

# Automatic Model Selection of the Mixtures of Gaussian Processes for Regression

Zhe Qiang and Jinwen Ma

Department of Information Science, School of Mathematical Sciences and LMAM  
Peking University, Beijing, 100871, China  
jwma@math.pku.edu.cn

**Abstract.** For the learning of mixtures of Gaussian processes, model selection is an important but difficult problem. In this paper, we develop an automatic model selection algorithm for mixtures of Gaussian processes in the light of the reversible jump Markov chain Monte Carlo framework for Gaussian mixtures. In this way, the component number and the parameters are updated according the five types of random moves and model selection can be made automatically. The key idea is that the moves of component splitting or merging preserve the zeroth, first and second moments of the components so that the covariance parameters of the new components can be related to the origin ones. It is demonstrated by the simulation experiments that this automatic model selection algorithm is feasible and effective.

**Keywords:** Mixtures of Gaussian processes, Reversible jump MCMC, Model selection, Regression, Split and merge moves.

## 1 Introduction

As a powerful statistical learning tool, Gaussian Process (GP) is widely used in machine learning and pattern recognition [1]. Since a single GP model cannot deal with the multimodality dataset, the Mixture of Gaussian Processes (MGP) [2] have been developed to model a general multimodality dataset. Obviously, MGP can also overcome a major disadvantage of the Gaussian process modeling that calculating the inversion of an  $n \times n$  covariance matrix requires the time complexity of  $O(n^3)$  for a training dataset with  $n$  points. Structurally, Shi et al. [3],[4] considered MGP as a hierarchical model and fit the data in two levels. Moreover, they proposed the hybrid Markove Chain Monte-Carlo algorithm to train the covariance parameters and then to utilize BIC to determine the number of GP components in the mixture. However, BIC is not so effective for the model selection on the mixtures of Gaussian processes.

On the other hand, Green [5] proposed the reversible jump Markov chain Monte Carlo (RJMCMC) framework to determine the dimension of parameters through the Markov chain Monte Carlo simulation. Later on, according to this theory, Richardson and Green [6] developed a RJMCMC approach to decide the number of actual components in the Gaussian mixture. Although this approach

is effective for Gaussian mixtures, its idea can also be used to design the learning algorithm for MGPs and solve the model selection problem in this case. However, the structure of the likelihood function of MGP is quite different from that of Gaussian mixture so that the split and merge moves in the algorithm cannot be implemented directly.

In this paper, we develop an automatic model selection algorithm for MGPs in the light of the RJMCMC framework by making the split and merge moves feasible and effective. As we find, the difficulty in the split moves focuses on that the overall dispersion should keep relatively constant when a covariance matrix is split into two matrices. In order to solve this difficulty, we can keep the first two moments of the components to be the same during a split or merge move. By mathematical analysis, we find that once the sampling interval is set small enough, these moments remain almost constant. In this way, our automatic model selection algorithm can be constructed effectively. It is demonstrated by simulation experiments that this automatic model selection algorithm is feasible and effective in the same way as the RJMCMC algorithm for Gaussian mixtures.

The rest of the paper is organized as follows. In Section 2, we revisit the mixture of Gaussian processes and give the Bayesian forms of the priors of the parameters. Section 3 presents the automatic model selection algorithm with five move types. Simulation experiments are conducted in Section 4. We conclude briefly in Section 5.

## 2 The Hierarchical Mixtures of Gaussian Processes

### 2.1 The Basic Model

For clarity, we consider MGP as the hierarchical mixture of Gaussian processes described by Shi et al. [3]. Clearly, it tries to model a large dataset with groups of repeated measurements. Actually, the lower-level model fits the measurements in the same group, while the higher-level model tries to describe the heterogeneity among different groups. For example, the dataset used in [4] is a set of repeated standing-up trajectories corresponding to different paraplegic patients. For this case, the lower-level model fitted those standing-up trajectories of paraplegic patients in the same type, while the higher-level model described the heterogeneity among different types of paraplegic patients.

Mathematically, in the higher-level model, we let  $\{(\mathbf{y}_{mn})_{n=1}^N\}_{m=1}^M$  be  $M$  curves belong to a number of Gaussian processes, which can represent all the points on the training curves. The mixture model can be given by

$$\mathbf{y}_m \sim \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}_m | \boldsymbol{\theta}_k). \quad (1)$$

In the lower-level model, each curve is assumed to be a function  $f_m(\cdot)$  plus a white noise, that is,

$$\mathbf{y}_{mn} = f_m(\mathbf{x}_{mn}) + \boldsymbol{\epsilon}_{mn}$$

where  $\boldsymbol{\epsilon}_{mn} \sim \mathcal{N}(0, \Sigma_k)$ . The  $n$ -th curve follows a Gaussian process if any finite subsequence or subset of the curve, say  $\mathcal{D}_m = \{\mathbf{X}_m, \mathbf{Y}_m\} = \{(\mathbf{x}_{mn}, \mathbf{y}_{mn})\}_{n \in \mathcal{D}_m}$

$\{1 \dots m\}$ , follows a Gaussian distribution if  $\mathbf{Y}_m \sim \mathcal{N}(\mathbf{0} \ \boldsymbol{\Sigma}(\mathbf{X}_m \ \mathbf{X}_m; \boldsymbol{\theta}_k))$ , where

$$\boldsymbol{\Sigma}(\mathbf{x}_i \ \mathbf{x}_j; \boldsymbol{\theta}_k) = (\mathbf{x}_i \ \mathbf{x}_j; \boldsymbol{\theta}_k) + \frac{2}{k} \delta_{ij}$$

where

$$(\mathbf{x}_i \ \mathbf{x}_j; \boldsymbol{\theta}_k) = \boldsymbol{\theta}_k \exp\left(-\frac{1}{2} \boldsymbol{\theta}_k (\mathbf{x}_i - \mathbf{x}_j)^2\right)$$

So, all the covariance parameters are  $\boldsymbol{\theta}_k = (\boldsymbol{\theta}_k \ \frac{2}{k})$ .

For a test input  $\mathbf{x}^*$ , if we assume it is on the  $k$ -th curve, and the  $k$ -th curve is belong to the  $k$ -th component, i.e., Gaussian process, the mean and variance of its prediction can be obtained as follows:

$$\begin{aligned} \mathbb{E}[y_m(\mathbf{x}^*)|\mathcal{D}_m] &= \boldsymbol{\sigma}^T(\mathbf{x}^*)^{-1}(\mathbf{X}_m \ \mathbf{X}_m; \boldsymbol{\theta}_k)\mathbf{Y}_m; \\ \text{var}[y_m(\mathbf{x}^*)|\mathcal{D}_m] &= (\mathbf{x}^* \ \mathbf{x}^*) - \boldsymbol{\sigma}^T(\mathbf{x}^*)^{-1}(\mathbf{X}_m \ \mathbf{X}_m; \boldsymbol{\theta}_k)\boldsymbol{\sigma}(\mathbf{x}^*) \end{aligned}$$

where  $\boldsymbol{\sigma}(\mathbf{x}^*) = (\sigma(\mathbf{x}^*, m, 1) \dots \sigma(\mathbf{x}^*, m, N_m))^T$ .

For the parameter learning of this hierarchical mixture model of Gaussian process, we introduce the latent variables  $z_m$  as follows:

$$y_m(\mathbf{X}_m)|z_m = \sum_k z_k y_m(\mathbf{x}_k)$$

and assume the mixing form of Eq.(1).

### 2.2 The Priors and its Bayesian Form

In the hierarchical mixtures of Gaussian processes, we adopt the forms of the priors as given in [4], that is,

$$\boldsymbol{\theta}_k \sim \text{Gamma}\left(\frac{1}{2} \ \frac{1}{2}\right) \ \boldsymbol{\theta}_k \sim \mathcal{LN}(-1 \ 1^2) \ \frac{2}{k} \sim \mathcal{LN}(-3 \ 3^2) \ \pi = 1 \dots$$

where  $\text{Gamma}$  represents the inverse gamma distribution, and

$$(\pi_1 \dots \pi_K) \sim \text{Dir}(1 \dots 1)$$

In the Bayesian analysis, the priors are as important as the posteriors and likelihood function. Supposing that  $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1 \dots \boldsymbol{\theta}_K)$ ,  $\boldsymbol{\pi} = (\pi_1 \dots \pi_K)$ ,  $\mathcal{D}$  is the set of training data, we then have

- (i) The posterior of the parameters:

$$(\boldsymbol{\Theta} \ \boldsymbol{\pi}|\mathcal{D}) \propto (\boldsymbol{\Theta} \ \boldsymbol{\pi}) (\mathcal{D}|\boldsymbol{\Theta} \ \boldsymbol{\pi});$$

- (ii) The likelihood of the mixture model:

$$(\mathcal{D}|\boldsymbol{\Theta} \ \boldsymbol{\pi}) = \prod_{m=1}^M \sum_{k=1}^K z_k \mathcal{N}(y_m|\mathbf{0} \ \boldsymbol{\Sigma}_k(\mathbf{x}_m));$$

- (iii) The prior distribution:

$$(\boldsymbol{\Theta} \ \boldsymbol{\pi}) = (\boldsymbol{\pi}) \prod_{k=1}^K (\boldsymbol{\theta}_k)$$

In the next section, we will sample parameters from  $(\boldsymbol{\Theta} \ \boldsymbol{\pi}|\mathcal{D})$ .

### 3 The Markov Chain Monte Carlo Algorithm for the Hierarchical Mixtures of Gaussian Processes

In this section, we construct the Markov Chain Monte Carlo (MCMC) algorithm for the hierarchical mixtures of Gaussian processes in a similar way as the MCMC algorithm for Gaussian mixtures in [6]. Actually, the mathematical framework of our MCMC algorithm keeps the same, including a number of moves of the components or their parameters with time, but the components become Gaussian processes.

#### 3.1 The Move Types

In [6], six types of moves were used. But here, since the adopted priors have no hyperparameters, there is no need for updating the hyperparameters. So, remaining five types of moves can be given as follows:

- (a)  $\boldsymbol{\pi} = (\pi_1 \cdots \pi_K)$ ;
- (b)  $\boldsymbol{\Theta} = (\boldsymbol{\theta}_1 \cdots \boldsymbol{\theta}_K)$ , where  $\boldsymbol{\theta}_k = (\alpha_k, \beta_k, \gamma_k)$ ;
- (c)  $\mathbf{z} = (z_1 \cdots z_M)$ ;
- (d) Split or merge;
- (e) Remove empty components.

For clarity, we refer to one process of implementing these five moves completely as a sweep, being set as the basic step of the MCMC algorithm.

#### 3.2 The Moves with the Component Number Fixed

The moves with fixed component number include the first three steps(a)(b)(c) in section 3.1. For this part, we adopt the algorithm in [4]:

For  $\boldsymbol{\pi}$  and  $\mathbf{z}$ , we use Gibbs sampling:

- (i) sample  $z_m$  from  $(z_m = k | \mathcal{D}, \boldsymbol{\Theta}, \boldsymbol{\pi}) \propto \pi_k \prod_{m: z_m = k} \mathcal{N}(\mathbf{y}_m | \boldsymbol{\theta}_k)$ ,  $k = 1 \cdots K$ ;
- (ii) sample  $\boldsymbol{\pi}$  from  $(\pi_1 \cdots \pi_K | \mathbf{z}) \sim \text{Dir}(\pi_1 + 1, \dots, \pi_K + 1)$ .

Here,  $\pi_k$  represents the element number in the set  $\{z_m = k | z_m = 1 \cdots K\}$  and  $\text{Dir}(\cdot)$  represents Dirichlet distribution.

And for  $\boldsymbol{\Theta}$ , we sample it from its posterior:

$$(\boldsymbol{\Theta} | \mathcal{D}, \mathbf{z}) \propto \prod_{k=1}^K (\boldsymbol{\theta}_k | \mathcal{D}_m, \mathbf{z})$$

Then, the posterior of  $\boldsymbol{\theta}_k$ ,  $k = 1 \cdots K$  are independent with each other and we can sample each  $\boldsymbol{\theta}_k$  separately. Here we adopt Hybrid Monte Carlo or Hamiltonian Monte Carlo to sample  $\boldsymbol{\theta}_k$ . Actually, this sampling method used to be adopted to simulate a physical system where a puck moves up and down along with a smooth surface and the total energy remains constant. However, in our application, we restate the potential energy as  $\mathcal{E}(\boldsymbol{\theta}_k) = -\log \pi(\boldsymbol{\theta}_k | \mathcal{D}_m, \mathbf{z})$  and the

kinetic energy as:  $\mathcal{K}(\phi_k) = \frac{1}{2} \sum_{k,i} \frac{2}{k,i} \cdot \phi_{k,i} \sim \mathcal{N}(0, 1)$   $\phi_k = (\phi_{k,1}, \phi_{k,2}, \phi_{k,3})$ . Then the total energy of the hamiltonian system for  $\theta_k$  is  $\mathcal{H}(\theta_k, \phi_k) = \mathcal{E}(\theta_k) + \mathcal{K}(\phi_k)$ .

For the convenience of calculation, we split  $\mathcal{H}(\theta_k, \phi_k)$  as follows:

$$\mathcal{H}(q, p) = \underbrace{-\frac{1}{2} \sum_{m \in \{z_m=k\}} \log p(\mathbf{y}_m | \theta_k)}_{\frac{U_0}{2}} + \underbrace{[-\log p(\theta_k) + \mathcal{K}(\phi_k)]}_{U_1} - \underbrace{\frac{1}{2} \sum_{m \in \{z_m=k\}} \log p(\mathbf{y}_m | \theta_k)}_{\frac{U_0}{2}}$$

thus the sample update step is:

- (i) From the current state  $(\theta_k, \phi_k)$ , we use a leapfrog step with step size  $\epsilon$  to calculate the new state  $(\theta_k^*, \phi_k^*)$ :

$$\phi_k^* = \phi_k - \frac{\epsilon}{2} \nabla_{\theta_k} \mathcal{H}(\theta_k, \phi_k)$$

$$\theta_k^* = \theta_k + \epsilon \nabla_{\theta_k} \mathcal{H}(\theta_k, \phi_k)$$

$$\phi_k^* = \phi_k^* - \frac{\epsilon}{2} \nabla_{\theta_k^*} \mathcal{H}(\theta_k^*, \phi_k^*)$$

$$\theta_k^* = \theta_k^* - \frac{\epsilon}{2} \nabla_{\theta_k^*} \mathcal{H}(\theta_k^*, \phi_k^*)$$

- (ii) Then new current state is:

$$(\theta_k^*, \phi_k^*) = \begin{cases} (\theta_k^*, \phi_k^*) & r < \exp(-\mathcal{H}(\theta_k^*, \phi_k^*) + \mathcal{H}(\theta_k, \phi_k)) \\ (\theta_k, \phi_k) & \text{otherwise} \end{cases}$$

- (iii) Finally, generate  $\phi_{k,i} \sim \mathcal{N}(0, 1)$ , and update  $\phi_k$  as:  $\phi_k^* = \phi_k + \sqrt{1 - \epsilon^2} \nu_k$ .

In addition, according to the advise in [7], we set  $\epsilon = 0.5 \cdot \frac{1}{m^{1/2}} = 0.95$ .

For  $k = 1 \dots K$ , by repeating (i)(ii)(iii)  $M = 20$  times, we finish one sweep of step(b). Note that  $\phi_{k,i} \sim \mathcal{N}(0, 1)$ , through the analysis of the method of handling constraints in [8], we reject the current state when  $\phi_{k,i} < 0$ .

### 3.3 The Moves with the Component Number Changed

The moves with the component number changed include the last two steps(d)(e). For move(e), all the components whose  $\phi_{k,i} < 1\%$  are considered to be empty components and we delete them directly.

**The Detailed Balance Framework.** For move(d), we construct a detailed balance framework for the move with the component number changed so that the covariance parameters of the new components can be related to the original ones. The key to success is that the first two moments remain almost constant. For convenience we shall denote the splitted component as  $k^*$ th component and the two new components as  $k_1, k_2$ .

- (I) Actually, for the zeroth moment of  $\mathbf{y}_m$ , by simple mathematics calculation we find it is  $\sum_{k=1}^K \phi_{k,i}$ , thus keeping the zeroth moment constant means  $\phi_{k^*} = \phi_{k_1} + \phi_{k_2}$ .

- (II) In the case of the first moment, it is  $\mathbf{0}$  since we assume that the Gaussian process is zero mean. Thus its first moment always keep constant.

(III) However, the second moment is somewhat complicated and it is  $\sum_{k=1}^K k \cdot k$ . Then keeping the second moment constant means  $k^* = k_1 + k_2$

(i) For the diagonal entry, let  $\cdot = \cdot$ , then we get a balance formula:

$$k^* \binom{k^*}{k^*} = k_1 \binom{k_1}{k_1} + k_2 \binom{k_2}{k_2}$$

(ii) For the off diagonal entry, let  $\cdot \neq \cdot$ , we have:

$$k^* \binom{k^*}{k^*} \exp\left(-\frac{\binom{k^*}{k^*} (i-j)^2}{2}\right) = k_1 \binom{k_1}{k_1} \exp\left(-\frac{\binom{k_1}{k_1} (i-j)^2}{2}\right) + k_2 \binom{k_2}{k_2} \exp\left(-\frac{\binom{k_2}{k_2} (i-j)^2}{2}\right) \quad (2)$$

For convenience, denote Eq.(2) as a function of  $\binom{i-j}{i-j}$ :

$$\cdot(\cdot) := k^* \binom{k^*}{k^*} \exp\left(-\frac{\binom{k^*}{k^*}}{2}\right) - k_1 \binom{k_1}{k_1} \exp\left(-\frac{\binom{k_1}{k_1}}{2}\right) - k_2 \binom{k_2}{k_2} \exp\left(-\frac{\binom{k_2}{k_2}}{2}\right) \equiv 0 \quad (3)$$

In fact,  $\cdot(\cdot)$  is exponential decreasing about  $\cdot$  and  $\cdot$  is in  $\{2, 2^2, 2^3, \dots, 2^m\}$ , thus we can conclude that when  $\cdot$  gets the minimum  $2^m$ ,  $\cdot(\cdot)$  reaches the maximum. Therefore, in order to keep  $\cdot(\cdot) \approx 0$  we just need keep  $\cdot(2^m) \approx 0$ . Do Taylor's expansion of  $\cdot(\cdot)$  at  $\cdot = 0$  and denote the n-th term as  $\cdot_n$  for convenience. What surprise us is that when  $\cdot = 2^m$  and if  $\cdot$  is small enough,  $\cdot_n$  will be monotonic decreasing, then  $\cdot_0$  will be the largest term and  $\cdot_1$  will be the second largest term. Therefore just let  $\cdot_0 = \cdot_1$  be zero, the second moment will keep almost constant.

Then, we have got our detailed balance framework:

$$k^* = k_1 + k_2 \quad (4a)$$

$$k^* \binom{k^*}{k^*} = k_1 \binom{k_1}{k_1} + k_2 \binom{k_2}{k_2} \quad (4b)$$

$$k^* \binom{k^*}{k^*} = k_1 \binom{k_1}{k_1} + k_2 \binom{k_2}{k_2} \quad (4c)$$

$$k^* \binom{k^*}{k^*} \binom{k^*}{k^*} = k_1 \binom{k_1}{k_1} \binom{k_1}{k_1} + k_2 \binom{k_2}{k_2} \binom{k_2}{k_2} \quad (4d)$$

**The Merge and Split Formula.** For merge move, we can use Eq.(4) to merge two components, so Eq.(4) is also our merge formula. For split move, reversible jump in [5] is needed. According to the theory above, 4 dimensions random variable  $\mathbf{u} = (u_1, u_2, u_3, u_4)$  need to be generated to decide these new parameters:  $u_i \sim \text{Beta}(2, 2), i = 1 \dots 4$ . By combining these random parameters, the balance formula Eq.(4) and the reversible jump theory, we can write the split formula as:

$$k_1 = u_1 k^* \quad k_2 = (1 - u_1) k^* \quad u_1 \in (0, 1) \quad (5a)$$

$$\frac{k_1}{k_1} = u_2 \frac{k^*}{k_1} \quad \frac{k_2}{k_2} = (1 - u_2) \frac{k^*}{k_2} \quad u_2 \in (0, 1) \quad (5b)$$

$$\binom{k_1}{k_1} = u_3 \binom{k^*}{k_1} \quad \binom{k_2}{k_2} = (1 - u_3) \binom{k^*}{k_2} \quad u_3 \in (0, 1) \quad (5c)$$

$$\binom{k_1}{k_1} \binom{k_1}{k_1} = \frac{1 - u_4}{u_3} \binom{k^*}{k_1} \binom{k_1}{k_1} \quad \binom{k_2}{k_2} \binom{k_2}{k_2} = \frac{u_4}{1 - u_3} \binom{k^*}{k_2} \binom{k_2}{k_2} \quad u_4 \in (0, 1) \quad (5d)$$

It is easy to validate that split formula Eq.(5) matches with merge formula Eq.(4).

The split and merge move is a Markov birth-death chain, and we set the split probability and merge probability as  $p_{k \rightarrow k+1} = 1 - p_{k+1 \rightarrow k}$  respectively, depending on  $k$ :  $p_{1 \rightarrow 0} = 0, p_{k_{max} \rightarrow k_{max}+1} = 0, p_{k \rightarrow k+1} = p_{k+1 \rightarrow k} = 0.5 \forall k = 2 \cdots k_{max} - 1$ , where  $k_{max}$  is the maximum component number that we set according to individual cases. Then based on the acceptance probability formula of reversible jump move in [5] the acceptance ratio for a split move is  $\min(1, \frac{p_{k+1 \rightarrow k} q_{k+1}(\theta_k)}{p_{k \rightarrow k+1} q_k(\theta_{k+1})})$  and a merge move is  $\min(1, \frac{p_{k \rightarrow k+1} q_k(\theta_{k+1})}{p_{k+1 \rightarrow k} q_{k+1}(\theta_k)})$  where

$$= \prod_{m=1}^M \frac{(\mathbf{Y}_m | \theta_{k+1})}{(\mathbf{Y}_m | \theta_k)} \times \frac{p_{k+1 \rightarrow k}}{p_{k \rightarrow k+1}} \times \underbrace{\frac{(\theta_{k+1})}{(\theta_k)}}_{T_1} \times \underbrace{\frac{1}{\text{Beta}(\mathbf{u} | \theta_{k+1}, \theta_k)}}_{T_2} \times \underbrace{\left| \frac{\theta_{k+1}}{(\theta_k, \mathbf{u})} \right|}_{T_3}$$

Here,  $\theta_k^*$   $\theta_{k+1}^*$  are chosen randomly from the  $k$  components.

### 4 Experimental Results

Our experiments are conducted on a set of simulated data generated from a mixture of three Gaussian processes with  $\mu_1 = (1.0, 0.2, 0.0025)$ ,  $\mu_2 = (0.5, 1.0, 0.001)$ ,  $\mu_3 = (10, 0.2, 0.0005)$ , respectively, each with 3 curves. The data points are with  $x = -4 : 0.08 : 4$ , being equally spaced. We have two kinds of prediction, type I prediction: choosing half data randomly from each of the 9 curves as training data, the rest as test data; type II prediction: generating a new curve from the first GP, choosing half data randomly on this curve as known data, using this half and training parameters from the type I prediction to simulate the other unknown half data. In the later two subsections, the two type of prediction is used to verify the component number fixed algorithm and the component number varied algorithm for regression and model selection. For both of the two algorithms, we run them for 20000 sweep. Here we discard the first 10000 iterations and select one sample from each 200 iterations, in order to have approximately independent draws.

#### 4.1 The Component Number Fixed

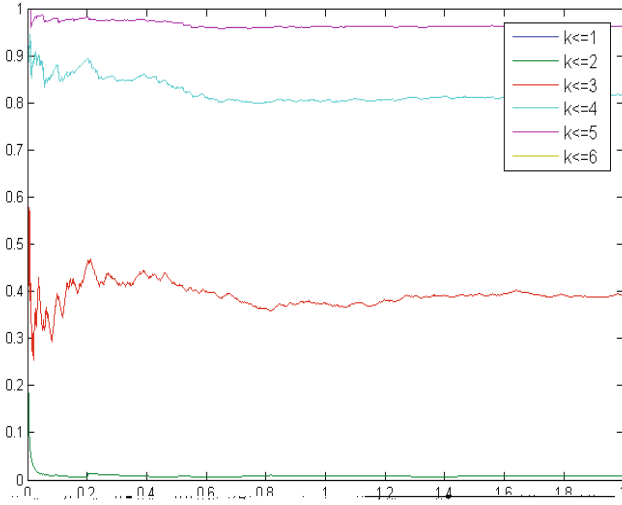
First, we fix  $k = 3$ , and use only the moves with the component number fixed for training. The log-likelihood tends to stabilize after about 1700 iteration. Table (1) presents the predictions. We will compare it with the component number varied result in the next subsection.

#### 4.2 The Component Number Varied

Then we use our component number varied method to train and predict the same curves as before. Fig.(1) is  $(\mu_1, \mu_2, \mu_3 | \mathcal{D}) = 1 \cdots 6$ , from which we can conclude

**Table 1.** RMSE and correlation coefficient( $r$ ) between true and predicted responses

Training data:half data on the first 9 curves		
model: fixed component number of GP mixture regression model		
Test data	RMSE	r
the first GP	0.2329	0.9485
the second GP	0.3612	0.8520
the third GP	0.2266	0.9024
the 10th curve	0.2037	0.7566
Training data:half data on the first 9 curves		
model: varied component number of GP mixture regression model		
Test data	RMSE	r
the first GP	0.0602	0.9900
the second GP	0.0253	0.9982
the third GP	0.0645	0.9867
the 10th curve	0.0493	0.9820



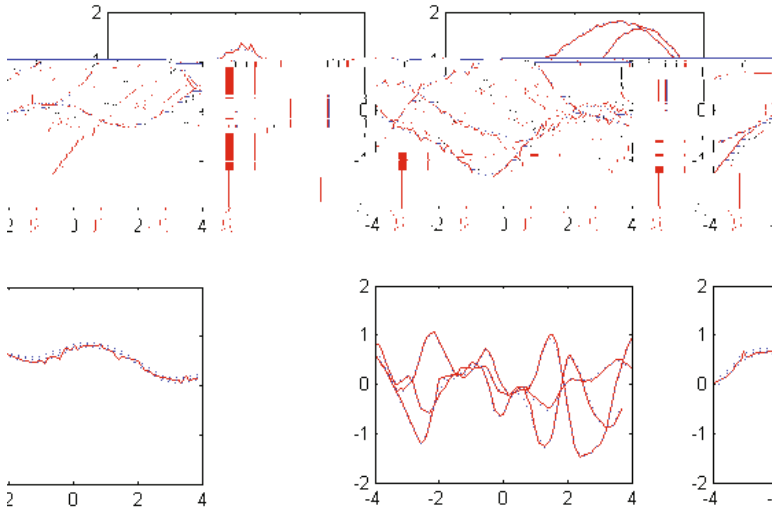
**Fig. 1.**  $p(k < j|\mathcal{D}), j = 1, \dots, 6$  for 40000 iterations

**Table 2.** posterior of  $k$  for Gaussian process mixture model

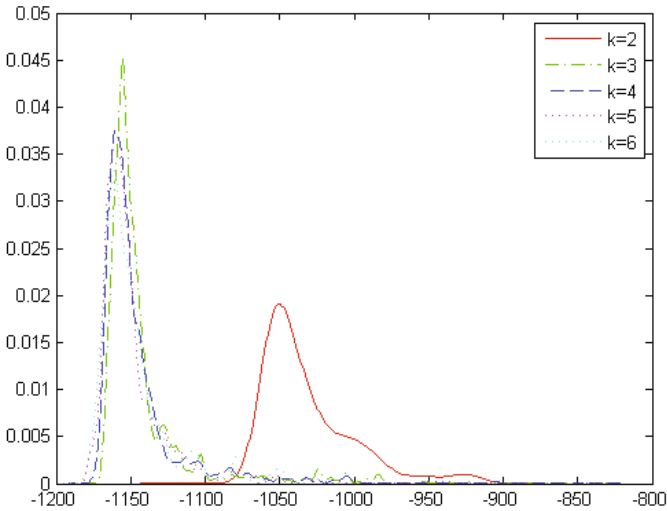
curve number	$p(k \mathcal{D})$			proportion (%) of moves accepted	
				split	merge
9	$p(1)=0.0001$	$p(2) = 0.0563$	$p(3)=0.3811$	14	2
	$p(4)=0.3677$	$p(5) = 0.1558$	$\sum_{k>6} p(k) = 0.0391$		

that the algorithm has converged after about 10000 iteration. Additionally, we also present the posterior of  $k$  for mixtures of Gaussian processes in Table (2). From this table, we can conclude this model favors 3 – 4 components. In this





**Fig. 2.** predict curve with the component number varied method on 9 curves, the solid line represent the real value, and the dash line represents the predict line with 95% confidence interval



**Fig. 3.** the deviance of the kernel smooth density posterior of  $k$

part, we initially add the birth and death moves in [6] to the algorithm in order to increase the opportunities for split and merge thus to speed up convergence of Markov chain, but finally it doesn't work, and we delete this type move.

The final predict results are presented in Fig.(2). We can see that predict curves almost overlap with the true curves. For type II prediction, the 95%

confidence interval is quite small. RMSE and correlation coefficient are in Table (1). The RMSE is about one tenth of that for the component number fixed case and the correlation coefficient  $\geq 0.9800$  which is also improved.

In addition, we analyze the deviance of the kernel smooth density of posterior of  $-2 * \log (|\mathcal{D}|)$  in Fig.(3). Since the deviance of  $k=2$  separate with  $k=3$  and from  $k=3$  to  $k=6$  the deviances are overlapping together, we can conclude that  $k=3$  can represent most information of train data.

Increasing the curve number to 30 and even to 90, we find that it does not improve the RMSE and correlation coefficient obviously. Therefore, this become a burden to waste time but of no use for improving algorithm accuracy and is unnecessary.

## 5 Conclusion

We have developed an automatic model selection algorithm for mixtures of Gaussian processes in the light of the reversible jump Markov chain Monte Carlo framework for Gaussian mixtures. The split and merge moves of Gaussian processes in the iteration keep the first two moments constant. In this way, the automatic model selection algorithm makes it possible to do the Bayesian analysis of both parameter estimation and model selection for the mixtures of Gaussian processes. Moreover, it is demonstrated that this developed algorithm is feasible and outperform the the hybrid MCMC algorithm.

**Acknowledgments.** This work was supported by the National Science Foundation of China under Grant 61171138.

## References

1. Rasmussen, C.E., Williams, C.K.I.: Gaussian process for machine learning. MIT Press, Cambridge (2006)
2. Tresp, V.: Mixtures of Gaussian processes. In: Proc. of the Conf. on Neural Information Processing Systems (NIPS), pp. 654–660 (2000)
3. Shi, J.Q., Murray-Smith, R., Titterington, D.M.: Bayesian regression and classification using mixtures of Gaussian processes. *International Journal of Adaptive Control and Signal Processing* 17(2), 149–161 (2003)
4. Shi, J.Q., Murray-Smith, R., Titterington, D.M.: Hierarchical Gaussian process mixtures for regression. *Statistics and Computing* 15(1), 31–41 (2005)
5. Green, P.J.: Reversible jump markov chain monte carlo computation and bayesian model determination. *Biometrika* 82(4), 711–732 (1995)
6. Richardson, S., Green, P.J.: On Bayesian Analysis of Mixtures with an Unknown Number of Components (with discussion). *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 59(4), 731–792 (1997)
7. Rasmussen, C.E.: Evaluation of Gaussian processes and other methods for non-linear regression. Diss. University of Toronto (1996)
8. Neal, R.M.: MCMC using Hamiltonian dynamics. In: *Handbook of Markov Chain Monte Carlo* (2011)