

Deconditional Downscaling with Gaussian processes



Siu Lun Chau^{*}, Shahine Bouabid^{*}, Dino Sejdinovic
35th Conference on Neural Information Processing, 2021

Background on Kernel Embeddings

Deconditional Mean Embedding

Deconditional Downscaling with Gaussian processes

Background on Kernel Embeddings



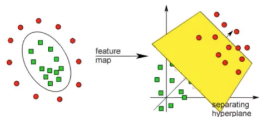
- | **Kernel method** is any method that endows a generic abstract domain X with an inner product structure induced by some feature transformation $\varphi : X \rightarrow H$.

- | **Kernel method** is any method that endows a generic abstract domain X with an inner product structure induced by some feature transformation $\varphi : X \rightarrow H$.
- | **Kernel function** is as an inner product of features: any function $k : X \times X \rightarrow \mathbb{R}$ for which there exists a **Hilbert space** H and a map $\varphi : X \rightarrow H$ s.t. $k(x, x^0) = \langle \varphi(x), \varphi(x^0) \rangle_H$ for all $x, x^0 \in X$.

- | **Kernel method** is any method that endows a generic abstract domain X with an inner product structure induced by some feature transformation $\varphi : X \rightarrow H$.
- | **Kernel function** is as an inner product of features: any function $k : X \times X \rightarrow \mathbb{R}$ for which there exists a **Hilbert space** H and a map $\varphi : X \rightarrow H$ s.t. $k(x, x^0) = \langle \varphi(x), \varphi(x^0) \rangle_H$ for all $x, x^0 \in X$.
- | There exists a canonical feature space H_k , called **reproducing kernel Hilbert space (RKHS)** with canonical feature map $x \mapsto k(\cdot, x)$, where:
 - | $\forall x \in X, k(\cdot, x) \in H_k$
 - | $\forall x \in X, \forall f \in H_k, \langle f, k(\cdot, x) \rangle_{H_k} = f(x)$.
 Thus also $k(x, x^0) = \langle k(\cdot, x), k(\cdot, x^0) \rangle_{H_k}$.

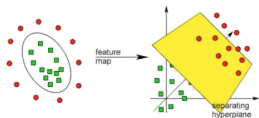
- | **Kernel method** is any method that endows a generic abstract domain X with an inner product structure induced by some feature transformation $\varphi : X \rightarrow H$.
- | **Kernel function** is as an inner product of features: any function $k : X \times X \rightarrow \mathbb{R}$ for which there exists a **Hilbert space** H and a map $\varphi : X \rightarrow H$ s.t. $k(x, x^0) = \langle \varphi(x), \varphi(x^0) \rangle_H$ for all $x, x^0 \in X$.
- | There exists a canonical feature space H_k , called **reproducing kernel Hilbert space (RKHS)** with canonical feature map $x \mapsto k(\cdot, x)$, where:
 - | $\forall x \in X, k(\cdot, x) \in H_k$
 - | $\forall x \in X, \forall f \in H_k, \langle f, k(\cdot, x) \rangle_{H_k} = f(x)$.
 Thus also $k(x, x^0) = \langle k(\cdot, x), k(\cdot, x^0) \rangle_{H_k}$.
- | **Moore-Aronszajn Theorem:** every positive semidefinite kernel is the kernel of a unique RKHS.

- | implicit feature map $x \mapsto k(\cdot, x) \in H_k$
replaces $x \mapsto [\phi_1(x), \dots, \phi_s(x)] \in \mathbb{R}^s$
- | $\langle k(\cdot, x), k(\cdot, y) \rangle_{H_k} = k(x, y)$
inner products readily available



Cortes & Vapnik, 1995; Schölkopf
& Smola, 2001

- | implicit feature map $x \mapsto k(\cdot, x) \in H_k$
replaces $x \mapsto [\phi_1(x), \dots, \phi_s(x)] \in \mathbb{R}^s$
- | $\langle k(\cdot, x), k(\cdot, y) \rangle_{H_k} = k(x, y)$
inner products readily available



Cortes & Vapnik, 1995; Schölkopf
& Smola, 2001

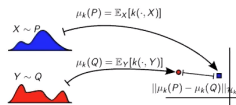
- | **RKHS embedding:** implicit feature mean

Smola et al, 2007; Sriperumbudur et al, 2010; Muandet et al, 2017

$P \mapsto \mu_k(P) = \mathbb{E}_X [k(\cdot, X)] \in H_k$
replaces

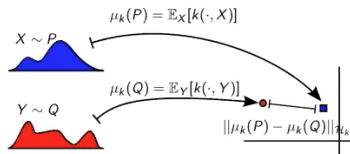
$P \mapsto [\mathbb{E}\phi_1(X), \dots, \mathbb{E}\phi_s(X)] \in \mathbb{R}^s$

- | $\langle \mu_k(P), \mu_k(Q) \rangle_{H_k} = \mathbb{E}_{X \sim P, Y \sim Q} k(X, Y)$
inner products easy to estimate



Gretton et al, 2005; Gretton et al, 2006; Fukumizu et al, 2007; DS et al, 2013; Muandet et al, 2012; Szabo et al, 2015

- Maximum Mean Discrepancy (MMD) (Borgwardt et al, 2006; Gretton et al, 2007) between P and Q :



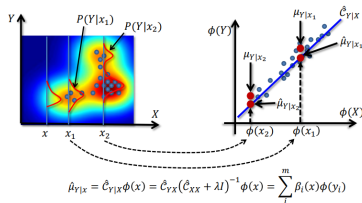
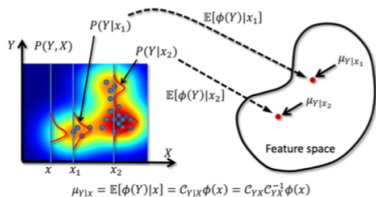
$$\text{MMD}_k(P, Q) = \sqrt{\| \mu_k(P) - \mu_k(Q) \|_{H_k}^2} = \sup_{f \in \mathcal{H}_k: \|f\|_{H_k} \leq 1} | \mathbb{E} f(X) - \mathbb{E} f(Y) |$$

- Characteristic** kernels: $\text{MMD}_k(P, Q) = 0$ iff $P = Q$ (Gaussian RBF $\exp(-\frac{1}{2\sigma^2} \|x - x'\|_2^2)$, Matérn family, inverse multiquadrics.).
- Can encode structural properties in the data: kernels on non-Euclidean domains, networks, images, text...

Consider a joint distribution P_{XY} over rvs (X, Y) taking values in $X \times Y$. Given a kernel $\ell : Y \times Y \rightarrow \mathbb{R}$, the conditional mean embedding (CME) of $Y|X = x$ is defined as:

$$\mu_{Y|X=x} := \mathbb{E}[\ell(\cdot, Y)|X = x] = \int_Y \ell(\cdot, y) dP_{Y|X=x}(y) \in H_\ell.$$

- | Allows to compute $\delta f \in H_\ell, \mathbb{E}[f(Y)|X = x] = \langle f, \mu_{Y|X=x} \rangle_{H_\ell}$.



To model conditional embeddings as functions of x , we associate them with a conditional mean operator (CMO) $C_{Y|X} : H_k \rightarrow H_\ell$ where $k : X \rightarrow X \rightarrow \mathbb{R}$, which satisfies

$$\mu_{Y|X=x} = C_{Y|X}k(\cdot, x).$$

This is essentially a feature-to-feature RKHS-valued ridge regression.

- | The same notion of a positive-definite kernel, but conceptual gaps between communities.
- | Orthogonal projection in RKHS , Conditioning in GPs.
- | 0/1 laws: GP sample paths with infinite-dimensional covariance kernel k lie a.s. outside of H_k . The space of sample paths can be thought of as an "outer shell" of H_k .

Deconditional Mean Embedding

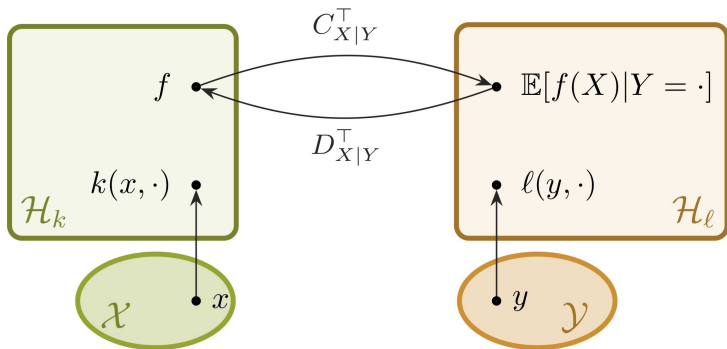


- | CMO $C_{Y|X} : H_k \rightarrow H_\ell$ allows us to reason about the conditional expectation of any $f \in H_\ell$ in the following way:

$$C_{Y|X}^\triangleright f = \mathbb{E}[f(Y) | X = \cdot]$$

- | Hsu and Ramos [2019] proposed Deconditional Mean Operator $D_{Y|X} : H_\ell \rightarrow H_k$ (DMO) as a natural counterpart to CMO.
- | Given the conditional mean function, we recover f by setting

$$D_{Y|X}^\triangleright \mathbb{E}[f(Y) | X = \cdot] = f$$



- | DMO can be rewritten in terms of CMOs and cross-covariance operators:

$$D_{Y|X} = (C_{Y|X}C_{XX})^\succ (C_{Y|X}C_{XX}C_{Y|X}^\succ)^{-1}$$

under regularity assumption.

- | In our work we further characterise their formulation and show that deconditioning can be seen as a **two-staged vector valued kernel ridge regression**.

- First introduced by Grunewalder et al. [2012], CME has a feature-to-feature regression interpretation:

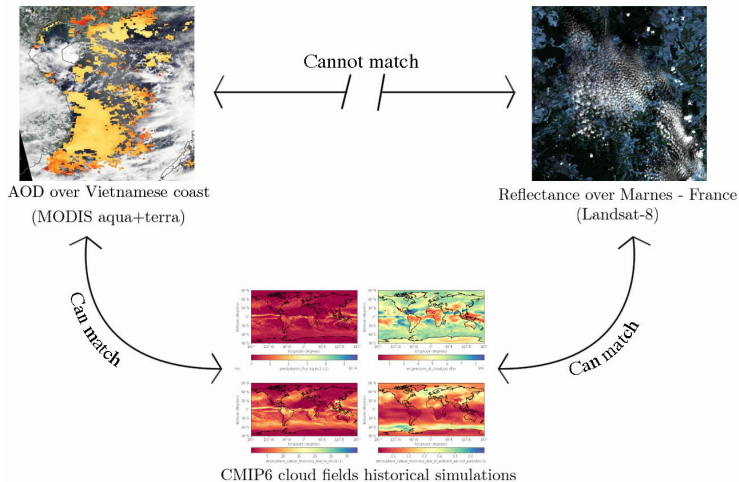
$$E_c(C) := \frac{1}{N} \sum_{i=1}^N jj\ell_{y_i} \quad C^>k_{x_i}jj^2_{H_\ell} + \lambda\Omega(C)$$

- Deconditioning has a similar interpretation, thus allowing one to apply regression convergence result to theoretically analyze DMO,

$$E_d(D) := \frac{1}{M} \sum_{j=1}^M jjk_{x_j} \quad DC_{YjX}^>k_{x_j}jj^2_{H_k} + \epsilon\Omega(D)$$

Deconditional Downscaling with Gaussian processes





Dataset

- | We have a dataset of N bags of high-resolution (HR) covariates ${}^b\mathbf{x}_j := \{x_j^{(1)}, \dots, x_j^{(n_j)}\}$ paired with a mediating variable y_j

$$D_1 = \{{}^b\mathbf{x}_j, y_j\}_{j=1}^N, \quad (1)$$

and a second separate dataset of M mediating variables \tilde{y}_j paired with a low-resolution (LR) response \tilde{z}_j .

$$D_2 = \{\tilde{y}_j, \tilde{z}_j\}_{j=1}^M. \quad (2)$$

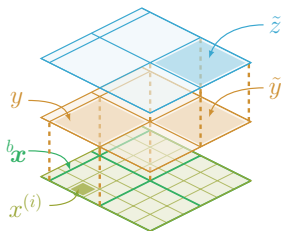


Figure 1: Representation of covariates indirect pairing

Objective

- Downscale the response \tilde{z} to the HR granularity level of $x^{(i)}$ covariates, i.e. find a function $f : \mathcal{X} \rightarrow \mathbb{R}$ which maps between HR covariates and HR responses (not observed!).

How?

- Assume there exists mapping $f : \mathcal{X} \rightarrow \mathbb{R}$.
- But only observe “aggregates” $\mathbb{E}_X[f(X) | Y = y]$.
- Recover underlying function f out of aggregated observations.

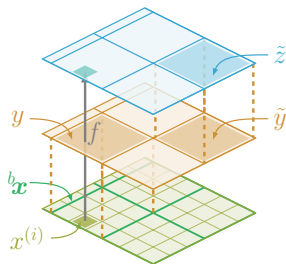


Figure 2: We wish to learn a map from HR covariates to HR estimate of the response.

Given dataset $D_2 = \{\tilde{y}_j, \tilde{z}_j\}_{j=1}^M$,

Observation Model

- | We suppose that the HR response f aggregates into the LR response \tilde{z}_j as

$$\tilde{z}_j = \mathbb{E}_X[f(X)|Y = \tilde{y}_j] + \varepsilon_j \quad (3)$$

with noise $\varepsilon_j \sim N(0, \sigma^2)$.

Recovering f corresponds to the *deconditioning* problem of Hsu and Ramos [1]:

- | Given $g : Y \rightarrow \mathbb{R}$, infer the function $f : X \rightarrow \mathbb{R}$ such that

$$g(y) = \mathbb{E}_X[f(X)|Y = y]. \quad (4)$$

f is called the *deconditional mean* of g w.r.t. $\mathbb{P}_{X|Y}$.

Given dataset $D_2 = \{\tilde{y}_j, \tilde{z}_j\}_{j=1}^M$,

Observation Model

- | We suppose that the HR response f aggregates into the LR response \tilde{z}_j as

$$\tilde{z}_j = \mathbb{E}_X[f(X)|Y = \tilde{y}_j] + \varepsilon_j \quad (5)$$

with noise $\varepsilon_j \sim N(0, \sigma^2)$.

Recovering f corresponds to the *deconditioning* problem of Hsu and Ramos [1]:

- | Given $g : Y \rightarrow \mathbb{R}$, infer the function $f : X \rightarrow \mathbb{R}$ such that

$$g(y) = \mathbb{E}_X[f(X)|Y = y]. \quad (6)$$

f is called the *deconditional mean* of g w.r.t. $\mathbb{P}_{X|Y}$.

Conditional Mean Process

- | By placing a GP prior on $f \sim \text{GP}(m, k)$, we can represent the LR field as:

$$g(y) = \mathbb{E}_X[f(X)|Y = y] = \int_X f(x) d\mathbb{P}_{X|Y=y}(x) \quad \text{GP}(\nu, q)$$

where by linearity of the integral,

$$\begin{aligned} \nu(y) &= \mathbb{E}_X[m(X)|Y = y] \\ q(y, y^0) &= \mathbb{E}_{X, X^0}[k(X, X^0)|Y = y, Y^0 = y^0] = \int \mu_{X|Y=y}, \mu_{X|Y=y^0} k. \end{aligned}$$

- | ν and q are estimated via conditional mean operator $C_{X|Y}$ using $D_1 = \{\mathbf{x}_j, y_j\}_{j=1}^N$ and a second kernel $\ell : Y \times Y \rightarrow \mathbb{R}$.

- | $f \sim \text{GP}(m, k)$ and $g \sim \text{GP}(\nu, q)$
- | The latent HR field $f(x)$ and the observed noisy LR field $\tilde{z} = g(\tilde{y}) + \epsilon$ are both normally distributed:

$$\begin{bmatrix} f(x) \\ \tilde{z} \end{bmatrix} \Big| \tilde{y} \sim \mathcal{N} \left(\begin{bmatrix} m(x) \\ \nu(\tilde{y}) \end{bmatrix}, \begin{bmatrix} k(x, x) & \text{Cov}(f(x), g(\tilde{y})) \\ \text{Cov}(g(\tilde{y}), f(x)) & q(\tilde{y}, \tilde{y}) + \sigma^2 \end{bmatrix} \right)$$

- | $f \sim \text{GP}(m, k)$ and $g \sim \text{GP}(\nu, q)$
- | The latent HR field $f(x)$ and the observed noisy LR field $\tilde{z} = g(\tilde{y}) + \epsilon$ are both normally distributed:

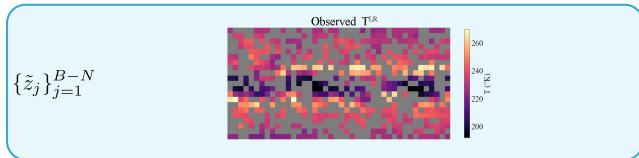
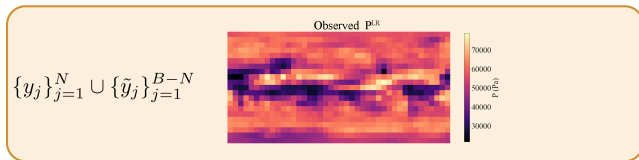
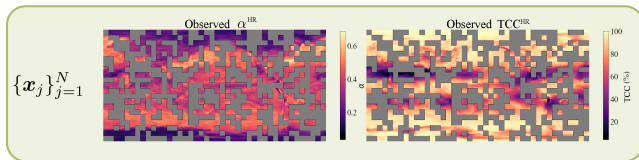
$$\begin{bmatrix} f(x) \\ \tilde{z} \end{bmatrix} | \tilde{y} \sim \mathcal{N} \left(\begin{bmatrix} m(x) \\ \nu(\tilde{y}) \end{bmatrix}, \begin{bmatrix} k(x, x) & \text{Cov}(f(x), g(\tilde{y})) \\ \text{Cov}(g(\tilde{y}), f(x)) & q(\tilde{y}, \tilde{y}) + \sigma^2 \end{bmatrix} \right)$$

- | Allows to directly obtain *deconditional posterior* $f | \tilde{z} \sim \text{GP}(m_d, k_d)$ from D_2 with:

$$\begin{aligned} \hat{m}_d(x) &= m(x) + k(x, \mathbf{x}) \mathbf{A} (\hat{\mathbf{Q}} + \sigma^2 \mathbf{I}_M)^{-1} (\tilde{\mathbf{z}} - \nu(\tilde{\mathbf{y}})) \\ \hat{k}_d(x, x^\theta) &= k(x, x^\theta) - k(x, \mathbf{x}) \mathbf{A} (\hat{\mathbf{Q}} + \sigma^2 \mathbf{I}_M)^{-1} \mathbf{A}^\top k(\mathbf{x}, x^\theta) \end{aligned}$$

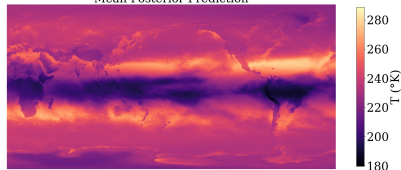
where $\mathbf{A} := (\ell(\mathbf{y}, \mathbf{y}) + N\lambda \mathbf{I}_N)^{-1} \ell(\mathbf{y}, \tilde{\mathbf{y}})$ with $\lambda > 0$, $\hat{\mathbf{Q}} := \hat{q}(\tilde{\mathbf{y}}, \tilde{\mathbf{y}})$.

Mediated Downscaling of Atmospheric Temperatures 23

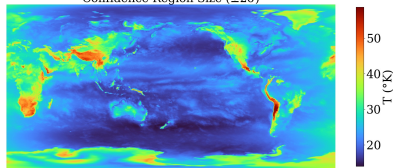


Mediated Downscaling of Atmospheric Temperatures 24

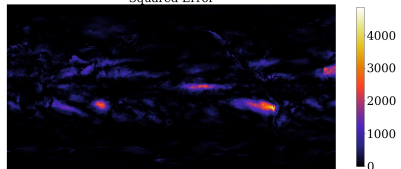
Mean Posterior Prediction



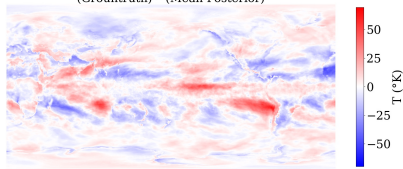
Confidence Region Size ($\pm 2\sigma$)

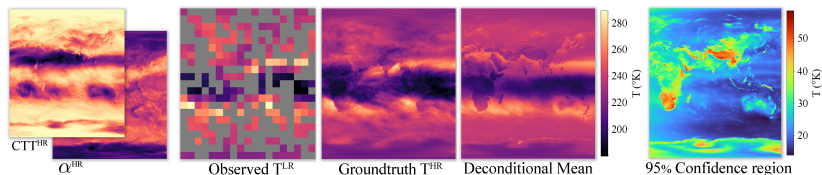


Squared Error



(Groundtruth) - (Mean Posterior)





Model	RMSE #		MAE #		Corr. "		SSIM "	
Krigging	8.02	0.28	5.55	0.17	0.831	0.012	0.212	0.011
VBAgg [2]	8.25	0.15	5.82	0.11	0.821	0.006	0.182	0.004
Our method	7.40	0.25	5.34	0.22	0.848	0.011	0.212	0.013

Table 1: Downscaling similarity scores of posterior mean against HR groundtruth; reports 1 s.d.

- [1] Kelvin Hsu and Fabio Ramos. Bayesian Deconditional Kernel Mean Embeddings. *Proceedings of Machine Learning Research*. PMLR, 2019.
- [2] Leon Ho Chung Law, Dino Sejdinovic, Ewan Cameron, Tim C.D. Lucas, Seth Flaxman, Katherine Battle, and Kenji Fukumizu. Variational learning on aggregate outputs with Gaussian processes. In *Advances in Neural Information Processing Systems*, 2018.

Usual CMO estimation

$$\begin{aligned}
 C_{XY} &= C_{XY}(C_{YY} + \lambda \text{Id}_{H_\ell})^{-1} \\
 &= E_{XY}[k(\cdot, X) \ell(Y, \cdot)](E_Y[\ell(\cdot, Y) \ell(Y, \cdot)] + \lambda \text{Id}_{H_\ell})^{-1} \\
 &= \frac{1}{N} k(\cdot, \mathbf{x}) \ell(\mathbf{y}, \cdot) \left(\frac{1}{N} \ell(\cdot, \mathbf{y}) \ell(\mathbf{y}, \cdot) + \lambda \text{Id}_{H_\ell} \right)^{-1} \\
 &= k(\cdot, \mathbf{x}) (\ell(\mathbf{y}, \cdot) \ell(\cdot, \mathbf{y}) + N \lambda \mathbf{I}_n)^{-1} \ell(\mathbf{y}, \cdot) \\
 &= k(\cdot, \mathbf{x}) (\mathbf{L}_{\mathbf{y}\mathbf{y}} + n \lambda \mathbf{I}_n)^{-1} \ell(\mathbf{y}, \cdot)
 \end{aligned}$$

where $k(\cdot, \mathbf{x}) = [k(\cdot, x_1), \dots, k(\cdot, x_N)]$, $\ell(\cdot, \mathbf{y}) = [\ell(\cdot, y_1), \dots, \ell(\cdot, y_N)]$.

Bagged CMO estimation

$$\begin{aligned}
 C_{X|Y} &= C_{XY}(C_{YY} + \lambda \text{Id}_{H_\ell})^{-1} \\
 &= \mathbb{E}_{X|Y}[k(\cdot, X) \ell(Y, \cdot)](C_{YY} + \lambda \text{Id}_{H_\ell})^{-1} \\
 &= \mathbb{E}_Y[\mathbb{E}_X[k(\cdot, X)|Y] \ell(Y, \cdot)](C_{YY} + \lambda \text{Id}_{H_\ell})^{-1} \\
 &= \mathbb{E}_Y[\mu_{X|Y} \ell(Y, \cdot)](C_{YY} + \lambda \text{Id}_{H_\ell})^{-1} \\
 &= \frac{1}{N} \hat{\mu}_{X|Y} \ell(\mathbf{y}, \cdot) \left(\frac{1}{N} \ell(\cdot, \mathbf{y}) \ell(\mathbf{y}, \cdot) + \lambda \text{Id}_{H_\ell} \right)^{-1} \\
 &= \hat{\mu}_{X|Y} (\mathbf{L}_{\mathbf{y}\mathbf{y}} + N\lambda \mathbf{I}_n)^{-1} \ell(\mathbf{y}, \cdot) \\
 &:= {}^S \hat{C}_{X|Y}
 \end{aligned}$$

where $\hat{\mu}_{X|Y} = [\hat{\mu}_{X|Y=y_1}, \dots, \hat{\mu}_{X|Y=y_N}]$ and

$$\hat{\mu}_{X|Y=y_j} = \frac{1}{n_j} \sum_{i=1}^{n_j} k(\cdot, x_j^{(i)}).$$

$$k(\cdot, \mathbf{x})\ell(\mathbf{y}, \cdot) = \sum_{i=1}^N k(\cdot, x_i) \ell(\cdot, y_i) \in H_k \otimes H_\ell$$

$$\begin{aligned} \ell(\mathbf{y}, \cdot)\ell(\cdot, \mathbf{y}) &= [\ell(\cdot, y_i), \ell(\cdot, y_j)]_{i,j} \in \mathbb{R}^{N \times N} \\ &= [\ell(y_i, y_j)]_{i,j} \in \mathbb{R}^{N \times N} \\ &= \mathbf{L}_{\mathbf{y}\mathbf{y}} \in \mathbb{R}^{N \times N} \end{aligned}$$