

Background

What is GP?

Gaussian process (GP) model is a class of important **Bayesian non-parametric models for machine learning**.

Application

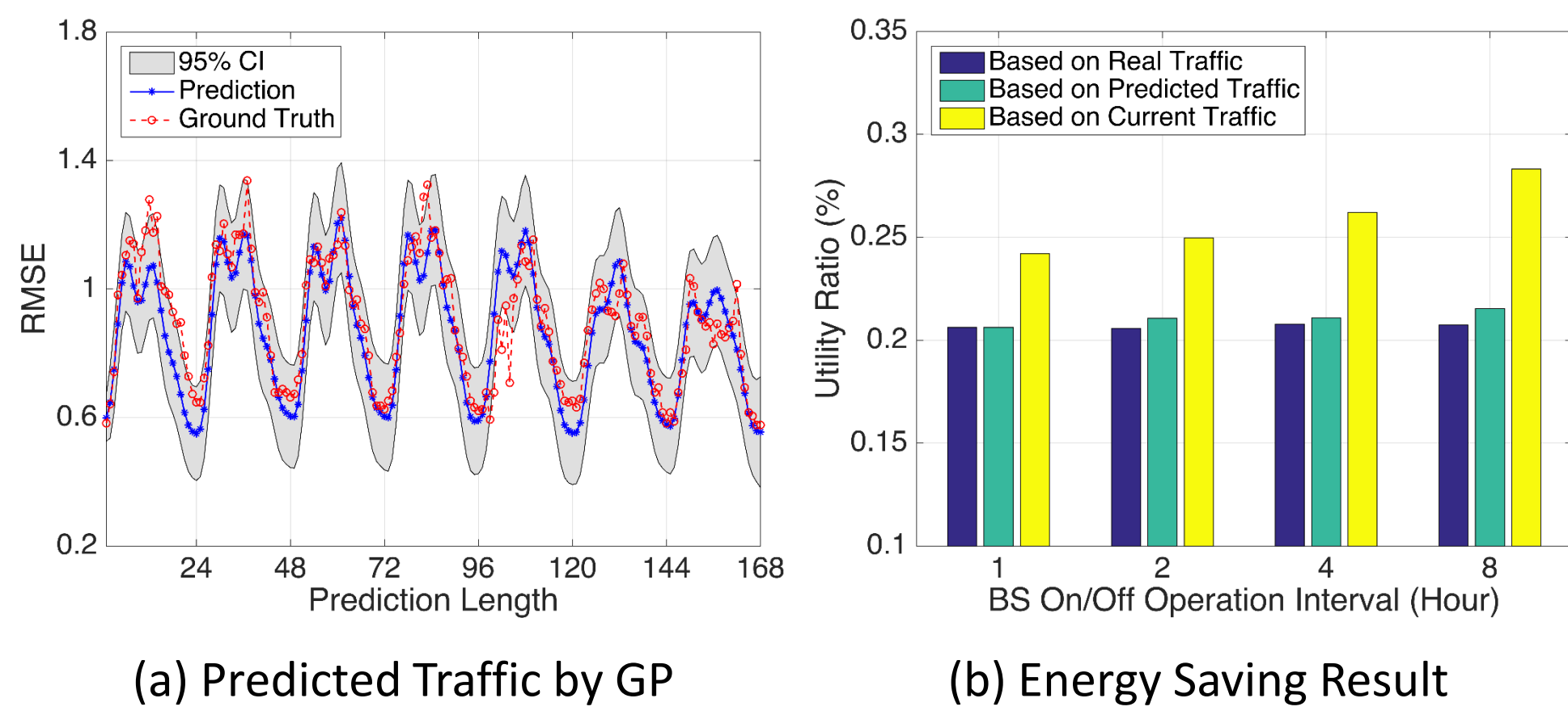


Figure 1. Performance of wireless traffic prediction based energy saving.

Challenge

- Standard GP suffers from the high complexity of **hyper-parameter optimization**, which scales as $\mathcal{O}(n^3)$ with n the number of training samples.
- Existing low-complex GP model reduces the complexity based on certain approximations, e.g., the subset-of-data (SOD) model based on sparse GP, the Bayesian committee machine (BCM).

Main Result

A Scalable GP Model for Processing Big Datasets

A novel **scalable GP regression model**, which is **parallelizable** over a large number of computation units and **does not involve any approximation** essentially.

Faster Hyper-parameter Optimization

A practical implementation with the Gauss-Seidel method, which reduces the complexity to $\mathcal{O}(n^3/k)$ with k the number of parallel computing units.

Simulation

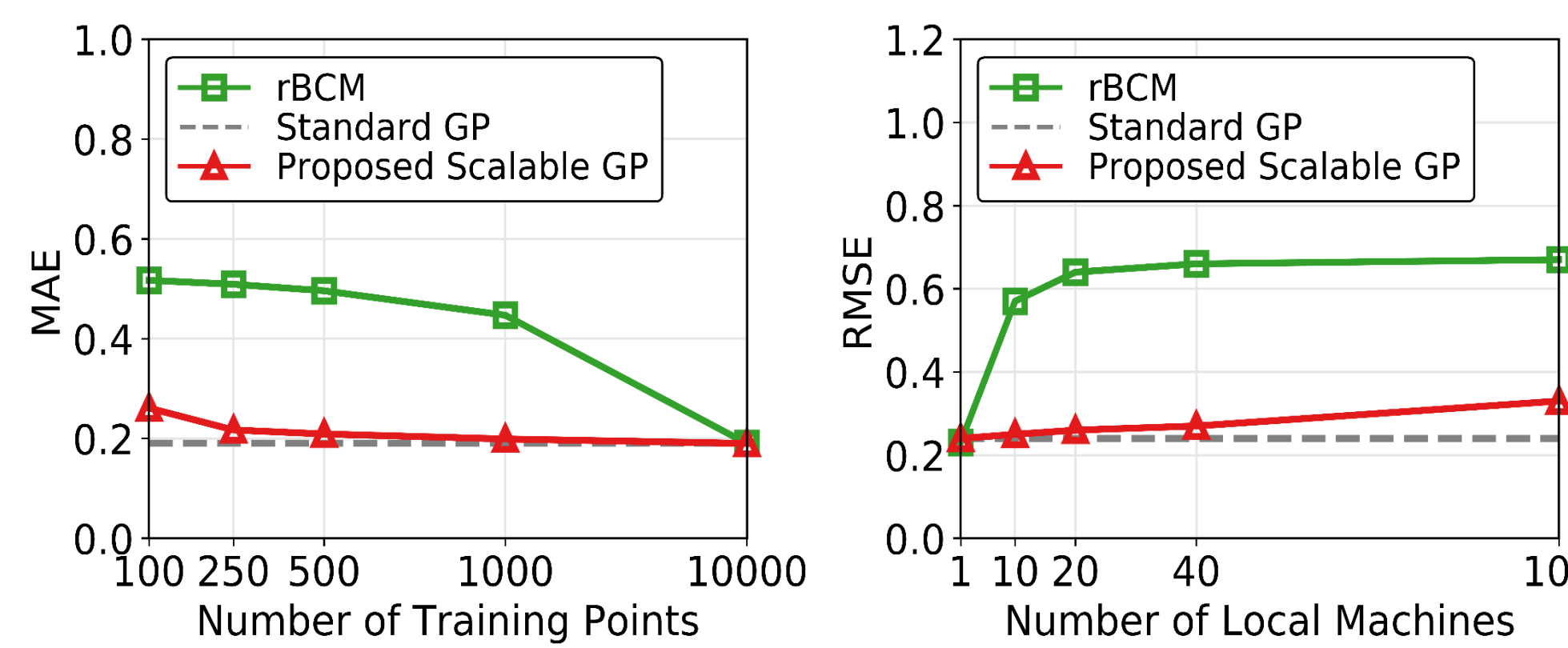


Figure 2: The Prediction Performance.

Baseline Model

- The state-of-art robust BCM (rBCM) model.
- The standard GP model (Optimizes the GP parameters as in \mathcal{P}_0)

Simulation Setting

- Artificially generated datasets with SE kernel.
- Using 10000 points as the training dataset to predict the next upcoming data point.
- Repeat 300 times (iteratively update the training set) to average the performance.
- Run at a workstation with eight Intel E3 Xeon CPU cores at 3.50 GHz.

Regression Model

GP Definition

A GP is a collection of random variables, any finite number of which follow a Gaussian distribution.

GP Function

$$f(x) \sim \mathcal{GP}(m(x), k(x, x'; \theta)),$$

where

- x : continuous-value input;
- $m(x)$: mean function (zero in practice);
- $k(x, x'; \theta)$: kernel function (e.g., SE, periodic).

GP-based Regression Model

$$y = f(x) + e,$$

where

- y : continuous-value output;
- e : noise (estimated independently).

Standard GP

Standard GP Hyper-parameter Optimization

$$\mathcal{P}_0: \arg \min_{\theta} g(\theta) = y^T C^{-1}(\theta) y + \log |C(\theta)|,$$

$$s.t. \quad \theta \in \Theta,$$

where

- $C(\theta) = K(X, X; \theta) + \sigma^2 I_n$: covariance matrix;
- $K(X, X; \theta)$: kernel matrix.

Gradient Decent (Benchmark Method)

$$\theta^{r+1} = \theta^r - \mu \cdot \nabla_{\theta} g(\theta)|_{\theta=\theta^r},$$

where

- $\frac{\partial g(\theta)}{\partial \theta_j} = \text{Tr} \left(C^{-1}(\theta) \frac{\partial C(\theta)}{\partial \theta_j} \right) - y^T C^{-1}(\theta) \frac{\partial C(\theta)}{\partial \theta_j} C^{-1}(\theta) y$;
- $C^{-1}(\theta)$: $\mathcal{O}(n^3)$ computational complexity.

A Motivating Example

1 Given a covariance matrix with four blocks:

$$C(\theta_1, \theta_2, \theta_3, \theta_4) = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix}.$$

2 Computation in the first unit:

$$C(\theta_1, z_{-1}^r) = \begin{bmatrix} a_{11} & a_{12} & a_{13}^r & a_{14}^r \\ a_{21} & a_{22} & a_{23}^r & a_{24}^r \\ a_{31}^r & a_{32}^r & a_{33}^r & a_{34}^r \\ a_{41}^r & a_{42}^r & a_{43}^r & a_{44}^r \end{bmatrix} = \begin{bmatrix} C_{11}(\theta_1) & C_{12}(\theta_2) \\ C_{21}(\theta_3) & C_{22}(\theta_4) \end{bmatrix}.$$

3 The block-wise matrix inverse satisfies:

$$\begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} I_{\frac{n}{2} \times \frac{n}{2}} & \mathbf{0}_{\frac{n}{2} \times \frac{n}{2}} \\ \mathbf{0}_{\frac{n}{2} \times \frac{n}{2}} & I_{\frac{n}{2} \times \frac{n}{2}} \end{bmatrix}.$$

4 The matrix derivative satisfies:

$$\frac{\partial C}{\partial \theta_{ij}} = \begin{bmatrix} \frac{\partial C_{11}(\theta)}{\partial \theta_{ij}} & \mathbf{0}_{\frac{n}{2} \times \frac{n}{2}} \\ \mathbf{0}_{\frac{n}{2} \times \frac{n}{2}} & \mathbf{0}_{\frac{n}{2} \times \frac{n}{2}} \end{bmatrix}.$$

5 Solution of the block-wise matrix inverse:

$$B_{11} = (C_{11} - C_{12} C_{22}^{-1} C_{21})^{-1},$$

$$B_{12} = B_{21}^T = -C_{22}^{-1} C_{21} B_{11},$$

$$B_{22} = (C_{22} - C_{21} C_{11}^{-1} C_{12})^{-1}.$$

Scalable GP

Goal

Break the problem \mathcal{P}_0 into smaller pieces that are easier to handle distributively without making any approximations.

ADMM-based GP Hyper-parameter Optimization

$$\mathcal{P}_1: \arg \min_{\{\theta_i\}} g(\{\theta_i\}),$$

$$s.t. \quad \theta_i - z = \mathbf{0}, \theta_i \in \Theta, i \in \mathcal{K},$$

where:

- $-g(\{\theta_i\}) = y^T C^{-1}(\{\theta_i\}) y + \log |y^T C^{-1}(\{\theta_i\})|$;
- $C^{-1}(\{\theta_i\})$: i -th block determined by θ_i .

Remark: \mathcal{P}_1 is equivalent to \mathcal{P}_0 .

Lagrangian Function

$$\mathcal{L}(\{\theta_i, z, \beta\}) \triangleq g(\{\theta_i\}) + \sum_{i=1}^k \beta_i^T (\theta_i - z) + \sum_{i=1}^k \frac{\rho}{2} \|\theta_i - z\|_2^2.$$

ADMM Iteration

$$\theta_i^{r+1} = \arg \min_{\theta_i} g(\theta_i, z_{-i}^r) + \beta_i^{r,T} (\theta_i - z^r) + \frac{\rho}{2} \|\theta_i - z^r\|_2^2,$$

$$z^{r+1} = \frac{1}{k} \sum_{i=1}^k \left(\theta_i^{r+1} + \frac{1}{\rho} \beta_i^r \right),$$

$$\beta_i^{r+1} = \beta_i^r + \rho (\theta_i^{r+1} - z^{r+1}).$$

where:

- $-g(\theta_i, z_{-i}^r) = y^T C^{-1}(\theta_i, z_{-i}^r) y + \log |y^T C^{-1}(\theta_i, z_{-i}^r)|$;
- $-z_{-i}^r = \{\theta_1^r, \dots, \theta_{i-1}^r, \theta_{i+1}^r, \dots, \theta_k^r\}$;
- $-C^{-1}(\theta_i, z_{-i}^r)$: one block determined by θ_i , the others by z_{-i}^r .

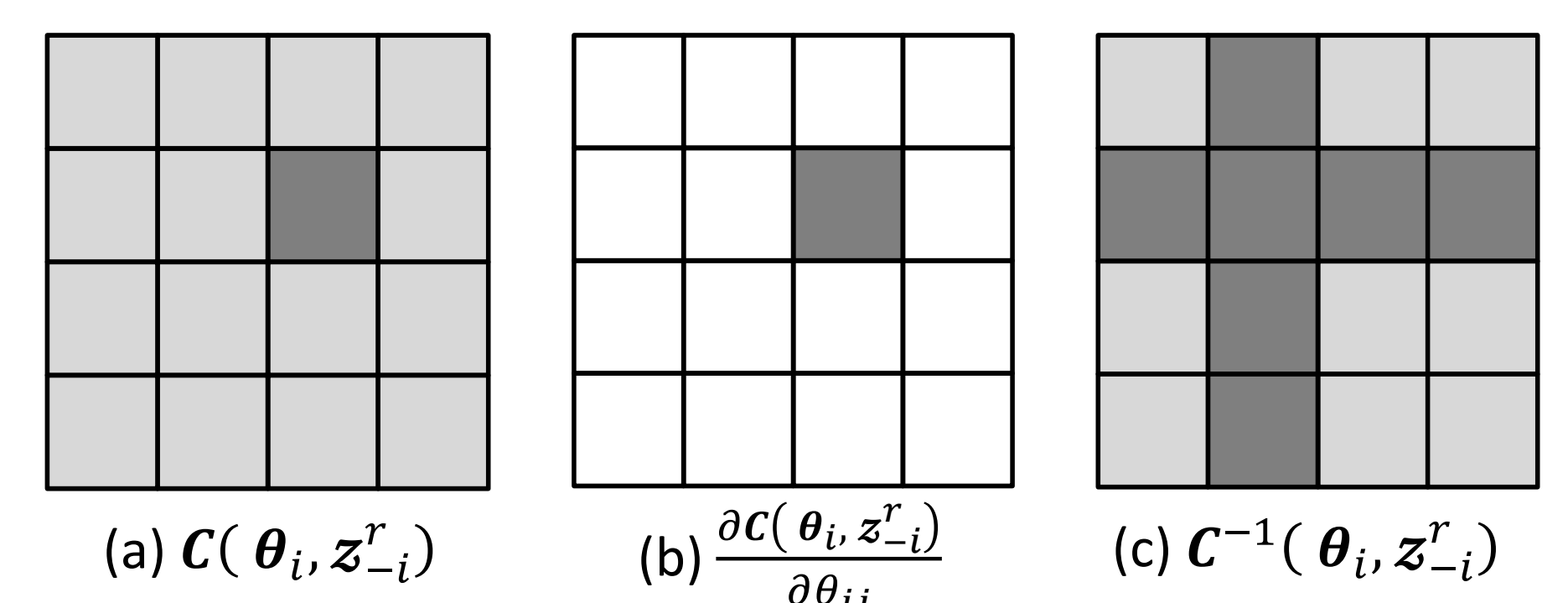


Figure 3. When updating the local hyper-parameter, say θ_i , the i -th local computing unit requires only one block of the full covariance matrix, e.g., the dark block in (a); only one block of the partial derivative matrix is non-zero, e.g., the dark block in (b); only one vertical and one horizontal slice of the full matrix inverse is needed, e.g., the two dark slices in (c), for gradient update.

A Practical Implementation (Gauss-Seidel Method)

$$\left[\nabla_{\theta_i} \mathcal{L}(\theta_i, z_{-i}^r, z^r, \beta^r) \right]_{\theta_i = \theta_i^r} = \frac{\partial}{\partial \theta_{ij}} \left(g(\theta_i, z_{-i}^r) + (\beta_i^r)^T (\theta_i - z) + \frac{\rho}{2} \|\theta_i - z^r\|_2^2 \right) \Bigg|_{\theta_i = \theta_i^r}$$

$$= \text{Tr} \left(C^{-1}(\theta_i, z_{-i}^r) \frac{\partial C(\theta_i, z_{-i}^r)}{\partial \theta_{ij}} \right) - y^T C^{-1}(\theta_i, z_{-i}^r) \frac{\partial C(\theta_i, z_{-i}^r)}{\partial \theta_{ij}} C^{-1}(\theta_i, z_{-i}^r) y + \beta_{ij}^r + \rho (\theta_{ij} - z_j^r) \Bigg|_{\theta_i = \theta_i^r}$$

$$\approx \text{Tr} \left(C^{-1}(z^r) \frac{\partial C(\theta_i, z_{-i}^r)}{\partial \theta_{ij}} \right) - y^T C^{-1}(z^r) \frac{\partial C(\theta_i, z_{-i}^r)}{\partial \theta_{ij}} C^{-1}(z^r) y + \beta_{ij}^r + \rho (\theta_{ij} - z_j^r).$$