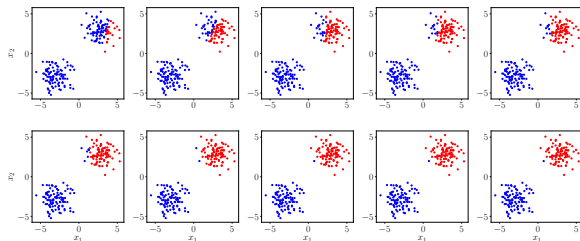
¹<https://github.com/andriygav/MixtureLib>



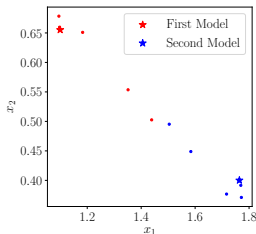
(a) Смесь Моделей



(b) Смесь Моделей



(c) Смесь Экспертов



(d) Смесь Экспертов

¹<https://github.com/andriygav/MixtureLib>

- ❶ Адуенко А. А., Курс лекций: Байесовский выбор моделей, 2018.
- ❷ Christopher Bishop, Pattern Recognition and Machine Learning, 2016.
- ❸ Yuksel Seniha Esen, Wilson Joseph N., Gader Paul D, Twenty Years of Mixture of Experts // IEEE Transactions on Neural Networks and Learning Systems. 2012. Issues. 23, No 8. pp. 1177–1193.
- ❹ Вывод вариационного ЕМ–алгоритма. <https://github.com/andriygav/EMprior/blob/master/Lecture/Grabovoy2019MeanField.pdf>
- ❺ Вывод смеси моделей с заданными априорными распределениями. <https://github.com/andriygav/EMprior/blob/master/paper/Grabovoy2019Draft.pdf>
- ❻ Вывод смеси экспертов с заданными априорными распределениями. <https://github.com/andriygav/EMprior/blob/master/paper/Grabovoy2019MixtureOfExpert.pdf>
- ❼ Програмная реализация мультимodelей. <https://github.com/andriygav/MixtureLib>