

---

## Appendix: Kernel Spectral Algorithm for HMM

---

**Le Song**

**Byron Boots**

School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA

LESONG@CS.CMU.EDU

BEB@CS.CMU.EDU

**Sajid M. Siddiqi**

Google, Pittsburgh, PA 15213, USA

SIDDIQI@GOOGLE.COM

**Geoffrey Gordon**

School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA

GGORDON@CS.CMU.EDU

**Alex Smola**

Yahoo! Research, Santa Clara, CA 95051, USA

ALEX@SMOLA.ORG

### 1. Sample Complexity

Let  $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$  and  $l(h, h') = \langle \phi(h), \phi(h') \rangle_{\mathcal{G}}$  be the kernels for the observation variables  $X_t$  and hidden states  $H_t$  respectively. We assume that the feature map is bounded, *i.e.*  $\|\varphi(x)\|_{\mathcal{F}} \leq 1$  and  $\|\phi(h)\|_{\mathcal{G}} \leq 1$ . In this case, the norms of the conditional embeddings of  $H_t$  are bounded by 1, *i.e.* by convexity

$$\begin{aligned} \|\mu_{H_t|x_{t-1:1}}\|_{\mathcal{G}} &= \|\mathbb{E}_{H_t|x_{t-1:1}}[\phi(H_t)]\|_{\mathcal{G}} \\ &\leq \mathbb{E}_{H_t|x_{t-1:1}}[\|\phi(H_t)\|_{\mathcal{G}}] \leq 1 \end{aligned} \quad (1)$$

We will use  $\|\cdot\|_2$  to denote the spectral norm of an operator. We assume that the spectral norms of  $\mathcal{O}$  and  $\mathcal{T}$  are finite, *i.e.*  $\|\mathcal{O}\|_2 \leq C_{\mathcal{O}}$  and  $\|\mathcal{T}\|_2 \leq C_{\mathcal{T}}$ . Let  $\hat{\mathcal{U}}$  be the column concatenation of the top  $N$  left singular vectors of  $\hat{\mathcal{C}}_{2,1}$ . We define

$$\tilde{\beta}_1 := \hat{\mathcal{U}}^\top \mu_1 \quad (2)$$

$$\tilde{\beta}_\infty := \mathcal{C}_{2,1}(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger \quad (3)$$

$$\tilde{\mathcal{B}}_x := \mathbb{P}(x) \left( \hat{\mathcal{U}}^\top (\mathcal{C}_{3,1|2} \varphi(x)) \right) \left( \hat{\mathcal{U}}^\top \mathcal{C}_{2,1} \right)^\dagger \quad (4)$$

We further define

$$\delta_1 := \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\tilde{\beta}_1 - \hat{\beta}_1)\|_{\mathcal{G}} \quad (5)$$

$$\delta_\infty := \|\tilde{\beta}_\infty - \hat{\beta}_\infty\|_2 \quad (6)$$

$$\Delta := \max_{x_t} \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\mathcal{B}}_{x_t} - \tilde{\mathcal{B}}_{x_t})(\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \quad (7)$$

$$\gamma := \max_{x_t} \|\mathcal{A}_{x_t}\|_2 \quad (\text{we assume } \gamma \leq 1). \quad (8)$$

**Lemma 1**  $\|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{G}} \leq (1 + \Delta)^t \delta_1 + (1 + \Delta)^t - 1$

**Proof** By induction on  $t$ .  $\|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\beta}_1 - \tilde{\beta}_1)\|_{\mathcal{G}} = \delta_1 = (1 + \Delta)^0 \delta_1 + (1 + \Delta)^0 - 1$  when  $t = 0$ . For the induction step define  $\hat{\beta}_t := \hat{\mathcal{B}}_{x_{t-1:1}} \hat{\beta}_1$  and  $\tilde{\beta}_t := \tilde{\mathcal{B}}_{x_{t-1:1}} \tilde{\beta}_1$ . Assume for  $t > 1$ :

$$\begin{aligned} &\|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\beta}_t - \tilde{\beta}_t)\|_{\mathcal{G}} \\ &\leq (1 + \Delta)^{t-1} \delta_1 + (1 + \Delta)^{t-1} - 1 \end{aligned} \quad (9)$$

Then, we have

$$\begin{aligned} &\|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{G}} \\ &\leq \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\mathcal{B}}_{x_t} - \tilde{\mathcal{B}}_{x_t})(\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \tilde{\beta}_t\|_{\mathcal{G}} \end{aligned} \quad (10)$$

$$+ \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\mathcal{B}}_{x_t} - \tilde{\mathcal{B}}_{x_t})(\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\beta}_t - \tilde{\beta}_t)\|_{\mathcal{G}} \quad (11)$$

$$+ \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \hat{\mathcal{B}}_{x_t} (\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\beta}_t - \tilde{\beta}_t)\|_{\mathcal{G}} \quad (12)$$

For (10), we have

$$\begin{aligned} &\|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\mathcal{B}}_{x_t} - \tilde{\mathcal{B}}_{x_t})(\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \tilde{\beta}_t\|_{\mathcal{G}} \\ &\leq \Delta \|\mu_{H_t|x_{t-1:1}}\|_{\mathcal{G}} \leq \Delta \end{aligned} \quad (13)$$

For (11), we have

$$\begin{aligned} &\|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\mathcal{B}}_{x_t} - \tilde{\mathcal{B}}_{x_t})(\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1}(\hat{\beta}_t - \tilde{\beta}_t)\|_{\mathcal{G}} \\ &\leq \Delta ((1 + \Delta)^{t-1} \delta_1 + (1 + \Delta)^{t-1} - 1) \end{aligned} \quad (14)$$

For (12), we have  $\blacksquare$

$$\begin{aligned} & \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \hat{\mathcal{B}}_{x_t} (\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\beta}_t - \tilde{\beta}_t)\|_{\mathcal{G}} \\ & \leq \|\mathcal{A}_{x_t}\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\beta}_t - \tilde{\beta}_t)\|_{\mathcal{G}} \\ & \leq (1 + \Delta)^{t-1} \delta_1 + (1 + \Delta)^{t-1} - 1 \end{aligned} \quad (15)$$

Summing the above three bounds, we have  $\blacksquare$

$$\begin{aligned} & (1 + \Delta)^t \delta_1 + (1 + \Delta)^t - (1 + \Delta) + \Delta \\ & = (1 + \Delta)^t \delta_1 + (1 + \Delta)^t - 1 \end{aligned} \quad (16)$$

Define:

$$\epsilon_1 := \|\mu_1 - \hat{\mu}_1\|_{\mathcal{F}} \quad (23)$$

$$\epsilon_2 := \|\mathcal{C}_{2,1} - \hat{\mathcal{C}}_{2,1}\|_{\mathcal{F} \otimes \mathcal{F}} \quad (24)$$

$$\epsilon_3 := \max_x \|\mu_{3,1|x} - \hat{\mu}_{3,1|x}\|_{\mathcal{F} \otimes \mathcal{F}} \quad (25)$$

**Lemma 2**  $\|\mu_{X_{t+1}|x_{t:1}} - \hat{\mu}_{X_{t+1}|x_{t:1}}\|_{\mathcal{F}} \leq C_{\mathcal{O}} \delta_{\infty} + (C_{\mathcal{O}} C_{\mathcal{T}} + C_{\mathcal{O}} \delta_{\infty}) ((1 + \Delta)^t \delta_1 + (1 + \Delta)^t - 1)$

**Proof** Using triangle inequality

$$\begin{aligned} & \|\mu_{X_{t+1}|x_{t:1}} - \hat{\mu}_{X_{t+1}|x_{t:1}}\|_{\mathcal{F}} \\ & = \|\tilde{\beta}_{\infty} \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1 - \hat{\beta}_{\infty} \hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1\|_{\mathcal{F}} \\ & \leq \|(\tilde{\beta}_{\infty} - \hat{\beta}_{\infty})(\hat{\mathcal{U}}^\top \mathcal{O})(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1\|_{\mathcal{F}} \end{aligned} \quad (17)$$

$$+ \|(\tilde{\beta}_{\infty} - \hat{\beta}_{\infty})(\hat{\mathcal{U}}^\top \mathcal{O})(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{F}} \quad (18)$$

$$+ \|\tilde{\beta}_{\infty} (\hat{\mathcal{U}}^\top \mathcal{O})(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{F}} \quad (19)$$

For (17), we have

$$\begin{aligned} & \|(\tilde{\beta}_{\infty} - \hat{\beta}_{\infty})(\hat{\mathcal{U}}^\top \mathcal{O})(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1\|_{\mathcal{F}} \\ & \leq \|(\tilde{\beta}_{\infty} - \hat{\beta}_{\infty})(\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1\|_{\mathcal{G}} \\ & \leq \|\tilde{\beta}_{\infty} - \hat{\beta}_{\infty}\|_2 \|\mathcal{O}\|_2 \|\mu_{H_{t+1}|x_{t:1}}\|_{\mathcal{G}} \\ & \leq C_{\mathcal{O}} \|\tilde{\beta}_{\infty} - \hat{\beta}_{\infty}\|_2 \leq C_{\mathcal{O}} \delta_{\infty} \end{aligned} \quad (20)$$

For (18), we have

$$\begin{aligned} & \|(\tilde{\beta}_{\infty} - \hat{\beta}_{\infty})(\hat{\mathcal{U}}^\top \mathcal{O})(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{F}} \\ & \leq \|(\tilde{\beta}_{\infty} - \hat{\beta}_{\infty})(\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{G}} \\ & \leq \|\tilde{\beta}_{\infty} - \hat{\beta}_{\infty}\|_2 \|\mathcal{O}\|_2 ((1 + \Delta)^t \delta_1 + (1 + \Delta)^t - \gamma) \\ & \leq C_{\mathcal{O}} \delta_{\infty} ((1 + \Delta)^t \delta_1 + (1 + \Delta)^t - 1) \end{aligned} \quad (21)$$

For (19), we have  $\blacksquare$

$$\begin{aligned} & \|\tilde{\beta}_{\infty} (\hat{\mathcal{U}}^\top \mathcal{O})(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{F}} \\ & \leq \|\tilde{\beta}_{\infty} (\hat{\mathcal{U}}^\top \mathcal{O})\|_2 \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_{t:1}} \hat{\beta}_1 - \tilde{\mathcal{B}}_{x_{t:1}} \tilde{\beta}_1)\|_{\mathcal{G}} \\ & \leq \|\mathcal{O}\|_2 ((1 + \Delta)^t \delta_1 + (1 + \Delta)^t - 1) \\ & \leq C_{\mathcal{O}} C_{\mathcal{T}} ((1 + \Delta)^t \delta_1 + (1 + \Delta)^t - 1) \end{aligned} \quad (22)$$

**Lemma 3** Suppose  $\epsilon_2 \leq \varepsilon \sigma_N(\mathcal{C}_{2,1})$  for some  $\varepsilon \leq 1/2$ . Let  $\varepsilon_0 = \epsilon_2^2 / ((1 - \varepsilon) \sigma_N(\mathcal{C}_{2,1}))^2$ . Then:

$$1. \varepsilon_0 < 1,$$

$$2. \sigma_n(\hat{\mathcal{U}}^\top \hat{\mathcal{C}}_{2,1}) \geq (1 - \varepsilon) \sigma_N(\mathcal{C}_{2,1}),$$

$$3. \sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1}) \geq \sqrt{1 - \varepsilon_0} \sigma_N(\mathcal{C}_{2,1}),$$

$$4. \sigma_N(\hat{\mathcal{U}}^\top \mathcal{O}) \geq \sqrt{1 - \varepsilon_0} \sigma_N(\mathcal{O}).$$

**Proof** This lemma can be proved by an extension of Lemma 9 in ?  $\blacksquare$

**Lemma 4**  $\delta_1 := \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\tilde{\beta}_1 - \hat{\beta}_1)\|_{\mathcal{G}} \leq \frac{\epsilon_1}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{O})}$

**Proof**

$$\begin{aligned} & \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\tilde{\beta}_1 - \hat{\beta}_1)\|_{\mathcal{G}} \\ & = \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{U}}^\top \mu_1 - \hat{\mathcal{U}}^\top \hat{\mu}_1)\|_{\mathcal{G}} \\ & \leq \|(\hat{\mathcal{U}}^\top \mathcal{O})^{-1} \hat{\mathcal{U}}^\top\|_2 \|\mu_1 - \hat{\mu}_1\|_{\mathcal{F}} \end{aligned} \quad (26)$$

$$\leq \frac{\epsilon_1}{\sigma_n(\hat{\mathcal{U}}^\top \mathcal{O})} \quad (27)$$

**Lemma 5**  $\delta_{\infty} := \|\tilde{\beta}_{\infty} - \beta_{\infty}\|_2 \leq \frac{1 + \sqrt{5}}{2} \frac{\epsilon_2}{\min\{\sigma_N(\hat{\mathcal{C}}_{2,1}), \sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})\}^2} + \frac{\epsilon_2}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})}$

**Proof**

$$\begin{aligned}
 & \left\| \tilde{\beta}_\infty - \beta_\infty \right\|_2 \\
 &= \left\| \mathcal{C}_{2,1} (\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger - \hat{\mathcal{C}}_{2,1} (\hat{\mathcal{U}}^\top \hat{\mathcal{C}}_{2,1})^\dagger \right\|_2 \\
 &\leq \left\| \hat{\mathcal{C}}_{2,1} ((\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger - (\hat{\mathcal{U}}^\top \hat{\mathcal{C}}_{2,1})^\dagger) \right\|_2 \\
 &+ \left\| (\mathcal{C}_{2,1} - \hat{\mathcal{C}}_{2,1}) (\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger \right\|_2 \\
 &\leq \left\| \hat{\mathcal{C}}_{2,1} \right\|_2 \left\| (\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger - (\hat{\mathcal{U}}^\top \hat{\mathcal{C}}_{2,1})^\dagger \right\|_2 \\
 &+ \left\| \mathcal{C}_{2,1} - \hat{\mathcal{C}}_{2,1} \right\|_2 \left\| (\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger \right\|_2 \\
 &\leq \frac{1 + \sqrt{5}}{2} \frac{\epsilon_2}{\min\{\sigma_N(\hat{\mathcal{C}}_{2,1}), \sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})\}^2} + \frac{\epsilon_2}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})}
 \end{aligned} \tag{28}$$

■

$$\begin{aligned}
 \textbf{Lemma 6} \quad \Delta := \max_{x_t} \left\| (\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_t} - \tilde{\mathcal{B}}_{x_t}) (\hat{\mathcal{U}}^\top \mathcal{O}) \right\|_2 \leq \\
 \frac{1}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{O})} \left( \frac{\epsilon_2}{\min\{\sigma_N(\hat{\mathcal{C}}_{2,1}), \sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})\}^2} + \frac{\epsilon_3}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})} \right)
 \end{aligned}$$

**Proof**

$$\begin{aligned}
 & \max_{x_t} \left\| (\hat{\mathcal{U}}^\top \mathcal{O})^{-1} (\hat{\mathcal{B}}_{x_t} - \tilde{\mathcal{B}}_{x_t}) (\hat{\mathcal{U}}^\top \mathcal{O}) \right\|_2 \\
 &\leq \max_{x_t} \frac{1}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{O})} \left( \left\| \hat{\mu}_{3,1|x_t} ((\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger - (\hat{\mathcal{U}}^\top \hat{\mathcal{C}}_{2,1})^\dagger) \right\|_2 \right. \\
 &\quad \left. + \left\| (\mu_{3,1|x_t} - \hat{\mu}_{3,1|x_t}) (\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})^\dagger \right\|_2 \right) \\
 &\leq \frac{1}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{O})} \left( \frac{\epsilon_2}{\min\{\sigma_N(\hat{\mathcal{C}}_{2,1}), \sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})\}^2} + \frac{\epsilon_3}{\sigma_N(\hat{\mathcal{U}}^\top \mathcal{C}_{2,1})} \right)
 \end{aligned} \tag{29}$$

■

Based on the properties of Hilbert space embeddings explained in the main text, we have

**Lemma 7** *With probability  $1 - \delta$ ,*

$$\epsilon_1 = \|\mu_1 - \hat{\mu}_1\|_{\mathcal{F}} = O_p(m^{-1/2}(\log(1/\delta))^{1/2}) \tag{30}$$

$$\epsilon_2 = \left\| \mathcal{C}_{2,1} - \hat{\mathcal{C}}_{2,1} \right\|_{\mathcal{F} \otimes \mathcal{F}} = O_p(m^{-1/2}(\log(1/\delta))^{1/2}) \tag{31}$$

$$\epsilon_3 = \max_{x_t} \left\| \mu_{3,1|x_t} - \hat{\mu}_{3,1|x_t} \right\|_{\mathcal{F} \otimes \mathcal{F}} = O_p(\lambda^{1/2} + (\lambda m)^{-1/2}(\log(1/\delta))^{1/2}). \tag{32}$$

Plugging these results, and we finally have

**Theorem 8** *Assume  $k(x, x')$  and  $l(h, h')$  bounded, and  $\max_x \|\mathcal{A}_x\|_2 \leq 1$ , then with probability*