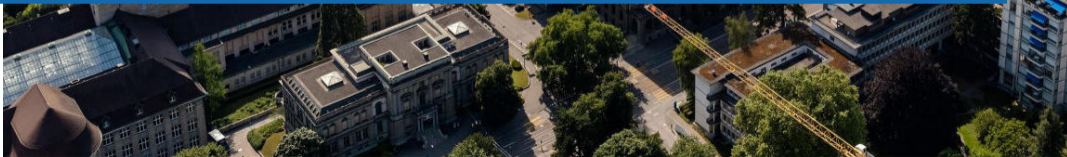


Error bounds for kernel-based linear system identification with unknown hyperparameters

Mingzhou Yin, Roy S. Smith

Sept 26, 2023, ERNSI 2023



System identification as a function learning problem

- Traditionally, SysID is studied as parameter estimation problems with known model structures

$$\min_{\theta} \sum_{k=1}^N (y_k - \hat{y}(k|\theta))^2 \quad (\text{PEM})$$

- Less accessible model structure non-parametric models

$$y_k = f(u_k, u_{k-1}, \dots, u_{-}) + v_k$$

- Restrict to causal linear systems

$$y_k = (g * u)_k + v_k, \quad g_k = 0, \quad k < 0, \quad : \text{discrete convolution}$$

- a function learning problem of g_k

More concretely...

- Consider finite impulse response model

$$G(q) = \sum_{l=0}^{n_g-1} g_l q^{-l}, \quad y_k = \sum_{l=0}^{n_g-1} g_l u_{k-l} + v_k$$

- Formulate data equation with collected input-output data

$$\begin{array}{ccccccc} y_1 & & u_1 & u_0 & \cdots & u_{2-n_g} & g_0 & & v_1 \\ y_2 & & u_2 & u_1 & \cdots & u_{3-n_g} & g_1 & & v_2 \\ \vdots & = & \vdots & \vdots & \ddots & \vdots & \vdots & + & \vdots \\ y_N & & u_N & u_{N-1} & \cdots & u_{N-n_g+1} & g_{n_g-1} & & v_N \\ \hline \mathbf{y} & & \mathbf{u} & & & \mathbf{u} & \mathbf{g} & & \mathbf{v} \end{array}$$

Kernel-based system identification

- One can try least-squares

$$\hat{\mathbf{g}}^{\text{LS}} = \arg\min_{\mathbf{g}} \|\mathbf{y} - \mathbf{g}\|_2^2 = (\mathbf{K}^{-1} \mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^{-1} \mathbf{y}, \text{ with covariance } \text{cov}(\hat{\mathbf{g}}^{\text{LS}}) = (\mathbf{K}^{-1} \mathbf{K} + \lambda \mathbf{I})^{-1}$$

- ... but usually leads to overfitting — too many unknowns
- The regularized version can be more effective

$$\hat{\mathbf{g}} = \arg\min_{\mathbf{g}} \|\mathbf{y} - \mathbf{g}\|_2^2 + \lambda \|\mathbf{g}\|_{\mathbf{K}^{-1}}^2 = (\mathbf{K}^{-1} \mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^{-1} \mathbf{y}$$

- ... by inducing prior assumptions on \mathbf{g}

Threefold interpretation

- **Ridge regression with basis expansion:** $\mathbf{g} = \sum_{i=1}^{n_b} \mathbf{g}_i = G$

$$\min_{\mathbf{g}} \|\mathbf{y} - G\mathbf{g}\|_2^2 + \lambda \|\mathbf{g}\|_2^2, \quad G^T K^{-1} G = I$$

- **Gaussian process:** Gaussian random design of $\mathbf{g} \sim N(\mathbf{0}, K)$

$$\begin{bmatrix} \mathbf{g} \\ \mathbf{y} \end{bmatrix} \sim N \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{y} \end{bmatrix}, \begin{bmatrix} K & K \\ K & K + \lambda I \end{bmatrix} \right)$$

Posterior distribution: $\mathbf{g}|\mathbf{y} \sim N(\hat{\mathbf{g}}, \Sigma)$, $\Sigma = \lambda^{-1} (K^{-1} + \lambda^{-1})^{-1}$

Threefold interpretation

- **Reproducing kernel Hilbert space:** g is sampled from CT function

$$g(t) \in H(k(\cdot, \cdot))$$

$$g(\cdot) = \arg \min_{g(\cdot) \in H} \|\mathbf{y} - \mathbf{g}\|_2^2 + \lambda \int_0^{n_g} g(\cdot)^2_H$$

$$\text{s.t. } \mathbf{g} = [g(0) \dots g(n_g - 1)]^T$$

- Representer theorem: $g(x) = \mathbf{k}_x^T (K + \lambda I)^{-1} \mathbf{y} = \hat{\mathbf{g}}^T \mathbf{g} = \hat{g}$

$$K_{l,l} = k(l, l), \quad \mathbf{k}_x = [k(x, 0) \dots k(x, n_g - 1)]$$

- Induced norm: $\|g(\cdot)\|_H^2 = \hat{\mathbf{g}}^T K^{-1} \hat{\mathbf{g}}$

How to choose K ?

Extensively studied, the common approach:

- Stable kernel structure:

$$K_{i,i}^{\text{DI}}(\cdot) = c^i, \quad K_{i,j}^{\text{DI}}(\cdot) = 0, \quad i \neq j \quad (\text{diagonal})$$

$$K_{i,j}^{\text{TC}}(\cdot) = c^{\max(i,j)} \quad (\text{tuned/correlated})$$

$$K_{i,j}^{\text{SS}}(\cdot) = c^{2 \max(i,j) - \frac{\min(i,j)}{2} - \frac{\max(i,j)}{6}} \quad (\text{stable spline})$$

- Maximum marginal likelihood to estimate hyperparameters :

$$\hat{\cdot} = \operatorname{argmin} -\log p(\mathbf{y}/\mathbf{u}, \cdot)$$

$$\text{Marginal likelihood: } p(\mathbf{y}/\mathbf{u}, \cdot) = \exp \left[-\frac{1}{2} \log \det (\cdot) - \frac{1}{2} \mathbf{y}^T (\cdot)^{-1} \mathbf{y} + \text{const.} \right]$$

- Certainty equivalence: $\hat{\cdot}$

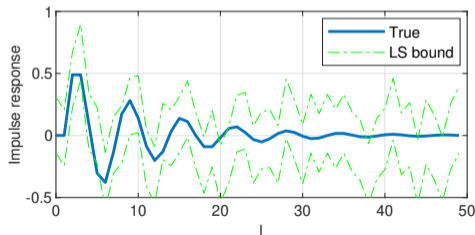
Error bound quantification

- For fixed design of g , LS gives unbiased estimator with minimum variance for i.i.d. Gaussian output noise
- Stochastic high-probability error bounds

$$P \left(\hat{g}_l^{\text{LS}} - g_l \leq \mu \sqrt{\frac{\text{LS}}{l,l}} \right) = 1 - \alpha, \quad F_N(\mu) = 1 - \alpha/2$$

- Still conservative due to overfitting

$$G_2(q) = \frac{0.0616}{q^2 - q + 0.9^2}, \quad \alpha = 0.5$$



Towards better error bounds

- Hope with random design of g : one of the main advantages of GP interpretation
- If $g \sim N(0, K(\cdot))$, stochastic bounds associated with posterior covariance

$$P \left[|\hat{g}_l - g_l| \leq \mu \sqrt{\frac{1}{l}} \right] = 1 - \epsilon$$

- Improvement is guaranteed

$$\sigma^2 = \sigma^2 + \sigma^2 K^{-1} \mathbf{1} \mathbf{1}^T \mathbf{K}^{-1} = \sigma^2 \mathbf{1}^T \mathbf{K}^{-1} \mathbf{1} = \text{LS}$$

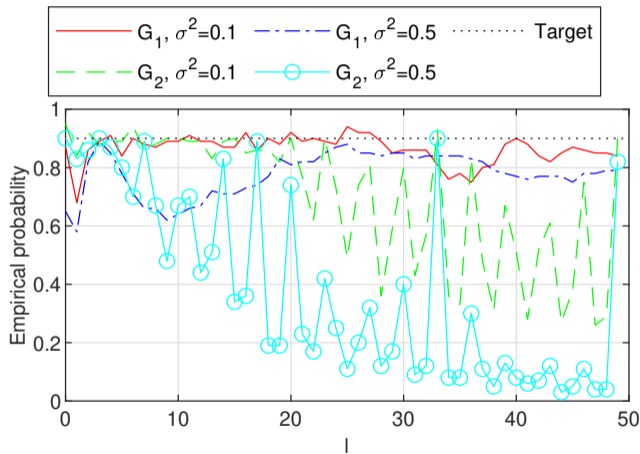
Are the bounds reliable?

$$G_1(q) = \frac{0.4888}{q^2 - 1.8q + 0.9^2}$$

$$G_2(q) = \frac{0.0616}{q^2 - q + 0.9^2}$$

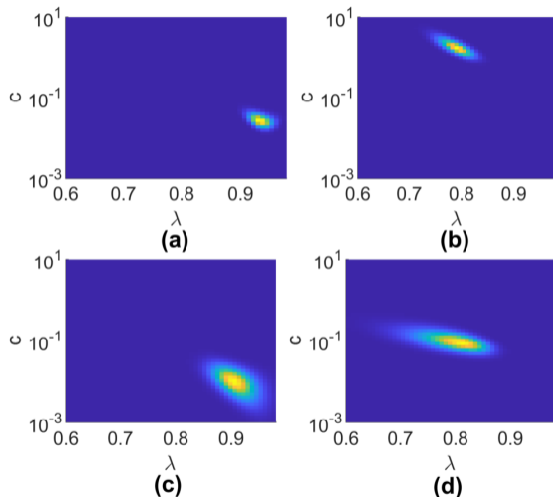
Target prob.: $1 - \quad = 0.9$

Too optimistic for lightly damped systems and low signal-to-noise ratio



What's the reason behind

- Certainty equivalence:
 $\hat{\theta}$
- ... but is it valid?
- Indirect evidence: how localized is the marginal likelihood function?
- $\hat{\theta}$ can be rather inaccurate in (b), (c), (d)



Toward more reliable error bounds

- Be more conservative in estimating
- Instead of using the maximum likelihood point $\hat{\theta}$, establishing a high-probability set for θ_0
- Assume a hyperprior of θ : $p(\theta)$ (uniform distribution if no prior knowledge)

$$\text{Posterior dist. of } \theta : p(\theta | \mathbf{u}, \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{u}, \theta) p(\theta)}{\int_{\mathcal{H}} p(\mathbf{y} | \mathbf{u}, \theta) p(\theta) d\theta}$$

$$\text{High-probability set: } P(\theta_0 \in [\theta_1, \theta_2]) = \frac{\int_{[\theta_1, \theta_2]} p(\mathbf{y} | \mathbf{u}, \theta) p(\theta) d\theta}{\int_{\mathcal{H}} p(\mathbf{y} | \mathbf{u}, \theta) p(\theta) d\theta} = 1 - \alpha$$

- $\alpha =$ Bounds robust to the whole set

Worst-case posterior covariance

- For general kernels, direct (non-convex) optimization for the worst case

$$\sigma_I^2 = \max_{[1, 2]} \sigma_{I,I}(\cdot).$$

- For DI & TC kernels, analytical results available

Lemma: Uniform worst-case covariance

The posterior covariance with true hyperparameters θ_0 can be bounded by

$$\sigma(\theta_0) \leq \frac{1}{4} \sigma^2 + \sigma^2 \frac{1}{2} K^{-1}(\theta_2)^{-1} =: \bar{\sigma}, \quad \sigma_I^2 = \bar{\sigma}_{I,I}$$

where $\bar{\sigma} = 0$ for DI kernels and $\bar{\sigma} = -1/\ln \theta_2 - 1$ for TC kernels.

Select the 'best' high-probability set

- DoF in choosing $\mathbf{u}_1, \mathbf{u}_2$ — only a feasibility problem

$$P(\mathbf{y} \in [\mathbf{u}_1, \mathbf{u}_2]) = \frac{\int_{[\mathbf{u}_1, \mathbf{u}_2]} p(\mathbf{y}/\mathbf{u}_1) p(\mathbf{y}) d\mathbf{y}}{\int_{\mathcal{H}} p(\mathbf{y}/\mathbf{u}_1) p(\mathbf{y}) d\mathbf{y}} \quad (1)$$

- Select $\mathbf{u}_1, \mathbf{u}_2$ that minimizes worst-case covariance — minimax problem

$$\sigma_{\mathbf{y}}^2 = \min_{\mathbf{u}_1, \mathbf{u}_2} \max_{[\mathbf{u}_1, \mathbf{u}_2]} \sigma_{\mathbf{y}}^2(\mathbf{y}) \quad \text{s.t. } (1)$$

- For DI & TC kernels, minimize the sum of uniform worst-case variances

$$\min_{\mathbf{u}_1, \mathbf{u}_2} \sum_{l=0}^{n_g-1} \sigma_l^2 = \text{tr}(\bar{\Sigma}) \quad \min_{\mathbf{u}_1, \mathbf{u}_2} \sum_{l=0}^{n_g-1} \sigma_l^2 \quad \text{s.t. } (1)$$

From worst-case covariance to stochastic bounds

Theorem: Stochastic error bounds

The regularized estimate \hat{g} admits stochastic error bounds:

$$P(|\hat{g}_I(\hat{\cdot}) - g_I| \leq \bar{\mu} \cdot \rho) = (1 - \alpha)(1 - \beta), \quad (1)$$

where $\bar{\mu} = \mu + \frac{\sigma}{S}$, $S = \frac{1}{\rho} \int_{\rho_1}^{\rho_2} \frac{1}{g_I(\rho)} d\rho$, if $\hat{\cdot} \in [\rho_1, \rho_2]$.

Proof sketch: decompose the error

$$|\hat{g}_I(\hat{\cdot}) - g_I| = \underbrace{|\hat{g}_I(\hat{\cdot}) - \hat{g}_I(\rho_0)|}_{\text{error in nominal estimate}} + \underbrace{|\hat{g}_I(\rho_0) - g_I|}_{\text{error with true hyperparam.}}$$

For $|\hat{g}_I(\rho_0) - g_I|$, we have bounded the worst-case covariance for ρ_0

$$|\hat{g}_I(\hat{\cdot}) - \hat{g}_I(\rho_0)| \leq \mu \frac{\rho_0}{g_I(\rho_0)} (1 - \alpha)(1 - \beta)$$

Still conservative. . .

- For $|\hat{g}_l(\hat{\cdot}) - \hat{g}_l(\mathbf{o})|$, no good bound yet. . .
- . . . a conservative bound: $|\hat{g}_l(\hat{\cdot}) - \hat{g}_l(\mathbf{o})| \leq |\hat{g}_l(\hat{\cdot})| + |\hat{g}_l(\mathbf{o})|$
- From RKHS theory,

$$|g(\mathbf{l})| \leq k^p(\mathbf{l}, \mathbf{l})^{\frac{1}{2}} \|g(\cdot)\|_{H^p} \leq \dots \leq \frac{2}{S} \|\mathbf{y}\|_S^2$$

$k^p(x, x)$: posterior kernel with $k^p(i, j) = \frac{1}{i, j}$

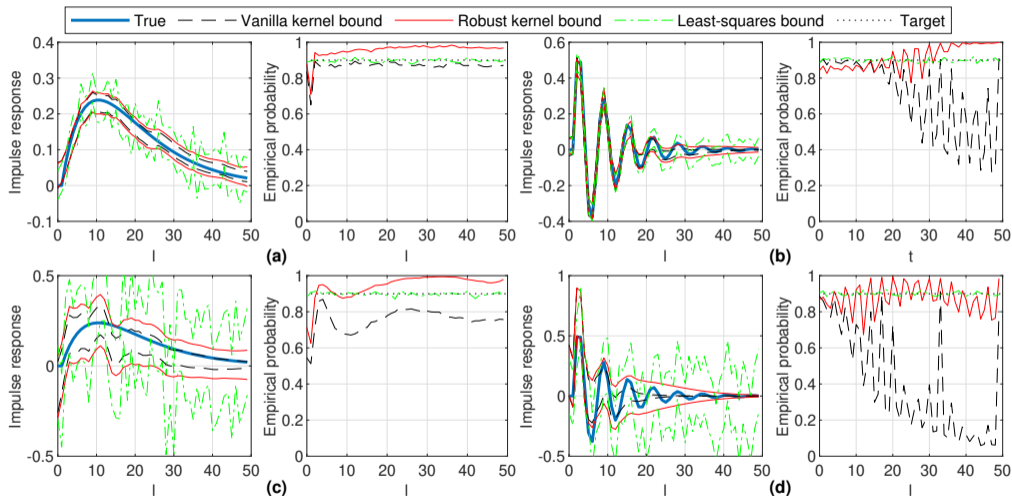
- True for all

$$|\hat{g}_l(\hat{\cdot})| + |\hat{g}_l(\mathbf{o})| \leq 2 \frac{2}{S} \|\mathbf{y}\|_S^2 = \frac{4}{S} \|\mathbf{y}\|_S^2$$

- Better than existing work in ML¹, but still not directly usable in practice

¹Capone, A., Lederer, A., & Hirche, S. (2022). Gaussian process uniform error bounds with unknown hyperparameters for safety-critical applications. In International Conference on Machine Learning (pp. 2609-2624).

Numerical verification



(a) $G_1, \sigma^2 = 0.1$, (b) $G_2, \sigma^2 = 0.1$, (c) $G_1, \sigma^2 = 0.5$, (d) $G_2, \sigma^2 = 0.5$

Error bounds for kernel-based linear system identification with unknown hyperparameters

- Posterior covariance error bounds are not reliable by default
- . . . when hyperparameters are not easy to identify
- Construct high-probability sets for true hyperparameters
- Robust error bounds from worst-case covariance in the set

