

Two-step Gaussian Process Regression Improving Performance of Training and Prediction

Wei Wang*, Santong Zhang, Wei Yang and Xiangbin Liu

School of Electronic and Information Engineering, Beijing Jiaotong University, China

*Corresponding author

Abstract—Since Gaussian process regression (GPR) cannot feasibly be applied to big and growing data sets, this paper introduces an integration algorithm called Two-step Gaussian Process Regression (TGPR) which speeds up both training and prediction to solve the problem. First, analyze the basics behind regular GPR. Then, introduce TGPR by using the inducing inputs to optimize the regular GPR algorithm. Last, apply TGPR to a three-dimension model, the experimental results compared with regular GPR show that TGPR is faster and more accurate.

Keywords—gaussian process; regression; inducing inputs; hyperparameter

I INTRODUCTION

Gaussian Process Regression is a newly developed method for machine learning based on Bayesian theory and statistical learning theory. It is first applied to machine learning by Rasmussen and Williams [1]. The prediction result of GPR has a clear probabilistic meaning compared with Artificial Neural Networks (ANN) [2] and Support Vector Machine (SVM) [3].

Although GPR has many advantages, its computational complexity increases sharply as the size of the measurements expands. Many sparse algorithms have been proposed to reduce computation. Williams presents the Nyström method for efficient GPR [4]. Snelson presents Sparse Pseudo-inputs Gaussian Process (SPGP) model whose covariance is parameterized by the locations of pseudo-input [5]. Tresp presents BCM (Bayesian Committee Machine) method to combine estimators that are trained on different data sets [6].

This paper is organized as follows. After providing the necessary background on GPR, the prediction and training of regular GPR is introduced in Section 2. In Section 3, a novel sparse algorithm TGPR which uses the inducing inputs is introduced. Section 4 presents the simulation on three-dimension model. The conclusions of TGPR algorithm is in Section 5.

II GAUSSIAN PROCESS REGRESSION

A. Prediction

A Gaussian Process (GP) is a collection of random variables, any finite numbers of which have joint Gaussian distributions. It can be expressed as follows:

$$\mathbf{f} \sim N(\mathbf{f} | \mathbf{m}, K) \quad (1)$$

Where, \mathbf{m} is mean function and $K=k(x, x)$ is covariance function, they are both called kernel function which describe the probabilistic meaning of GPR [7].

Regression is that, first get the function relation between measurement inputs X_m and measurement outputs \mathbf{f}_m , being $\mathbf{f}_m = f(X_m)$. Then predict the unknown test outputs \mathbf{f}_* given test inputs X_* . The first step is called training while the second step is called prediction.

Considering the measurement input set $X_m = \{x_{m_1}, x_{m_2}, \dots, x_{m_{n_m}}\}$ that consists of n_m points, and the test input set $X_* = \{x_{*_1}, x_{*_2}, \dots, x_{*_n_*}\}$ that consists of n_* points. The prior distribution of \mathbf{f}_m and \mathbf{f}_* is given by:

$$\begin{bmatrix} \mathbf{f}_m \\ \mathbf{f}_* \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{m}_m \\ \mathbf{m}_* \end{bmatrix}, \begin{bmatrix} K_{mm} & K_{m*} \\ K_{*m} & K_{**} \end{bmatrix} \right) \quad (2)$$

The measurement outputs are generally affected by Gaussian white noise $\mathbf{v} \sim N(0, \hat{\Sigma}_{f_m})$ where $\hat{\Sigma}_{f_m} = \text{diag}(\hat{\sigma}_{f_{m1}}^2, \hat{\sigma}_{f_{m2}}^2, \dots, \hat{\sigma}_{f_{m_{n_m}}}^2)$. As a result, the measurement output becomes $\hat{\mathbf{f}}_m = \mathbf{f}_m + \mathbf{v}$. The posterior distribution of \mathbf{f}_* can be calculated from the conditional probability formula:

$$\mathbf{f}_* \sim N(\mathbf{m}_* + K_{*m} K_{mm}^{-1} (\mathbf{f}_m - \mathbf{m}_m), K_{**} - K_{*m} (K_{mm} + \hat{\Sigma}_{f_m})^{-1} K_{m*}) \quad (3)$$

B. Training

In the training step, the mean function generally is initialized to zero mean function, and covariance function is chosen as the well-known Squared Exponential covariance function (SE):

$$k(\mathbf{x}, \mathbf{x}') = \lambda_f^2 \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \Lambda_x^{-1}(\mathbf{x} - \mathbf{x}')\right). \quad (4)$$

Where, $\lambda_f^2 = k(\mathbf{x}, \mathbf{x})$,
 $\Lambda_x = \lambda_x^2 = \text{diag}(\lambda_{x_{n_1}}^2, \lambda_{x_{n_2}}^2, \dots, \lambda_{x_{n_m}}^2)$ while
 $\lambda_{x_{n_1}}^2, \lambda_{x_{n_2}}^2, \dots, \lambda_{x_{n_m}}^2$ is length scale for each input.

Then there is a set $\theta = \{\lambda_x, \lambda_f, \hat{\Sigma}_{f_m}\}$, in which all parameters are called hyperparameters [8]. Only with the true set of hyperparameters called the most likely hyperparameters, the model can be explained exactly by the function relation between X_m and f_m . The most likely hyperparameters can be chosen using Maximum Likelihood method. According to Bayesian theory, the posterior hyperparameter distribution is found as follows:

$$p(\theta | \hat{f}_m, X_m) = \frac{p(\hat{f}_m | \theta, X_m) p(\theta | X_m)}{p(\hat{f}_m | X_m)}. \quad (5)$$

The key now is to find the maximum of (5). Remember that the covariance function is SE, now take the logarithm of (5), and then consider the derivative with respect to the hyperparameters:

$$\frac{\partial \log(p)}{\partial \theta} = -\frac{1}{2} \text{tr}\left(P^{-1} \frac{\partial P}{\partial \theta}\right) + \frac{1}{2} \text{tr}\left(\alpha^T \frac{\partial P}{\partial \theta} \alpha\right) = \frac{1}{2} \text{tr}\left((\alpha \alpha^T - P^{-1}) \frac{\partial P}{\partial \theta}\right). \quad (6)$$

Where $P = (K_{mm} + \hat{\Sigma}_{f_m})$, $\alpha = P^{-1}(\hat{f}_m - m_m)$. Note that $\theta = \{\lambda_x, \lambda_f, \hat{\Sigma}_{f_m}\}$, it means there will be three derivatives. Applying gradient ascent algorithm to all these derivatives, the most likely hyperparameters can be found [9].

C. Downsides of GPR

In the training step of GPR, the main computation focus on inverting the matrix $K_{mm} + \hat{\Sigma}_{f_m}$, resulting $O(n_m^3)$ runtime requirement. In the prediction step of GPR, it takes $O(n_m^2 n_*)$ time to calculate f^* . When the number of measurements n_m grows big, GPR runs into computational problems, requiring amounts of time.

III TWO-STEP GAUSSIAN PROCESS REGRESSION

A. Speeding up Prediction: Using Inducing Inputs

The inducing inputs set $X_u = \{x_{u_1}, x_{u_2}, \dots, x_{u_{n_u}}\}$ that consists of n_u inducing inputs [10] is introduced to speed up prediction.

STEP 1. Calculating the posterior distribution of the inducing outputs f_u given measurements (X_m, f_m) .

Assuming that f_m and f^* are conditionally independent given f_u . The prior distribution of f_m and f_u is given by:

$$\begin{bmatrix} f_m \\ f_u \end{bmatrix} \sim N\left(\begin{bmatrix} m_m \\ m_u \end{bmatrix}, \begin{bmatrix} K_{mm} & K_{mu} \\ K_{um} & K_{uu} \end{bmatrix}\right). \quad (7)$$

The outputs f_m are still affected by Gaussian white noise $v \sim N(0, \hat{\Sigma}_{f_m})$ and can be written as:

$$f_m = \hat{f}_m - v \sim N(\hat{f}_m, \hat{\Sigma}_{f_m}). \quad (8)$$

There are two distributions for f_m , (7) and (8). Here the idea is to combine these two distributions. The operation called "combine" is introduced as follows:

Assuming that, for n independent measurements, the random variable \mathbf{x} can be explained by n distributions, being $N(m_1, K_1), \dots, N(m_n, K_n)$. Then the posterior distribution for \mathbf{x} can be calculated by combining these n distributions. Here use the operator \oplus as "combine", the posterior distribution for \mathbf{x} is now given by:

$$f_x(\mathbf{x}) \sim N(\mu_x, \Sigma_x) = N(m_1, K_1) \oplus \dots \oplus N(m_n, K_n). \quad (9)$$

$$\mu_x = \Sigma_x (K_1^{-1} m_1 + K_2^{-1} m_2 + \dots + K_n^{-1} m_n). \quad (10)$$

$$\Sigma_x = (K_1^{-1} + K_2^{-1} + \dots + K_n^{-1})^{-1}. \quad (11)$$

Then the posterior distribution of f_m and f_u can be calculated by combining (7) and (8):

$$\begin{bmatrix} f_m \\ f_u \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_m \\ \mu_u \end{bmatrix}, \begin{bmatrix} \Sigma_{mm} & \Sigma_{mu} \\ \Sigma_{um} & \Sigma_{uu} \end{bmatrix}\right). \quad (12)$$

$$\begin{bmatrix} \mu_m \\ \mu_u \end{bmatrix} = \begin{bmatrix} \Sigma_{mm} (K_{mm}^{-1} m_m + \hat{\Sigma}_{f_m}^{-1} \hat{f}_m) \\ m_u + K_{um} (K_{mm} + \hat{\Sigma}_{f_m})^{-1} (\hat{f}_m - m_m) \end{bmatrix}. \quad (13)$$

$$\begin{bmatrix} \Sigma_{mm} & \Sigma_{mu} \\ \Sigma_{um} & \Sigma_{uu} \end{bmatrix} = \begin{bmatrix} K_{mm} (K_{mm} + \hat{\Sigma}_{f_m})^{-1} \hat{\Sigma}_{f_m} & \hat{\Sigma}_{f_m} (K_{mm} + \hat{\Sigma}_{f_m})^{-1} K_{mu} \\ K_{um} (K_{mm} + \hat{\Sigma}_{f_m})^{-1} \hat{\Sigma}_{f_m} & K_{uu} - K_{um} (K_{mm} + \hat{\Sigma}_{f_m})^{-1} K_{mu} \end{bmatrix}. \quad (14)$$

From (12), not only the posterior distribution of inducing outputs \mathbf{f}_u is found, but also that of the measurement outputs \mathbf{f}_m which are noiseless. It is one of the advantages of TGPR algorithm.

STEP 2. Just using the posterior distribution of \mathbf{f}_u to predict the test outputs \mathbf{f}^* .

The prior distribution of \mathbf{f}_u and \mathbf{f}^* is given by:

$$\begin{bmatrix} \mathbf{f}_u \\ \mathbf{f}^* \end{bmatrix} \sim N\left(\begin{bmatrix} \mathbf{m}_u \\ \mathbf{m}^* \end{bmatrix}, \begin{bmatrix} K_{uu} & K_{u^*} \\ K_{*u} & K_{**} \end{bmatrix}\right). \quad (15)$$

Another distribution of \mathbf{f}_u is from (12), it is:

$$\mathbf{f}_u \sim N(\boldsymbol{\mu}_u, \Sigma_{uu}). \quad (16)$$

The idea is still to combine (15) and (16). But note that, the prior knowledge of \mathbf{f}_u , being $\mathbf{f}_u \sim N(\mathbf{m}_u, K_{uu})$, has been used in both of these two distribution of \mathbf{f}_u . Combining (15) and (16) will use the prior knowledge twice, which is unacceptable, it is necessary to separate that prior for one time. Here use the operator $!$ as “separate”, so the operation called “separate” is introduced as follows:

$$N(\boldsymbol{\mu}_u, \Sigma_{uu})! N(\mathbf{m}_u, K_{uu}) = N(\boldsymbol{\mu}_u, \Sigma_{uu}) \oplus N(\mathbf{m}_u, -K_{uu}). \quad (17)$$

Then the posterior distribution of \mathbf{f}_u and \mathbf{f}^* can be calculated by combining (15) and (16), and separating $N(\mathbf{m}_u, K_{uu})$ from the result for one time:

$$\begin{bmatrix} \mathbf{f}_u \\ \mathbf{f}^* \end{bmatrix} \sim N\left(\begin{bmatrix} \boldsymbol{\mu}_u \\ \boldsymbol{\mu}^* \end{bmatrix}, \begin{bmatrix} \Sigma_{uu} & \Sigma_{u^*} \\ \Sigma_{*u} & \Sigma_{**} \end{bmatrix}\right). \quad (18)$$

$$\begin{bmatrix} \boldsymbol{\mu}_u \\ \boldsymbol{\mu}^* \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_u \\ \mathbf{m}^* + K_{*u} K_{uu}^{-1} (\boldsymbol{\mu}_u - \mathbf{m}_u) \end{bmatrix}. \quad (19)$$

$$\begin{bmatrix} \Sigma_{uu} & \Sigma_{u^*} \\ \Sigma_{*u} & \Sigma_{**} \end{bmatrix} = \begin{bmatrix} \Sigma_{uu} & \Sigma_{uu} K_{uu}^{-1} K_{u^*} \\ K_{*u} K_{uu}^{-1} \Sigma_{uu} & K_{**} - K_{*u} K_{uu}^{-1} (K_{uu} - \Sigma_{uu}) K_{uu}^{-1} K_{u^*} \end{bmatrix}. \quad (20)$$

Performance of TGPR.

Assuming K_{uu}^{-1} has been calculated in STEP 1., calculating the posterior distribution of \mathbf{f}^* only takes

$O(n_u n_u^*)$ time. Compared with $O(n_m^2 n_u^*)$ time regular GPR takes, TGPR algorithm has speeded up prediction.

B. Speeding up Training: Train n_m Times Individually

STEP 1. Calculating \mathbf{f}_u individually using each of n_m measurements.

Assuming that each test output $f_{m_1}, \dots, f_{m_{n_m}}$ is fully independent given \mathbf{f}_u . Firstly use the first measurement $(\mathbf{x}_{m_1}, \hat{f}_{m_1})$ to calculate the distribution of \mathbf{f}_u , resulting in $\mathbf{f}_u^1 \sim N(\boldsymbol{\mu}_u^1, \Sigma_{uu}^1)$. Then use the second measurement $(\mathbf{x}_{m_2}, \hat{f}_{m_2})$ to find $\mathbf{f}_u^2 \sim N(\boldsymbol{\mu}_u^2, \Sigma_{uu}^2)$, and so on. At last n_m distributions of \mathbf{f}_u are calculated.

STEP 2. Combining n_m distributions of \mathbf{f}_u .

The next thing to do is to combine all these n_m distributions together. Note that, each of these distributions contains the prior knowledge of \mathbf{f}_u . Combining n_m distributions will use the prior knowledge n_m times, so separating that prior distribution for n_m times is indispensable:

$$\mathbf{f}_u \sim N(\boldsymbol{\mu}_u, \Sigma_{uu}) = N(\boldsymbol{\mu}_u^1, \Sigma_{uu}^1) \oplus N(\boldsymbol{\mu}_u^2, \Sigma_{uu}^2) \oplus \dots \oplus N(\boldsymbol{\mu}_u^{n_m}, \Sigma_{uu}^{n_m})! \frac{N(\mathbf{m}_u, K_{uu})}{(n_m - 1) \text{ times}}. \quad (21)$$

Then the posterior distribution of \mathbf{f}_m and \mathbf{f}_u becomes:

$$\begin{bmatrix} \mathbf{f}_m \\ \mathbf{f}_u \end{bmatrix} \sim N\left(\begin{bmatrix} \boldsymbol{\mu}_m \\ \boldsymbol{\mu}_u \end{bmatrix}, \begin{bmatrix} \Sigma_{mm} & \Sigma_{mu} \\ \Sigma_{um} & \Sigma_{uu} \end{bmatrix}\right). \quad (22)$$

$$\begin{bmatrix} \boldsymbol{\mu}_m \\ \boldsymbol{\mu}_u \end{bmatrix} = \begin{bmatrix} \mathbf{m}_m + \Sigma_{mm} \hat{\Sigma}_{f_m}^{-1} (\hat{f}_m - \mathbf{m}_m) \\ \mathbf{m}_u + \Sigma_{um} \hat{\Sigma}_{f_m}^{-1} (\hat{f}_m - \mathbf{m}_m) \end{bmatrix}. \quad (23)$$

$$\begin{bmatrix} \Sigma_{mm} & \Sigma_{mu} \\ \Sigma_{um} & \Sigma_{uu} \end{bmatrix} = \begin{bmatrix} \Sigma_{mm} & \hat{\Sigma}_{f_m} (\Lambda_{mm} + \hat{\Sigma}_{f_m})^{-1} K_{mu} \Delta_{uu}^{-1} K_{um} \\ K_{um} \Delta_{uu}^{-1} K_{um} (\Lambda_{mm} + \hat{\Sigma}_{f_m})^{-1} \hat{\Sigma}_{f_m} & K_{um} \Delta_{uu}^{-1} K_{uu} \end{bmatrix}. \quad (24)$$

$$\Sigma_{mm} = (\Lambda_{mm}^{-1} + \hat{\Sigma}_{f_m}^{-1})^{-1} + \hat{\Sigma}_{f_m} (\Lambda_{mm} + \hat{\Sigma}_{f_m})^{-1} K_{mu} \Delta_{uu}^{-1} K_{um} (\Lambda_{mm} + \hat{\Sigma}_{f_m})^{-1} \hat{\Sigma}_{f_m}. \quad (25)$$

Where $\Lambda_{mm} = \text{diag}(K_{mm} - K_{mu} K_{uu}^{-1} K_{um})$,

$$\Delta_{uu} = K_{uu} + K_{um} (\Lambda_{mm} + \hat{\Sigma}_{f_m})^{-1} K_{mu}$$

Performance of TGPR

The combined posterior distribution of f_u can be obtained from (22), by which the training only takes $O(n_m n_u^2)$ time. Compared with $O(n_m^3)$ time regular GPR takes, TGPR algorithm has speeded up training.

C. Choosing the Inducing Inputs

The assumptions of TGPR algorithm cause K_{mm} to become:

$$K_{mm} = K_{mu} K_{uu}^{-1} K_{um} + \text{diag}(K_{mm} - K_{mu} K_{uu}^{-1} K_{um}) = K_{mu} K_{uu}^{-1} K_{um} + \Lambda_{mm}. \quad (26)$$

This means that K_{mm} now depends on K_{uu} . The idea to choose inducing inputs is to treat inducing inputs set X_u as another hyperparameter. That is, maximize the posterior

hyperparameter distribution $p(\theta | \hat{f}_m, X_m)$ with respect to K_{uu} . $p(\hat{f}_m | X_m, X_u)$ in (5) now becomes:

$$p(\hat{f}_m | X_m, X_u) = \mathcal{N}(\hat{f}_m | m_m, K_{mm} + \hat{\Sigma}_{f_m}). \quad (27)$$

Then by applying the Maximum Likelihood method in Section 2.2 to (26) and (27), the most likely inducing inputs will be found.

IV SIMULATION

The simulations are performed on a personal computer. And the Gaussian function used in the simulation is the semi-randomly chosen function as follows:

$$f(x_1, x_2) = \left(\frac{x_1}{4}\right)^2 + \left(\frac{x_2}{4}\right)^2 - \sin\left(2\pi \frac{x_1}{5}\right) \left(1 - \frac{1}{2} \cos\left(2\pi \frac{x_2}{4}\right)\right) + \sin\left(2\pi \frac{x_2}{5}\right) - 1. \quad (28)$$

The three-dimension model of (28), which lies on the interval $x_1, x_2 \in [-2, 2]$, is shown in Figure I.

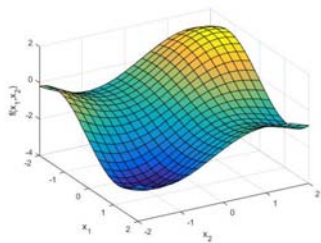


FIGURE I. THE THREE-DIMENSION MODEL

In simulations, the mean is zero mean function, while the covariance is SE with $\lambda_f = \lambda_{x_1} = \lambda_{x_2} = 1$. The function is affected by Gaussian white noise with $\hat{\sigma}_{f_m} = 0.5$.

One of the prediction results for above model using TGPR with $n_m = 500$ is shown in Figure II(a), while the error is shown in Figure II(b). The red dots are measurements; the blue plane is the exact function while the green plane is the prediction result; the translucent gray planes indicate 95% confidence intervals (CIs); the purple plane indicates error, being the prediction result minus the exact function.

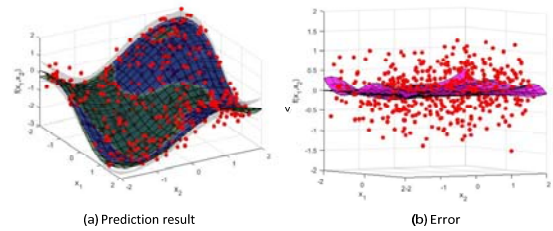


FIGURE II. THE PREDICTION RESULT AND ERROR

Five simulations for further analyses and comparisons are performed in the case of $n_u = 25$ and $n_u = 100$ respectively. The results are the averages of 100 simulation runs and are shown in Table I. The root mean square error (RMSE) is calculated over inputs to analyze the accuracy of the algorithms. A lower RMSE means a better accuracy.

TABLE I. SIMULATION RESULTS

Simulation	$n_u = 25$			$n_u = 100$		
	n_m	Runtime	RMS E	n_m	Run time	RM SE
1. Regular GPR with $n_m = 5000$	5000	2.09s	0.048	5000	2.01 s	0.049
2. Regular GPR with $n_m = 10000$	10000	10.11s	0.038	10000	9.95 s	0.038
3. TGPR with $n_m = 5000$	5000	0.028s	0.056	5000	0.065s	0.050
4. TGPR with $n_m = 10000$	10000	0.046s	0.047	10000	0.109s	0.038
5. TGPR with the same runtime as 1	10000	2.09s	0.042	25000	2.01 s	0.0086

The analyses for Table I are as follows:

(1) No matter which GPR algorithm is applied, the more measurements used in training step, the more accurate the prediction becomes, and more time the prediction will take.

(2) In TGPR algorithm, the more inducing inputs used, the more accurate the prediction becomes, and more time the prediction will take.

(3) In regular GPR, as the number of measurements grows big, there is a significant increase of runtime (Simulation 1 and 2).

(4) Given the same number of measurements, compared with regular GPR, TGPR reduces the regression time with a slight decrease in accuracy (Simulation 1 and 3).

(5) As the number of measurements grows from 5000 to 10000, the improvement of regression performance that gains from TGPR algorithm is more significant. TGPR reduces the regression time without any decrease in accuracy (Simulation 2 and 4).

(6) Given the same amount of runtime, compared with regular GPR, TGPR algorithm noticeably improves the accuracy of prediction because it can incorporate more measurements (Simulation 1 and 5).

V CONCLUSION

For reducing the computational complexity in regular GPR, Two-step Gaussian Process Regression algorithm is proposed in this paper. By using the inducing inputs, TGPR reduces the prediction time from $O(n_m^2 n_u)$ to $O(n_u n_m^2)$. By training n_m times individually, TGPR reduces the training time from $O(n_m^3)$ to $O(n_m n_u^2)$. The performance of TGPR algorithm and regular GPR are compared based on several simulations. The results show that, when working on a large number of measurements, TGPR algorithm significantly speeds up training and prediction. And due to that more measurements can be incorporated in less time, TGPR algorithm also provides a more accurate prediction.

ACKNOWLEDGEMENTS

This research is supported by Key Projects in the National Science & Technology Support Program under No. 2013BAK06B03 and the Natural Science Foundation of China (61673050).

REFERENCE

- [1] Rasmussen, Carl Edward, and Christopher KI Williams. Gaussian processes for machine learning. Vol. 1. Cambridge: MIT press, 2006.
- [2] Zhang, Guoqiang, B. Eddy Patuwo, and Michael Y. Hu. Forecasting with artificial neural networks: The state of the art. *International journal of forecasting* 14.1 (1998): 35-62.
- [3] Cristianini, Nello, and John Shawe-Taylor. An introduction to support vector machines and other kernel-based learning methods. Cambridge university press, 2000.
- [4] Williams, Christopher KI, and Matthias Seeger. Using the Nyström method to speed up kernel machines. *Advances in neural information processing systems*. 2001.
- [5] Snelson, Edward, and Zoubin Ghahramani. Sparse Gaussian processes using pseudo-inputs. *Advances in neural information processing systems*. 2006.
- [6] Tresp, Volker. A Bayesian committee machine. *Neural computation* 12.11 (2000): 2719-2741.
- [7] Titsias, Michalis K., and Neil D. Lawrence. Bayesian Gaussian process latent variable model. *International Conference on Artificial Intelligence and Statistics*. 2010.
- [8] Blum, Manuel, and Martin A. Riedmiller. Optimization of Gaussian process hyperparameters using Rprop. *ESANN*. 2013.
- [9] Bottou, Léon. Large-scale machine learning with stochastic gradient descent. *Proceedings of COMPSTAT2010*. Physica-Verlag HD, 2010. 177-186.
- [10] Quíñero-Candela, Joaquin, and Carl Edward Rasmussen. A unifying view of sparse approximate Gaussian process regression. *Journal of Machine Learning Research* 6.Dec (2005): 1939-1959.