# Probabilistic Approximate Least-Squares (APPENDIX)

**Simon Bartels**  **Philipp Hennig**
Max Planck Institute for Intelligent Systems
Tübingen Germany

**An Upper Bound on the Approximation Error**

Let $\boldsymbol{a}, \boldsymbol{b} \in \mathbb{R}^N$ be arbitrary vectors. The Gaussian measure (11) on $\boldsymbol{H}$ implies that the scalar $\hat{\mu} := \boldsymbol{a}^\mathsf{T} \boldsymbol{H} \boldsymbol{b}$ is Gaussian distributed as well, with mean $\boldsymbol{a}^\mathsf{T} \boldsymbol{H}_M \boldsymbol{b}$ and a variance we denote with $\hat{\epsilon}^2$ (derived in the proof below).

**Theorem 1.** *The absolute error* $|\hat{\mu} - \boldsymbol{a}^\mathsf{T} \boldsymbol{H} \boldsymbol{b}|$ *divided by the standard deviation* $\hat{\epsilon}$ *is always less than 1:*

$$\frac{|\boldsymbol{a}^\mathsf{T} \boldsymbol{H}_M \boldsymbol{b} - \boldsymbol{a}^\mathsf{T} \boldsymbol{H} \boldsymbol{b}|}{\hat{\epsilon}} < 1 \qquad (32)$$

*Proof.* If $\boldsymbol{v} \in \mathbb{R}^N$ is a Gaussian random vector $\mathcal{N}(\boldsymbol{v}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ and $\boldsymbol{A} \in \mathbb{R}^{M \times N}$ is of rank $M$ then $\boldsymbol{A}\boldsymbol{v}$ is also Gaussian with $\mathcal{N}(\boldsymbol{A}\boldsymbol{v}; \boldsymbol{A}\boldsymbol{\mu}, \boldsymbol{A}^\mathsf{T} \boldsymbol{\Sigma} \boldsymbol{A})$. We can rewrite $\boldsymbol{a}^\mathsf{T} \boldsymbol{H}_M \boldsymbol{b}$ as $\overrightarrow{\boldsymbol{a}^\mathsf{T} \boldsymbol{H}_M \boldsymbol{b}} = (\boldsymbol{a}^\mathsf{T} \otimes \boldsymbol{b}^\mathsf{T}) \overrightarrow{\boldsymbol{H}_M}$. and it follows that $\hat{\epsilon}$ has the form

$$\hat{\epsilon}^2 = (\boldsymbol{a}^\mathsf{T} \otimes \boldsymbol{b}^\mathsf{T}) (\boldsymbol{W}_M \otimes \boldsymbol{W}_M) (\boldsymbol{a} \otimes \boldsymbol{b}) \qquad (33)$$

To simplify this expression we use that $\boldsymbol{a} \otimes \boldsymbol{b}$ is an $N^2$ dimensional vector and thus $\boldsymbol{a} \otimes \boldsymbol{b} = (\boldsymbol{a} \otimes \boldsymbol{b}) \overrightarrow{1} = \overrightarrow{\boldsymbol{a} 1 \boldsymbol{b}^\mathsf{T}} = \overrightarrow{\boldsymbol{a}\boldsymbol{b}^\mathsf{T}}$. Therefore $\hat{\epsilon}^2$ reduces to

$$\hat{\epsilon}^2 = (\boldsymbol{a}^\mathsf{T} \otimes \boldsymbol{b}^\mathsf{T}) \boldsymbol{\Gamma} (\boldsymbol{W}_M \otimes \boldsymbol{W}_M) \boldsymbol{\Gamma} \overrightarrow{\boldsymbol{a}\boldsymbol{b}^\mathsf{T}} \qquad (34)$$

$$= \frac{1}{2} (\boldsymbol{a}^\mathsf{T} \otimes \boldsymbol{b}^\mathsf{T}) \boldsymbol{\Gamma} (\boldsymbol{W}_M \otimes \boldsymbol{W}_M) \overrightarrow{\boldsymbol{a}\boldsymbol{b}^\mathsf{T} + \boldsymbol{b}\boldsymbol{a}^\mathsf{T}} \qquad (35)$$

$$= \frac{1}{2} (\boldsymbol{a}^\mathsf{T} \otimes \boldsymbol{b}^\mathsf{T}) \boldsymbol{\Gamma} \overrightarrow{\boldsymbol{W}_M \boldsymbol{a}\boldsymbol{b}^\mathsf{T} \boldsymbol{W}_M + \boldsymbol{W}_M \boldsymbol{b}\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M} \qquad (36)$$

$$= \frac{1}{2} (\boldsymbol{a}^\mathsf{T} \otimes \boldsymbol{b}^\mathsf{T}) \overrightarrow{\boldsymbol{W}_M \boldsymbol{a}\boldsymbol{b}^\mathsf{T} \boldsymbol{W}_M + \boldsymbol{W}_M \boldsymbol{b}\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M} \qquad (37)$$

$$= \frac{1}{2} \overrightarrow{\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{a}\boldsymbol{b}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b} + \boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b}\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b}} \qquad (38)$$

$$= \frac{1}{2} \left( \boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{a}\boldsymbol{b}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b} + (\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b})^2 \right) \qquad (39)$$

If we now choose $\boldsymbol{W} = \sqrt{2}\boldsymbol{H}$, and therefore $\boldsymbol{W}_M = \sqrt{2}(\boldsymbol{H} - \boldsymbol{H}_M)$, then

$$\frac{|\boldsymbol{a}^\mathsf{T} (\boldsymbol{H} - \boldsymbol{H}_M) \boldsymbol{b}|}{\sqrt{\hat{\epsilon}^2}} \qquad (40)$$

$$= \frac{|\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b}|}{\sqrt{2} \sqrt{\frac{1}{2} \boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{a} \cdot \boldsymbol{b}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b} + \frac{1}{2} (\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b})^2}} \qquad (41)$$

$$= \frac{|\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b}|}{\sqrt{\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{a} \cdot \boldsymbol{b}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b} + (\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b})^2}} \qquad (42)$$

$$< \frac{|\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b}|}{\sqrt{(\boldsymbol{a}^\mathsf{T} \boldsymbol{W}_M \boldsymbol{b})^2}} = 1 \qquad \text{(as } \boldsymbol{W}_M \text{ is s.p.d.)} \qquad (43)$$
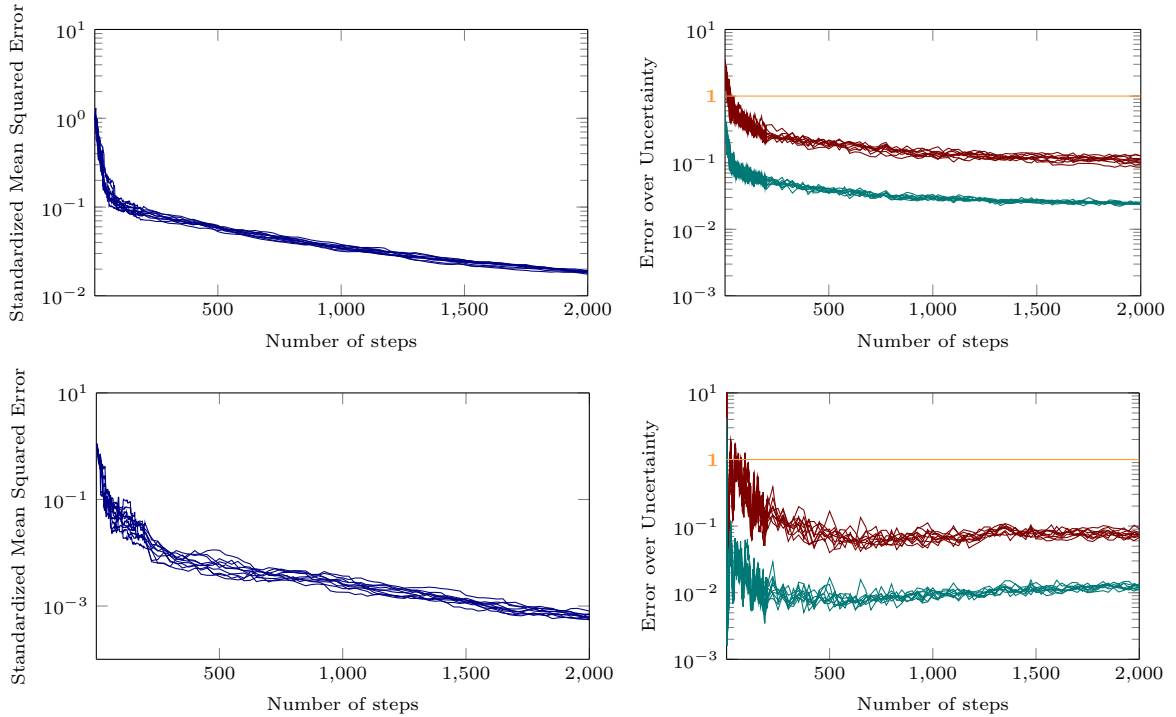
$\square$

Figure 1: Ten Farthest Point Clustering initializations of the probabilistic subset of data approximation on the PUMADYN (**top row**) and CPU (**bottom row**) data sets, using the ARD Squared Exponential kernel. **Left:** standardized mean squared error for Subset of Data. **Right:** ratio between absolute error and uncertainty. The upper lines are the maximum, the lower lines the average over all test inputs. The horizontal line shows the theoretical bound at 1 that would be guaranteed if $\boldsymbol{W} = \sqrt{2}\boldsymbol{H}$ where estimated exactly.
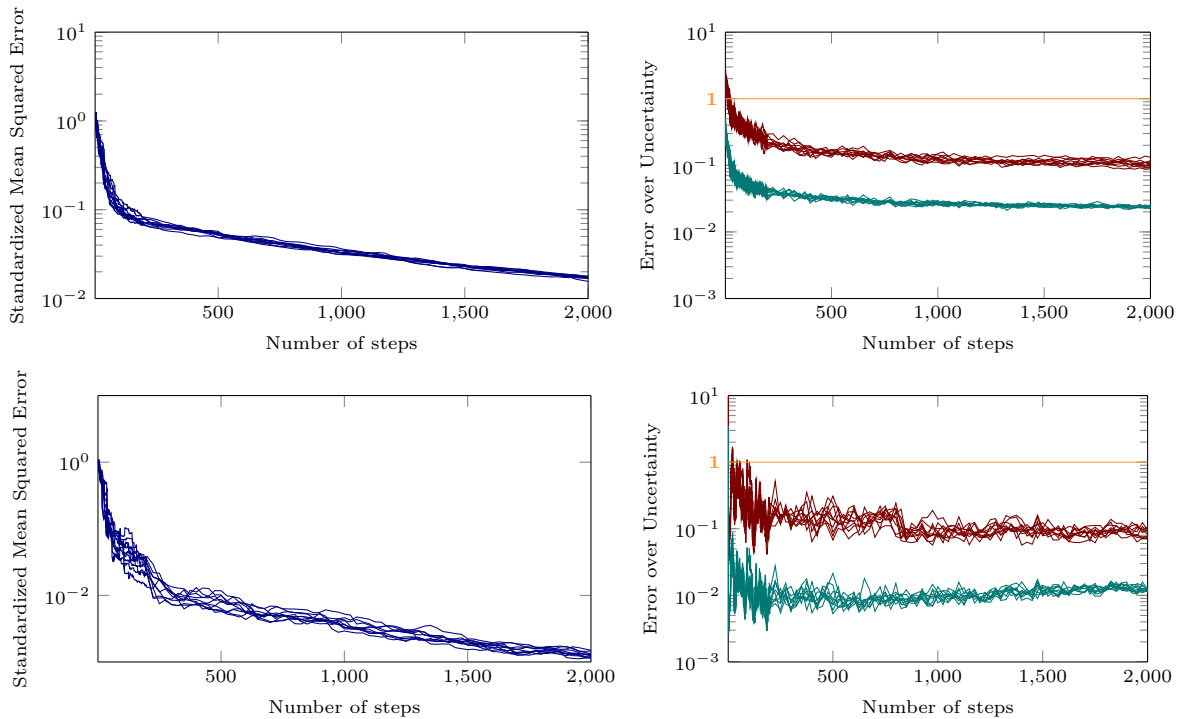


Figure 2: Same setup as Figure 1, but using the ARD Matérn $^5/_2$ kernel.