

Gaussian Process Regression with Grid Spectral Mixture Kernel: Distributed Learning for Multidimensional Data

Presenter: Richard Cornelius Suwandi

The Chinese University of Hong Kong, Shenzhen

July 7, 2022

FUS  ON 2022

Contents

- 1 Introduction
- 2 Background
- 3 Problem Statement
- 4 Proposed Algorithms
- 5 Results
- 6 Conclusion

GPR with
GSM Kernel

R.C. Suwandi

Introduction

Background

Problem
Statement

Proposed
Algorithms

Results

Conclusion

Introduction

- Kernel design for Gaussian processes (GPs) along with the associated hyper-parameter optimization is a **challenging problem**
- The computational complexity for training the model hyper-parameters can be very **demanding** and even **prohibitive for large data sets**
- Large amount of labeled training data are usually aggregated from a large number of local agents or mobile devices, which may cause **severe data privacy issues**

- We propose a novel grid spectral mixture (GSM) kernel [1, 2] design for GPs that can **automatically fit multidimensional data**
- Two **efficient distributed learning** algorithms are proposed to alleviate the computational complexity owing to the curse of dimensionality in the kernel hyper-parameter optimization
- The proposed algorithms can help with **preserving data privacy** during the learning process

Background

- We consider the following **GP regression** model

$$y = f(\mathbf{x}) + e, \quad e \sim \mathcal{N}(0, \sigma_e^2) \quad (1)$$

where $f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}_h))$ is a real-valued, scalar Gaussian process with mean function $\mu(\mathbf{x})$ and covariance function $k(\mathbf{x}, \mathbf{x}'; \boldsymbol{\theta}_h)$

- The set of **unknown hyper-parameters** that needs to be tuned is denoted by $\boldsymbol{\theta} \triangleq [\boldsymbol{\theta}_h^T, \sigma_e^2]^T$

Grid Spectral Mixture (GSM) Kernel

- The basic idea behind the GSM kernel is to undertake an approximation of the underlying stationary kernel using the fact that any stationary kernel and its spectral density are **Fourier duals** [3]

Theorem

If the spectral density exists, then the stationary kernel function, $k(\boldsymbol{\tau})$, and its spectral density, $S(\boldsymbol{\omega})$, are Fourier duals of each other

$$k(\boldsymbol{\tau}) = \int_{\mathbb{R}^{d_x}} S(\boldsymbol{\omega}) \exp \left[j2\pi \boldsymbol{\tau}^\top \boldsymbol{\omega} \right] d\boldsymbol{\omega}, \quad (2a)$$

$$S(\boldsymbol{\omega}) = \int_{\mathbb{R}^{d_x}} k(\boldsymbol{\tau}) \exp \left[-j2\pi \boldsymbol{\tau}^\top \boldsymbol{\omega} \right] d\boldsymbol{\tau}. \quad (2b)$$

Grid Spectral Mixture (GSM) Kernel

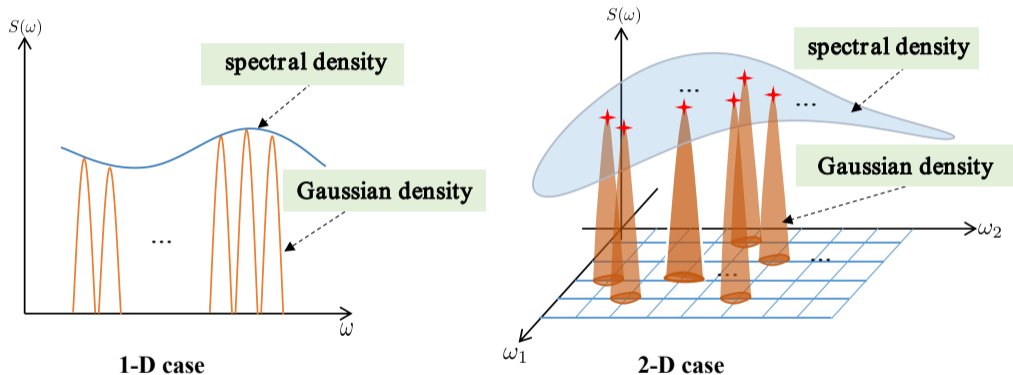


Figure: Basis kernel spectral density space under different input dimensions.

- The GSM kernel approximates the spectral density of the underlying kernel function in the frequency domain by a **Gaussian mixture**,

$$S(\omega) = \frac{1}{2} \sum_{q=1}^Q \theta_q [\mathcal{N}(\omega | \mu_q, v_q) + \mathcal{N}(\omega | -\mu_q, v_q)], \quad (3)$$

where $\{\mu_q\}_{q=1}^Q$ and $\{v_q\}_{q=1}^Q$ are fixed to **preselected grid points** and $\{\theta_q\}_{q=1}^Q$ are weights to be optimized

- Taking the inverse Fourier transform of $S(\omega)$, yields the **original GSM kernel** as

$$k(\tau) = \sum_{q=1}^Q \theta_q \underbrace{\cos(2\pi\tau\mu_q) \exp[-2\pi^2\tau^2v_q]}_{k_q(\tau)}, \quad (4)$$

where $\tau = |x - x'| \in \mathbb{R}$

- Learning the model hyper-parameters θ in GPR model typically resorts to the **type-II maximum likelihood**,

$$\hat{\theta} = \arg \min_{\theta} \underbrace{\mathbf{y}^{\top} [\mathbf{C}(\theta)]^{-1} \mathbf{y} + \log \det [\mathbf{C}(\theta)]}_{\triangleq l(\theta)} + \text{constant}, \quad (5)$$

where $\mathbf{C}(\theta) \triangleq \mathbf{K}_{XX} + \sigma_e^2 \mathbf{I}_n$

- The optimization problem in Eq. (5) is a well-known **difference-of-convex programming (DCP) problem** [4], where $g(\theta) \triangleq \mathbf{y}^{\top} [\mathbf{C}(\theta)]^{-1} \mathbf{y}$ and $h(\theta) \triangleq -\log \det [\mathbf{C}(\theta)]$ are convex functions w.r.t. θ
- This DCP problem can be efficiently solved using the **successive convex approximation (SCA) algorithm** [5]

- The vanilla SCA algorithm generates a sequence of feasible points $\boldsymbol{\theta}^t, t \in \mathbb{N}$ by solving

$$\boldsymbol{\theta}^{t+1} = \arg \min_{\boldsymbol{\theta}} \tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) \quad (6)$$

where $\tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) : \Theta \times \Theta \mapsto \mathbb{R}$ is called the **surrogate function**

Assumption

The surrogate function $\tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) : \Theta \times \Theta \mapsto \mathbb{R}$ satisfies the following conditions:

- $\tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$ is strongly convex on space Θ ;
- $\tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$ is differentiable with $\nabla_{\boldsymbol{\theta}} \tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = \nabla_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) \big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^t}$

- By performing the **first-order Taylor expansion**, we can make the convex function $h(\boldsymbol{\theta})$ affine and construct $\tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t)$:

$$\tilde{l}(\boldsymbol{\theta}, \boldsymbol{\theta}^t) = g(\boldsymbol{\theta}) - h(\boldsymbol{\theta}^t) - \nabla_{\boldsymbol{\theta}} h(\boldsymbol{\theta}^t)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^t) \quad (7)$$

- The problem in Eq. (6) becomes a **convex optimization problem**, and can be solved effectively by using the commercial solver MOSEK [6, 2]
- The computational complexity in each iteration scales as $\mathcal{O}(Qn^3)$, where n is the number of training samples and Q is the number of basis kernels

Problem Statement

GPR with
GSM Kernel

R.C. Suwandi

Introduction

Background

Problem
Statement

Proposed
Algorithms

Results

Conclusion

- We can similarly use a Gaussian mixture to approximate the spectral density of the underlying kernel function,

$$S(\boldsymbol{\omega}) = \sum_{q=1}^Q \theta_q [\mathcal{N}(\boldsymbol{\omega}; \boldsymbol{\mu}_q, \mathbf{V}_q) + \mathcal{N}(\boldsymbol{\omega}; -\boldsymbol{\mu}_q, \mathbf{V}_q)], \quad (8)$$

where $\boldsymbol{\mu}_q = [\mu_q^{(1)}, \dots, \mu_q^{(d_x)}]^\top$ and $\mathbf{V}_q = \text{diag}(v_q^{(1)}, \dots, v_q^{(d_x)})$

- Taking the inverse Fourier transform of $S(\boldsymbol{\omega})$ yields the **GSM kernel with multidimensional input** as

$$k(\boldsymbol{\tau}) = \sum_{q=1}^Q \theta_q \cos(2\pi \boldsymbol{\tau}^\top \boldsymbol{\mu}_q) \prod_{p=1}^{d_x} \exp\left\{-2\pi^2 \tau_p^2 v_q^{(p)}\right\} \quad (9)$$

- By assuming that each input dimension is **independent**, we empirically sample Q' frequencies, either uniformly or randomly, from the frequency region $[0, \mu_u^{(p)})$ for the p -th dimension and obtain $[\mu_1^{(p)}, \dots, \mu_{Q'}^{(p)}]$
- The highest frequency $\mu_u^{(p)}$ is set to be equal to $1/2$ over the **minimum input spacing** between two adjacent training data points in the dimension p
- Using the sampled frequencies, we can generate $Q = Q'^{d_x}$ **grid points** in the \mathbb{R}^{d_x} space
- Finally, using the generated grid points, we can construct Q **isotropic multivariate Gaussian densities** to approximate the underlying spectral density in the frequency domain

Proposed Algorithms

GPR with
GSM Kernel

R.C. Suwandi

Introduction

Background

Problem
Statement

Proposed
Algorithms

Results

Conclusion

- To alleviate the computational burden owing to the curse of dimensionality, we leverage a **multicore computing environment** to optimize θ in **parallel**
- The feasible set Θ in the GSM kernel admits a Cartesian product structure, i.e., $\Theta = \Theta_1 \times \Theta_2 \times \dots \times \Theta_s$ with $\Theta_i \subseteq \mathbb{R}^{Q/s}$
- We can **partition** the optimization variable into s blocks, $\theta = [\theta_1, \theta_2, \dots, \theta_s]^\top$, and construct a surrogate function that is **additively separable** in the blocks:

$$\tilde{l}(\theta, \theta^t) = \sum_{i=1}^s \tilde{l}_i(\theta_i, \theta^t). \quad (10)$$

where

$$\tilde{l}_i(\theta_i, \theta^t) = g(\theta_i, \theta_{-i}^t) - h(\theta^t) - \nabla_{\theta_i} h(\theta^t)^\top (\theta_i - \theta_i^t), \quad (11)$$

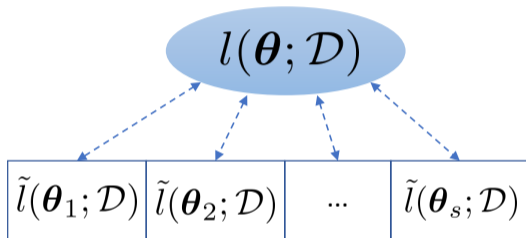


Figure: **Distributed SCA (DSCA)** for linear multiple kernel learning.

- The computational complexity of each computing core scales as $\mathcal{O}(\frac{Q}{s}n^3)$, where n is the number of training samples

Doubly Distributed SCA (D²SCA)

- DSCA is **impractical for big data** and **prone to data privacy issues**
- We propose a doubly distributed algorithm based on the **alternating direction method of multipliers (ADMM)** [7, 8] which enables N multicore agents to collaboratively learn the global hyper-parameters while preserving the data privacy of the local agents
- Each agent optimizes the hyper-parameters using its **local data** and then exchanges the hyper-parameters with a central agent to reach a **global consensus**

Doubly Distributed SCA (D^2SCA)

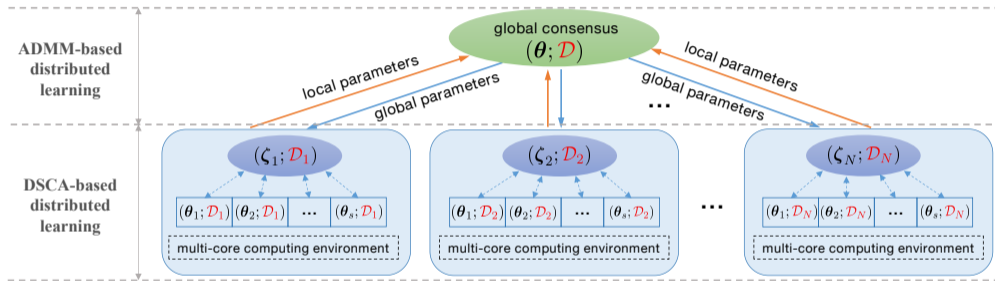


Figure: Doubly Distributed SCA (D^2SCA) for linear multiple kernel learning.

- The overall computational complexity of the D^2SCA algorithm scales as $O\left(\frac{Qn^3}{sN^3}\right)$

Results

GPR with
GSM Kernel

R.C. Suwandi

Introduction

Background

Problem
Statement

Proposed
Algorithms

Results

Conclusion

- We investigate the training and prediction performance of the proposed learning algorithms, **DSCA** and **D²SCA**, on various data sets
- We have selected 8 **one-dimensional input data sets** and 4 **multidimensional input data sets** as our benchmark
- We compare the proposed GSM kernel-based GP (GSMGP) with the SM kernel-based GP (SMGP) proposed by Wilson *et al.* in [9] and the squared-exponential kernel-based GP (SEGP) in terms of the **prediction mean squared error (MSE)**

1-D Case: Prediction Performance

| Data Set | GSMGP DSCA | GSMGP D ² SCA | SMGP | SEGP | LSTM | ARIMA |
|--------------|----------------|-----------------------------|---------|---------|---------|---------|
| ECG | 1.1E-02 | 1.2E-02 | 1.9E-02 | 1.6E-01 | 1.6E-01 | 1.8E-01 |
| CO2 | 9.2E-01 | 1.4E+00 | 1.1E+00 | 1.5E+03 | 2.9E+02 | 4.9E+00 |
| Electricity | 4.3E+03 | 3.6E+03 | 7.5E+03 | 8.3E+03 | 8.0E+03 | 1.2E+04 |
| Employment | 5.4E+01 | 7.0E+01 | 7.0E+02 | 8.4E+03 | 1.9E+03 | 3.9E+02 |
| Hotel | 4.2E+02 | 1.5E+03 | 2.8E+03 | 5.6E+04 | 5.0E+04 | 1.7E+04 |
| Passenger | 6.9E+01 | 1.1E+02 | 1.6E+02 | 8.8E+02 | 7.0E+02 | 4.5E+03 |
| Clay | 8.5E+01 | 2.4E+02 | 3.3E+02 | 1.5E+03 | 3.6E+02 | 3.3E+02 |
| Unemployment | 2.0E+03 | 3.1E+03 | 1.4E+04 | 5.6E+05 | 1.7E+05 | 1.5E+04 |

Table: Performance comparison between the proposed GSMGP and its competitors in terms of the prediction MSE.

1-D Case: Prediction Performance

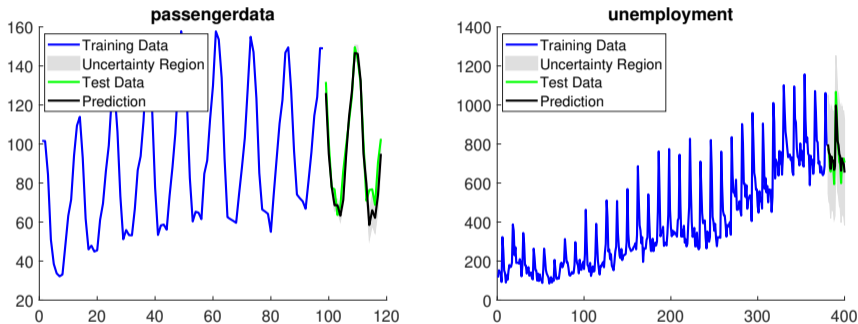


Figure: Training and prediction performance of the GSMGP with $\sigma = 0.001$ and $Q = 500$ uniformly generated grids. The optimal weights are solved via the distributed SCA (DSCA) algorithm.

- GPR with GSM Kernel
- R.C. Suwandi
- Introduction
- Background
- Problem Statement
- Proposed Algorithms
- Results**
- Conclusion

1-D Case: Training Performance

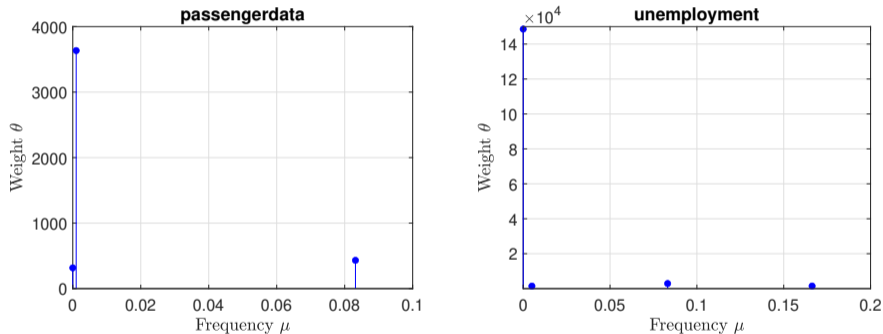


Figure: Estimated weights and frequencies generated by the distributed SCA (DSCA) algorithm for selected data sets.

- GPR with GSM Kernel
- R.C. Suwandi
- Introduction
- Background
- Problem Statement
- Proposed Algorithms
- Results
- Conclusion

| Data Set | GSMGP DSCA | GSMGP D ² SCA | SMGP | SEGP | LSTM |
|----------|----------------|-----------------------------|---------|---------|---------|
| ALE | 2.4E-02 | 2.3E-02 | 3.8E-01 | 3.7E-02 | 3.4E-02 |
| CCCP | 1.9E+01 | 1.6E+01 | 2.1E+05 | 1.7E+01 | 2.8E+02 |
| Airfoil | 1.7E+01 | 5.2E+01 | 6.9E+01 | 7.7E+01 | 7.3E+01 |
| Concrete | 6.7E+01 | 4.0E+01 | 1.7E+03 | 1.3E+02 | 1.4E+02 |

Table: Performance comparison between the proposed GSMGP and its competitors, SMGP and SEGP, in terms of the **prediction MSE**.

| Data Set | GSMGP-DSCA # of Iter. | GSMGP-D ² SCA # of Iter. | SMGP # of Iter. | SEGP # of Iter. |
|----------|--------------------------|--|--------------------|--------------------|
| ALE | 4 | 3 | 63 | 122 |
| CCCP | 5 | 7 | 82 | 141 |
| Airfoil | 7 | 5 | 500 | 137 |
| Concrete | 3 | 3 | 500 | 154 |

¹ The number of iterations in GSMGP-D²SCA is the number of global iterations.

Table: Performance comparison between the proposed GSMGP and its competitors, SMGP and SEGP, in terms of the total number of iterations.

Conclusion

GPR with
GSM Kernel

R.C. Suwandi

Introduction

Background





Problem
Statement





Proposed
Algorithms

Results

Conclusion

- We propose an extension of the grid spectral mixture (GSM) kernel for **multidimensional input**
- Two efficient distributed learning algorithms, namely DSCA and D²SCA, are proposed to **alleviate the computational complexity** owing to the curse of dimensionality in the kernel hyper-parameter optimization
- The proposed algorithms can learn the global hyper-parameters with **lower computational complexity** and **preserve data privacy** during the learning process
- Experimental results verify that the proposed GSM kernel and the associated learning algorithms are **superior in terms of training and prediction performance** compared to their competitors

-  F. Yin, X. He, L. Pan, T. Chen, Z.-Q. Luo, and S. Theodoridis, “Sparse structure enabled grid spectral mixture kernel for temporal Gaussian process regression,” in *Proc. Int. Conf. Inf. Fusion (FUSION)*, (Cambridge, UK), pp. 47–54, July 2018.
-  F. Yin, L. Pan, T. Chen, S. Theodoridis, Z.-Q. Luo, and A. M. Zoubir, “Linear multiple low-rank kernel based stationary Gaussian processes regression for time series,” *IEEE Trans. Signal Process.*, vol. 68, pp. 5260–5275, 2020.
-  C. E. Rasmussen and C. I. K. Williams, *Gaussian Processes for Machine Learning*. MIT Press, 2006.
-  S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.

-  G. Scutari and Y. Sun, “Parallel and distributed successive convex approximation methods for big-data optimization,” in *Multi-agent Optimization*, pp. 141–308, Springer, 2018.
-  MOSEK ApS, “MOSEK optimization toolbox for Matlab,” *User’s Guide and Reference Manual, Version*, vol. 4, 2019.
-  S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed optimization and statistical learning via the alternating direction method of multipliers,” *Found. Trends Mach. Learn.*, vol. 3, pp. 1–122, January 2011.
-  Y. Xu, F. Yin, W. Xu, J. Lin, and S. Cui, “Wireless traffic prediction with scalable gaussian process: Framework, algorithms, and verification,” *IEEE J. Sel. Areas Commun.*, vol. 37, no. 6, pp. 1291–1306, 2019.



A. Wilson and R. P. Adams, “Gaussian process kernels for pattern discovery and extrapolation,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, (Atlanta, USA), pp. 1067–1075, 2013.

GPR with
GSM Kernel

R.C. Suwandi

Introduction

Background

Problem
Statement

Proposed
Algorithms

Results

Conclusion