

# A Dual Formulation for Probabilistic Principal Component Analysis

## Context & Contributions

- ▶ We characterize Probabilistic Principal Component Analysis (Tipping & Bishop, 1999) in Hilbert spaces and demonstrate how the **optimal solution admits a representation in dual space**.
- ▶ We develop a new extension of **KPCA** (Mika et al., 1998; Schölkopf et al., 1998) incorporating a **noise assumption** on the feature map.
- ▶ We give a **probabilistic** interpretation to the **generation in KPCA** (Schreurs & Suykens, 2018).
- ▶ We illustrate how the dual model works on a toy and a real dataset and **show its connections to KPCA**.

## Definitions

As some kernels lead to **infinite dimensional feature maps**, we need to carefully define **finite subspaces** to allow the proper definition of **probability distributions**. Given a set of observations  $\{\varphi_i \in \mathcal{H}\}_{i=1}^N$  and a kernel space  $\mathcal{E}$  with basis  $\{e_i\}_{i=1}^N$  and a finite feature space  $\mathcal{H}_{\mathcal{E}} = \text{span}\{\varphi_1, \dots, \varphi_N\}$ . We define the mapping  $\Phi : \mathcal{E} \rightarrow \mathcal{H}_{\mathcal{E}} : \sum_{i=1}^N \varphi_i e_i^*$ . This defines the covariance  $\Phi \circ \Phi^*$  and the kernels  $\Phi^* \circ \Phi$ . We consider a latent space  $\mathcal{L} \subset \mathcal{E}$  of dimension  $q$  and define the interconnection operator  $W : \mathcal{L} \rightarrow \mathcal{H}_{\mathcal{E}}$ , where  $\mathcal{H}_{\mathcal{L}} \subset \mathcal{H}_{\mathcal{E}}$ . We refer to the feature spaces  $\mathcal{H}$ ,  $\mathcal{H}_{\mathcal{E}}$  and  $\mathcal{H}_{\mathcal{L}}$  as the **primal spaces** and  $\mathcal{E}$  and  $\mathcal{L}$  as the **dual spaces**.

Distribution	Interpretation	Primal (features)	Dual (kernels)
latent   observation	latent projection	$h \phi \sim \mathcal{N}(\Sigma_{h \phi}^{-1} \circ W^* \phi, \sigma^2 \Sigma_{h \phi}^{-1})$	$h k \sim \mathcal{N}(\Sigma_{h k}^{-1} \circ A k, \Sigma_{h k}^{-1})$
observation   latent	latent-based generation	$\phi h \sim \mathcal{N}(W h, \sigma^2 I_{\mathcal{H}_{\mathcal{E}}})$	$k h \sim \mathcal{N}((\Phi^* \circ \Phi) \circ A h, \sigma^2 \Phi^* \circ \Phi)$
latent	latent prior	$h \sim \mathcal{N}(\mathbf{0}, I_{\mathcal{L}})$	$h \sim \mathcal{N}(\mathbf{0}, I_{\mathcal{L}})$
observation	absolute generation	$\phi \sim \mathcal{N}(\mu, W \circ W^* + \sigma^2 I_{\mathcal{H}_{\mathcal{E}}})$	$k \sim \mathcal{N}(\mathbf{0}, A^* \circ A + \sigma^2 (\Phi^* \circ \Phi)^{-1})$

Table: Interpretation of the different distributions of the Prob. PCA framework after training, in both primal and dual formulations. The covariance operators are given by  $\Sigma_{h|\phi} = (W^* \circ W + \sigma^2 I_{\mathcal{L}})^{-1}$  and  $\Sigma_{h|k} = (A^* \circ (\Phi^* \circ \Phi) \circ A + \sigma^2 I_{\mathcal{L}})^{-1}$ . For the simplicity of this presentation, we do not consider the centering of kernels of feature maps and refer to the paper for these considerations.

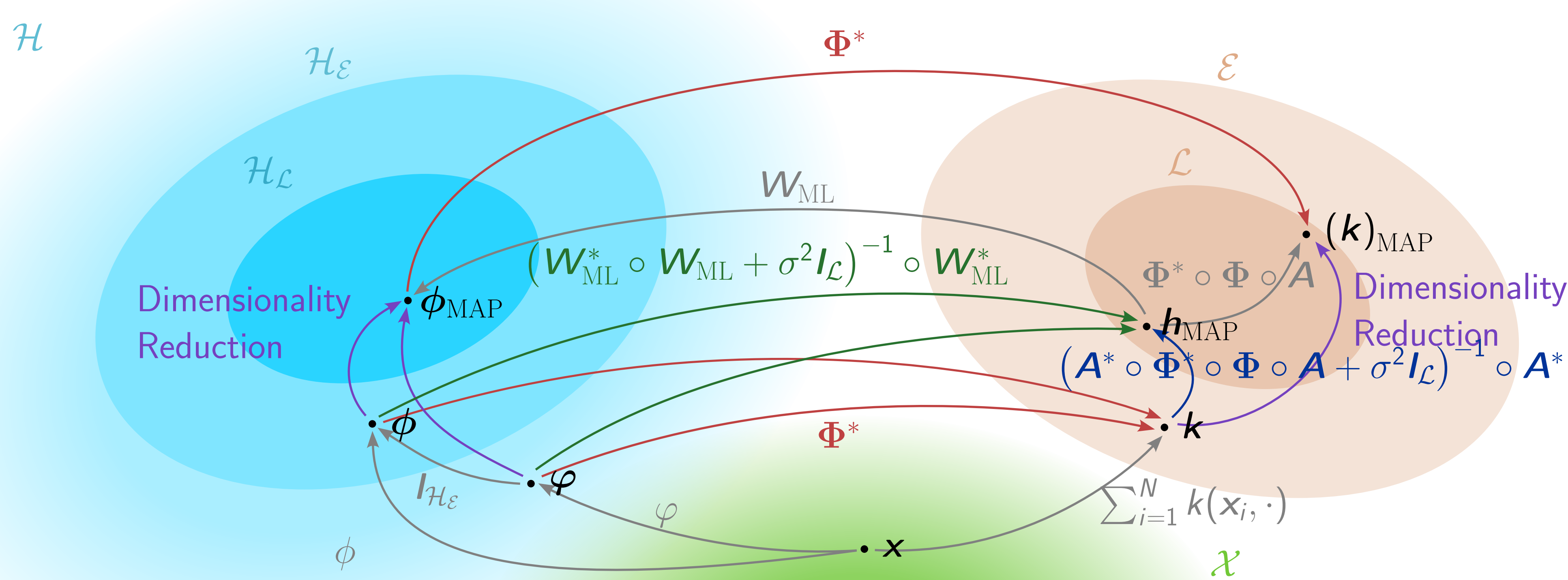


Figure: Global overview of the Probabilistic Principal Component Analysis in both primal and dual formulations. The primal spaces, or feature  $\mathcal{H}$ ,  $\mathcal{H}_{\mathcal{E}}$  and  $\mathcal{H}_{\mathcal{L}}$  are in blue. The dual, or kernel and latent spaces  $\mathcal{E}$  and  $\mathcal{L}$  are in brown. The input space  $\mathcal{X}$  is in green. The color or the applications (arrows) is just for the readability and has nothing to do with the color of the spaces.

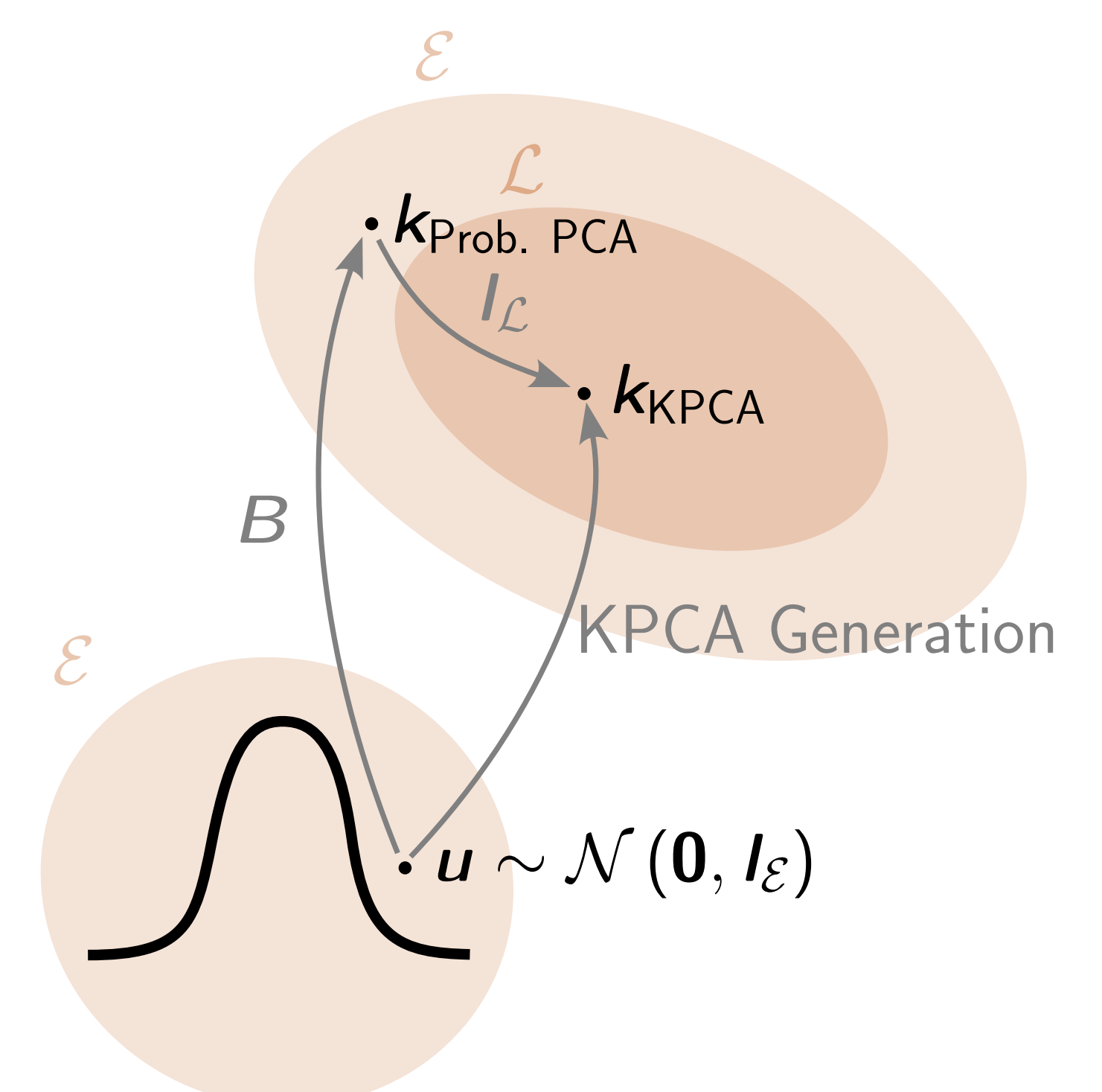


Figure: Schematic overview of the dual sampling in Prob. PCA compared to the generation in KPCA, with  $B : \mathcal{E} \rightarrow \mathcal{E} : N^{-1/2} \sum_{p=1}^q \lambda_p \epsilon_p r_p^* + \sum_{p=q+1}^N \sigma \lambda_p^{1/2} \epsilon_p r_p^*$ .

## Maximum Likelihood

If we consider the eigendecomposition of the covariance  $\Phi \circ \Phi^* = \sum_{i=1}^N \lambda_i v_i v_i^*$  and take the  $q$  dominant eigenpairs ( $\lambda_1 \geq \dots \geq \lambda_q \geq \dots \geq \lambda_N$ ), we have

$$W_{ML} = \sum_{p=1}^q \sqrt{\lambda_p / N - \sigma_{ML}^2} v_p r_p^*,$$

$$\sigma_{ML}^2 = \frac{1}{N(N-q)} \sum_{p=q+1}^N \lambda_p,$$

with  $\{r_p\}_{p=1}^q$  and arbitrary orthonormal base of the latent space  $\mathcal{L}$ . We note that  $W_{ML}$  is not unique as it is rotational invariant.

## Representation

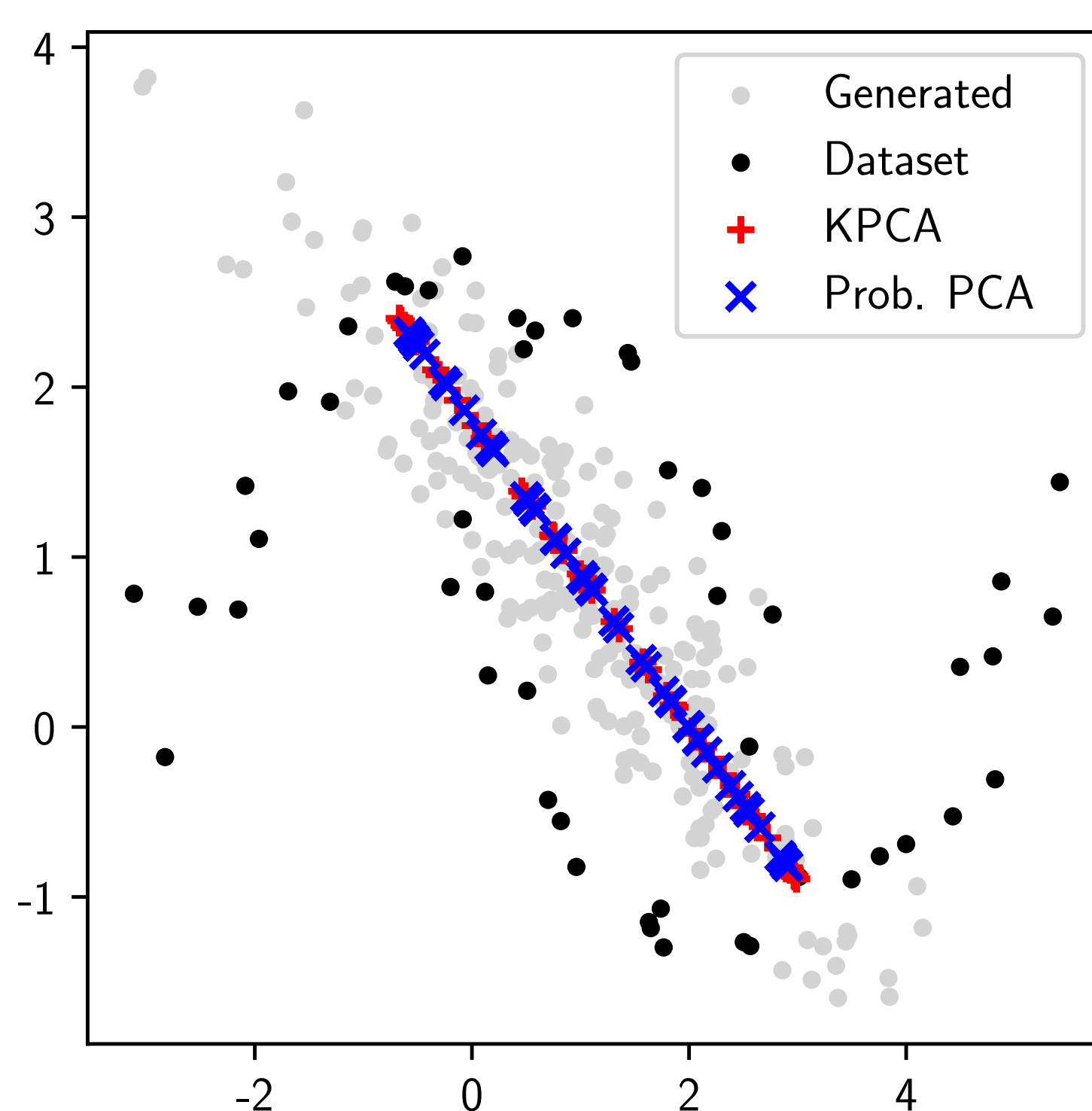
As a consequence of the choice of  $\mathcal{H}_{\mathcal{E}}$  as  $\text{span}\{\varphi_1, \dots, \varphi_N\}$ , we have

$$W = \Phi \circ A, \quad \text{with } A : \mathcal{L} \rightarrow \mathcal{L},$$

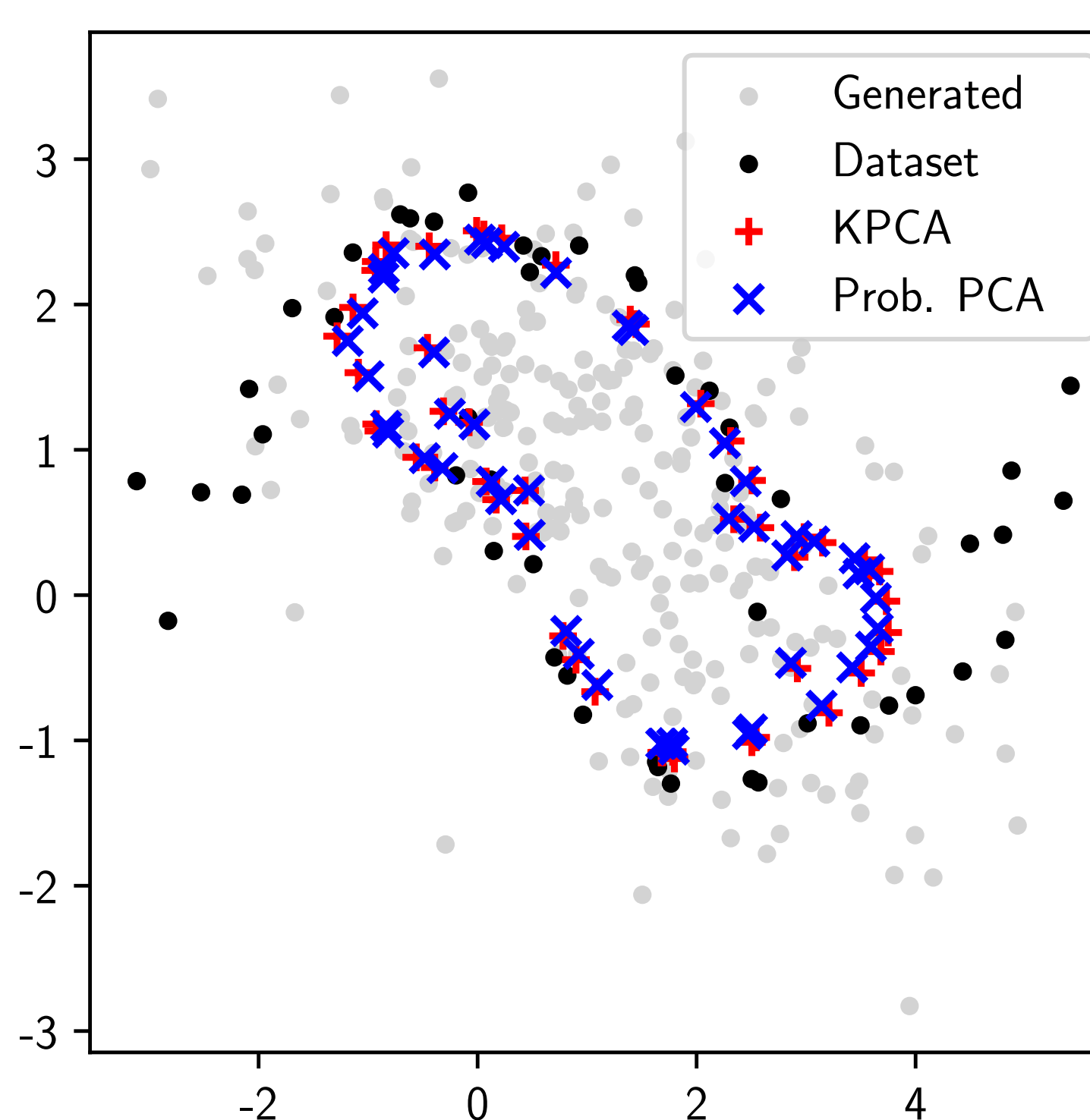
and in particular, as  $\Phi^* \circ \Phi$  and  $\Phi \circ \Phi^*$  share the same spectrum:

$$A_{ML} = \sum_{p=1}^q \sqrt{1/N - \sigma_{ML}^2 \lambda_p^{-1}} \epsilon_p r_p^*,$$

with  $\{(\lambda_p, \epsilon_p)\}_{p=1}^q$  the  $q$  dominant eigenpairs of  $\Phi^* \circ \Phi$ .



(a) With  $q = 1$  component, the explained variance is 31.23% and  $\sigma_{ML}^2 = 1.40\%$ .



(b) With  $q = 3$  components, the explained variance is 54.03% and  $\sigma_{ML}^2 = 0.98\%$ .

Figure: Visualisation of the Probabilistic PCA reconstruction (in blue) the classical KPCA (in red). Samples generated by are also given (in grey). The dataset contains  $N = 20$  points (in black).

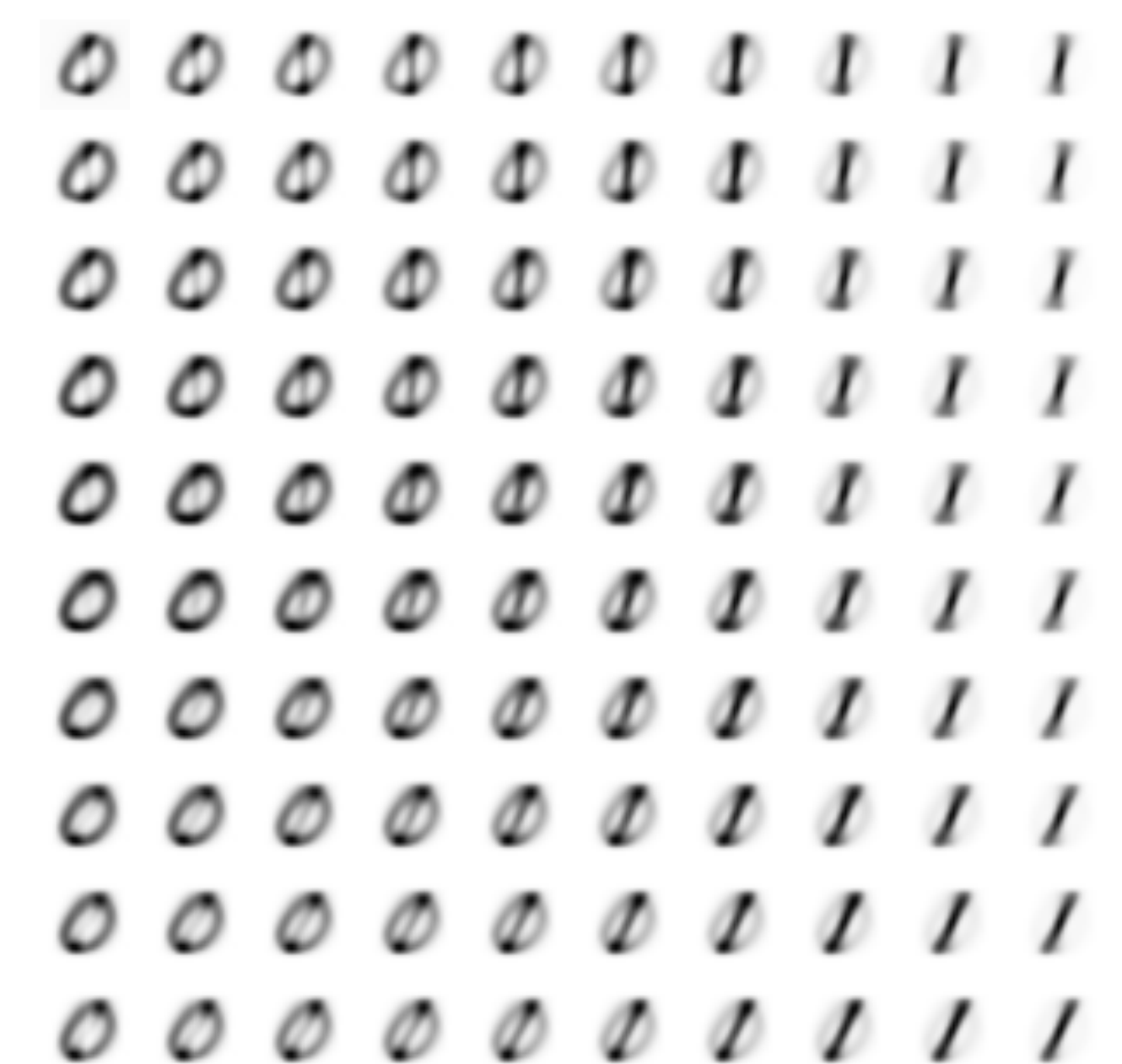


Figure: Generated datapoints on the MNIST dataset restricted to 0's and 1's, with  $N = 500$  datapoints, with  $q = 2$  components. The sample  $u$  is uniform on  $[-1, 1]$  for the two first components and zero for the others. The horizontal axis varies in the first component and the vertical one in the second component. The explained variance is 27.97% and  $\sigma_{ML}^2 = 0.14\%$ .