

# Improving the Gaussian Process Sparse Spectrum Approximation by Representing Uncertainty in Frequency Inputs

Yarin Gal, Richard Turner

University of Cambridge  
{yg279, ret26}@cam.ac.uk



## In short:

- Gaussian process sparse pseudo-input approximations **cannot handle complex functions well**.
- Alternatively, sparse spectrum approximations are **known to over-fit**.

We develop **variational inference for the sparse spectrum approximation** avoiding both issues, **treating the covariance function as a random variable**.

## Background

### What is a Gaussian process (GP)?

- A powerful tool in statistics, robust to over-fitting.
- Models distributions over functions.
- Supervised/unsupervised, regression/classification.
- Offers uncertainty estimates over the function values (in blue).

- Given training inputs  $\mathbf{X} = \{\mathbf{x}_i\}_{i=1}^N \in \mathbb{R}^{N \times Q}$  and outputs  $\mathbf{Y} = \{\mathbf{y}_i\}_{i=1}^N \in \mathbb{R}^{N \times D}$ , estimate a function  $\mathbf{y} = \mathbf{f}(\mathbf{x})$  that is **likely to have generated Y**.

- We place a joint Gaussian distribution over all function values:

$$p(\mathbf{Y} | \mathbf{X}) = \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X}) + \tau^{-1} \mathbf{I}_N)$$

with precision hyper-parameter  $\tau$  and covariance function  $\mathbf{K}(\mathbf{X}, \mathbf{X})$ .

### What is variational inference?

- Condition the model on a finite set of random variables  $\omega$ .

- The predictive distribution for a new input point  $\mathbf{x}^*$

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^* | \mathbf{x}^*, \omega) p(\omega | \mathbf{X}, \mathbf{Y}) d\omega,$$

- The distribution  $p(\omega | \mathbf{X}, \mathbf{Y})$  cannot be evaluated analytically — define an “easier” approximating *variational* distribution  $q(\omega)$ .

- Minimise the Kullback–Leibler (KL) divergence:  $\text{KL}(q(\omega) || p(\omega | \mathbf{X}, \mathbf{Y}))$ .

- Minimising the KL divergence = maximising *log evidence lower bound*,

$$\mathcal{L}_{VI} := \int q(\omega) \log p(\mathbf{Y} | \mathbf{X}, \omega) d\omega - \text{KL}(q(\omega) || p(\omega))$$

with respect to the variational parameters defining  $q(\omega)$ .

## Variational Sparse Spectrum Approximation

1. We are given the Fourier transform of the squared exponential covariance function

$$\mathbf{K}(\mathbf{x} - \mathbf{y}) = \int \sigma^2 \mathcal{N}(\mathbf{w}; 0, \mathbf{I}_Q) \cos(2\pi \mathbf{w}^T (\mathbf{x} - \mathbf{y})) d\mathbf{w}$$

2. Introduce auxiliary variable  $b$ :

$$= \int 2\sigma^2 \mathcal{N}(\mathbf{w}; 0, \mathbf{I}_Q) \text{Unif}[0, 2\pi] \cos(2\pi \mathbf{w}^T \mathbf{x} + b) \cos(2\pi \mathbf{w}^T \mathbf{y} + b) d\mathbf{w} db$$

3. Approximate with Monte Carlo integration with  $K$  terms:

$$\hat{\mathbf{K}}(\mathbf{x} - \mathbf{y}) = \frac{\sigma^2}{K} \sum_{k=1}^K \sqrt{2} \cos(2\pi \mathbf{w}_k^T (\mathbf{x} - \mathbf{z}_k) + b_k) \sqrt{2} \cos(2\pi \mathbf{w}_k^T (\mathbf{y} - \mathbf{z}_k) + b_k)$$

with  $\mathbf{w}_k \sim \mathcal{N}(0, \mathbf{I}_Q)$ ,  $b_k \sim \text{Unif}[0, 2\pi]$ , and variational parameters  $\mathbf{z}_k$ . This is a **random covariance function**.

4. The GP is re-parametrised as

$$\mathbf{w}_k \sim \mathcal{N}(0, \mathbf{I}_Q), \quad b_k \sim \text{Unif}[0, 2\pi], \quad \omega = \{\mathbf{w}_k, b_k\}_{k=1}^K$$

$$\Phi_{n,k}(\omega) = \sqrt{\frac{2\sigma^2}{K}} \cos(2\pi \mathbf{w}_k^T (\mathbf{x}_n - \mathbf{z}_k) + b_k), \quad \Phi \in \mathbb{R}^{N \times K}$$

$$p(\mathbf{Y} | \mathbf{X}, \omega) = \mathcal{N}(\mathbf{Y}; \mathbf{0}, \Phi(\omega) \Phi(\omega)^T + \tau^{-1} \mathbf{I}_N)$$

$$p(\mathbf{Y} | \mathbf{X}) = \int p(\mathbf{Y} | \mathbf{X}, \omega) p(\omega) d\omega$$

5. Use variational distribution  $q(\omega) = \prod_{k=1}^K q(\mathbf{w}_k) q(b_k)$  to approximate posterior  $p(\omega | \mathbf{X}, \mathbf{Y})$ :

$$q(\mathbf{w}_k) = \mathcal{N}(\mu_k, \Sigma_K), \quad q(b_k) = \text{Unif}(\alpha_k, \beta_k),$$

with  $\Sigma_K$  diagonal.

6. We maximise the log evidence lower bound:

$$\mathcal{L}_{VSSGP} = \sum_{d=1}^D \left( -\frac{N}{2} \log(2\pi\tau^{-1}) - \frac{\tau}{2} \mathbf{y}_d^T \mathbf{y}_d + \frac{1}{2} \log(|\tau^{-1} \Sigma|) \right.$$

$$\left. + \frac{1}{2} \tau \mathbf{y}_d^T E_{q(\omega)}(\Phi) \Sigma E_{q(\omega)}(\Phi^T) \mathbf{y}_d \right) - \text{KL}(q(\omega) || p(\omega))$$

with  $\Sigma = (E_{q(\omega)}(\Phi^T \Phi) + \tau^{-1} \mathbf{I})^{-1}$ . We can evaluate the KL and the expectations analytically using the identity

$$E_{q(\mathbf{w})}(\cos(\mathbf{w}^T \mathbf{x} + b)) = e^{-\frac{1}{2} \mathbf{x}^T \Sigma \mathbf{x}} \cos(\mu^T \mathbf{x} + b).$$

7. This requires  $\mathcal{O}(NK^2 + K^3)$  time complexity. Using parallel inference with  $T$  workers this reduces to  $\mathcal{O}(\frac{NK^2}{T} + K^3)$ .

## Factorised Variational Sparse Spectrum Approximation

We often use large  $K$ , for which the above is still slow. We can do better.

1. We introduce auxiliary random variables  $\mathbf{a}_d$ . The GP is re-parametrised as

$$\mathbf{a}_d \sim \mathcal{N}(0, \mathbf{I}_K), \quad \mathbf{A} = [\mathbf{a}_d]_{d=1}^D \in \mathbb{R}^{K \times D}$$

$$p(\mathbf{Y} | \mathbf{X}, \mathbf{A}, \omega) = \mathcal{N}(\mathbf{Y}; \Phi(\omega) \mathbf{A}, \tau^{-1} \mathbf{I}_N)$$

$$p(\mathbf{Y} | \mathbf{X}) = \int p(\mathbf{Y} | \mathbf{X}, \mathbf{A}, \omega) p(\omega) d\omega d\mathbf{A}$$

2. Use variational distribution  $q(\omega) = \prod_{k=1}^K q(\mathbf{w}_k) q(b_k) \prod_{d=1}^D q(\mathbf{a}_d)$  to approximate posterior  $p(\omega, \mathbf{A} | \mathbf{X}, \mathbf{Y})$ :

$$q(\mathbf{a}_d) = \mathcal{N}(\mathbf{m}_d, \mathbf{s}_d)$$

with  $\mathbf{s}_K$  diagonal.

3. This requires  $\mathcal{O}(NK^2)$  time complexity — no matrix inverse. Using parallel inference with  $T$  workers this reduces to  $\mathcal{O}(\frac{NK^2}{T})$ .

## Stochastic Factorised Variational Sparse Spectrum Approximation

With stochastic optimisation, we can do even better.

$$\mathcal{L}_{sfVSSGP} \approx \frac{N}{|S|} \sum_{n \in S} \sum_{d=1}^D \mathcal{L}_{nd} - \text{KL}(q(\mathbf{A}) || p(\mathbf{A})) - \text{KL}(q(\omega) || p(\omega))$$

with random data subset  $S$  is an unbiased estimator of  $\mathcal{L}_{fVSSGP}$ .

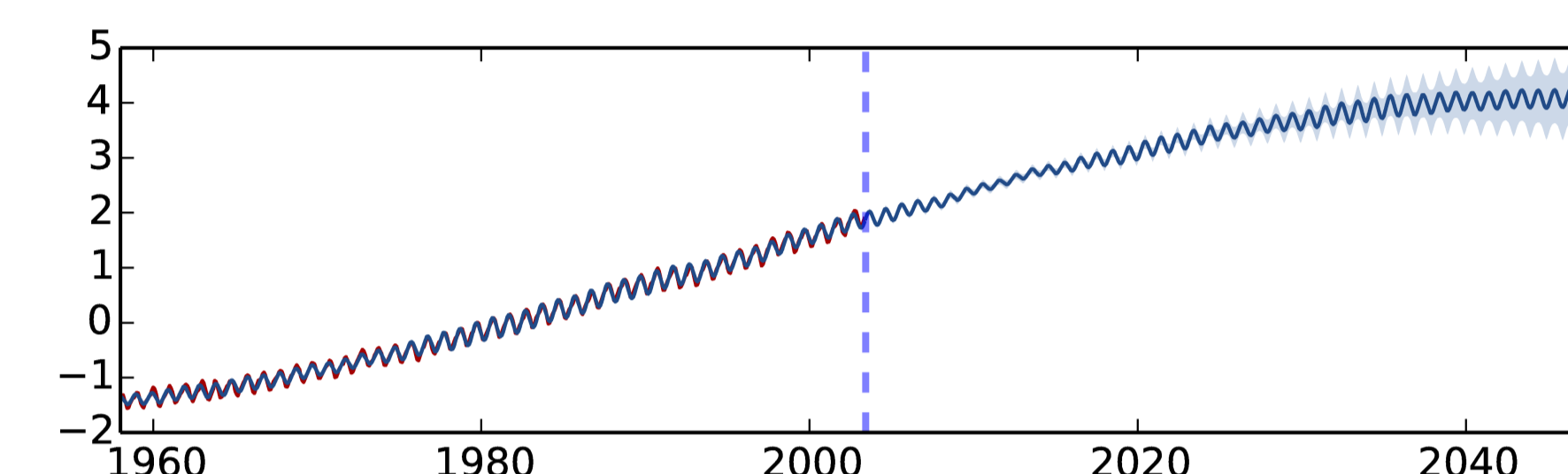
- Requires  $\mathcal{O}(SK^2)$  time complexity with  $S \ll N$  size of random subset, compared to  $\mathcal{O}(SK^2 + K^3)$  of GP SVI using sparse pseudo-input approximation.

## Approximation Properties

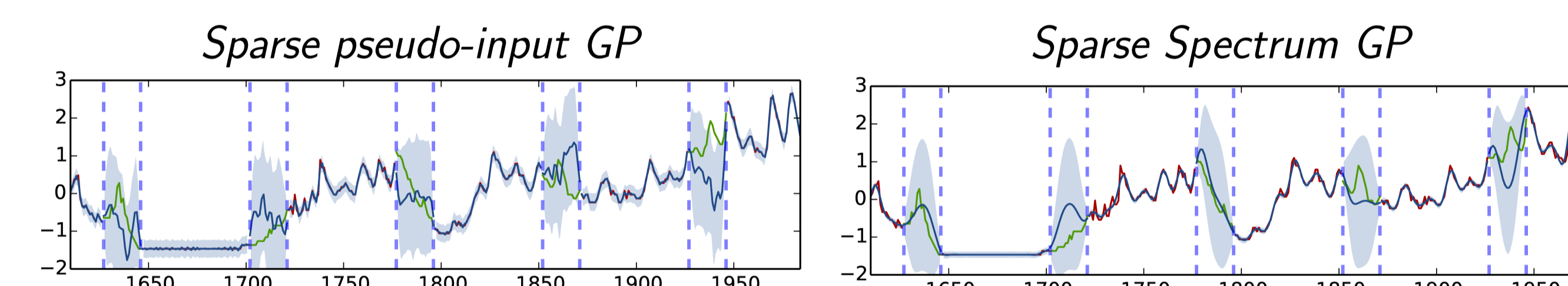
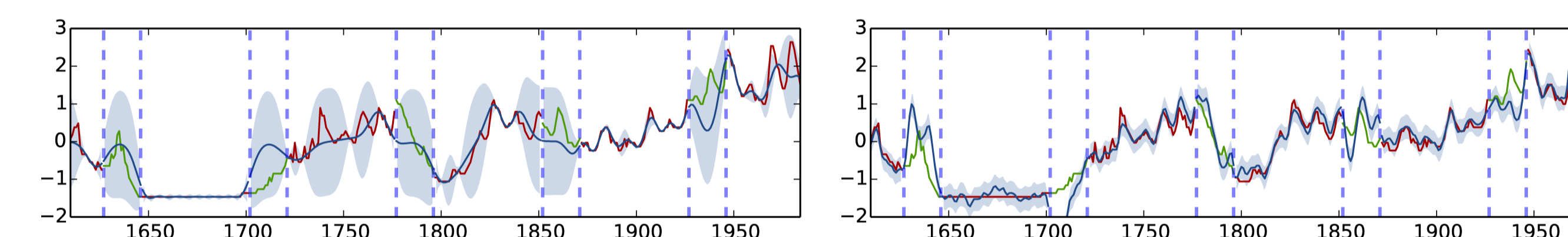
### Model uncertainty

- Model uncertainty in interpolation and extrapolation tasks

Extrapolation on the Mauna Loa CO<sub>2</sub> concentrations dataset (sum of two SE covariance functions,  $K = 10$  inducing inputs):



Interpolation on the reconstructed solar irradiance dataset (SE covariance function,  $K = 50$  inducing inputs):



	Solar	SPGP	SSGP	RP <sub>1</sub>	RP <sub>2</sub>	GP	VSSGP
Train	0.23	0.15	0.32	0.04	0.08	<b>0.13</b>	
Test	0.61	0.63	0.65	0.76	0.50	<b>0.41</b>	

Interpolation RMSE on train / test sets

### Variational Sparse Spectrum GP

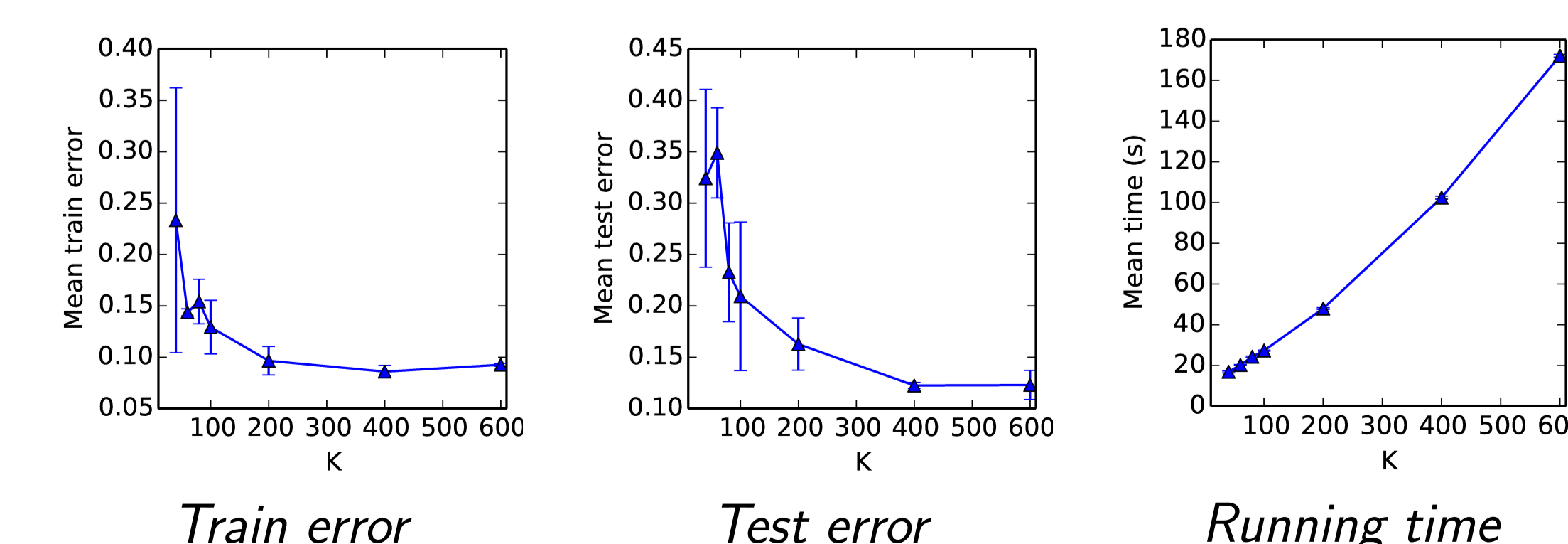
### Imputation accuracy

- Imputation RMSE on train / test sets, for a speech signal of length 1K ( $K = 100$ ):

Audio 1K	SSGP	VSSGP	fVSSGP	sfVSSGP
Train	0.0091 ± 0.0042	<b>0.0062 ± 0.00048</b>	0.0054 ± 0.00083	0.005 ± 0.003
Test	0.088 ± 0.033	<b>0.034 ± 0.0043</b>	0.038 ± 0.0049	0.04 ± 0.0066

### Speed-Accuracy trade-off

- Train error, test error, and running time as functions of number of inducing points ( $K$ ) for a speech signal of length 4K:



Code available at <http://github.com/yaringal/VSSGP>.