

SUPPLEMENTARY MATERIALS: A Rigorous Theory of Conditional Mean Embeddings*

Ilja Klebanov[†], Ingmar Schuster[‡], and T. J. Sullivan[§]

SM1. Technical Results. This section contains several technical results used in the proofs of the theorems given in the article. The following well-known result due to [SM10, Theorem 1] (see also [SM11, Theorem 2.1]) is used several times:

Theorem SM1.1. *Let \mathcal{H} , \mathcal{H}_1 and \mathcal{H}_2 be Hilbert spaces and let $A: \mathcal{H}_1 \rightarrow \mathcal{H}$ and $B: \mathcal{H}_2 \rightarrow \mathcal{H}$ be bounded linear operators with $\text{ran } A \subseteq \text{ran } B$. Then $Q := B^\dagger A: \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is a well-defined bounded linear operator, where B^\dagger denotes the Moore–Penrose pseudo-inverse of B . It is the unique operator that satisfies the conditions*

$$(SM1.1) \quad A = BQ, \quad \ker Q = \ker A, \quad \text{ran } Q \subseteq \overline{\text{ran } B^*}.$$

Remark SM1.2. In the original work of [SM10] only the existence of a bounded operator Q such that $A = BQ$ was shown. However, the construction of Q in the proof is identical to that of B^\dagger (multiplied by A). This connection has been observed before by [SM1, Corollary 2.2 and Remark 2.3], where it was proven in the case of closed range operators, leaving the proof of the general case to the reader.

The following result partially generalises [SM9, Proposition 4.1]:

Lemma SM1.3. *Let \mathcal{H} be a separable Hilbert space, let \mathcal{G} be an RKHS over \mathcal{Y} with canonical feature map ψ , and suppose that \mathcal{G} is a subset of $\mathcal{L}^2(\nu)$, where ν is a σ -finite measure on \mathcal{Y} . Then any bounded linear operator $A: \mathcal{H} \rightarrow \mathcal{G}$ is Hilbert–Schmidt as an operator $A: \mathcal{H} \rightarrow L^2(\nu)$.*

Proof. Let $h \in \mathcal{H}$ and $y \in \mathcal{Y}$. Then $(Ah)(y) = \langle \psi(y), Ah \rangle_{\mathcal{G}} = \langle A^* \psi(y), h \rangle_{\mathcal{H}}$. Thus A is a Carleman operator and the claim follows from [SM20, Theorem 6.15]. ■

The following results are used in the proofs of Sections 3, 4, and 7. Note that Lemma SM1.4 is essentially one direction of Proposition 5 in [SM13], but does not require k to be bounded, which makes a separate proof necessary.

Lemma SM1.4. *Under Assumption 2.1, if k is a characteristic kernel, then \mathcal{H}_C is dense in $L_C^2(\mathbb{P}_X)$.*

*Submitted to the editors December 6, 2019.

Funding: The research presented here was supported in part by the German Research Foundation (Deutsche Forschungsgemeinschaft) through project TrU-2 “Demand modelling and control for e-commerce using RKHS transfer operator approaches” of the Excellence Cluster “MATH+ The Berlin Mathematics Research Centre” (EXC-2046/1, project ID: 390685689). The authors also wish to thank S. Klus, H. C. Lie, M. Mollenhauer, and B. Sprungk for helpful and collegial discussions.

[†]Zuse Institute Berlin, Takustraße 7, 14195 Berlin, Germany (klebanov@zib.de).

[‡]Zalando SE, 11501 Berlin, Germany (ingmar.schuster@zalando.de).

[§]Freie Universität Berlin, Arnimallee 6, 14195 Berlin, Germany (t.j.sullivan@fu-berlin.de); Mathematics Institute and School of Engineering, The University of Warwick, Coventry, CV4 7AL, United Kingdom (t.j.sullivan@warwick.ac.uk); and Zuse Institute Berlin, Takustraße 7, 14195 Berlin, Germany (sullivan@zib.de).

29 *Proof.* Suppose that \mathcal{H}_C is not dense in $L^2_C(\mathbb{P}_X)$. Then there exists $f \in L^2(\mathbb{P}_X)$ that is
 30 not \mathbb{P}_X -a.e. constant such that $[f] \perp_{L^2_C(\mathbb{P}_X)} \mathcal{H}_C$. Choose $\tilde{f} := f - \mathbb{E}[f(X)]$ and set

$$31 \quad Q_1(E) := \int_E |\tilde{f}| \, d\mathbb{P}_X, \quad Q_2(E) := \int_E (|\tilde{f}| - \tilde{f}) \, d\mathbb{P}_X$$

32 for every Borel-measurable subset $E \subseteq \mathcal{X}$. Since $\|\tilde{f}\|_{L^1(\mathbb{P}_X)} \neq 0$, we may assume without loss
 33 of generality that $\|\tilde{f}\|_{L^1(\mathbb{P}_X)} = 1$, making Q_1 and Q_2 two *distinct* probability distributions.
 34 Since, for every $h \in \mathcal{H}$,

$$35 \quad \langle \tilde{f}, h \rangle_{L^2(\mathbb{P}_X)} = \langle f - \mathbb{E}[f(X)], h \rangle_{L^2(\mathbb{P}_X)} \stackrel{[f] \perp_{\mathcal{H}_C}}{=} \langle f - \mathbb{E}[f(X)], \mathbb{E}[h(X)] \rangle_{L^2(\mathbb{P}_X)} = 0,$$

36 it follows that $\tilde{f} \perp_{L^2(\mathbb{P}_X)} \mathcal{H}$. Let $Z_1 \sim Q_1$ and $Z_2 \sim Q_2$ and $x \in \mathcal{X}$. Since $\varphi(x) \in \mathcal{H}$,

$$37 \quad (\mathbb{E}[\varphi(Z_1)] - \mathbb{E}[\varphi(Z_2)])(x) = \langle \tilde{f}, \varphi(x) \rangle_{L^2(\mathbb{P}_X)} = 0,$$

38 which contradicts the assumption that k is characteristic. Note that, by [Assumption 2.1](#),
 39 $\mathbb{E}[\varphi(Z_1)]$ and $\mathbb{E}[\varphi(Z_2)]$ are well defined. In fact, by the Cauchy–Schwarz inequality,

$$40 \quad \mathbb{E}[\|\varphi(Z_1)\|_{\mathcal{H}}] = \int_{\mathcal{X}} \|\varphi(x)\|_{\mathcal{H}} |\tilde{f}(x)| \, d\mathbb{P}_X(x) \leq \mathbb{E}[\|\varphi(X)\|_{\mathcal{H}}^2]^{1/2} \mathbb{E}[\tilde{f}(X)^2]^{1/2} < \infty$$

41 and similarly for Z_2 . ■

42 **Lemma SM1.5.** Under [Assumption 2.1](#), $\ker C_X = \{h \in \mathcal{H} \mid h \text{ is } \mathbb{P}_X\text{-a.e. constant in } \mathcal{X}\}$
 43 and $\ker {}^u C_X = \{h \in \mathcal{H} \mid h = 0 \text{ } \mathbb{P}_X\text{-a.e. in } \mathcal{X}\}$.

44 *Proof.* This is a direct consequence of the facts that $\langle h, C_X h \rangle_{\mathcal{H}} = \mathbb{V}[h(X)]$ and that
 45 $\langle h, {}^u C_X h \rangle_{\mathcal{H}} = \|h\|_{L^2(\mathbb{P}_X)}^2$. ■

46 **Lemma SM1.6.** Under [Assumption 2.1](#), for all $h \in \mathcal{H}$ and $g \in \mathcal{G}$,

$$47 \quad \text{Cov}[h(X), f_g(X)] = \langle h, C_{XY} g \rangle_{\mathcal{H}}, \quad {}^u \text{Cov}[h(X), f_g(X)] = \langle h, {}^u C_{XY} g \rangle_{\mathcal{H}}.$$

48 *Proof.* Let $h \in \mathcal{H}$ and $g \in \mathcal{G}$ be arbitrary. Then

$$\begin{aligned} 49 \quad \text{Cov}[h(X), f_g(X)] &= \mathbb{E}[h(X)\mathbb{E}[g(Y)|X]] - \mathbb{E}[h(X)]\mathbb{E}[\mathbb{E}[g(Y)|X]] \\ 50 &= \mathbb{E}[h(X)g(Y)] - \mathbb{E}[h(X)]\mathbb{E}[g(Y)] \\ 51 &= \text{Cov}[h(X), g(Y)] \\ 52 &= \langle h, C_{XY} g \rangle_{\mathcal{H}}, \end{aligned}$$

54 as required. The second statement is proved analogously using uncentred covariance operators
 55 and without subtracting the (products of) expected values. ■

56 **Lemma SM1.7.** Under [Assumption 2.1](#), let $A: \mathcal{G} \rightarrow \mathcal{H}$ be a bounded linear operator. Then,
 57 for all $x \in \mathcal{X}$ and $y \in \mathcal{Y}$,

$$58 \quad (A\psi(y))(x) = (A^* \varphi(x))(y), \quad \mathbb{E}[(A\psi(y))(X)] = (A^* \mu_X)(y).$$

59 *Proof.* By the reproducing properties of ψ , φ , and μ_X ,

$$\begin{aligned} 60 \quad (A\psi(y))(x) &= \langle A\psi(y), \varphi(x) \rangle_{\mathcal{H}} = \langle \psi(y), A^*\varphi(x) \rangle_{\mathcal{G}} = (A^*\varphi(x))(y), \\ 61 \quad \mathbb{E}[(A\psi(y))(X)] &= \langle A\psi(y), \mu_X \rangle_{\mathcal{H}} = \langle \psi(y), A^*\mu_X \rangle_{\mathcal{G}} = (A^*\mu_X)(y), \end{aligned}$$

63 as claimed. ■

64 **Lemma SM1.8.** *Let V be a Hilbert space, let $U_1 \subseteq U_2 \subseteq \dots$ be an increasing sequence of*
 65 *closed subspaces $U_n \subseteq V$, $n \in \mathbb{N}$, and let $U := \bigcup_{n \in \mathbb{N}} U_n$. Further, let $P_{U_n}: V \rightarrow U_n$ denote the
 66 *orthogonal projection onto U_n . Then, for all $v \in \bar{U}$,**

$$67 \quad P_{U_n} v \xrightarrow{n \rightarrow \infty} v.$$

68 *Proof.* Let $v \in \bar{U}$ and $\varepsilon > 0$. Then there exists $u \in U$ such that $\|u - v\| < \varepsilon$. Since the
 69 sequence $(U_n)_{n \in \mathbb{N}}$ is increasing and U is its union, there exists an $n_0 \in \mathbb{N}$ such that $u \in U_{n_0}$
 70 and thereby $P_{U_n} u = u$ for all $n \geq n_0$. We therefore obtain, for $n \geq n_0$,

$$71 \quad \|P_{U_n} v - v\| \leq \|P_{U_n} v - P_{U_n} u\| + \|P_{U_n} u - u\| + \|u - v\| \leq \|P_{U_n}\| \|v - u\| + \|u - v\| < 2\varepsilon,$$

72 by the triangle inequality and non-expansivity of orthogonal projection. ■

73 **Lemma SM1.9.** *Under [Assumption 2.1](#), with $\mathcal{H}^{(n)}$, $C^{(n)}$, $h_g^{(n)}$ as in [Theorem 4.4](#) and ${}^u C^{(n)}$,*
 74 *${}^u h_g^{(n)}$ as in [Theorem 5.4](#),*

75 (a) $\text{Cov}[h(X), f_g(X)] = \lim_{n \rightarrow \infty} \text{Cov}[h(X), h_g^{(n)}(X)]$ for all $h \in \mathcal{H}$;

76 (b) $[h_g^{(n)}]$ is the $L^2_{\mathcal{C}}$ -orthogonal projection of $[f_g]$ onto $\mathcal{H}_{\mathcal{C}}^{(n)}$ for all $g \in \mathcal{G}$;

77 (c) ${}^u \text{Cov}[h(X), f_g(X)] = \lim_{n \rightarrow \infty} {}^u \text{Cov}[h(X), {}^u h_g^{(n)}(X)]$ for all $h \in \mathcal{H}$;

78 (d) ${}^u h_g^{(n)}$ is the L^2 -orthogonal projection of f_g onto $\mathcal{H}^{(n)}$ for all $g \in \mathcal{G}$.

79 *Proof.* We only give the proofs of (a) and (b); (c) and (d) can be proven similarly. It is
 80 clear that that $C^{(n)} \rightarrow C$ (in the strong and thereby in the weak sense) as $n \rightarrow \infty$ and that
 81 C_X and $C_X^{(n)}$ agree on $\mathcal{H}^{(n)} \ni h_g^{(n)}$. Using [Lemma SM1.6](#) we obtain, for all $h \in \mathcal{H}$,

$$\begin{aligned} 82 \quad \text{Cov}[h(X), f_g(X)] &= \langle h, C_{XY} g \rangle_{\mathcal{H}} \\ 83 \quad &= \lim_{n \rightarrow \infty} \langle h, C_{XY}^{(n)} g \rangle_{\mathcal{H}} \\ 84 \quad &= \lim_{n \rightarrow \infty} \langle h, C_X^{(n)} h_g^{(n)} \rangle_{\mathcal{H}} \\ 85 \quad &= \lim_{n \rightarrow \infty} \langle h, C_X h_g^{(n)} \rangle_{\mathcal{H}} \\ 86 \quad &= \lim_{n \rightarrow \infty} \text{Cov}[h(X), h_g^{(n)}(X)], \\ 87 \end{aligned}$$

88 which yields (a). Also, for arbitrary $h^{(n)} \in \mathcal{H}^{(n)}$, Lemma SM1.6 yields

$$\begin{aligned}
89 \quad \langle [h^{(n)}], [f_g] \rangle_{L^2_{\mathcal{C}}} &= \text{Cov}[h^{(n)}(X), f_g(X)] \\
90 &= \langle h^{(n)}, C_{XY}g \rangle_{\mathcal{H}} \\
91 &= \langle C_{YX}h^{(n)}, g \rangle_{\mathcal{G}} \\
92 &= \langle C_{YX}^{(n)}h^{(n)}, g \rangle_{\mathcal{G}} \\
93 &= \langle h^{(n)}, C_{XY}^{(n)}g \rangle_{\mathcal{H}} \\
94 &= \langle h^{(n)}, C_X^{(n)}h_g^{(n)} \rangle_{\mathcal{H}} \\
95 &= \langle h^{(n)}, C_X h_g^{(n)} \rangle_{\mathcal{H}} \\
96 &= \text{Cov}[h^{(n)}(X), h_g^{(n)}(X)] \\
97 &= \langle [h^{(n)}], [h_g^{(n)}] \rangle_{L^2_{\mathcal{C}}}, \\
98
\end{aligned}$$

99 which yields (b). ■

100 **Lemma SM1.10.** *Let Assumption 2.1 hold and $\mathcal{H}^{(1)} \subseteq \mathcal{H}^{(2)} \subseteq \dots$ be an increasing sequence*
101 *of closed subspaces $\mathcal{H}^{(n)}$ of $L^2(\mathbb{P}_X)$, $n \in \mathbb{N}$. Further, let $\mathbf{m}, \mathbf{m}^{(n)} \in L^2(\mathbb{P}_X; \mathcal{G}) \simeq L^2(\mathbb{P}_X) \otimes \mathcal{G}$*
102 *and denote $\bar{f} := f - \mathbb{E}[f(X)]$ for $f \in L^2(\mathbb{P}_X)$ and $\bar{\mathbf{f}} := \mathbf{f} - \mathbb{E}[\mathbf{f}(X)]$ for $\mathbf{f} \in L^2(\mathbb{P}_X; \mathcal{G})$.*

- 103 (a) *If $([\mathbf{m}^{(n)}](\cdot))(y)$ is the orthogonal projection in $L^2_{\mathcal{C}}(\mathbb{P}_X)$ of $([\mathbf{m}](\cdot))(y)$ onto $\mathcal{H}^{(n)}_{\mathcal{C}}$ for*
104 *each $y \in \mathcal{Y}$, then $[\mathbf{m}^{(n)}]$ is the orthogonal projection in $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})$ of $[\mathbf{m}]$ onto $(\mathcal{H}^{(n)} \otimes$
105 $\mathcal{G})_{\mathcal{C}}$.*
106 (b) *If, in addition to the assumption in (a), $([\mathbf{m}^{(n)}](\cdot))(y) \rightarrow ([\mathbf{m}](\cdot))(y)$ in $L^2_{\mathcal{C}}(\mathbb{P}_X)$ as*
107 *$n \rightarrow \infty$ for each $y \in \mathcal{Y}$, then $[\mathbf{m}^{(n)}] \rightarrow [\mathbf{m}]$ in $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})$, or, in other words, $\bar{\mathbf{m}}^{(n)}(X) \rightarrow$
108 $\bar{\mathbf{m}}(X)$ in $L^2(\mathbb{P}; \mathcal{G})$.*
109 (c) *If $(\mathbf{m}^{(n)}(\cdot))(y)$ is the orthogonal projection in $L^2(\mathbb{P}_X; \mathbb{R})$ of $(\mathbf{m}(\cdot))(y)$ onto $\mathcal{H}^{(n)}$ for*
110 *each $y \in \mathcal{Y}$, then $\mathbf{m}^{(n)}$ is the orthogonal projection in $L^2(\mathbb{P}_X; \mathcal{G})$ of \mathbf{m} onto $\mathcal{H}^{(n)} \otimes \mathcal{G}$.*
111 (d) *If, in addition to the assumption in (c), $(\mathbf{m}^{(n)}(\cdot))(y) \rightarrow (\mathbf{m}(\cdot))(y)$ in $L^2(\mathbb{P}_X; \mathbb{R})$ as*
112 *$n \rightarrow \infty$ for each $y \in \mathcal{Y}$, then $\mathbf{m}^{(n)} \rightarrow \mathbf{m}$ in $L^2(\mathbb{P}_X; \mathcal{G})$, or, in other words, $\mathbf{m}^{(n)}(X) \rightarrow$
113 $\mathbf{m}(X)$ in $L^2(\mathbb{P}; \mathcal{G})$.*

114 *Proof.* We only give the proofs of (a) and (b); (c) and (d) can be proven similarly with
115 fewer technicalities. Let $h \in \mathcal{H}^{(n)}$ and $y \in \mathcal{Y}$. Then

$$\begin{aligned}
116 \quad \langle [\mathbf{m}^{(n)}] - [\mathbf{m}], [h \otimes \psi(y)] \rangle_{L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})} &= \mathbb{E}[\langle \bar{\mathbf{m}}^{(n)}(X) - \bar{\mathbf{m}}(X), \psi(y) \rangle_{\mathcal{G}} \bar{h}(X)] \\
117 &= \mathbb{E} \left[\left((\bar{\mathbf{m}}^{(n)}(X))(y) - (\bar{\mathbf{m}}(X))(y) \right) \bar{h}(X) \right] \\
118 &= \langle ([\mathbf{m}^{(n)}](\cdot))(y) - ([\mathbf{m}](\cdot))(y), [h] \rangle_{L^2_{\mathcal{C}}(\mathbb{P}_X; \mathbb{R})} \\
119 &= 0, \\
120
\end{aligned}$$

121 which proves (a). Hence, by Lemma SM1.8, $[\mathbf{m}^{(n)}]$ converges in $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})$ to some limit $[\mathbf{m}']$.

122 This implies pointwise convergence for each $y \in \mathcal{Y}$ in the following sense:

$$\begin{aligned}
 123 \quad & \|([\mathbf{m}^{(n)}](\cdot))(y) - ([\mathbf{m}'](\cdot))(y)\|_{L^2_{\mathcal{G}}(\mathbb{P}_X; \mathbb{R})}^2 = \mathbb{E}[\langle \psi(y), \bar{\mathbf{m}}^{(n)}(X) - \bar{\mathbf{m}}'(X) \rangle_{\mathcal{G}}^2] \\
 124 \quad & \leq \|\psi(y)\|_{\mathcal{G}}^2 \mathbb{E}[\|\bar{\mathbf{m}}^{(n)}(X) - \bar{\mathbf{m}}'(X)\|_{\mathcal{G}}^2] \\
 125 \quad & = \|\psi(y)\|_{\mathcal{G}}^2 \|\mathbf{m}^{(n)} - \mathbf{m}'\|_{L^2_{\mathcal{G}}(\mathbb{P}_X; \mathcal{G})}^2 \\
 126 \quad & \xrightarrow{n \rightarrow \infty} 0. \\
 127 \quad &
 \end{aligned}$$

128 Therefore, by assumption, $[\mathbf{m}']$ agrees with $[\mathbf{m}]$ \mathbb{P}_X -a.e., proving (b). ■

129 **Lemma SM1.11.** *Under the assumptions and notation of Theorem 4.4, $(\mu^{(n)}(X, \cdot))_{n \in \mathbb{N}}$ is*
 130 *a martingale in $L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$ with respect to the filtration $(\sigma(V^{(n)}))_{n \in \mathbb{N}}$ of Σ , where $V^{(n)} :=$*
 131 *$P_{\mathcal{H}^{(n)}}(\varphi(X))$ and $P_{\mathcal{H}^{(n)}}: \mathcal{H} \rightarrow \mathcal{H}^{(n)}$ denotes the orthogonal projection in \mathcal{H} onto $\mathcal{H}^{(n)}$.*

132 *Proof.* Consider the Karhunen–Loève expansion of $\varphi(X)$,

$$133 \quad \varphi(X) = \mu_X + \sum_{i \in \mathbb{N}} Z_i h_i,$$

134 where $Z_i: (\Omega, \Sigma, \mathbb{P}) \rightarrow \mathbb{R}$ are uncorrelated real-valued random variables with $\mathbb{E}[Z_i] = 0$ and
 135 $\mathbb{V}[Z_i] = \sigma_i$ for all $i \in \mathbb{N}$, $\sigma_i \geq 0$ denoting the eigenvalue of C_X corresponding to the
 136 eigenvector h_i . We observe that $\sigma(V^{(n)}) = \sigma(Z^{(n)})$, where $Z^{(n)} := (Z_1, \dots, Z_n)$. Now let
 137 $A^{(n)} = (C_X^{(n)\dagger} C_{XY}^{(n)})^*$, $n \in \mathbb{N}$, and observe that $A^{(n)}v = A^{(n)}P_{\mathcal{H}^{(n)}}v$ and that $A^{(n+1)}$ and $A^{(n)}$
 138 agree on $\mathcal{H}^{(n)}$. Hence, for $n \in \mathbb{N}$,

$$\begin{aligned}
 139 \quad & \mathbb{E}[\mu^{(n+1)}(X, \cdot) | V^{(n)}] = \mu_Y + A^{(n+1)} \mathbb{E}[V - \mu_X | Z^{(n)}] \\
 140 \quad & = \mu_Y + A^{(n+1)} \mathbb{E}\left[\sum_{i \in \mathbb{N}} Z_i h_i \mid Z^{(n)}\right] \\
 141 \quad & = \mu_Y + A^{(n)} \sum_{i=1}^n Z_i h_i \\
 142 \quad & = \mu_Y + A^{(n)} (\varphi(X) - \mu_X) \\
 143 \quad & = \mu^{(n)}(X, \cdot),
 \end{aligned}$$

145 proving the martingale property. ■

146 **Lemma SM1.12.** *Let Assumptions 2.1 and B* hold. Then $\mathbb{E}[C_{Y|X}] = \int_{\mathcal{X}} C_{Y|X=x} d\mathbb{P}_X(x)$*
 147 *is well defined as a strong (Bochner) integral, i.e. $\int_{\mathcal{X}} \|C_{Y|X=x}\| d\mathbb{P}_X(x) < \infty$.*

148 *Proof.* The Cauchy–Schwarz inequality and (2.2) imply that, for $x \in \mathcal{X}_Y$,

$$\begin{aligned}
 149 \quad & \|C_{Y|X=x}\| = \sup_{\|g\|_{\mathcal{G}} \leq 1, \|\tilde{g}\|_{\mathcal{G}} \leq 1} \langle g, C_{Y|X=x} \tilde{g} \rangle_{\mathcal{G}} \\
 150 \quad & = \sup_{\|g\|_{\mathcal{G}} \leq 1, \|\tilde{g}\|_{\mathcal{G}} \leq 1} \langle g, \tilde{g} \rangle_{\mathcal{L}^2_{\mathcal{G}}(\mathbb{P}_{Y|X=x})} \\
 151 \quad & \leq \sup_{\|g\|_{\mathcal{G}} \leq 1, \|\tilde{g}\|_{\mathcal{G}} \leq 1} \|g\|_{\mathcal{L}^2(\mathbb{P}_{Y|X=x})} \|\tilde{g}\|_{\mathcal{L}^2(\mathbb{P}_{Y|X=x})} \\
 152 \quad & \leq \mathbb{E}[\|\psi(Y)\|_{\mathcal{G}}^2 | X = x], \\
 153 \quad &
 \end{aligned}$$

154 which, by the law of total expectation and (2.1), yields that

$$155 \quad \mathbb{E}[\|C_{Y|X}\|] \leq \mathbb{E}[\mathbb{E}[\|\psi(Y)\|_{\mathcal{G}}^2 \mid X]] = \mathbb{E}[\|\psi(Y)\|_{\mathcal{G}}^2] < \infty,$$

156 as claimed. ■

157 **SM2. Empirical Estimates for CMEs.** In practice, the kernel mean embeddings and
 158 kernel (cross-)covariance operators will often be estimated empirically from observed data, and
 159 so empirical versions of the CME, along with convergence guarantees, are of great importance.
 160 As mentioned in Remark 1.2, this topic is beyond the scope of this paper. However, we wish
 161 to point out why this is a complex problem and briefly address the main difficulties.

162 In the simplest setting, given $J \in \mathbb{N}$ independent samples $(X_1, Y_1), \dots, (X_J, Y_J) \sim \mathbb{P}_{XY}$,
 163 we have the empirical estimators

$$164 \quad \mu_X \approx \hat{\mu}_X := \frac{1}{J} \sum_{j=1}^J \varphi(X_j), \quad C_{XY} \approx \hat{C}_{XY} := \frac{1}{J} \sum_{j=1}^J (\varphi(X_j) - \hat{\mu}_X) \otimes (\psi(Y_j) - \hat{\mu}_Y),$$

166 and so on. (To simplify the notation, we suppress the obvious J -dependence of these estima-
 167 tors.) Laws of large numbers for these empirical estimators have already been established —
 168 see e.g. [SM19, Theorem 2] and [SM17, Lemma 5.8] — but the impact of this approximation
 169 error upon conditioning is, to the best of our knowledge, not yet fully quantified. One natural
 170 approach to approximate the CME $\mu_{Y|X=x}$ is the regularisation of \hat{C}_X or ${}^u\hat{C}_X$,

$$171 \quad \mu_{Y|X=x} \approx \left(({}^u\hat{C}_X + \varepsilon \text{Id}_{\mathcal{H}})^\dagger {}^u\hat{C}_{XY} \right)^* \varphi(x) = {}^u\hat{C}_{YX} ({}^u\hat{C}_X + \varepsilon \text{Id}_{\mathcal{H}})^{-1} \varphi(x),$$

172 where $\varepsilon > 0$ is a regularisation parameter which may depend on J . Note that such a
 173 regularisation can be viewed as an approximation both to the new CME formula derived
 174 in Theorem 5.3, $\mu_{Y|X=x} = ({}^u C_X^\dagger {}^u C_{XY})^* \varphi(x)$, as well as to the original (uncentred) one,
 175 $\mu_{Y|X=x} = {}^u C_{YX} {}^u C_X^{-1} \varphi(x)$. Therefore, this approach is rather well studied and convergence
 176 rates for this strategy have been established under certain conditions [SM12, SM15, SM18].

177 However, the new formulae (4.3), (4.6), and (5.3) relying on the Moore–Penrose pseudo-
 178 inverse suggest another type of approximation, where we will focus on the centred case from
 179 now on. The naïve estimate would be

$$180 \quad (\text{SM2.1}) \quad \mu_{Y|X=x} \approx \hat{\mu}_Y + (\hat{C}_X^\dagger \hat{C}_{XY})^* (\varphi(x) - \hat{\mu}_X).$$

181 Note that $\text{ran } \hat{C}_{XY} \subseteq \text{ran } \hat{C}_X$ and so (SM2.1) is well defined. However, the convergence of
 182 \hat{C}_X to C_X (e.g. in the Hilbert–Schmidt norm, as $J \rightarrow \infty$) translates badly to the convergence
 183 of \hat{C}_X^\dagger to the pseudo-inverse C_X^\dagger . One problem is that small eigenvalues of C_X might be
 184 approximated by eigenvalues of \hat{C}_X that are orders of magnitude smaller, causing \hat{C}_X^\dagger to
 185 “blow up”. So, in addition to the convergence of \hat{C}_X in the classical norms (such as the
 186 Hilbert–Schmidt norm or operator norm), we need to control the the smallest eigenvalue of
 187 \hat{C}_X .

188 A natural workaround, inspired by the finite-rank approximation in [Theorem 4.4](#), is to
 189 truncate¹ the (cross-)covariance operators to a subspace $\mathcal{H}^{(n)} = \text{span}\{h_1, \dots, h_n\}$ of \mathcal{H} with
 190 $\dim \mathcal{H}^{(n)} = n = n(J) \ll J$. One might thus hope to approximate the dominant n eigenvalues
 191 of C_X well while artificially setting the others to zero and preventing the blow-up of \widehat{C}_X^\dagger .
 192 There are several results from random matrix theory that control the behaviour of the n^{th}
 193 eigenvalue of (truncated) empirical covariance matrices for growing J and $n = n(J)$ [[SM2](#),
 194 [SM3](#), [SM4](#), [SM16](#)]. Most of these results are formulated for the case where the true mean is
 195 known to be zero and the true covariance matrix is the identity matrix and are typically of
 196 the following form, where $\lambda_{\max}(M)$ and $\lambda_{\min}(M)$ denote the largest and smallest eigenvalues
 197 of a matrix M , respectively:

198 **Theorem SM2.1** ([\[SM4, Theorem 2\]](#)). *Let $(\xi_{ij})_{i,j \in \mathbb{N}}$ be a double array of independent and*
 199 *identically distributed random variables with zero mean and unit variance. For $J \in \mathbb{N}$, let*
 200 *$n = n(J)$ be such that $n(J) \rightarrow \infty$ and $n(J)/J \rightarrow \gamma \in (0, 1)$ for $J \rightarrow \infty$ and let*

$$201 \text{ (SM2.2)} \quad A_J = (\xi_{ij})_{i=1, \dots, n(J), j=1, \dots, J}, \quad S_J = \frac{1}{J} A_J A_J^\top.$$

202 *Then, if $\mathbb{E}[\xi_{11}^4] < \infty$,*

$$203 \quad \lambda_{\max}(S_J) \xrightarrow[J \rightarrow \infty]{a.e.} (1 + \sqrt{\gamma})^2, \quad \lambda_{\min}(S_J) \xrightarrow[J \rightarrow \infty]{a.e.} (1 - \sqrt{\gamma})^2.$$

204 In our case the covariance operator C_X is not the identity; however, our case follows
 205 partially from [Theorem SM2.1](#) using the eigendecomposition of C_X ,

$$206 \quad C_X = \sum_{i \in \mathbb{N}} \sigma_i h_i \otimes h_i, \quad \sigma_1 \geq \sigma_2 \geq \dots \geq 0,$$

207 and the Karhunen–Loève expansion of $V = \varphi(X) - \mu_X$,

$$208 \quad V = \sum_{i \in \mathbb{N}} \sqrt{\sigma_i} \xi_i h_i,$$

209 where $\xi_n: \Omega \rightarrow \mathbb{R}$ are uncorrelated random variables with $\mathbb{E}[\xi_n] = 0$ and $\mathbb{V}[\xi_i] = 1$ for each
 210 $i \in \mathbb{N}$ and $(h_i)_{i \in \mathbb{N}}$ is an orthonormal eigenbasis of \mathcal{H} . To this end, let $V_j := \varphi(X_j) - \mu_X$
 211 be i.i.d. copies of V and let $V_j^{(n)}$ be their respective orthogonal projections onto $\mathcal{H}^{(n)} :=$
 212 $\text{span}\{h_1, \dots, h_n\}$, i.e.

$$213 \quad V_j = \sum_{i \in \mathbb{N}} \sqrt{\sigma_i} \xi_{ij} h_i \text{ with } \xi_{ij} \stackrel{\text{i.i.d.}}{\sim} \xi_i \text{ for all } i, j, \quad V_j^{(n)} = \sum_{i=1}^n \sqrt{\sigma_i} \xi_{ij} h_i.$$

214 To simplify notation, let us work with $(n \times n)$ -matrices instead of operators, expressed in the
 215 basis (h_1, \dots, h_n) of $\mathcal{H}^{(n)}$. Then

$$216 \text{ (SM2.3)} \quad \widehat{C}_X^{(n)} = \frac{1}{J} \sum_{j=1}^J V_j \otimes V_j = (C_X^{(n)})^{1/2} S_J (C_X^{(n)})^{1/2},$$

¹Naturally, truncation can be viewed as another form of regularisation. For further regularised estimates of large covariance and precision matrices by tapering, banding, sparsifying or similar see e.g. [[SM5](#), [SM6](#), [SM7](#), [SM21](#)] and references therein.

217 where S_J is defined by (SM2.2). So, we are *nearly* in the setup of Theorem SM2.1 and ready
 218 to conclude

$$219 \quad \lim_{J \rightarrow \infty} \lambda_{\min}(\widehat{C}_X^{(n)}) \geq \lim_{J \rightarrow \infty} \sigma_n^{1/2} \lambda_{\min}(S_J) \sigma_n^{1/2} \geq \sigma_n (1 - \sqrt{\gamma})^2,$$

220 with $n = n(J)$ and γ as in Theorem SM2.1. (Here we make use of the fact that $\|Ax\| \geq$
 221 $\lambda_{\min}(A)\|x\|$ for all x when A is positive semi-definite.) However, some obstacles remain:

- 222 • Since the random variables V_j are independent copies of V , the distributions of ξ_{ij} and
 223 $\xi_{ij'}$ agree for all j, j' , but we cannot expect the distributions of ξ_{ij} and $\xi_{i'j}$ to agree
 224 for all i, i' . Hence, ξ_{ij} do not fulfil the requirement of being identically distributed.
 225 However, there are generalisations of Theorem SM2.1 to this case under technical
 226 assumptions; see e.g. [SM3, Theorem 2.8].
- 227 • It is unclear when the condition $\mathbb{E}[\xi_{11}^4] < \infty$ is fulfilled. There are, however, some
 228 results in the case of infinite fourth moments; see [SM16] and references therein.
- 229 • The random variables ξ_i are uncorrelated but, in general, not independent. We are
 230 not aware of a result similar to Theorem SM2.1 for the uncorrelated case.²

231 Even if we manage to formulate a version of Theorem SM2.1 which suits our needs, note
 232 that the considerations so far, if executed rigorously, only guarantee that C_X^\dagger does not blow
 233 up. This is still some way short of establishing the convergence of the corresponding CME
 234 estimator

$$235 \quad (\text{SM2.4}) \quad \mu_{Y|X=x} \approx \widehat{\mu}_Y + (\widehat{C}_X^{(n)\dagger} \widehat{C}_{XY}^{(n)})^* (\varphi(x) - \widehat{\mu}_X), \quad n = n(J).$$

236 We conclude here by summarising some of the necessary steps:

- 237 • In the above considerations we assumed μ_X to be known. In practice, however, we
 238 may only employ its empirical estimate $\widehat{\mu}_X$. Since $\mu_X \otimes \mu_X$ is a rank-one estimator,
 239 this issue might be partially resolved by [SM14, Corollary 2.1], which implies that the
 240 eigenvalues of ${}^u C_X$ and $C_X = {}^u C_X - \mu_X \otimes \mu_X$ have a similar decay rate.
- 241 • In order to project onto $\mathcal{H}^{(n)} = \text{span}\{h_1, \dots, h_n\}$, we require the eigenvectors h_i of
 242 C_X . While the n dominant eigenvectors of \widehat{C}_X can be used as estimates of h_i , it is
 243 unclear how the approximation error affects the theoretical results presented above.
- 244 • Controlling the eigenvalues alone is insufficient. For a reasonable approximation of
 245 $C_X^{(n)}$ in (SM2.3), we would need a result which tells us that S_J becomes close to
 246 the identity matrix for large J , which requires us to understand the behaviour of
 247 its eigenvectors as well. The investigation of the eigenvectors of S_J turns out to be
 248 extremely challenging; see [SM2, Chapter 10] for a survey of existing results.

249

REFERENCES

- 250 [1] M. L. ARIAS, G. CORACH, AND M. C. GONZALEZ, *Generalized inverses and Douglas equations*, Proc.
 251 Amer. Math. Soc., 136 (2008), pp. 3177–3183, <https://doi.org/10.1090/S0002-9939-08-09298-8>.
 252 [2] Z. BAI AND J. W. SILVERSTEIN, *Spectral Analysis of Large Dimensional Random Matrices*, Springer Series
 253 in Statistics, Springer, New York, second ed., 2010, <https://doi.org/10.1007/978-1-4419-0661-8>.

²[SM8] discuss the estimation of covariance and precision matrices for time series data where the random variables V_j are not assumed to be independent. However, we need to drop independence in the rows of A_J , not in its columns.

- 254 [3] Z. D. BAI, *Methodologies in spectral analysis of large-dimensional random matrices, a review*, Statist.
 255 Sinica, 9 (1999), pp. 611–677, https://doi.org/10.1142/9789812793096_0015. With comments by G.
 256 J. Rodgers and Jack W. Silverstein; and a rejoinder by the author.
- 257 [4] Z. D. BAI AND Y. Q. YIN, *Limit of the smallest eigenvalue of a large-dimensional sample covariance*
 258 *matrix*, Ann. Probab., 21 (1993), pp. 1275–1294, <https://doi.org/10.1214/aop/1176989118>.
- 259 [5] P. J. BICKEL AND E. LEVINA, *Covariance regularization by thresholding*, Ann. Statist., 36 (2008),
 260 pp. 2577–2604, <https://doi.org/10.1214/08-AOS600>.
- 261 [6] P. J. BICKEL AND E. LEVINA, *Regularized estimation of large covariance matrices*, Ann. Statist., 36
 262 (2008), pp. 199–227, <https://doi.org/10.1214/009053607000000758>.
- 263 [7] T. T. CAI, C.-H. ZHANG, AND H. H. ZHOU, *Optimal rates of convergence for covariance matrix estima-*
 264 *tion*, Ann. Statist., 38 (2010), pp. 2118–2144, <https://doi.org/10.1214/09-AOS752>.
- 265 [8] X. CHEN, M. XU, AND W. B. WU, *Covariance and precision matrix estimation for high-dimensional*
 266 *time series*, Ann. Statist., 41 (2013), pp. 2994–3021, <https://doi.org/10.1214/13-AOS1182>.
- 267 [9] E. DE VITO, L. ROSASCO, AND A. CAPONNETTO, *Discretization error analysis for Tikhonov regulariza-*
 268 *tion*, Analysis and Applications, 04 (2006), pp. 81–99, <https://doi.org/10.1142/S0219530506000711>.
- 269 [10] R. G. DOUGLAS, *On majorization, factorization, and range inclusion of operators on Hilbert space*, Proc.
 270 Amer. Math. Soc., 17 (1966), pp. 413–415, <https://doi.org/10.2307/2035178>.
- 271 [11] P. A. FILLMORE AND J. P. WILLIAMS, *On operator ranges*, Adv. Math., 7 (1971), pp. 254–281, [https://doi.org/10.1016/S0001-8708\(71\)80006-3](https://doi.org/10.1016/S0001-8708(71)80006-3).
- 272 [12] K. FUKUMIZU, *Nonparametric Bayesian inference with kernel mean embedding*, in Modern Methodology
 273 and Applications in Spatial-Temporal Modeling, Springer Japan, Tokyo, 2015, pp. 1–24, https://doi.org/10.1007/978-4-431-55339-7_1.
- 274 [13] K. FUKUMIZU, F. R. BACH, AND M. I. JORDAN, *Kernel dimension reduction in regression*, Ann. Statist.,
 275 37 (2009), pp. 1871–1905, <https://doi.org/10.1214/08-AOS637>.
- 276 [14] I. C. GOHBERG AND M. G. KREĬN, *Introduction to the Theory of Linear Nonselfadjoint Operators*, vol. 18
 277 of Translations of Mathematical Monographs, American Mathematical Society, Providence, R.I., 1969.
 278 Transl. A. Feinstein.
- 279 [15] S. GRÜNEWÄLDER, G. LEVER, L. BALDASSARRE, S. PATTERSON, A. GRETTON, AND M. PONTIL, *Condi-*
 280 *tional mean embeddings as regressors*, in Proceedings of the 29th International Conference on Machine
 281 Learning, 2012, pp. 1823–1830, <https://icml.cc/2012/papers/898.pdf>.
- 282 [16] J. HEINY AND T. MIKOSCH, *Almost sure convergence of the largest and smallest eigenvalues of high-*
 283 *dimensional sample correlation matrices*, Stochastic Process. Appl., 128 (2018), pp. 2779–2815, <https://doi.org/10.1016/j.spa.2017.10.002>.
- 284 [17] M. MOLLENHAUER, *Singular value decomposition of operators on reproducing kernel Hilbert spaces*, mas-
 285 ter’s thesis, Freie Universität Berlin, 2018.
- 286 [18] J. PARK AND K. MUANDET, *A measure-theoretic approach to kernel conditional mean embeddings*, 2020,
 287 <https://arxiv.org/abs/2002.03689>.
- 288 [19] A. J. SMOLA, A. GRETTON, L. SONG, AND B. SCHÖLKOPF, *A Hilbert space embedding for distributions*, in
 289 Proceedings of the 18th International Conference on Algorithmic Learning Theory, Berlin, Heidelberg,
 290 2007, Springer, pp. 13–31, https://doi.org/10.1007/978-3-540-75225-7_5.
- 291 [20] J. WEIDMANN, *Linear Operators in Hilbert Spaces*, vol. 68 of Graduate Texts in Mathematics, Springer-
 292 Verlag, New York-Berlin, 1980, <https://doi.org/10.1007/978-1-4612-6027-1>. Transl. J. Szücs.
- 293 [21] M. YUAN, *High dimensional inverse covariance matrix estimation via linear programming*, J. Mach. Learn.
 294 Res., 11 (2010), pp. 2261–2286, <http://www.jmlr.org/papers/volume11/yuan10b/yuan10b.pdf>.