



Automatic Model Selection Algorithm Based on BYY Harmonic Learning for Mixture of Gaussian Process Functional Regressions Models

Xiaogang Guo, Taohong Guo, and Jiwei Ma^(✉)

Department of Informatics and Computational Science, School of Mathematical Science and LMAM, Peking University, Beijing 100871, China

Abstract. For the mixture model, determining the number of components is a difficult problem of model selection. This paper proposes an automatic model selection algorithm based on Bayesian Ying-Yang (BYY) harmonic learning for mixture of Gaussian process functional regression (mix-GPFR) model. BYY harmonic learning has been successfully applied to the model selection of Gaussian mixture model (GMM), but cannot be directly used for mix-GPFR model. We deal with the case of this problem and propose a new mechanism for cross-correlation based Gaussian process (GP) model, through which we transform mix-GPFR model into a GMM. Then, we can make model selection for mix-GPFR model via BYY harmonic learning. Experiments show that the proposed automatic model selection algorithm can find the optimal number of components in simulated data set.

Keywords: Mixture of Gaussian Process Functional Regression Model, Selection, Bayesian Ying-Yang Harmonic Learning, Cross-Correlation

1 Introduction

Gaussian process (GP) model is an effective tool for Bayesian nonlinear classification and regression, e.g., classification of handwritten digits and image classification [1]. However, the current algorithm for model selection in data set is still limited. To overcome this limitation, mixture of Gaussian process functional regression (mix-GPFR) model was proposed [2, 3] and the extensive research has been devoted to estimate their parameters, although their effectiveness and algorithmic scalability is still unclear [4–8].

Like the mixture model, mix-GPFR model also faces the problem of model selection, and determining the number of Gaussian process functional regression (GPFR) components. Since a large number of GPFR components will inevitably lead to overgeneralization, model selection is a great importance. In addition, making model selection is a difficult task. In this paper, we propose

al-*de*ig-*a*t-*m*atic-*m*del-*e*lecti-*a*lg-*i*th-*m*. The traditi-*a*l-*m*eth-*d*i-*t*ch-*e* the-*t*imal-*n*umber-*f*-*G*PFR-*c*m-*e*t-*t*h-*r*ough-*c*er-*t*ain-*t*ati-*c*al-*e*lecti-*c*rite-*r*i-. For-*e*x-*a*m-*p*le, Qia-*g*-et-*a*l. [6] s-*e*lect-*e*d-*t*he-*l*iti-*g*-*e*ct-*i*-*m*ax-*i*mi-*a*ti-*(*SEM $)$ -*a*lg-*i*th-*m* ba-*e*d-*o*n-*t*he-*B*a-*e*i-*a*-*i*n-*f*-*o*r-*m*al-*c*rite-*r*i-*(*BIC $)$ [9]. H-*e*-*e*-*x*, all-*t*he-*e*-*i*-*t*-*i*g-*n*ati-*c*al-*e*lecti-*c*rite-*r*i-*a*re-*c*o-*n*s-*i*der-*e*d-*a*s-*a*-*h*igh-*t*ime-*c*m-*e*-*t* and-*t*he-*e*-*f*-*a*-*t*ati-*c*al-*e*lecti-*c*rite-*r*i-*i*s-*s*-*a*-*h*igh-*t*ime-*c*m-*e*-*t*, i-*c*e-*e*-*s*e-*e*ed-*t*o-*r*e-*e*-at-*t*he-*v*-*h*-*l*e-*a*-*r*am-*e*t-*e*-*r*im-*a*t-*i*g-*-*c-*e*-*f*-*-*d-*i*-f-*-*f-*e*-*r*e-*n*t-*n*umber-*f*-*G*PFR-*c*m-*e*-*t*. More-*o*-*v*er, st-*a*t-*i*st-*i*c-*i*n-*f*-*o*r-*m*al-*m*eth-*d*, i-*n*-*c*h-*a*s-*e*-*e*-*x*-*i*ble-*j*-*m* Mark-*-*chai-*-*M-*-*te-*-*Carl [10] and-*-*Dirichlet-*-*c-*-*ce-*-*e [11], ha-*-*e-*-*al-*-*bee-*-*ed-*-*t-*-*deal-*-*with-*-*the-*-*m-*-*del-*-*electi-*-*-*-*blem-*-*f-*-*mix-*-*-G-*-*PFR-*-*m-*-*del [5, 7, 8]. H-*-*e-*-*e-*-*x, the-*-*e-*-*meth-*-*d-*-*re-*-*lect-*-*gal-*-*age-*-*number-*-*f-*-*am-*-*le-*-*, *-*h-*-*ich-*-*e-*-*l-*-*i-*-*a-*-*high-*-*c-*-*m-*-*ati-*-*al-*-*c-*-*t. For-*-*G-*-*au-*-*ss-*-*ia-*-*mix-*-*-c-*-*m-*-*del (GMM), the-*-*a-*-*t-*-*m-*-*atic-*-*m-*-*del-*-*electi-*-*alg-*-*i-*-*th-*-*m-*-* ba-*-*e-*-*d-*-* Ba-*-*e-*-*ia-*-*Yi-*-*g-*-*Ya-*-*g (BYY) ha-*-*m-*-* bee-*-*ed-*-*lea-*-*rn-*-*ing [12, 13] ha-*-*e-*-*ac-*-*chie-*-*ved-*-*better-*-*re-*-*s-*-*ult-*-*s-*-*and-*-*high-*-*c-*-*m-*-*ati-*-*eed-*-*tha-*-*n-*-*the-*-*ba-*-*e-*-*d-*-*ati-*-*tical-*-*electi-*-*crite-*-*ria-*-* and-*-*st-*-*a-*-*t-*-*i-*-*c-*-*in-*-*f-*-*o-*-*r-*-*m-*-*al-*-* meth-*-*d [14-20].

Here, we indicate the variable z is called the GMM , we establish the following BYY theorem: $q(z = g) = \pi_g$; $q(z | z = g) = \mathcal{N}(z | \mu_g, \Sigma_g)$; $p(z) = \frac{1}{I} \sum_{i=1}^I \delta(z - z_i)$, i.e. the empirical distribution;

$$p(z = g | \Theta) = \frac{\pi_g \mathcal{N}(z | \mu_g, \Sigma_g)}{\sum_{s=1}^G \pi_s \mathcal{N}(z | \mu_s, \Sigma_s)}. \tag{3}$$

Moreover, we require the regularity condition, i.e. $\text{rank}(\Theta) = 1$. Then, we have

$$H(p||q) = J(\Theta) = \frac{1}{I} \sum_{i=1}^I \sum_{g=1}^G \frac{\pi_g \mathcal{N}(z_i | \mu_g, \Sigma_g)}{\sum_{s=1}^G \pi_s \mathcal{N}(z_i | \mu_s, \Sigma_s)} \ln \left(\frac{\pi_g \mathcal{N}(z_i | \mu_g, \Sigma_g)}{\sum_{s=1}^G \pi_s \mathcal{N}(z_i | \mu_s, \Sigma_s)} \right), \tag{4}$$

where $J(\Theta)$ is called harmonic function and $\Theta = \{\pi_g, \mu_g, \Sigma_g\}_{g=1}^G$.

According to BYY harmonic learning, the maximum of $J(\Theta)$ corresponds to the optimal number of Gaussian components and the best parameter [14–20]. Hence, we can make model selection to determine the parameter by maximizing $J(\Theta)$. In the process of maximizing $J(\Theta)$, the mixing coefficients of the selected Gaussian components converge. Combined with the automatic model selection algorithm based on variational Bayes method, the based BYY harmonic learning has a more better result and higher computational speed [14–20].

3 Automatic Model Selection Algorithm Based on BYY Harmonic Learning

First, we briefly introduce the mix-GPFR model. A GP is a collection of random variables, a joint probability distribution over a Gaussian distribution [1]. Specifically, a GP $\{f(x) | x \in \mathcal{X} \subseteq \mathbb{R}^D\}$, we need to determine its mean function $m(x)$ and covariance function $c(x, x')$, where

$$m(x) = \mathbb{E}[f(x)] \text{ and } c(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x')))]. \tag{5}$$

where the GP is denoted as

$$f(x) \sim \mathcal{GP}(m(x), c(x, x')). \tag{6}$$

In mix-GPFR model, since $D = 1$, we denote the input variable x . The mix-GPFR model with G GPFR components can be established through the following function:

$$q(z = g) = \pi_g, \text{ where } \pi_g \geq 0 \text{ and } \sum_{g=1}^G \pi_g = 1; \tag{7}$$

$$q(y(x)|z = g) = \mathcal{GPFR}(x|\mathbf{b}_g, \theta_g, r_g) = \mathcal{GP}(m(x|\mathbf{b}_g), c(x, x'|\theta_g) + r_g^{-1} \delta(x, x')). \tag{8}$$

In E. (8), $\delta(x, x')$ is the Kronecker delta function,

$$m(x|\mathbf{b}_g) = \varphi(x)^T \mathbf{b}_g \text{ and } c(x, x'|\boldsymbol{\theta}_g) = \theta_{g0}^2 \boldsymbol{\alpha} \cdot \left\{ -\frac{(x - x')^2}{2\theta_{g1}^2} \right\}, \tag{9}$$

where $\varphi(x) = [\varphi_1(x), \varphi_2(x), \dots, \varphi_P(x)]^T$ is a column vector of B-spline [21] and $c(x, x'|\boldsymbol{\theta}_g)$ is defined as the scaled exponential correlation function. θ_{g0}, θ_{g1} , and r_g are the parameters.

The Yilmazchi of the mix-GPFR model is

$$q(z = g, y(x)) = q(z = g)q(y(x)|z = g) = \pi_g \mathcal{GPFR}(x|\mathbf{b}_g, \boldsymbol{\theta}_g, r_g) \tag{10}$$

and Yamachi is

$$p(z = g, y(x)) = p(y(x))p(z = g|y(x)) = p(y(x)) \frac{\pi_g \mathcal{GPFR}(x|\mathbf{b}_g, \boldsymbol{\theta}_g, r_g)}{\sum_{s=1}^G \pi_s \mathcal{GPFR}(x|\mathbf{b}_s, \boldsymbol{\theta}_s, r_s)}. \tag{11}$$

We denote a training set $\mathcal{D} = \{\mathcal{C}_i\}_{i=1}^I$, where $\mathcal{C}_i = \{(x_{in}, y_{in})\}_{n=1}^{N_i}$ is a set of training samples of length N_i . It is generally assumed that x_{i1}, \dots, x_{iN_i} are iid and distributed in the interval $[x_{\min}, x_{\max}]$ ($i = 1, \dots, I$). Let $\mathbf{x}_i = [x_{i1}, \dots, x_{iN_i}]^T$, $\mathbf{y}_i = [y_{i1}, \dots, y_{iN_i}]^T$, and $\Theta = \{\pi_g, \mathbf{b}_g, \boldsymbol{\theta}_g, r_g\}_{g=1}^G$. For the mix-GPFR model,

$$H(p||q) = \sum_{g=1}^G \int p(y(x))p(z = g|y(x)) \ln(q(z = g)q(y(x)|z = g)) dy(x) \tag{12}$$

can be approximated by

$$J(\Theta) = \frac{1}{I} \sum_{i=1}^I \sum_{g=1}^G \frac{\pi_g \mathcal{N}(\mathbf{x}_i|\mathbf{m}_{ig}, \mathbf{C}_{ig})}{\sum_{s=1}^G \pi_s \mathcal{N}(\mathbf{x}_i|\mathbf{m}_{is}, \mathbf{C}_{is})} \ln(\pi_g \mathcal{N}(\mathbf{x}_i|\mathbf{m}_{ig}, \mathbf{C}_{ig})) \tag{13}$$

with $\mathbf{m}_{ig} = m(\mathbf{x}_i|\mathbf{b}_g)$ and $\mathbf{C}_{ig} = c(\mathbf{x}_i, \mathbf{x}_i|\boldsymbol{\theta}_g) + r_g^{-1} \mathbf{I}_{N_i}$, where \mathbf{I}_{N_i} is the $N_i \times N_i$ identity matrix.

$\hat{C}_i = \{(x_n, \hat{y}_{in})\}_{n=1}^N$ is an empirical covariance matrix with $x_n = x_{ni} + (n-1)\Delta$. Denoting the variance of Gaussian noise

$$\sigma_1^2 = \frac{1}{N_i} \sum_{n=1}^{N_i} (y_{in} - \hat{f}_i(x_{in}))^2. \tag{14}$$

It is clear that

$$\sigma_2^2 = \frac{1}{N_i} \sum_{n=1}^{N_i} (y_{in} - f_i(x_{in}))^2 \tag{15}$$

is a biased estimate of the variance of Gaussian noise of the function C_i for $f_i(x)$. Hence, σ_1^2 is a good estimate of the variance of the function that there are significant differences between $f_i(x)$ and $\hat{f}_i(x)$. In addition, \hat{C}_i is a good estimate of C_i . That is to say, the difference between the two covariance matrices is neglected in \hat{C}_i and is small, which will be validated through experiments in Sect. 4.

Let $\hat{D} = \left\{ \hat{C}_i \right\}_{i=1}^I$, $\mathbf{x} = [x_1, x_2, \dots, x_N]^T$, and $\hat{\mathbf{y}}_i = [\hat{y}_{i1}, \hat{y}_{i2}, \dots, \hat{y}_{iN}]^T$. $\hat{\mathbf{y}}_i$ can be regarded as a sample of the full vector GMM:

$$q(z = g) = \pi_{g\mathbf{v}}, \text{ here } \pi_g \geq 0 \text{ and } \sum_{g=1}^G \pi_g = 1; q(\hat{\mathbf{y}} | z = g) = \mathcal{N}(\hat{\mathbf{y}} | \mathbf{m}_g, \mathbf{C}_g), \tag{16}$$

where $\mathbf{m}_g = m(\mathbf{b}_g)$ and $\mathbf{C}_g = c(\mathbf{b}_g, \mathbf{r}_g^{-1} \mathbf{I}_N)$. In Bayesian machine learning

$$q(z = g, \hat{\mathbf{y}}) = q(z = g)q(\hat{\mathbf{y}} | z = g) = \pi_g \mathcal{N}(\hat{\mathbf{y}} | \mathbf{m}_g, \mathbf{C}_g) \tag{17}$$

and the Bayesian machine learning

$$p(z = g, \hat{\mathbf{y}}) = p(\hat{\mathbf{y}})p(z = g | \hat{\mathbf{y}}) = p(\hat{\mathbf{y}}) \frac{\pi_g \mathcal{N}(\hat{\mathbf{y}} | \mathbf{m}_g, \mathbf{C}_g)}{\sum_{s=1}^G \pi_s \mathcal{N}(\hat{\mathbf{y}} | \mathbf{m}_s, \mathbf{C}_s)}. \tag{18}$$

Then, it is convenient to define

$$J(\Theta) = \frac{1}{I} \sum_{i=1}^I \sum_{g=1}^G \frac{\pi_g \mathcal{N}(\hat{\mathbf{y}}_i | \mathbf{m}_g, \mathbf{C}_g)}{\sum_{s=1}^G \pi_s \mathcal{N}(\hat{\mathbf{y}}_i | \mathbf{m}_s, \mathbf{C}_s)} \mathbf{1}(\pi_g \mathcal{N}(\hat{\mathbf{y}}_i | \mathbf{m}_g, \mathbf{C}_g)). \tag{19}$$

As the covariance matrix GMM, the maximum of $J(\Theta)$ corresponds to the final number of GPFR components and the best parameters. Therefore, we can make model selection and learn the parameters by maximizing $J(\Theta)$ through numerical optimization methods.

After the training process, we can determine the classification accuracy according to the maximum a posteriori probability, i.e. let

$$z_i = \underset{g \in \{1, 2, \dots, G\}}{\text{argmax}} \frac{\pi_g \mathcal{N}(\hat{\mathbf{y}}_i | \mathbf{m}_g, \mathbf{C}_g)}{\sum_{s=1}^G \pi_s \mathcal{N}(\hat{\mathbf{y}}_i | \mathbf{m}_s, \mathbf{C}_s)} \quad (i = 1, 2, \dots, I). \tag{20}$$

The test data GPFR components do not get a training case and their small mixing coefficients. The classification accuracy can be determined through a Bayesian method. Besides, we can predict the test data by calculating their conditional distributions given the test data. The details are referred to [4-8].

4 Experimental Results

In this section, we evaluate the predictive performance of the proposed methods on synthetic data and real-world data to compare the effectiveness of the proposed automatic model selection algorithm. We compare the mix-GPFR model trained via the proposed algorithm with GP model, mix-GP model, GPFR model, and mix-GPFR model trained through the traditional EM algorithm [2, 3] and the SEM algorithm [6].

Since we are mainly concerned with the predictive ability of mix-GPFR model, the standard measurement is the RMSE in which is a standard metric. It is assumed that there are T test cases and the test data of the t th ($t = 1, 2, \dots, T$) test case are $y_{t1}, y_{t2}, \dots, y_{tM}$, which correspond to the predicted value $\hat{y}_{t1}, \hat{y}_{t2}, \dots, \hat{y}_{tM}$, respectively. If $\|y_{tm} - \hat{y}_{tm}\|$

$$\text{RMSE} = \sqrt{\frac{1}{TM} \sum_{t=1}^T \sum_{m=1}^M (y_{tm} - \hat{y}_{tm})^2}. \tag{21}$$

As a result, a smaller RMSE indicates a better predictive result.

4.1 On Synthetic Datasets

The synthetic data are generated as S_2, S_3, \dots, S_{10} , respectively, where the subscripts represent the number of features. For each feature, we sample 20 training cases and 10 test cases from a GP with a normal mean function. The mean function and a parameter of the Gaussian covariance function generate the synthetic data are listed in Table 1, where S_l ($l = 2, 3, \dots, 10$) are generated by the l GP. Each case consists of 100 input, which is a random distribution in $[-3, 3]$. The 60 input on the left side of test cases are k and the 40 on the right side are used for testing.

First, we demonstrate that the effectiveness of the constructed GP model through α -experiment. A training case is a random choice from each component S_9 . Figure 1 shows the reconstruction error of the 10 training cases. Figure 1 is composed of 9 sub-graphs, each figure shows a training case, its reconstruction error, and their test mean function. As can be seen from the graphs, although there are significant differences between training cases and their reconstruction errors, their test mean functions are similar, which implies that the proposed constructed GP model is effective.

When testing the proposed algorithm, G is initialized as $l+3$ for S_l . To illustrate the bad effect of a large number of GPFR components on predictive ability, we train mix-GPFR model consisting of $l-1$ and $l+1$ GPFR components via the EM algorithm [2, 3], which are denoted as $\text{mix-GPFR}(-1)$ and $\text{mix-GPFR}(+1)$, respectively. Similarly, mix-GP model with $l-1$ and $l+1$ GP components are denoted as $\text{mix-GP}(-1)$ and $\text{mix-GP}(+1)$, respectively. Besides, P is set to be 20. Table 2 shows the α -experimental results.

From Table 2, we see that the RMSE of the GPFR model and the mix-GPFR model are smaller than those of the GP model and the mix-GP model, respectively, which

Table 1. Mean function and parameter of the Gaussian process used to generate the synthetic data set.

Mean function	θ^T	$\sqrt{r^{-1}}$
x^2	[0.5, 0.5]	0.15
$(-4(x + 1.5)^2 + 9)1_{\{x < 0\}} + (4(x - 1.5)^2 - 9)1_{\{x \geq 0\}}$	[0.528, 0.4]	0.144
$8 \sin(1.5x - 1)$	[0.556, 0.3]	0.139
$\sin(1.5x) + 2x - 5$	[0.583, 0.2]	0.133
$\sin(4x) - 0.5x^2 - 2x$	[0.611, 0.1]	0.128
$-x^2$	[0.639, 0.1]	0.122
$(4(x + 1.5)^2 - 9)1_{\{x < 0\}} + (-4(x - 1.5)^2 + 9)1_{\{x \geq 0\}}$	[0.667, 0.2]	0.117
$5 \cos(3x + 2)$	[0.694, 0.3]	0.111
$\cos(1.5x) - 2x + 5$	[0.722, 0.4]	0.106
$\cos(4x) + 0.5x^2 + 2x$	[0.75, 0.5]	0.1

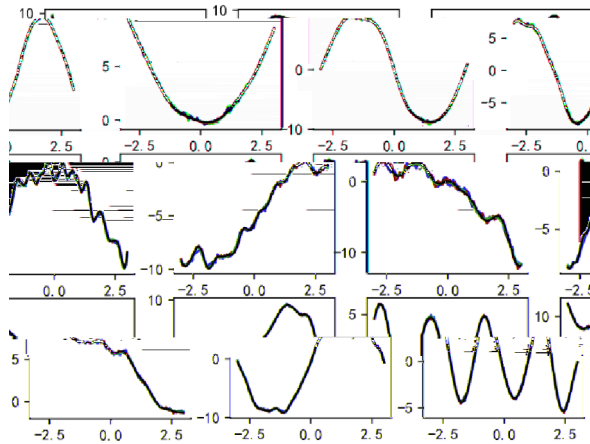


Fig. 1. The result of reconstruction of the synthetic data set. The red, green, blue, and black color represent the original case, the reconstructed case, the true mean function of the original case, and the true mean function of the reconstructed case, respectively.

demonstrate the effectiveness of modeling the mean function as a linear combination of B-spline. Because of the mix-GP (mix-GPFR) model and the GP (GPFR) model, the need for fitting the mixture test case is demonstrated. For the same reason, we can see that a large number of GPFR cannot affect the prediction result badly. For S_2, S_3, \dots, S_9 , both the SEM algorithm and the proposed algorithm find the correct number of GPFR cannot affect their RMSE and the time complexity. However, the time complexity of the SEM algorithm is higher than that of the proposed algorithm. On the other hand,

the SEM algorithm need to treat the whole parameter learning case for different number of GPFR components. On the other hand, since different training cases have different initial values, the likelihood of the algorithm is different. This is the main reason why the SEM algorithm has a high time complexity. For S_{10} , since the SEM algorithm failed to find the number of GPFR components, its RMSE is larger than that for the proposed algorithm.

Taking S_9 for example, we see that the clustering result for the proposed algorithm in Fig. 2, where different colors represent different components. On the left and right side of Fig. 2 are the clustering result for the proposed algorithm on the training and test data set, respectively. It is clear that the proposed algorithm correctly found all the components.

Table 2. RMSE and running time for all the methods on the synthetic data set.

	S_2		S_3		S_4	
	RMSE	Time (mi)	RMSE	Time (mi)	RMSE	Time (mi)
GP	5.5831	6.87	4.7878	9.88	4.7798	13.09
mix-GP (-1)	5.5239	6.12	4.6125	8.90	4.3580	15.84
mix-GP (+1)	4.8240	7.39	4.6035	17.40	4.3488	23.31
GPFR	5.0759	6.68	4.6864	12.42	4.3051	17.72
mix-GPFR (-1)	5.0214	8.07	0.9416	15.95	0.9510	18.93
mix-GPFR (+1)	1.6846	14.20	1.0680	22.66	0.9319	25.79
mix-GPFR (SEM)	0.4312	20.63	0.4856	41.59	0.5469	58.97
mix-GPFR (BYY)	0.4401	9.46	0.4746	15.03	0.5403	18.64
	S_5		S_6		S_7	
	RMSE	Time (mi)	RMSE	Time (mi)	RMSE	Time (mi)
GP	4.9213	17.85	4.9897	15.07	5.3096	20.47
mix-GP (-1)	4.4775	15.03	4.3834	20.14	4.3025	30.66
mix-GP (+1)	4.5205	28.59	4.3813	29.19	4.3082	30.53
GPFR	4.8079	26.73	4.8649	24.25	4.9871	21.45
mix-GPFR (-1)	0.8756	31.66	1.0776	29.67	1.3295	32.09
mix-GPFR (+1)	0.9252	30.58	1.0270	35.37	1.0281	38.44
mix-GPFR (SEM)	0.5638	81.34	0.6057	87.52	0.6540	92.78
mix-GPFR (BYY)	0.5573	25.49	0.6137	23.66	0.6571	27.82
	S_8		S_9		S_{10}	
	RMSE	Time (mi)	RMSE	Time (mi)	RMSE	Time (mi)
GP	4.8180	19.67	4.4758	17.82	4.8438	21.67
mix-GP (-1)	4.4904	24.13	4.1223	32.36	4.5730	32.78

(continued)

Table 2. (continued)

	S_2		S_3		S_4	
	RMSE	Time (mi)	RMSE	Time (mi)	RMSE	Time (mi)
mix.-GP (+1)	4.4818	23.70	4.1214	30.95	4.5878	28.14
GPFR	4.6871	26.73	4.3686	21.67	4.6865	20.97
mix.-GPFR (-1)	1.5325	33.78	1.0585	40.27	1.5279	41.56
mix.-GPFR (+1)	1.1891	35.96	0.9789	39.38	1.0343	49.78
mix.-GPFR (SEM)	0.6448	99.49	0.6233	116.85	1.4379	130.65
mix.-GPFR (BYY)	0.6421	28.91	0.6199	28.62	0.6317	32.46

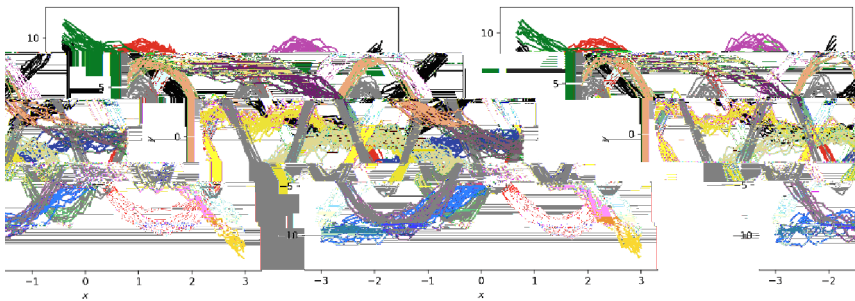


Fig. 2. Clustering results for datasets S_7 and S_9 .

4.2 On Real-World Datasets

Here, we utilize the electricity load data provided by the North China Grid Company [8], which records electricity load every 15 minutes in 2009 and 2010. Hence, daily electricity load can be regarded as a sequence with 96 points. We divide the data into training and testing data according to the early, which are respectively \mathcal{R}_1 and \mathcal{R}_2 , sequentially. Each training data consists of 200 training sequences for a total of 165 test cases. Moreover, the 56 points in the left side of the test cases are used for the left side of the test cases and the 40 points in the right side are used for the right side of the test cases.

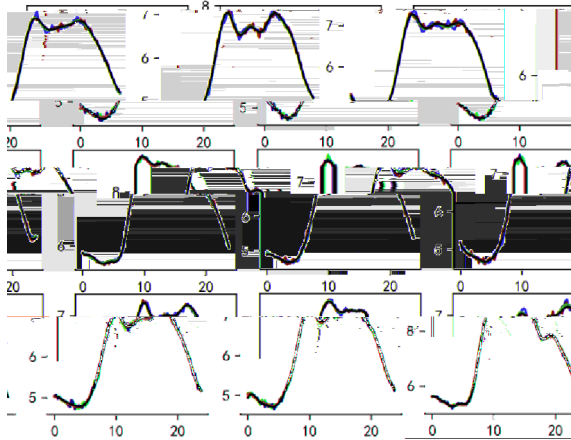


Fig. 3. The reconstruction results for the electricity load data set. The red, green, blue, and black colors represent the original data, the reconstructed data, the text mean function of the original data, and the text mean function of the reconstructed data, respectively.

Although all the cases have the same input, we treat them as if the data have the same input. Like the synthetic data set, we also add nine training cases for \mathcal{R}_1 , whose reconstruction cases are presented in Fig. 3. As can be seen from the green, red, and blue colors, the reconstruction based on GP model is effective for the electricity load data set.

Since the number of components in \mathcal{R}_1 and \mathcal{R}_2 are k_{v_1}, k_{v_2} and $G = 3, 6, 9, 12, 15$ for the mix-GP and mix-GPFR model trained using the EM algorithm. For the proposed algorithm and the SEM algorithm, G is set to be 15. Besides, P is set to be 30. The experimental results are described in Table 3. For \mathcal{R}_1 and \mathcal{R}_2 , the RMSE of the proposed algorithm is smaller than that of the SEM algorithm in the number of components given by the SEM algorithm is smaller than the minimal. The clustering results are presented in Fig. 4. On the left and right side of Fig. 4 are the clustering results for the proposed algorithm on the training and test data set, respectively. For \mathcal{R}_1 and \mathcal{R}_2 , the number of components given by the proposed algorithm are 13 and 11, respectively. As can be seen from Fig. 4, the belonging of different components are basically different in a certain interval, that is to say, the clustering results given by the algorithm are reasonable.

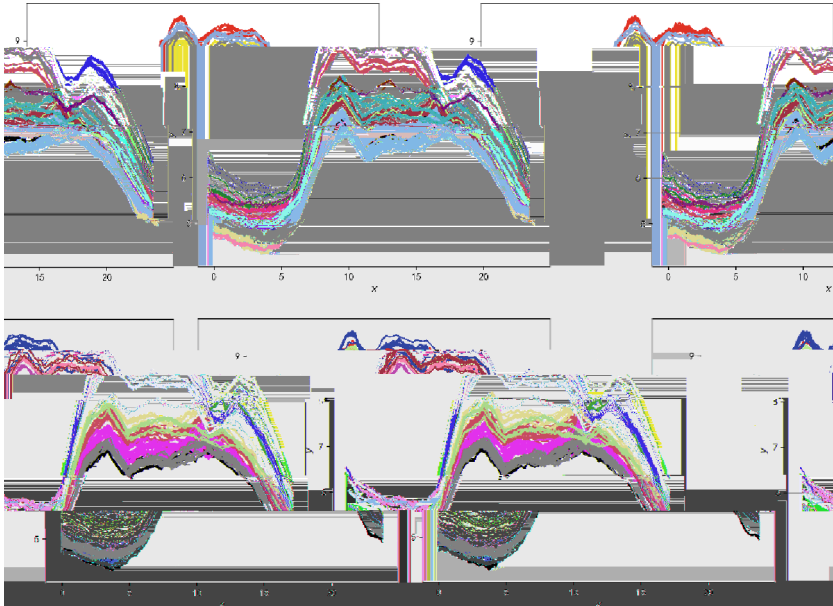


Fig. 4. Clustering results for the automatic model selection algorithm \mathcal{R}_1 and \mathcal{R}_2 .

Table 3. RMSE and clustering time for all the methods \mathcal{R}_1 and \mathcal{R}_2 .

	\mathcal{R}_1		\mathcal{R}_2	
	RMSE	Time (min)	RMSE	Time (min)
GP	0.9599	19.43	0.8977	20.39
mix-GP (3)	0.9390	20.33	0.8846	21.42
mix-GP (6)	0.9387	22.54	0.8854	22.66
mix-GP (9)	0.9380	25.99	0.8853	26.09
mix-GP (12)	0.9395	29.83	0.8847	31.23
mix-GP (15)	0.9401	34.76	0.8872	36.91
GPFR	0.5584	21.30	0.5499	21.59
mix-GPFR (3)	0.2089	24.45	0.2133	20.64
mix-GPFR (6)	0.1701	25.76	0.1731	24.77
mix-GPFR (9)	0.1356	28.93	0.1455	29.45
mix-GPFR (12)	0.1248	34.65	0.1314	33.63
mix-GPFR (15)	0.1178	35.88	0.1301	36.78
mix-GPFR (SEM)	0.1323	150.76	0.1377	170.17
mix-GPFR (BYY)	0.1109	33.97	0.1201	34.58

5 Conclusion

In this paper, we compare automatic del electi alg with ba ed BYY ha m lea i g f m k-GPFR m del . Si ce differ t tra i g c e ha e differ t i t , BYY ha m lea i g ca t be directl a lied t the m del electi c blom f m k-GPFR m del . T tackle thi , e c e c e c t c i ba ed GP m del , th gh , hich , e if the i t f all the tra i g c e . The , e ca make m del electi f m k-GPFR m del ia BYY ha m lea i g . Ex e i m e tal ce l t thetic a d ce al- c ld data et h , that c c ed a t m a i c m del electi alg c i t m ca d the t i m a l m b e r f c m e t i a m li- c ce c e data et a d i t i m e c m l e x i t i l , e t h a t h a t f the SEM alg c i t m .

Acknowledgement. Thi s k i s s u c t e d b y the Nat i o n a l Ke R & D P r o g r a m f Chi a (2018AAA0100205).

References

1. R a m e , C.E., W i l l i a m , C.K.I.: G a i a P r o c e e f s : M a c h i n e L e a r i n g . M I T P r e s s , C a m b r i d g e (2 0 0 6)
2. S h i , J.Q., W a n g , B.: C r e e r e d i c t i a d c l t e r i g i t h m i k t c e f G a i a c c e f c t i a l c e g e i m del . S t a t i s t . C o m m . t . **18**, 267-283 (2008)
3. S h i , J.Q., C h i , T.: G a i a P r o c e R e g e i A a l i f s : F c t i a l D a t a . C R C P r e s s , B e a R a t (2 0 1 1)
4. W u , D., M a , J.: A D A E M a l g c i t m f s m i k t c e f G a i a c c e f c t i a l c e g e i . I : H a g , D.-S., H a , K., H u a i , A. (e d .) I C I C 2 0 1 6 . L N C S (L N A I) , **1**. 9773, . 294-303. S c i e n c e , C h a m (2 0 1 6) . [h t t p s : / / d o i . o r g / 1 0 . 1 0 0 7 / 9 7 8 - 3 - 3 1 9 - 4 2 2 9 7 - 8 _ 2 8](https://doi.org/10.1007/978-3-319-42297-8_28)
5. Q i a g , Z., M a , J.: A t m a i c m del electi f the m i k t c e f G a i a c c e e f c c e g e i . I : H u , X., X i a , Y., Z h a g , Y., Z h a , D. (e d .) I S N N 2 0 1 5 . L N C S , **1**. 9377, . 335-344. S c i e n c e , C h a m (2 0 1 5) . [h t t p s : / / d o i . o r g / 1 0 . 1 0 0 7 / 9 7 8 - 3 - 3 1 9 - 2 5 3 9 3 - 0 _ 3 7](https://doi.org/10.1007/978-3-319-25393-0_37)
6. Q i a g , Z., L i , J., M a , J.: C r e e c l t e r i g i a the l i t l e a r i g f m i k t c e f G a i a c c e e . I : 2 0 1 6 I E E E 1 3 t h I t e r a t i a l C f e r e c e S i g a l P r o c e e i n g (I C S P) , . 1089-1094 (2016)
7. Q i a g , Z., M a , J.: M del electi c r e d i c t i f t h e m i k t c e f G a i a c c e e i t h R J M C M C . I : S h i , Z., P e a r t , C., H a g , T. (e d .) I C I S 2 0 1 8 . I A I C T , **1**. 539, . 310-317. S c i e n c e , C h a m (2 0 1 8) . [h t t p s : / / d o i . o r g / 1 0 . 1 0 0 7 / 9 7 8 - 3 - 0 3 0 - 0 1 3 1 3 - 4 _ 3 3](https://doi.org/10.1007/978-3-030-01313-4_33)
8. L i , T., M a , J.: D i r i c h l e t c e m i k t c e f G a i a c c e f c t i a l c e g e i a d i t a r i a t i a l E M a l g c i t m . P a t t e r n R e c o g . **134**, 109-129 (2023)
9. S c h a r , G.: E t i m a t i n g the d i m e n s i o n a l i t y . S t a t i s t . **6**(2), 461-464 (1978)
10. R i c h a r d s , S., G r e e , P.J.: O B a e i a a a l i f m i k t c e i t h a . k . m b e r f c m e t i . J R a l S t a t i s t . S c . (S e r . B) **59**(4), 731-792 (1997)
11. E c b a r , M.D., W e t t , M.: B a e i a d e i t e t i m a t i a d i f e r e c e . i g m i k t c e . J . A m . S t a t i s t . A s s o c . **90**(430), 577-588 (1995)
12. X u , L.: Y i g - Y a g m a c h i e : a B a e i a - K l l b a c k c h e m e f c . i e d l e a r i g a d c e c e l t e c t c a t i a t i . I : P r o c e e d i n g f t h e 1 9 9 5 I t e r a t i a l C f e r e c e N e c a l I f f m a t i P r o c e e i n g , **1**. 2, . 977-988 (1995)
13. X u , L.: B e t h a m , i e d R P C L a d a t m a t e d m del electi f . . e i e d a d e i e d l e a r i g . G a i a m i k t c e , t h r e e - l a e r a d M E - R B F - S V M m del . I . I t . J . N e c a l S t a t . **11**(1), 43-69 (2001)

14. Cheng, G., Li, L., Ma, J.: A gradient Bayesian learning algorithm for trajectory detection. In: Song, F., Zhang, J., Tang, Y., Cao, J., Yin, W. (ed.) ISSN 2008. LNCS, vol. 5263, pp. 618–626. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-87732-5_69
15. Li, L., Ma, J.: A Bayesian calibration EM algorithm for Gaussian mixture learning. *Acta Mathematica Sinica* **205**(2), 832–840 (2008)
16. Li, L., Ma, J.: A Bayesian likelihood EM algorithm for Gaussian mixture learning. In: Song, F., Zhang, J., Tang, Y., Cao, J., Yin, W. (ed.) ISSN 2008. LNCS, vol. 5263, pp. 600–609. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-87732-5_67
17. Ma, J., Wang, T., Xu, L.: A gradient Bayesian learning algorithm for Gaussian mixture learning. *Neurocomputing* **56**, 481–487 (2004)
18. Ma, J., Gao, B., Wang, Y., Cheng, Q.: Conjugate gradient algorithm for Bayesian learning of Gaussian mixture learning. *Int. J. Pattern Recogn. Artif. Intell.* **19**(5), 701–713 (2005)
19. Ma, J., Li, J.: The Bayesian algorithm for Gaussian mixture learning. *Pattern Recogn.* **40**(7), 2029–2037 (2007)
20. Ma, J., He, X.: A fast algorithm for Bayesian learning of Gaussian mixture learning. *Pattern Recogn. Lett.* **29**(6), 701–711 (2008)
21. Berger, C.D.: On calculating the Bayes risk. *J. Am. Statist. Assoc.* **6**, 50–62 (1972)