

# A variational hardcut EM algorithm for the mixtures of Gaussian processes

Tao LI & Jinwen MA\*

Department of Information and Computational Sciences, School of Mathematical Sciences and LMAM,  
Peking University, Beijing 100871, China

Received 16 June 2021/Revised 9 December 2021/Accepted 26 January 2022/Published online 10 February 2023

**Citation** Li T, Ma J W. A variational hardcut EM algorithm for the mixtures of Gaussian processes. Sci China Inf Sci, 2023, 66(3): 139103, https://doi.org/10.1007/s11432-021-3477-3

Gaussian process (GP) is the dominant model in the non-parametric Bayesian community, but it is not flexible enough to model non-stationary/multi-modal data, because a conventionally trained Gaussian process is stationary. To overcome this problem, the mixture of Gaussian processes (MGP) [1] was proposed. Compared with GP, the MGP model is more flexible, but its parameter learning is rather challenging. When applying the EM algorithm to the case of MGPs, the E-step is intractable due to two reasons. First, the posterior distributions of latent indicators are intractable because the samples are correlated. Second, even though the posterior of latent indicators is given, taking expectation with respect to these correlated latent indicators leads to exponential many summation terms. This study presents a variational hardcut EM (VHEM) algorithm to tackle the computational difficulties in the E-step, and theoretical analysis reveals that the VHEM algorithm plays the intermediate role between the hardcut EM algorithm [2] and MCMC-EM algorithm [3].

*Mixture of Gaussian processes.* Given a dataset  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  where  $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{R}$  and let  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]^T, \mathbf{y} = [y_1, y_2, \dots, y_N]^T$  for brevity. If  $y$  and  $\mathbf{x}$  are linked by a Gaussian process, then  $\mathbf{y}|\mathbf{X} \sim \mathcal{N}(\mathbf{m}, \mathbf{C})$ . For simplicity, we generally assume  $\mathbf{m}$  to be  $\mathbf{0}$ . Given a covariance function  $c(\cdot, \cdot; \boldsymbol{\theta})$  with parameters  $\boldsymbol{\theta}$ , the  $(i, j)$ -th element of covariance matrix  $\mathbf{C}$  is  $c(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta})$ . Here, we use the squared exponential covariance function  $c(\mathbf{x}_i, \mathbf{x}_j; \boldsymbol{\theta}) = \theta_0^2 \exp(-\sum_{l=1}^d \theta_l^2 \frac{(\mathbf{x}_{il} - \mathbf{x}_{jl})^2}{2}) + \sigma^2 \mathbb{I}(\mathbf{x}_i = \mathbf{x}_j)$  where  $\mathbb{I}$  is the indicator function, but extensions to the other covariance functions are straightforward. The parameters can be learned via the Type-II maximum likelihood estimation [4].

The mixture of Gaussian processes assumes there are  $K$  independent GP components. Latent variable  $z_i$  indicates the GP component that generates  $(\mathbf{x}_i, y_i)$ . The information flow can be characterized as  $z \rightarrow \mathbf{x} \rightarrow y$ . Given the mixing proportions  $\{\pi_k\}_{k=1}^K$ , we first sample  $z_i$  according to  $p(z_i = k) = \pi_k$ . Conditioned on  $z_i = k$ , the input  $\mathbf{x}_i$  is generated from the  $k$ -th multivariate normal distribution,  $\mathbf{x}_i|z_i = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ , where  $\boldsymbol{\mu}_k$  and  $\boldsymbol{\Sigma}_k$  are the mean

vector and the covariance matrix of the  $k$ -th multivariate normal distribution. Given  $\mathbf{z} = \{z_1, \dots, z_N\}$ , we can divide the samples according to the component labels. For each component, we assume  $\mathbf{x} \rightarrow y$  is generated by a Gaussian process. Let  $\mathbf{X}_k(\mathbf{z}) = \{\mathbf{x}_i|z_i = k, i = 1, \dots, N\}, \mathbf{y}_k(\mathbf{z}) = \{y_i|z_i = k, i = 1, \dots, N\}$ , and  $\mathbf{C}_k(\mathbf{z})$  be the covariance matrix of the  $k$ -th Gaussian process parameterized by  $\boldsymbol{\theta}_k$ , and then we have  $\mathbf{y}_k(\mathbf{z})|\mathbf{X}_k(\mathbf{z}) \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_k(\mathbf{z}))$ . These variables are dependent on  $\mathbf{z}$ , which is the main difficulty that makes the inference of MGP rather challenging.

The parameters  $\boldsymbol{\Theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k, \boldsymbol{\theta}_k\}_{k=1}^K$  are learned by the EM algorithm. The complete data log-likelihood is

$$\begin{aligned} \mathcal{L}(\boldsymbol{\Theta}, \mathbf{z}) &= \log p(\mathbf{X}, \mathbf{y}, \mathbf{z}; \boldsymbol{\Theta}) \\ &= \sum_{k=1}^K \left( \sum_{i=1}^N \mathbb{I}(z_i = k) [\log \pi_k + \log p(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)] \right. \\ &\quad \left. + \log p(\mathbf{y}_k(\mathbf{z})|\mathbf{X}_k(\mathbf{z}); \boldsymbol{\theta}_k) \right). \end{aligned} \quad (1)$$

In the E-step, we need to calculate the expectation of  $\mathcal{L}(\boldsymbol{\Theta}, \mathbf{z})$  with respect to the posterior  $p(\mathbf{z}|\mathbf{X}, \mathbf{y}; \boldsymbol{\Theta}^{\text{old}})$  to obtain the Q-function  $\mathcal{Q}(\boldsymbol{\Theta}|\boldsymbol{\Theta}^{\text{old}}) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\Theta}^{\text{old}}, \mathcal{D})}[\mathcal{L}(\boldsymbol{\Theta}, \mathbf{z})]$ . In the M-step, we optimize  $\mathcal{Q}(\boldsymbol{\Theta}|\boldsymbol{\Theta}^{\text{old}})$  to get new estimates of the parameters. The M-step is relatively easy, but the E-step is very challenging. Since the samples are not i.i.d. but correlated, the posterior  $p(\mathbf{z}|\mathbf{X}, \mathbf{y}; \boldsymbol{\Theta}^{\text{old}})$  is difficult to calculate and  $\mathcal{L}(\boldsymbol{\Theta}, \mathbf{z})$  involves  $K^N$  summation terms in the expectation, which is prohibitively large for computation.

Two kinds of methods [2,3] have been developed to tackle this problem. The hardcut EM algorithm [2] approximates the latent variables  $\mathbf{z}$  deterministically and ignores the dependency between the samples. Instead of approximating the posterior, the MCMC-EM algorithm [3] generates samples from the posterior using the Gibbs sampling technique and then approximates the expectation using these samples.

*VHEM algorithm.* To overcome the difficulties in the E-step, we employ the variational inference to approximate the posterior. We use a factorized distribution  $q(\mathbf{z}; \boldsymbol{\Lambda}) = \prod_{i=1}^N q(z_i; \boldsymbol{\lambda}_i)$  to approximate the  $p(\mathbf{z}|\mathbf{X}, \mathbf{y}; \boldsymbol{\Theta}^{\text{old}})$ . Here,

\* Corresponding author (email: jwma@math.pku.edu.cn)

$\Lambda = \{\lambda_i\}_{i=1}^N$ ,  $\lambda_i = \{\lambda_{i,k}\}_{k=1}^K$  and  $q(z_i; \lambda_i)$  is the probability mass function of a categorical distribution, i.e.,  $q(z_i = k) = \lambda_{i,k}$ . According to the mean-field variational inference theory [5], the optimal  $q(\mathbf{z}; \Lambda)$  satisfies

$$\lambda_{i,k} \propto \exp(\mathbb{E}_{q(\mathbf{z}_{-i}; \Lambda)}[\log p(\mathbf{X}, \mathbf{y}, \mathbf{z}_{-i} \cup \{z_i = k\}; \Theta)]), \quad (2)$$

where  $\mathbf{z}_{-i} = \mathbf{z} - \{z_i\}$ . Therefore, we can perform the fixed point iteration based on (2) to find the optimal  $q(\mathbf{z}; \Lambda)$ .

Although the components of  $\mathbf{z}_{-i}$  are independent, the complete log-likelihood is not separable with respect to  $\mathbf{z}_{-i}$ , thus the expectation inside (2) is intractable. We apply the hardcut approximation to side-step this problem. For  $j \neq i$ , we approximate  $q(z_j; \lambda_j)$  via a deterministic allocation  $\tilde{q}(z_j; \lambda_j) = \mathbb{I}(z_j = \arg \max_{k=1,2,\dots,K} \lambda_{j,k})$ . According to  $\tilde{q}(\mathbf{z}_{-i}; \Lambda)$ , the latent variables  $\mathbf{z}_{-i}$  are deterministic, and we write them as  $\tilde{\mathbf{z}}_{-i}$ . We can approximate the intractable expectation in (2) by  $\log p(\mathbf{X}, \mathbf{y}, \tilde{\mathbf{z}}_{-i} \cup \{z_i = k\}; \Theta)$ . Some calculation reveals Eq. (2) can be directly approximated by

$$\lambda_{i,k} \propto \pi_k p(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \frac{p(\mathbf{y}_k(\tilde{\mathbf{z}}_{-i} \cup \{z_i = k\}) | \mathbf{X}_k(\tilde{\mathbf{z}}_{-i} \cup \{z_i = k\}); \boldsymbol{\theta}_k)}{p(\mathbf{y}_{-i,k}(\tilde{\mathbf{z}}_{-i}) | \mathbf{X}_{-i,k}(\tilde{\mathbf{z}}_{-i}); \boldsymbol{\theta}_k)}, \quad (3)$$

where

$$\begin{aligned} \mathbf{X}_{-i,k}(\mathbf{z}) &= \{\mathbf{x}_j | z_j = k, j = 1, \dots, N \text{ and } j \neq i\}, \\ \mathbf{y}_{-i,k}(\mathbf{z}) &= \{y_j | z_j = k, k = 1, \dots, N \text{ and } j \neq i\}. \end{aligned}$$

Finally, we find out that although  $\lambda_i$  is important in deriving the algorithm, but it is not essential for implementation purpose. Instead, we can perform iterations on  $\tilde{\mathbf{z}}$  directly,

$$\begin{aligned} \tilde{z}_i &= \arg \max_{k=1,2,\dots,K} \pi_k p(\mathbf{x}_i; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &\frac{p(\mathbf{y}_k(\tilde{\mathbf{z}}_{-i} \cup \{z_i = k\}) | \mathbf{X}_k(\tilde{\mathbf{z}}_{-i} \cup \{z_i = k\}); \boldsymbol{\theta}_k)}{p(\mathbf{y}_{-i,k}(\tilde{\mathbf{z}}_{-i}) | \mathbf{X}_{-i,k}(\tilde{\mathbf{z}}_{-i}); \boldsymbol{\theta}_k)}. \quad (4) \end{aligned}$$

Once the above variational E-step iteration converges, we obtain an approximate posterior  $\tilde{q}(\mathbf{z}) = \mathbb{I}(\mathbf{z} = \tilde{\mathbf{z}})$  and we can calculate the approximate Q-function  $\tilde{Q}(\Theta; \Theta^{\text{old}}) = \mathbb{E}_{\tilde{q}(\mathbf{z})}[\mathcal{L}(\Theta, \mathbf{z})] = \mathcal{L}(\Theta, \tilde{\mathbf{z}})$ . Then we maximize  $\tilde{Q}(\Theta; \Theta^{\text{old}})$  with respect to  $\Theta$  to estimate the parameters. The entire algorithm is summarized in Algorithm A1 (see Appendix A).

*Comparisons with the other algorithms on the learning of MGPs.* The iteration formula (4) is very similar to the Gibbs sampling step in the MCMC-EM algorithm. The only difference is that the MCMC-EM algorithm samples  $z_i$  according to the probability, while the VHEM algorithm assigns  $z_i$  to be the class label with the highest probability.

In the hardcut EM algorithm, dependences among samples in the same Gaussian process component are ignored, and we only use  $p(y_i | \mathbf{x}_i; \boldsymbol{\theta}_k)$  to measure the probability that the  $i$ -th sample belonging to the  $k$ -th Gaussian process component. In this way, we do not need to perform iterations since  $\{z_i\}_{i=1}^N$  are independent. From (4), we can see that in the VHEM algorithm, when we calculate the probability that the  $i$ -th sample coming from the  $k$ -th component, other samples temporarily assigned with label  $k$  are also taken into consideration. Let

$$\begin{aligned} \mathbf{c} &= c(\mathbf{X}_{-i,k}(\tilde{\mathbf{z}}_{-i}), \mathbf{x}_i; \boldsymbol{\theta}_k), \\ \mathbf{C}_- &= c(\mathbf{X}_{-i,k}(\tilde{\mathbf{z}}_{-i}), \mathbf{X}_{-i,k}(\tilde{\mathbf{z}}_{-i}); \boldsymbol{\theta}_k), \end{aligned}$$

then the last term of (4) is given by

$$\log \frac{p(\mathbf{y}_k(\tilde{\mathbf{z}}_{-i} \cup \{z_i = k\}) | \mathbf{X}_k(\tilde{\mathbf{z}}_{-i} \cup \{z_i = k\}); \boldsymbol{\theta}_k)}{p(\mathbf{y}_{-i,k}(\tilde{\mathbf{z}}_{-i}) | \mathbf{X}_{-i,k}(\tilde{\mathbf{z}}_{-i}); \boldsymbol{\theta}_k)}$$

$$\begin{aligned} &= -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\theta_{k,0}^2 + \sigma_k^2 - \mathbf{c}^T \mathbf{C}_-^{-1} \mathbf{c}) \\ &\quad - \frac{(\mathbf{y}_{-i,k}(\tilde{\mathbf{z}}_{-i})^T \mathbf{C}_-^{-1} \mathbf{c} - y_i)^2}{2(\theta_{k,0}^2 + \sigma_k^2 - \mathbf{c}^T \mathbf{C}_-^{-1} \mathbf{c})}. \quad (5) \end{aligned}$$

The hardcut EM algorithm can be regarded as a further approximation of the VHEM algorithm when  $\mathbf{c}^T \mathbf{C}_-^{-1}$  is replaced by  $\mathbf{0}$ . This is not realizable because  $\mathbf{C}_-^{-1}$  is always invertible, and  $\mathbf{c}^T \mathbf{C}_-^{-1} = \mathbf{0}$  implies  $\mathbf{c} = \mathbf{0}$ , which cannot hold for general covariance functions. However, when  $\mathbf{x}_i$  is far from the points in  $\mathbf{X}_-$ , we may expect  $\mathbf{c}^T \mathbf{C}_-^{-1} \approx \mathbf{0}$ . For the case  $\mathbf{x}_i$  lies in the high-probability region of the  $l$ -th component, the difference between the VHEM algorithm and the hardcut EM algorithm for computing  $\lambda_{i,k}, k \neq l$  is negligible. If  $\mathbf{x}_i$  lies in the overlapping region of the  $k$ -th component and the  $l$ -th component, the difference would be significant.

Furthermore, Eq. (5) also presents an intuitive interpretation on the VHEM algorithm from the perspective of leave-one-out cross validation [6]. See Appendix B for more detailed discussion.

*Experiments.* See Appendix C.

*Conclusion.* We have proposed a new kind of variational inference based learning algorithm (VHEM) for the generative MGP model. The main advantages of the VHEM algorithm are three folds. First, compared with the hardcut EM algorithm, its derivation relies on variational inference; thus its theory is more solid and sound. Second, it connects existing learning algorithms for MGP, including the hardcut EM algorithm, the MCMC-EM algorithm, and the LOOCV algorithm. Last, it is remarkably faster than the MCMC-EM algorithm and significantly more accurate than the hardcut EM algorithm. The VHEM algorithm is able to achieve comparable performances with the MCMC-EM algorithm, with the cost of a little longer running times compared with the hardcut EM algorithm. To balance performance and computational cost, the VHEM algorithm is a good choice for real applications.

**Acknowledgements** This work was supported by National Key Research and Development Program of China (Grant No. 2018AAA0100205).

**Supporting information** Appendixes A–C. The supporting information is available online at [info.scichina.com](http://info.scichina.com) and [link.springer.com](http://link.springer.com). The supporting materials are published as submitted, without typesetting or editing. The responsibility for scientific accuracy and content remains entirely with the authors.

**References**

- 1 Tresp V. Mixtures of Gaussian processes. In: Proceedings of Advances in Neural Information Processing Systems, 2001. 654–660
- 2 Chen Z Y, Ma J W, Zhou Y T. A precise hard-cut EM algorithm for mixtures of Gaussian processes. In: Proceedings of International Conference on Intelligent Computing. Berlin: Springer, 2014. 68–75
- 3 Wu D, Ma J W. An effective EM algorithm for mixtures of Gaussian processes via the MCMC sampling and approximation. *Neurocomputing*, 2019, 331: 366–374
- 4 Williams C K, Rasmussen C E. *Gaussian Processes for Machine Learning*. Cambridge: MIT Press, 2006
- 5 Bishop C M. *Pattern Recognition and Machine Learning*. Berlin: Springer, 2006
- 6 Yang Y, Ma J W. An efficient EM approach to parameter learning of the mixture of Gaussian processes. In: Proceedings of International Symposium on Neural Networks. Berlin: Springer, 2011. 165–174