

Geometric Rates of Convergence for Kernel-based Sampling Algorithms – Supplementary material

Rajiv Khanna¹

Liam Hodgkinson^{1,2}

Michael W. Mahoney^{1,2}

¹Department of Statistics, UC Berkeley

²International Computer Science Institute (ICSI), Berkeley

A APPENDIX

A.1 PROOF OF THEOREM 2

We begin by first proving Theorem 2, since the additional assumption of realizability makes it an easier read. For further ease of exposition, instead of directly working with $g(\cdot)$, we translate the function to remove any constants not dependent on the variable. We write,

$$l(\mathbf{S}) := \|\mu_\pi\|_k^2 - g(\mathbf{S}) = \mathbf{z}^\top \mathbf{K}^{-1} \mathbf{z}.$$

Some auxiliary Lemmas are proved later in this section. We use $Z(\mathbf{S}_j) := \sum_j w_j \phi(\mathbf{x}_j)$. Further, note that the Assumption 2, when applied for $h(\cdot)$, ensures that for any iterates considered in this proof we have that

$$\begin{aligned} & -\frac{m_\omega}{2} \|Z(\mathbf{S}_i) - Z(\mathbf{S}_j)\|_k^2 \\ & \geq l(\mathbf{S}_i) - l(\mathbf{S}_j) - \langle \nabla l(\mathbf{S}_i), Z(\mathbf{S}_i) - Z(\mathbf{S}_j) \rangle_k \\ & \geq -\frac{M_\Omega}{2} \|Z(\mathbf{S}_i) - Z(\mathbf{S}_j)\|_k^2. \end{aligned}$$

Proof. Say $(i - 1)$ steps of the Algorithm 1 have been performed to select the set \mathbf{S} . Let $\mathbf{w} \in \mathbb{R}^{(i-1)}$ be the corresponding weight vector. Let $h(\mathbf{S}, \mathbf{u}) := \|\mu_\pi\|_k^2 - \|\mu_\pi - \sum_j u_j \phi(\mathbf{x}_j)\|_k^2$, so that $l(\mathbf{S}) = \min_{\mathbf{u}} h(\mathbf{S}, \mathbf{u})$ (as per Lemma 1). Set weight vector $\mathbf{u} \in \mathbb{R}^i$ as follows. For $j \in [0, i - 1]$, $u_i = w_i$. Set $u_i = \alpha$, where α is an arbitrary scalar.

From weight optimality proved in Lemma 1,

$$l(\mathbf{S} \cup \{\mathbf{x}_i\}) - l(\mathbf{S}) \geq h(\mathbf{S} \cup \{\mathbf{x}_i\}, \mathbf{u}) - l(\mathbf{S}),$$

for an arbitrary $\alpha \in \mathbb{R}$. From Assumption 2 (smoothness),

$$l(\mathbf{S} \cup \{\mathbf{x}_i\}) - l(\mathbf{S}) \geq \alpha \langle \nabla l(\mathbf{S}), \phi(\mathbf{x}_i) \rangle_k - \alpha^2 \frac{M_\Omega}{2}.$$

Let $\gamma_{\mathbf{S}}$ be the optimum value of the solution of the inner LMO problem. Since \mathbf{x}_i is the optimizing atom,

$$l(\mathbf{S} \cup \{\mathbf{x}_i\}) - l(\mathbf{S}) \geq \alpha \gamma_{\mathbf{S}} - \alpha^2 \frac{M_\Omega}{2}.$$

Let \mathbf{S}_\perp^* be the set obtained by orthogonalizing \mathbf{S}_r^* with respect to \mathbf{S} using the Gram-Schmidt procedure. Putting in $\alpha = \frac{\gamma_{\mathbf{S}}}{M_\Omega}$, we get,

$$\begin{aligned} l(\mathbf{S} \cup \{\mathbf{x}_i\}) - l(\mathbf{S}) & \geq \frac{1}{2M_\Omega} \gamma_{\mathbf{S}} & (1) \\ & \geq \frac{1}{2rM_\Omega} \sum_{\mathbf{x}_j \in \mathbf{S}_\perp^*} \langle \phi(\mathbf{x}_j), \nabla l(\mathbf{S}) \rangle_k^2 \\ & \geq \frac{m_\omega}{rM_\Omega} (l(\mathbf{S} \cup \mathbf{S}_\perp^*) - l(\mathbf{S})) & (2) \\ & \geq \frac{m_\omega}{rM_\Omega} (l(\mathbf{S}_r^*) - l(\mathbf{S})) \\ & = \frac{m_\omega}{rM_\Omega} (\|\mu_\pi\|_k^2 - l(\mathbf{S})). \end{aligned}$$

The second inequality is true because $\gamma_{\mathbf{S}} = \langle \nabla l(\mathbf{S}), \mathbf{x}_i \rangle_k$ is the optimum value of the inner LMO problem in the i^{th} iteration. The third inequality follows from Lemma 2. The fourth inequality is true because of monotonicity of $l(\cdot)$, and the final equality is true because of Assumption 1 (realizability).

Let $C := \frac{m_\omega}{rM_\Omega}$. We have $l(\mathbf{S} \cup \{\mathbf{x}_i\}) - l(\mathbf{S}) = g(\mathbf{S}) - g(\mathbf{S} \cup \{\mathbf{x}_i\}) \geq Cg(\mathbf{S}) \implies g(\mathbf{S} \cup \{\mathbf{x}_i\}) \leq (1 - C)g(\mathbf{S})$. The result now follows. \square

A.2 PROOF OF THEOREM 1

Proof. We proceed as in the proof of Theorem 2, but by replacing \mathbf{S}_r^* with \mathbf{T}_r . From (2),

$$l(\mathbf{S} \cup \{\mathbf{x}_i\}) - l(\mathbf{S}) \geq \frac{m_\omega}{rM_\Omega} (l(\mathbf{T}_r) - l(\mathbf{S})).$$

Adding and subtracting $l(\mathbf{T}_r)$ on the LHS and rearranging,

$$l(\mathbb{T}_r) - l(\mathbb{S} \cup \{\mathbf{x}_i\}) \leq \left(1 - \frac{m_\omega}{rM_\Omega}\right)(l(\mathbb{T}_r) - l(\mathbb{S})).$$

Thus after k iterations,

$$l(\mathbb{T}_r) - l(\mathbb{S}_k) \leq \left(1 - \frac{m_\omega}{rM_\Omega}\right)^k (l(\mathbb{T}_r) - l(\emptyset)).$$

Rearranging,

$$\begin{aligned} l(\mathbb{S}_k) &\geq \left(1 - \left(1 - \frac{m_\omega}{rM_\Omega}\right)^k\right) l(\mathbb{T}_r) \\ &\geq \left(1 - \exp\left(-\frac{km_\omega}{rM_\Omega}\right)\right) l(\mathbb{T}_r). \end{aligned}$$

With $k = \left(r \frac{M_\Omega}{m_\omega} \log \frac{1}{\epsilon}\right)$, we get,

$$l(\mathbb{S}_k) \geq (1 - \epsilon)l(\mathbb{T}_r).$$

The result now follows. \square

A.3 AUXILIARY LEMMAS

The following Lemma proves that the weights w_i in $g(\mathbb{S})$ obtained using the posterior inference are an optimum choice that minimize the distance to μ_π in the RKHS over any set of weights [Khanna et al., 2019].

Lemma 1. *The residual $\mu_\pi - \sum_j w_j \phi(\mathbf{x}_j)$ is orthogonal to $\mathbf{x}_i \in \mathbb{S} \forall i$. In other words, for any set of samples \mathbb{S} , $g(\mathbb{S}) = \min_{\mathbf{u}} \|\mu_\pi - \sum_i u_i \phi(\mathbf{x}_i)\|_k$.*

Proof. Recall that $w_i = \sum_j [\mathbf{K}^{-1}]_{ij} \mathbf{z}_j$, and $\mathbf{z}_i = \int k(\mathbf{x}, \mathbf{x}_i) d\pi(\mathbf{x})$. For an arbitrary index i ,

$$\begin{aligned} &\langle \mu_\pi - \sum_j w_j \phi(\mathbf{x}_j), \phi(\mathbf{x}_i) \rangle_k \\ &= \int k(\mathbf{x}, \mathbf{x}_i) d\pi(\mathbf{x}) - \langle \sum_j w_j \phi(\mathbf{x}_j), \phi(\mathbf{x}_i) \rangle_k \\ &= \mathbf{z}_i - \langle \sum_j w_j \phi(\mathbf{x}_j), \phi(\mathbf{x}_i) \rangle_k \\ &= \mathbf{z}_i - \sum_j w_j k(\mathbf{x}_j, \mathbf{x}_i) \\ &= \mathbf{z}_i - \sum_j \sum_t [\mathbf{K}^{-1}]_{tj} \mathbf{z}_t k(\mathbf{x}_j, \mathbf{x}_i) \\ &= \mathbf{z}_i - \sum_t \mathbf{z}_t \sum_j \mathbf{K}_{ji} [\mathbf{K}^{-1}]_{tj} \\ &= \mathbf{z}_i - \mathbf{z}_i, \end{aligned}$$

where the last equality follows by noting that $\sum_j \mathbf{K}_{ji} [\mathbf{K}^{-1}]_{tj}$ is inner product of row i of \mathbf{K} and row t of \mathbf{K}^{-1} , which is 1 if $t = i$ and 0 otherwise. This completes the proof. \square

Lemma 2. *For any set of chosen samples $\mathbb{S}_1, \mathbb{S}_2$, let \mathcal{P} be the operator of projection onto $\text{span}(\mathbb{S}_1 \cup \mathbb{S}_2)$. Then, $l(\mathbb{S}_1 \cup \mathbb{S}_2) - l(\mathbb{S}_1) \leq \frac{\mathcal{P}(\nabla l(\mathbb{S}_1))}{2m_\omega}$.*

Proof. Observe that

$$\begin{aligned} 0 &\leq l(\mathbb{S}_1 \cup \mathbb{S}_2) - l(\mathbb{S}_1) \\ &\leq \langle \nabla l(\mathbb{S}_1), Z(\mathbb{S}_1 \cup \mathbb{S}_2) - Z(\mathbb{S}_1) \rangle_k \\ &\quad - \frac{m_\omega}{2} \|Z(\mathbb{S}_1 \cup \mathbb{S}_2) - Z(\mathbb{S}_1)\|_k^2 \\ &\leq \arg \max_{X \in \text{span}(\mathbb{S}_1 \cup \mathbb{S}_2)} \langle \nabla l(\mathbb{S}_1), X - Z(\mathbb{S}_1) \rangle_k - \frac{m_\omega}{2} \|X - Z(\mathbb{S}_1)\|_k^2 \\ &= \arg \max_X \langle \mathcal{P}(\nabla l(\mathbb{S}_1)), X - Z(\mathbb{S}_1) \rangle_k - \frac{m_\omega}{2} \|X - Z(\mathbb{S}_1)\|_k^2. \end{aligned}$$

Solving the argmax problem on the RHS for X , we get the required result. \square

A.4 PROOF OF THEOREM 3

We next present some notation and few lemmas that lead up to the main result of this section (Theorem 3). The domain of candidate atoms \mathcal{X} is split into $\{\mathcal{X}_j, j \in [s]\}$ over s machines, where machine j runs WKH on \mathcal{X}_j . Let \mathbb{G}_j be the k -sized solution returned by running Algorithm 1 on \mathcal{X}_j , i.e., $\mathbb{G}_j = \text{WKH}(\mathcal{X}_j, k)$. Note that each \mathcal{X}_j induces a partition onto the optimal r -sized solution \mathbb{S}_r^* as follows ($r = 1$ for this theorem):

$$\begin{aligned} \mathbb{T}_j &:= \{x \in \mathbb{S}_1^* : x \notin \text{WKH}(\mathcal{X}_j \cup x, k)\}, \\ \mathbb{T}_j^c &:= \{x \in \mathbb{S}_1^* : x \in \text{WKH}(\mathcal{X}_j \cup x, k)\}. \end{aligned}$$

In other words, $\mathbb{T}_j = \mathbb{S}_1^*$ if the j^{th} machine running WKH on $\mathcal{X}_j \cup \mathbb{S}_1^*$ will not select it as among its output, and it is empty otherwise, since \mathbb{S}_1^* is a singleton. We re-use the definition of $l(\cdot)$ used in Appendix A.1.

Before moving to the proof of the main theorem, we prove two prerequisites. Recall \mathbb{G}_j is the set of iterates selected by machine j . In this mini-result, we lower bound the expected improvement in the loss at the aggregator machine.

Lemma 3. *For the aggregator machine that runs WKH over $\cup_j \mathbb{G}_j$ (step 6 of Algorithm 2), we have, at selection of next sample point \mathbf{x}_i after having selected \mathbb{S} , \exists machine j such that*

$$\mathbb{E}[l(\mathbb{S} \cup \{\mathbf{x}_i\}) - l(\mathbb{S})] \geq \frac{m_\omega}{M_\Omega} \mathbb{E}(l(\mathbb{T}_j^c) - l(\mathbb{S})).$$

Proof. The expectation is over the random split of \mathcal{X} into \mathcal{X}_j for $j \in [s]$. We denote \mathbb{T}_j^c to be the complement of \mathbb{T}_j . Then, we have that

$$\begin{aligned}
& \mathbb{E}[l(\mathbb{S} \cup \{\mathbf{x}_i\}) - l(\mathbb{S})] \\
& \geq \mathbb{E}\left[\frac{1}{2M_\Omega} \gamma_{\mathbb{S}}\right] \\
& \geq \frac{1}{2M_\Omega} \sum_{\mathbf{x} \in \mathbb{S}_1^*} \mathbb{P}(\mathbf{x} \in \cup_j \mathbb{G}_j) \mathbb{E}\langle \phi(\mathbf{x}), \nabla l(\mathbb{S}) \rangle_k^2 \\
& = \frac{1}{2sM_\Omega} \sum_{\mathbf{x} \in \mathbb{S}_1^*} \left[\sum_{b=1}^s \mathbb{P}(\mathbf{x} \in \mathbb{T}_b^c) \right] \mathbb{E}\langle \phi(\mathbf{x}), \nabla l(\mathbb{S}) \rangle_k^2 \\
& = \frac{1}{2sM_\Omega} \sum_{b=1}^s \sum_{\mathbf{x} \in \mathbb{T}_b^c} \mathbb{E}\langle \phi(\mathbf{x}), \nabla l(\mathbb{S}) \rangle_k^2 \\
& \geq \frac{m_\omega}{sM_\Omega} \sum_{b=1}^s \mathbb{E}(l(\mathbb{S} \cup \mathbb{T}_b^c) - l(\mathbb{S})) \\
& \geq \frac{m_\omega}{sM_\Omega} \sum_{b=1}^s \mathbb{E}(l(\mathbb{T}_b^c) - l(\mathbb{S})) \\
& \geq \frac{m_\omega}{M_\Omega} \min_{b \in [s]} \mathbb{E}(l(\mathbb{T}_b^c) - l(\mathbb{S})).
\end{aligned}$$

The equality in step 3 above is because of Lemma 5. \square

In the following lemma, we lower bound the greedy improvement in the loss on each machine.

Lemma 4. *For any individual worker machine j running local WKH, if \mathbb{S} is the set of $(i-1)$ iterates already chosen, then at the selection of next sample point \mathbf{x}_i , $l(\mathbb{S} \cup \{\mathbf{x}_i\}) \geq (l(\mathbb{T}_j) - l(\mathbb{S}))$.*

Proof. We proceed as in proof of Theorem 2 in Appendix A.1. From (1), we have,

$$\begin{aligned}
l(\mathbb{S} \cup \{\mathbf{x}\}) - l(\mathbb{S}) & \geq \frac{1}{2M_\Omega} \gamma_{\mathbb{S}} \\
& \geq \frac{1}{2M_\Omega} \sum_{\mathbf{x}_j \in \mathbb{T}_j} \langle \phi(\mathbf{x}_j), \nabla l(\mathbb{S}) \rangle_k^2 \\
& \geq \frac{m_\omega}{M_\Omega} (l(\mathbb{S} \cup \mathbb{T}_j) - l(\mathbb{S})) \\
& \geq \frac{m_\omega}{M_\Omega} (l(\mathbb{T}_j) - l(\mathbb{S})).
\end{aligned}$$

\square

We are now ready to prove Theorem 3.

Proof of Theorem 3. If, for a random split of \mathcal{X} , for any $j \in [s]$, $\mathbb{T}_j = \mathbb{S}_1^*$, then the given rate follows from Lemma 4, after following the straightforward steps covered in proof of Theorem 2 for proving the rate from the given condition

on $l(\cdot)$. On the other hand, if none of $j \in [s]$, $\mathbb{T}_j = \mathbb{S}_1^*$, then $\forall j \in [s], \mathbb{T}_j = \emptyset \implies \mathbb{T}_j^c = \mathbb{S}_1^*$. In this case, the given rate follows from Lemma 3. \square

Finally, here is the statement and proof of an auxiliary lemma that was used above.

Lemma 5. *For any $x \in \mathcal{X}$, $\mathbb{P}(x \in \cup_j \mathbb{G}_j) = \frac{1}{s} \sum_j \mathbb{P}(x \in \mathbb{T}_j^c)$.*

Proof. We have

$$\begin{aligned}
& \mathbb{P}(x \in \cup_j \mathbb{G}_j) \\
& = \sum_j \mathbb{P}(x \in \mathcal{X}_j \cap x \in \text{WKH}(\mathcal{X}_j, k)) \\
& = \sum_j \mathbb{P}(x \in \mathcal{X}_j) \mathbb{P}(x \in \text{WKH}(\mathcal{X}_j, k) | x \in \mathcal{X}_j) \\
& = \sum_j \mathbb{P}(x \in \mathcal{X}_j) \mathbb{P}(x \in \mathbb{T}_j^c) \\
& = \frac{1}{s} \mathbb{P}(x \in \mathbb{T}_j^c).
\end{aligned}$$

\square