

Modular Dimensionality Reduction [★]

Henry W J Reeve¹, Tingting Mu², and Gavin Brown²

¹ University of Birmingham, Edgbaston, Birmingham,
B15 2TT, United Kingdom

² University of Manchester, Oxford Rd, Manchester,
M13 9PL, United Kingdom

Abstract. We introduce an approach to modular dimensionality reduction, allowing efficient learning of multiple complementary representations of the same object. Modules are trained by optimising an unsupervised cost function which balances two competing goals: Maintaining the inner product structure within the original space, and encouraging structural diversity between complementary representations. We derive an efficient learning algorithm which outperforms gradient based approaches without the need to choose a learning rate. We also demonstrate an intriguing connection with Dropout. Empirical results demonstrate the efficacy of the method for image retrieval and classification.

Keywords: Ensemble learning · Dimensionality reduction · Dropout · Kernel principal components analysis.

1 Introduction

High dimensional data is a widespread challenge in machine learning applications, from computer vision through to bioinformatics and natural language processing. A natural solution is to find a structure-preserving mapping to a low dimensional space, many techniques for which can be found in the literature, such as kernel PCA, Isomap, LLE and Laplacian Eigenmaps [6, 23]. This paper provides a meta-level tool for modular dimensionality reduction, applicable to each of the aforementioned approaches.

We start from the observation that *multiple* abstractions of the same concept can be taken, and may provide complementary views on a task of interest. We therefore propose a *modular* approach to unsupervised dimensionality reduction, in which we learn a diverse collection of low-dimensional representations of the data. Once a modular representation is learned, each module may be used independently – with their respective predictions combined at test time. This procedure is naturally parallelisable in a distributed computing architecture; and, since each representation is low-dimensional, processing for each module is fast and efficient.

[★] Acknowledgments: H. Reeve was supported by the EPSRC through the Centre for Doctoral Training Grant [EP/1038099/1]. G. Brown was supported by the EPSRC LAMBDA project [EP/N035127/1].

In the context of supervised learning, successful ensemble performance emanates from a fruitful trade-off between the accuracy of the individual members of the ensemble and the degree of diversity [15],[4]. We carry this insight across to the domain of unsupervised dimensionality reduction, by demonstrating the importance of diversity for a set of representation modules. We introduce an unsupervised loss function for training a *set* of dimensionality reduction modules, which balances two competing objectives. The first objective is for each module to preserve relational structure within the original feature space; the second is for modules to exhibit a *diversity* of relational structures.

The contributions of this paper are as follows:

1. An unsupervised loss function for modular dimensionality reduction.
2. A bespoke optimisation procedure which outperforms gradient based methods such as stochastic gradient descent in our setting.
3. A detailed empirical comparison with competitors.
4. An intriguing connection to the dropout algorithm from deep learning [13].

2 Background

We first review work on dimensionality reduction and ensemble learning.

2.1 Unsupervised Dimensionality Reduction

The canonical approach to unsupervised dimensionality reduction is PCA, and its kernelised generalisation, KPCA [21]. KPCA is a general approach which may be applied to a wide variety of application domains through an appropriate choice of kernel [19,25]. Several manifold learning techniques have also been shown to be special cases of KPCA, with a data-dependent kernel [12].

Classically, KPCA has been viewed as the orthogonal projection which maximises the preserved variance [21]. We shall adopt an alternative perspective in which we view KPCA as a form of *unsupervised similarity learning*, whereby a mapping is chosen so that inner-products in the low dimensional space approximate the kernel. To make this precise we require some notation. We let X denote our original feature space and let $k : X \times X \rightarrow \mathbb{R}$ denote a symmetric positive semi-definite kernel function. Take an unlabelled dataset $D = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset X$. For simplicity, we assume throughout that k is centred with respect to D [21]. Let H_k denote the associated reproducing kernel Hilbert space (RKHS) of real-valued functions. For each $H \in \mathbb{N}$, we let H_k^H denote the class of H -dimensional mappings $\varphi : X \rightarrow \mathbb{R}^H$ with coordinate functions taken from the RKHS H_k . That is, for each $\varphi \in H_k^H$ there exists $\varphi^1, \dots, \varphi^H \in H_k$ such that for all $x \in X$, $\varphi(x) = (\varphi^h(x))_{h=1}^H$.

Definition 1 (Inner product loss function). *Given an unsupervised data set D and a mapping $\varphi \in H_k^H$, the inner product loss is given by*

$$L_k(\varphi, D) = \frac{1}{N^2} \sum_{i,j \in N} (\langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j))^2.$$

We can interpret KPCA as a form of unsupervised similarity learning which minimises the inner product loss. Let $\xi : X \rightarrow H_k$ denote the canonical embedding given by $\xi(x)(y) = k(x, y)$.

Proposition 1. *The inner product loss $L_k(\varphi, D)$ is minimised by taking φ to be the member of H_k^H obtained by embedding D into H_k via ξ and projecting onto the top H kernel principal components.*

The proofs of all results within the main text are given in the appendices.

2.2 Ensembles and Diversity

Combining the outputs of multiple predictors often brings both statistical advantages, such as bias or variance reduction, and computational advantages, through parallelism. In order to outperform an individual model, ensembles promote a level of diversity or disagreement between the predictions the constituent models [10, 15]. Whilst methods such as bagging and boosting encourage diversity through a manipulation of the training data, a more direct approach is the *Negative Correlation Learning* (NCL) algorithm of Liu and Yao [18] in which diversity is targeted explicitly.

Suppose we have a supervised regression ensemble $H = fh_m g_{m=1}^M$ consisting of predictors h_m . In the previous section we used an unsupervised dataset $D = f\mathbf{x}_1, \dots, \mathbf{x}_N g$. To distinguish this we use notation $T = f(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) g$ for a *supervised* dataset. We let $\mathbb{V}(\cdot)$ denote the empirical variance of a finite sequence. The NCL algorithm can be understood in terms of the following *modular loss function*.

Definition 2 (Modular loss function). *The modular loss E_λ is defined by*

$$E_\lambda(H, T) = \frac{1}{NM} \sum_{n=1}^N \sum_{m=1}^M (h_m(\mathbf{x}_n) - y_n)^2 + \lambda \frac{1}{N} \sum_{n=1}^N \mathbb{V} \left(fh_m(\mathbf{x}_n) g_{m=1}^M \right).$$

The modular loss function consists of two terms: A squared loss term which targets the average individual accuracy of the predictors h_m , combined with a diversity term which encourages disagreement between the predictors. The hyper-parameter λ controls the degree of emphasis placed on the diversity. This has the special property that when $\lambda = 1$, $E_\lambda(H, T)$ is exactly the squared loss for the ensemble predictor $\frac{1}{M} \sum_m h_m(\mathbf{x})$ from the target y .

The NCL algorithm is equivalent to stochastic gradient descent applied to the modular loss. This perspective differs from original formulation of the NCL algorithm first introduced by Liu and Yao which utilises a multiplicity of interacting cost functions [18]. However, the updates of the two formulations are equal up to a factor of $1/M$ applied to the learning rate.

3 The Modular Inner Product Loss

Our goal is to train a collection of M distinct but complementary representations of the data. With this goal in mind, we introduce the modular inner product loss which combines two contrasting objectives. On the one hand, we seek high quality representations which faithfully preserve the relational structure encoded by the kernel. On the other hand, we would like the relational structure encoded in our different representations to be diverse. Let $F(H, M)$ denote the class of all M -tuples $\Phi = f\varphi_m g_{m=1}^M$ with each $\varphi_m \in \mathbb{H}_k^H$. Recall that $\mathbb{V}(\cdot)$ denotes the empirical variance.

Definition 3 (The modular inner product loss). *Suppose we have an unlabelled data set $D \subseteq X$ and a kernel k . Given $\Phi \in F(H, M)$, the modular inner product loss is given by*

$$L_k^\lambda(\Phi, D) = \frac{1}{M} \sum_{m=1}^M L_k(\varphi_m, D) + \lambda \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{V}\left(f\varphi_m(\mathbf{x}_i) - \varphi_m(\mathbf{x}_j) g_{m=1}^M\right). \quad (1)$$

The modular inner product loss is an analogue of the supervised modular loss function (Definition 2), with inner products between a pair of examples in a representation module replacing predictions for a single example, and the target replaced by an unsupervised inner product.

An equivalent reformulation of the modular inner product loss is as a convex combination between the average inner product loss of the individual modules and the inner product loss of a composite representation. Given $\Phi \in F(H, M)$ we define $\bar{\Phi} \in (\mathbb{H}_k)^{H \times M}$ by $\bar{\Phi}(\mathbf{x}) = \left(1/\sqrt{M}\right) [\varphi_1(\mathbf{x})^T, \dots, \varphi_M(\mathbf{x})^T]^T$. Proposition 2 is proved in Appendix 9.

Proposition 2. $L_k^\lambda(\Phi, D) = (1 - \lambda) \frac{1}{M} \sum_{m=1}^M L_k(\varphi_m, D) + \lambda L_k(\bar{\Phi}, D)$.

When $\lambda = 0$ the loss $L_k^\lambda(\Phi, D)$ is minimised by taking each φ_m to be a projection onto the top H kernel principal components, whilst for $\lambda = 1$, $L_k^\lambda(\Phi, D)$ is minimised by taking $\bar{\Phi}$ to be the projection onto the top M/H kernel principal components. Hence, $L_k^\lambda(\Phi, D)$ blends smoothly between training representation modules as individuals and targeting the composite representation.

4 Efficient Optimization

We now introduce the *module-by-module* (MBM) algorithm, which is a form of alternating optimisation designed to minimise the modular inner product loss without the need to choose a learning rate. Our objective is to minimise $L_\lambda(\Phi, D)$ over $\Phi \in F(H, M)$. We require an empirical kernel map.

Definition 4. *A rank R empirical kernel map is a function $\psi \in \mathbb{H}_k^R$ such that $\psi(\mathbf{x}_i)^T \psi(\mathbf{x}_j) = k(\mathbf{x}_i, \mathbf{x}_j)$ for all pairs $(\mathbf{x}_i, \mathbf{x}_j) \in D^2$.*

One can always construct an empirical kernel map of rank N by taking $\psi(\mathbf{x}) = \mathbf{K}(D)^{\frac{1}{2}} [k(\mathbf{x}, \mathbf{x}_1), \dots, k(\mathbf{x}, \mathbf{x}_N)]^T$, where $\mathbf{K}(D) = (k(\mathbf{x}_i, \mathbf{x}_j))_{i,j}$ denotes the kernel gram matrix. Moreover, given a kernel k we can often obtain a low rank empirical kernel map ψ for a kernel \tilde{k} which closely approximates k by employing a method such as random Fourier features [11] or the Nyström method [27]. By reasoning analogous to [20] we have the following useful proposition.

Proposition 3. *Given a rank R empirical kernel map ψ , the minimum for $L_k^\lambda(\Phi, D)$ is attained by $\Phi = \sum_{m=1}^M \varphi_m \mathcal{G}_{m=1}^M$ with each φ_m of the form $\varphi_m(\mathbf{x}) = \mathbf{W}_m \psi(\mathbf{x})$ for some matrix $\mathbf{W}_m \in \mathbb{R}^{H \times R}$.*

Hence, our objective reduces to the following matrix optimisation problem:
Minimise

$$C^\lambda(W, \Psi) := \sum_{m=1}^M \left\| \mathbf{F}_m^T \mathbf{F}_m - \Psi^T \Psi \right\|^2 + \lambda \sum_{m=1}^M \left\| \mathbf{F}_m^T \mathbf{F}_m - \frac{1}{M} \sum_{q=1}^M \mathbf{F}_q^T \mathbf{F}_q \right\|^2$$

$$\quad / L_k^\lambda(\Phi_{\mathcal{W}}, D),$$

where $\Psi = [\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_n)] \in \mathbb{R}^{n \times N}$, $\mathbf{F}_m = \mathbf{W}_m \Psi$ and $\Phi_{\mathcal{W}} = \sum_{m=1}^M \varphi_m \mathcal{G}_{m=1}^M$. We make use the concept of the rank-constrained approximate square root of a symmetric matrix.

Definition 5. Define $RT_r : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{r \times d}$ by

$$RT_r(\mathbf{M}) = \operatorname{argmin}_{\mathbf{F} \in \mathbb{R}^{r \times d}} \left\{ \left\| \mathbf{F}^T \mathbf{F} - \mathbf{M} \right\|^2 \right\}.$$

Dax has shown that the rank-constrained approximate square root $RT_r(\mathbf{M})$ of any $d \times d$ symmetric matrix \mathbf{M} (not necessarily positive semi-definite) may be computed via the singular value decomposition in $O(d^2 - r)$ time and $O(d^2)$ space complexity [7]. The following proposition allows us to optimise the weights of a single module φ_m , whilst leaving the remaining modules fixed.

Proposition 4. *Suppose we take $m \geq 1$, \dots, M , fix \mathbf{W}_q for $q \neq m$, and let*

$$\mathbf{T}_m = \frac{M}{M - \lambda + (M - 1)} \Psi^T \left(\mathbf{I}_D - \frac{\lambda}{M} \sum_{q \neq m} \mathbf{W}_q^T \mathbf{W}_q \right) \Psi.$$

Take $\mathbf{F}_m = RT_H(\mathbf{T}_m)$. Setting $\mathbf{W}_m = \mathbf{F}_m \Psi^\dagger$ minimises $C^\lambda(W, \Psi)$ with respect to \mathbf{W}_m , under the constraint that \mathbf{W}_q remains fixed for $q \neq m$, where Ψ^\dagger denotes the pseudo-inverse of Ψ .

Unfortunately, computing \mathbf{F}_m via Proposition 4 is $O(N^2 - H)$ which is intractable for large N . The following proposition enables us to reduce the complexity of this optimisation whenever we have access to an empirical kernel map of rank $R \ll N$.

Proposition 5. *Suppose that ψ is an empirical kernel map of rank R . Take $\tilde{\Psi} = (RT_R(\Psi\Psi^T))^T \in \mathbb{R}^{R \times R}$. For all $W = \sum_{m=1}^M \mathbf{W}_m \mathcal{G}_{m=1}^M$ with $\mathbf{W}_m \in \mathbb{R}^{H \times R}$ we have $C_\lambda(W, \tilde{\Psi}) = C_\lambda(W, \Psi)$. Moreover, computing $\tilde{\Psi}$ is $O(R^2 N)$ in time complexity and $O(R^2)$ in space complexity.*

Combining Propositions 3, 4 and 5 gives rise to the *module-by-module* algorithm (MBM, Algorithm 1), which is $O(NR^2 + EHR^2)$ in time and $O(NR)$ in space complexity, and has the advantage of reducing the modular inner product loss at every iteration until a critical point is reached.

Inputs: A data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, a rank R empirical kernel map ψ , a number of modules M , a number of dimensions per module H , a diversity parameter λ and $\epsilon > 0$.
 Compute $\Psi = [\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_N)]$;
 Update $\tilde{\Psi} = (RT_\rho(\Psi\Psi^T))^T$;
 Randomly initialise $\mathbf{F}_m \in \mathbb{R}^{H \times R}$ for $m = 1, \dots, M$;
 Compute $\mathbf{Q} = \tilde{\Psi}^T \tilde{\Psi}$ and $\mathbf{S} = \sum_{m=1}^M \mathbf{F}_m^T \mathbf{F}_m$;
 Compute $c = ((1 - \lambda) + \lambda/M + \epsilon)^{-1}$;
for $e = 1, \dots, E$ **do**
 for $m = 1, \dots, M$ **do**
 Compute $\mathbf{S}_m = \mathbf{S} - \mathbf{F}_m^T \mathbf{F}_m$;
 Compute $\mathbf{T} = c \cdot (\mathbf{Q} - (\lambda/M) \mathbf{S}_m + \epsilon \cdot \mathbf{F}_m^T \mathbf{F}_m)$;
 Update $\mathbf{F}_m = RT_H(\mathbf{T})$;
 Update $\mathbf{S} = \mathbf{S}_m + \mathbf{F}_m^T \mathbf{F}_m$;
 end
end
 Compute $\mathbf{W}_m = \mathbf{F}_m \Psi^y$ for $m = 1, \dots, M$;
Output: $\Phi = \{\mathbf{W}_m \cdot \psi\}_{m=1}^M$.

Algorithm 1: The module-by-module (MBM) algorithm.

The following theorem justifies the use of the MBM algorithm - it is guaranteed to reduce the modular inner product loss at every epoch until a critical point is reached.

Theorem 1. *Given $E \in \mathbb{N}$, let $\Phi^E \in F(H, M)$ denote the set obtained by training with Algorithm 1, for E epochs. Then for all $E \in \mathbb{N}$, $L_k^\lambda(\Phi^{E+1}, D) < L_k^\lambda(\Phi^E, D)$, unless Φ^E is a critical point of $L_k^\lambda(\Phi, D)$, in which case $L_k^\lambda(\Phi^{E+1}, D) = L_k^\lambda(\Phi^E, D)$.*

5 The Dropout Connection

In this section we introduce a surprising connection between the modular inner product loss and the dropout algorithm [13, 22]. Dropout is a state of the art approach to regularising deep neural networks in which a random collection of hidden neurons is “dropped out” at each stochastic gradient update. The dropout algorithm can be understood as implicitly minimising the expectation of a stochastic loss function based on predictions from a random sub-network [22, 26]. There is a natural analogue of this, in our setting: to minimise the expectation of a stochastic variant of the inner product loss, based on inner products computed from a random subset of modules. We refer to this analogue as the drop-module (DM) algorithm. To be precise, given an ensemble of feature mappings $\Phi \in F(H, M)$ each binary vector $\eta = (\eta_1, \dots, \eta_M) \in \{0, 1\}^M$ corresponds to a ‘noisy’ representation Φ_η given by $\Phi_\eta(\mathbf{x}) = (1/\sqrt{M}) (\eta_1 \varphi_1(\mathbf{x}), \dots, \eta_M \varphi_M(\mathbf{x}))$.

Fix a probability $p \in [0, 1]$ and let $B(p)$ denote the probability measure on $\{0, 1\}^M$ with $\mathbb{E}_{B(p)}(\eta) = p$. Let Θ denote the parameters of Φ . The DM algorithm proceeds by randomly sampling $\mathbf{x}_i, \mathbf{x}_j \in D$ and $\eta_m \in B(p)$ and updating

$$\Theta \leftarrow \Theta + \alpha \frac{\partial}{\partial \Theta} (\langle \Phi_\eta(\mathbf{x}_i), \Phi_\eta(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j))^2.$$

The DM algorithm implicitly minimises the following stochastic loss function

$$L_{k,p}^{\text{drop}}(\Phi, D) = \mathbb{E}_{\eta \in B(p)^M} [L_k(\Phi_\eta, D)].$$

Previously, Baldi et al. demonstrated that dropout may be understood as training an exponentially large ensemble with shared weights [1]. In our setting this corresponds to shared weight a ensemble of size 2^M with an ensemble member for each $\eta \in \{0, 1\}^M$. We demonstrate that the DM algorithm can be related to an ensemble of size M , trained via the modular inner product loss (see Definition 3). We emphasise that unlike the shared weight ensembles considered by Baldi et al. [1], here we consider ensembles with separate weights in which the interaction takes place purely via the diversity term in the modular inner product loss.

Theorem 2. *The drop-module inner product loss at p is equivalent to the modular inner product loss at $\lambda = Mp/(1+p(M-1))$. To be precise, for $\Phi_0 \in F(H, M)$ we have*

$$\left. \frac{\partial L_{k,p}^{\text{drop}}(\Phi, D)}{\partial \Phi} \right|_{\Phi=\Phi_0} = p \left. \frac{\partial L_k^\lambda(\Phi, D)}{\partial \Phi} \right|_{\Phi=(\frac{p}{1+p(M-1)})\Phi_0}.$$

Theorem 2 implies that if we take $\lambda = Mp/(1+p(M-1))$ then the minima of $L_k^\lambda(\Phi, D)$ are equal to the minima of $L_{k,p}^{\text{drop}}(\Phi, D)$, up to a constant scaling factor of $\sqrt{p/\lambda}$. In this sense, the minima for the two loss functions are *representationally equivalent*. The relationship between the diversity parameter λ in the MBM algorithm and the probability of keeping module p in the DM algorithm is illustrated in Figure 1.

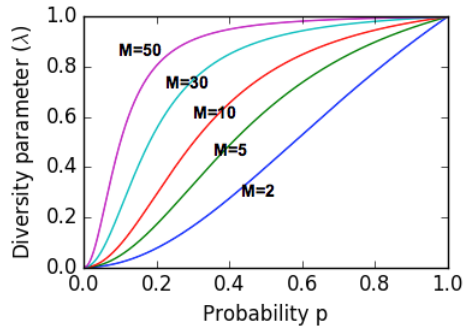


Fig. 1. The diversity parameter λ in MBM vs. the corresponding p in DM.

6 Experimental Results

In this section we first demonstrate the optimisation performance of the MBM algorithm before comparing our method for other natural approaches for training multiple kernelised representations. The data sets used in all experiments are described in Section 12.1.

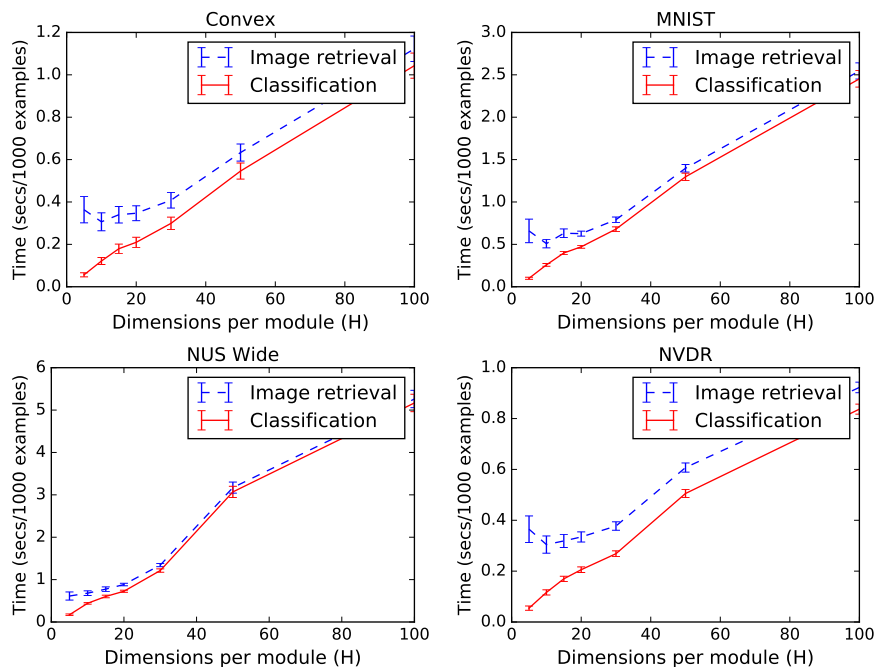
6.1 Optimisation performance of the MBM algorithm

In this section we assess the MBM algorithm (Algorithm 1) in terms of its efficiency at optimising the modular inner product loss. We compare with two gradient based approaches. As a baseline we consider stochastic gradient descent (SGD) applied directly to the modular inner product loss. This is expected to perform poorly in our setting since the modular inner product loss sums over all pairs of examples. We also consider the state of the art Adam optimiser of Kingma et al. [14] (Adam). The Adam optimiser is applied in batch mode, first compressing the data by applying Proposition 5. We set $M = 10$, $H = 10$, $\lambda = 0.9$ and let k to be the Gaussian kernel with γ set using the heuristic of Kwok and Tsang [16, Section 4]. In addition, we employ a rank $R = 1000$ Nyström approximation [27]. For SGD and Adam we consider learning rates in the range $\ell 10^{-6}$, $\ell 10^{-5}$, $\ell 10^{-4}$, $\ell 10^{-3}$, $\ell 10^{-2}$, $\ell 10^{-1}$, $\ell 10^0$, $\ell 10^1$, $\ell 10^2$. We evaluated the algorithms by training for one hour and evaluating the minimum value of the loss function attained during training and the convergence time - the time taken for the loss function to fall within 1% of its minimum. For SGD and Adam we report results corresponding to the learning rate which achieves the lowest minimum loss. The results are shown in Table 1. The SGD method was extremely slow and typically failed to converge within one hour. The compressed Adam method performed much better and typically converged within 30 minutes. However, the bespoke MBM algorithm achieved the same minimum loss in at most twice the speed on each of the data sets. The MBM algorithm also has the advantage of not requiring the user to set a learning rate.

Table 1. A comparison of the MBM algorithm with gradient based methods.

Data set	Loss function minimum			Convergence time (seconds)		
	SGD	Adam	MBM	SGD	Adam	MBM
Convex	110.0±21.2	13.4±0.0	13.3±0.0	3369.8±163.8	963.9±368.9	334.7±151.2
MNIST	427.6±16.3	59.9±0.0	59.9±0.0	3484.7±63.2	1609.1±69.2	400.6±69.4
NUS Wide	598.6±8.4	73.1±0.0	73.1±0.0	3376.3±185.2	1495.1±149.5	422.7±47.7
NVDR	110.0±18.5	10.7±0.0	10.6±0.0	1891.3±578.8	1157.2±227.0	402.4±35.1
Rectangles	29.4±0.4	10.3±0.0	10.3±0.0	1796.1±687.6	746.9±162.0	256.2±36.5

6.2 Image retrieval & classification performance of MBM modules


Fig. 2. Test time as vs. H. See Appendix 12.2 for discussion & other data sets.

We compare four unsupervised approaches to training multiple kernelised feature mappings:

Partition We compute the top HM KPCAs and randomly partition these into M sets of H , so that each mapping $\varphi_m \in \mathbb{H}_k^H$ is a projection onto a disjoint subset of the top H/M KPCAs.

Bootstrap Bagging [3] applied to KPCA. For each $m \in \{1, \dots, M\}$ take a bootstrap sample \tilde{D}_m , of size N , and let $\varphi_m \in \mathbb{H}_k^H$ be the KPCA projection mapping onto H dimensions for \tilde{D}_m .

Random A kernelised variant of the widely used technique of random projections [2, 8]. For each $m = 1, \dots, M$ we sample a random matrix $\mathbf{R}_m \in \mathbb{R}^{H \times D}$ from an $H \times D$ standard normal distribution, and normalise each row so that it has unit norm. In order to kernelise this technique the feature space for the random matrices is the output of the empirical kernel map ψ (see Section 4).

MBM Our proposed approach in which Φ is trained to minimise the modular inner product loss via MBM algorithm. The diversity parameter λ is set based upon performance on a validation set. We shall consider $H \in \{5, 10, 15, 20, 30, 50, 100\}$ and take M so that $H \times M = 300$. We also compare with the following non-modular base line.

Monolithic A single mapping $\varphi \in \mathbb{H}_k^H$ - the projection onto the top 300 KPCAs.

In each case we take k to be the Gaussian kernel with the γ parameter set via the heuristic of Kwok and Tsang [16, Section 4]. For computational efficiency we employ a rank 1000 Nyström approximation [27] in each case. We shall consider two distinct tasks:

Image retrieval: We shall consider the modular low dimensional representation’s capability for efficiently retrieving a set of κ close-by images. Let $\Phi = \{\varphi_m\}_{m=1}^M \in F(H, M)$ be a modular representation and D an unlabelled training set. Given a test point $\mathbf{x} \in X$, for each module, we compute the set $I_{\kappa, n}^{\varphi_m}(\mathbf{x}) \subset D$ of κ -nearest neighbours of \mathbf{x} based the distance $k\varphi_m(\mathbf{x}_q) - \varphi_m(\mathbf{x})k_2$. We then extract subset of size κ , $I_{\kappa, n}^{\Phi}(\mathbf{x}) = \bigcup_{m=1}^M I_{\kappa, n}^{\varphi_m}(\mathbf{x})$ so that the elements $\mathbf{x}_q \in I_{\kappa, n}^{\Phi}(\mathbf{x})$ minimise the average squared distance from the test point \mathbf{x} over the low dimensional spaces ie. $(1/M) \sum_{m=1}^M k\varphi_m(\mathbf{x}_q) - \varphi_m(\mathbf{x})k_2^2$ is minimised. Let $I_{\kappa, n}(\mathbf{x})$ denote the set of κ nearest neighbours as computed in the original space X . To assess performance we compute the precision: the average value of $(1/\kappa) \#(I_{\kappa, n}^{\Phi}(\mathbf{x}) \cap I_{\kappa, n}(\mathbf{x}))$. This procedure is based upon the method of [24] and gives a quantitative assessment of the representation’s ability to preserve structural information. The results of the image retrieval task for $\kappa = 10$, $H = 20$ and $M = 15$ are shown in Table 2. On each of the eight data sets the precision attained by the MBM method significantly exceeds the precision attained by the other modular methods: partition, bootstrap and random. Table 3 compares MBM with the monolithic method in which we simply compute the 10 nearest neighbours in the φ -projected space, where φ is the projection onto the top 300 KPCAs. For a relatively modest reduction in performance the MBM method obtains a significant speed up at test time. The speed up is due to the fact that each set of nearest neighbours $I_{\kappa, n}^{\varphi_m}(\mathbf{x})$ may be computed in parallel on a low-dimensional space (see Figure 2 and Appendix 12.2). Figure 3 demonstrates the precision as a function of the number of dimensions per module (H)

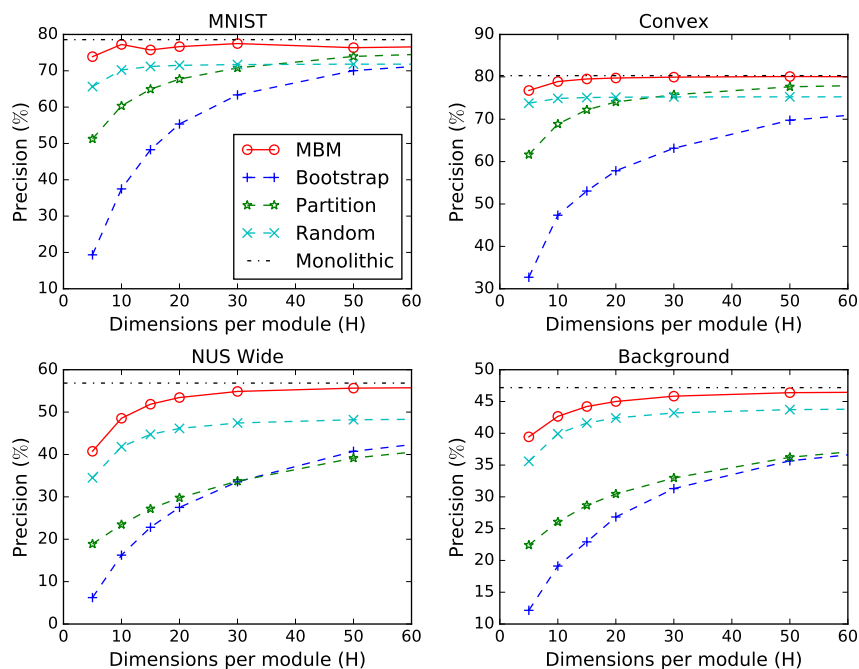


Fig. 3. Precision as a function of H (See Section 6.2).

with $\kappa = 10$ and $H = M = 300$. As H increases, the precision approaches the precision attained by the 300-dimensional Monolithic approach. The precision attained by the MBM approach typically exceeds that attained by the other modular approaches (Bootstrap, Partition, Random) across a range of values of H . Corresponding figures for other data sets are given in Appendix 12.3.

Classification We compare the methods in terms of their capacity for extracting multiple sets of features for use in a classification ensemble. Given a modular representation $\Phi = \{ \varphi_m \}_{m=1}^M \subseteq F(H, M)$, for each m we train a classifier f_m based on the features extracted by φ_m . Given a test point \mathbf{x} we combine the outputs of $f_m(\varphi_m(\mathbf{x}))$ by taking a modal average. Table 2 shows the classification accuracy for ensembles consisting of 15 5-nearest neighbour classifiers trained on 20-dimensional spaces. The MBM approach significantly outperforms the other approaches on five out of eight data sets, and performs comparably or better than the alternatives on every data set. Table 3 compares with the monolithic approach - a single 5-nearest neighbour classifier on 300 KPCAs. The MBM approach is both faster and more accurate than the monolithic method on all but one data set.

Table 2. A comparison of methods for modular dimensionality reduction. See Section 6.2 for details.

Data set	Image retrieval precision (%)				Ensemble classification accuracy (%)			
	Partition	Bootstrap	Random	MBM	Partition	Bootstrap	Random	MBM
NUS Wide	29.8±0.3	27.5±0.3	46.1±0.3	53.4±0.3	32.2±0.7	37.7±0.7	42.0±0.7	43.7±0.7
NVDR	68.0±0.5	48.1±0.5	73.4±0.5	77.6±0.4	44.9±1.0	39.6±1.0	43.9±1.0	44.8±1.0
MNIST	67.7±0.1	55.4±0.1	71.5±0.1	76.6±0.1	89.2±0.2	95.0±0.2	94.5±0.2	95.8±0.1
Background	30.5±0.2	26.8±0.2	42.4±0.2	45.0±0.2	41.0±0.4	47.3±0.4	53.8±0.4	57.3±0.4
Random	7.0±0.1	10.9±0.1	6.1±0.1	17.2±0.1	40.9±0.4	89.5±0.2	53.1±0.4	90.3±0.2
Rotations	59.6±0.2	54.3±0.1	68.4±0.1	75.3±0.1	72.0±0.3	85.5±0.3	83.8±0.3	87.4±0.2
Convex	74.1±0.1	57.8±0.2	75.2±0.1	79.7±0.1	55.7±0.4	65.5±0.3	60.3±0.4	65.6±0.3
Rectangles	63.3±0.1	44.6±0.1	46.7±0.1	71.6±0.1	93.0±0.2	95.5±0.2	93.9±0.2	98.1±0.1

Table 3. Comparing the MBM & the Monolithic approach (See Section 6.2).

Data set	Image retrieval				Ensemble classification			
	Monolithic	MBM Δ	Speed	λ	Monolithic	MBM Δ	Speed	λ
NUS Wide	56.8±0.3	-3.4±0.5	10.8×	0.999	40.6±0.7	+3.1±1.0	11.7×	0.990
NVDR	79.0±0.5	-1.4±0.6	7.7×	0.999	40.8±1.0	+3.9±1.4	12.6×	0.999
MNIST	78.6±0.1	-1.9±0.1	9.1×	0.990	95.7±0.1	+0.1±0.2	12.9×	0.900
Background	47.2±0.2	-2.2±0.2	5.8×	0.999	52.5±0.4	+4.9±0.5	8.1×	0.999
Random	23.1±0.1	-5.9±0.1	5.7×	0.950	83.9±0.3	+6.4±0.3	7.1×	0.500
Rotations	76.5±0.1	-1.2±0.1	7.9×	0.990	86.3±0.3	+1.0±0.4	10.9×	0.800
Convex	80.3±0.1	-0.6±0.1	7.2×	0.999	57.6±0.4	+8.0±0.5	23.0×	0.200
Rectangles	70.1±0.1	+1.5±0.2	2.4×	0.990	95.8±0.1	-2.3±0.2	10.7×	0.990

The two Monolithic columns show the image retrieval precision (%) and the classification accuracy (%) of the monolithic method. The MBM columns show the corresponding change in performance due to using the MBM method for each task. The Speed columns show the corresponding speed ups ie. the test time for the Monolithic method divided by the test time for the MBM method.

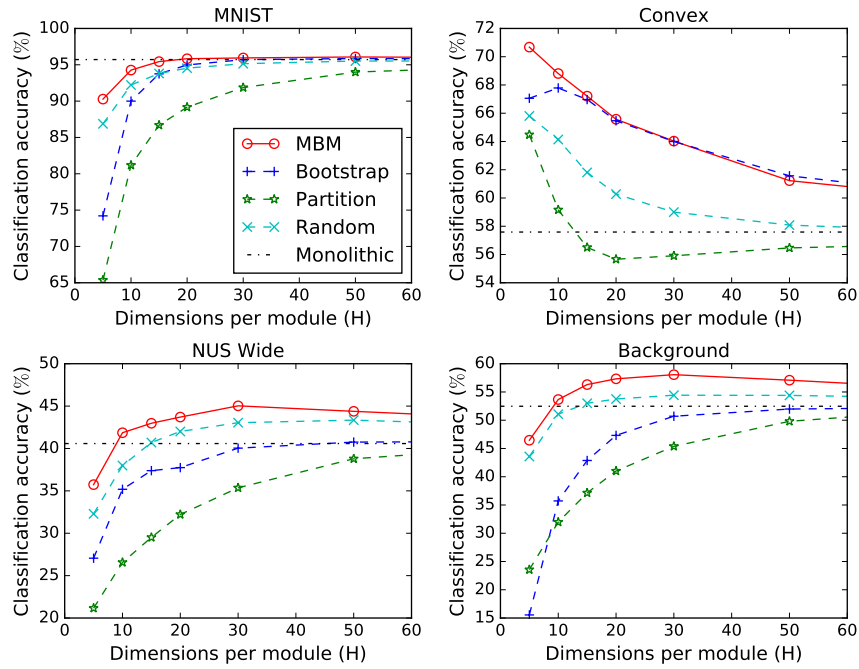


Fig. 4. Classification accuracy as a function of H (See Section 6.2).

The diversity parameter The diversity parameter λ in the MBM algorithm controls the level of emphasis placed upon encouraging a diversity of representations. We found that the optimal performance (both in terms of information retrieval and classification) was typically attained with λ just below 1, with performance declining sharply by taking $\lambda = 1$ (see Figure 5, cols 1& 2). It is interesting to observe that the dropout algorithm often performs well with $p = 0.5$ and this corresponds a value of λ just below 1, when M is large (see Figure 1). However, whilst this pattern was observed on all data sets for image retrieval (see Appendix 12.6), for some data sets the best classification performance was attained by taking much lower values of λ (see Figure 5, col 3 and Appendix 12.6). Ultimately, the optimal value of λ is data dependent and must be set based on validation performance.

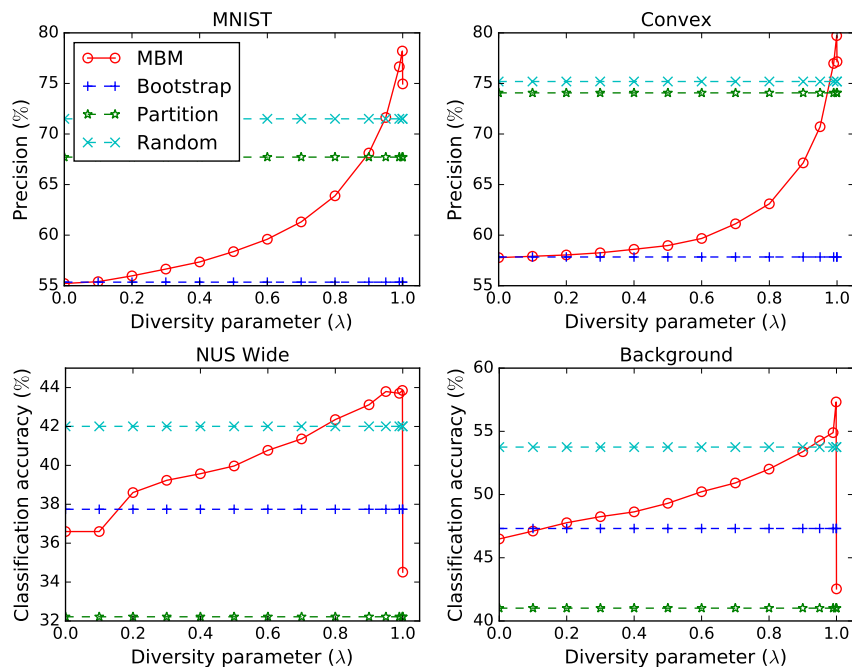


Fig. 5. Performance as a function of the diversity parameter (λ) (see Section 6.2).

7 Discussion

We have investigated a method for *modular* unsupervised dimensionality reduction. Our method is based upon the *modular inner product loss* (Definition 3), an adaptation of concepts from both negative correlation learning [4, 18] and kernel principal components analysis [21]. Whilst the modular loss could be optimised by gradient based methods we introduced a novel *module-by-module* algorithm, which converges at least twice as fast as a state of the art gradient based optimiser [14] without the need to tune the learning rate.

Modular representations have the potential to be applied on range of tasks. Empirical results on both image retrieval and classification tasks confirm that the MBM algorithm is superior to a range of competitors including random projections and bootstrapping, whilst providing a parallelisation advantage over “monolithic” dimensionality reduction. We also demonstrated an intriguing equivalency between our proposal and an analogue of the dropout algorithm - drop module, which deserves further attention.

In summary, this work has shown the potential of explicitly managing diversity in unsupervised representation learning.

References

1. Pierre Baldi and Peter J Sadowski. Understanding dropout. In *Advances in Neural Information Processing Systems*, pages 2814–2822, 2013.
2. Ella Bingham and Heikki Mannila. Random projection in dimensionality reduction: applications to image and text data. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 245–250. ACM, 2001.
3. Leo Breiman. Bagging predictors. *Machine learning*, 24(2):123–140, 1996.
4. Gavin Brown, Jeremy L Wyatt, and Peter Tiño. Managing diversity in regression ensembles. *The Journal of Machine Learning Research*, 6:1621–1650, 2005.
5. Tat-Seng Chua, Jinhui Tang, Richang Hong, Haojie Li, Zhiping Luo, and Yantao Zheng. Nus-wide: a real-world web image database from national university of singapore. In *Proceedings of the ACM international conference on image and video retrieval*, page 48. ACM, 2009.
6. John P Cunningham and Zoubin Ghahramani. Linear dimensionality reduction: Survey, insights, and generalizations. *Journal of Machine Learning Research*, 16:2859–2900, 2015.
7. Achiya Dax. Low-rank positive approximants of symmetric matrices. *Advances in Linear Algebra & Matrix Theory*, 4(03):172, 2014.
8. Robert J Durrant and Ata Kabán. Random projections as regularizers: Learning a linear discriminant ensemble from fewer observations than dimensions.
9. Carl Eckart and Gale Young. The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
10. Pascal Germain, Alexandre Lacasse, Francois Laviolette, Mario Marchand, and Jean-Francis Roy. Risk bounds for the majority vote: From a pac-bayesian analysis to a learning algorithm. *The Journal of Machine Learning Research*, 16(1):787–860, 2015.
11. Nathan Halko, Per-Gunnar Martinsson, and Joel A Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 53(2):217–288, 2011.
12. Jihun Ham, Daniel D Lee, Sebastian Mika, and Bernhard Schölkopf. A kernel view of the dimensionality reduction of manifolds. In *Proceedings of the twenty-first international conference on Machine learning*, page 47. ACM, 2004.
13. Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
14. Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
15. Anders Krogh, Jesper Vedelsby, et al. Neural network ensembles, cross validation, and active learning. *Advances in neural information processing systems*, pages 231–238, 1995.
16. James Tin-Yau Kwok and Ivor Wai-Hung Tsang. The pre-image problem in kernel methods. *Neural Networks, IEEE Transactions on*, 15(6):1517–1525, 2004.
17. Hugo Larochelle, Dumitru Erhan, Aaron Courville, James Bergstra, and Yoshua Bengio. An empirical evaluation of deep architectures on problems with many factors of variation. In *Proceedings of the 24th international conference on Machine learning*, pages 473–480. ACM, 2007.
18. Yong Liu and Xin Yao. Ensemble learning via negative correlation. *Neural Networks*, 12(10):1399–1404, 1999.

19. Huma Lodhi, Craig Saunders, John Shawe-Taylor, Nello Cristianini, and Chris Watkins. Text classification using string kernels. *The Journal of Machine Learning Research*, 2:419–444, 2002.
20. Bernhard Schölkopf, Ralf Herbrich, and Alex J Smola. A generalized representer theorem. In *International Conference on Computational Learning Theory*, pages 416–426. Springer, 2001.
21. Bernhard Schölkopf, Alexander Smola, and Klaus-Robert Müller. Kernel principal component analysis. In *Artificial Neural Networks/ICANN'97*, pages 583–588. Springer, 1997.
22. Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
23. Dmitry Storcheus, Afshin Rostamizadeh, and Sanjiv Kumar. A survey of modern questions and challenges in feature extraction. In *Proceedings of The 1st International Workshop on Feature Extraction: Modern Questions and Challenges, NIPS*, pages 1–18, 2015.
24. Jarkko Venna, Jaakko Peltonen, Kristian Nybo, Helena Aidos, and Samuel Kaski. Information retrieval perspective to nonlinear dimensionality reduction for data visualization. *Journal of Machine Learning Research*, 11(Feb):451–490, 2010.
25. S Vichy N Vishwanathan, Nicol N Schraudolph, Risi Kondor, and Karsten M Borgwardt. Graph kernels. *The Journal of Machine Learning Research*, 11:1201–1242, 2010.
26. Sida I Wang and Christopher D Manning. Fast dropout training.
27. Christopher Williams and Matthias Seeger. Using the nystrom method to speed up kernel machines. In *Proceedings of the 14th Annual Conference on Neural Information Processing Systems*, number EPFL-CONF-161322, pages 682–688, 2001.
28. Xiao Wu, Alexander G Hauptmann, and Chong-Wah Ngo. Practical elimination of near-duplicates from web video search. In *Proceedings of the 15th ACM international conference on Multimedia*, pages 218–227. ACM, 2007.

Appendix for Modular Dimensionality Reduction

8 Proof of Proposition 1

In this section we prove Proposition 1. We will require some notation. Recall that $\xi : X \rightarrow H_k$ denotes the canonical embedding given by $\xi(x)(y) = k(x, y)$. We also let $P_D^H : H_k \rightarrow \mathbb{R}^H$ denote corresponding Hotelling transform, i.e. the projection onto the top H kernel principal components embedded isometrically into \mathbb{R}^H . Given a pair of Hilbert spaces H^1, H^2 with $a^i \in H^i$ for $i \in \{1, 2\}$ we define a bounded linear operator $a^2 - a^1 : H^1 \rightarrow H^2$ by $(a^2 - a^1)(z) = \langle a^1, z \rangle a^2$ for all $z \in H^1$. Note that $(a^2 - a^1)^* = a^2 - a^1$, i.e. given $z_i \in H^i$ for $i \in \{1, 2\}$ we have $\langle (a^2 - a^1)(z_1), z_2 \rangle = \langle z_1, (a^1 - a^2)(z_2) \rangle = \langle a_1, z_1 \rangle \langle a_2, z_2 \rangle$. Finally, we let $\{e_1, \dots, e_N\}$ denote the canonical orthonormal basis for \mathbb{R}^N .

Proposition 1. *The inner product loss $L_k(\varphi, D)$ is minimised by taking φ to be the member of H_k^H obtained by embedding D into H_k via ξ and projecting onto the top H kernel principal components.*

Proof. We must show that the inner product loss $L_k(\varphi, D)$ is minimised over $\varphi \in X \rightarrow \mathbb{R}^H$ by taking $\varphi(x) = P_D^H(\xi(x))$. We shall define a data dependent operator $\Psi : \mathbb{R}^N \rightarrow H_k$ by $\Psi = \sum_{n=1}^N \xi(x_n) e_n$. Clearly Ψ has rank at most N , so by the singular value theorem for finite rank operators there exists an orthonormal basis $\{v_1, \dots, v_N\} \subset \mathbb{R}^N$, a set of orthonormal vectors $\{u_1, \dots, u_N\} \subset H_k$ and a set of singular values $\sigma_1, \dots, \sigma_N \in \mathbb{R}^+$ with $\sum_{j=1}^N \sigma_j^2 = \sum_{j=1}^N \sigma_{Nj}^2$ such that $\Psi = \sum_{q=1}^N \sigma_q (u_q - v_q)$. Now let C denote the covariance operator given by $C = (1/N) \sum_{n=1}^N \xi(x_n) \xi(x_n)$, where we use the fact that k is centred with respect to D . Since $\{e_1, \dots, e_N\}$ is orthonormal we have

$$\begin{aligned} N C &= \sum_{n=1}^N \xi(x_n) \xi(x_n) = \left(\sum_{n=1}^N \xi(x_n) e_n \right) \left(\sum_{n=1}^N e_n \xi(x_n) \right) = \Psi \Psi \\ &= \left(\sum_{q=1}^N \sigma_q (u_q - v_q) \right) \left(\sum_{\bar{q}=1}^N \sigma_{\bar{q}} (v_{\bar{q}} - u_{\bar{q}}) \right) = \sum_{q=1}^N \sigma_q^2 (u_q - u_q). \end{aligned}$$

Thus, u_1, \dots, u_H are the top H kernel principal components, i.e. the eigenvectors of C with maximal eigenvalues. Consequently, $P_D^H = \sum_{h=1}^H e_h u_h$. Let $A : H_k \rightarrow \tilde{H}$ be an arbitrary linear operator into a Hilbert space \tilde{H} . We have $A\Psi = \sum_{n=1}^N (A\xi(x_n) e_n)$, so $\Psi^T A A\Psi$ corresponds to an $N \times N$ matrix with entries

$$\begin{aligned} (\Psi^T A A\Psi)_{n_1, n_2} &= e_{n_1}^T (\Psi^T A A\Psi) e_{n_2} = \langle e_{n_1}, (A\Psi) (A\Psi) e_{n_2} \rangle \\ &= \langle (A\Psi) e_{n_1}, (A\Psi) e_{n_2} \rangle = \langle A\xi(x_{n_1}), A\xi(x_{n_2}) \rangle. \end{aligned}$$

In particular, taking A to be the identity $I_{N \times N}$ we have

$$(\Psi^T \Psi)_{n_1, n_2} = \langle \xi(x_{n_1}), \xi(x_{n_2}) \rangle = k(x_{n_1}, x_{n_2})$$

Hence, if let \mathbf{K} denote the matrix with entries $k(\mathbf{x}_i, \mathbf{x}_j)$, we have

$$\mathbf{K} = \Psi \Psi^T = \left(\sum_{q=1}^N \sigma_q \begin{pmatrix} \mathbf{v}_q & \mathbf{u}_q \end{pmatrix} \right) \left(\sum_{q=1}^N \sigma_q \begin{pmatrix} \mathbf{u}_q & \mathbf{v}_q \end{pmatrix} \right) = \sum_{q=1}^N \sigma_q^2 \begin{pmatrix} \mathbf{v}_q & \mathbf{v}_q \end{pmatrix}.$$

Thus, by the Eckart-Young theorem [9] the rank H matrix \mathbf{M} which minimises $\|\mathbf{K} - \mathbf{M}\|_F$ (where $\|\cdot\|_F$ denotes the Frobenius norm) is $\mathbf{M}_H = \sum_{q=1}^H \sigma_q^2 \begin{pmatrix} \mathbf{v}_q & \mathbf{v}_q \end{pmatrix}$. Given any matrix $\varphi : X \rightarrow \mathbb{R}^H$ we shall let $\mathbf{Z}(\varphi)$ denote the $H \times N$ matrix with entries $\varphi_h(x_n)$ for $h \in \{1, \dots, H\}$ and $n \in \{1, \dots, N\}$. It follows that $\mathbf{Z}(\varphi) \mathbf{Z}(\varphi)^T$ is the $N \times N$ matrices with entries $\langle \varphi(x_{n_1}), \varphi(x_{n_2}) \rangle$, so

$$\|\mathbf{Z}(\varphi) \mathbf{Z}(\varphi)^T - \mathbf{K}\|_F^2 = \sum_{i,j \in \{1, \dots, N\}} (\langle \varphi(x_i), \varphi(x_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j))^2 = N^2 L_k(\varphi, D).$$

Given any function $\varphi : X \rightarrow \mathbb{R}^H$, $\mathbf{Z}(\varphi) \mathbf{Z}(\varphi)^T$ is of rank at most H , so if we can choose φ so that $\mathbf{Z}(\varphi) \mathbf{Z}(\varphi)^T = \mathbf{M}_H$ then $L_k(\varphi, D)$ will have attained its minimum. Now consider $\varphi \in \mathbb{H}_k^H$ defined by $\varphi(x) = \mathbf{P}_D^H \xi(x)$ for $x \in X$. Given $n_1, n_2 \in \{1, \dots, N\}$ we have

$$(\mathbf{Z}(\varphi) \mathbf{Z}(\varphi)^T)_{n_1, n_2} = \langle \varphi(x_{n_1}), \varphi(x_{n_2}) \rangle = \langle \mathbf{P}_D^H \xi(x_{n_1}), \mathbf{P}_D^H \xi(x_{n_2}) \rangle = \left(\Psi \begin{pmatrix} \mathbf{P}_D^H & \mathbf{P}_D^H \Psi \end{pmatrix} \right)_{n_1, n_2}.$$

Hence, we have

$$\begin{aligned} \mathbf{Z}(\varphi) \mathbf{Z}(\varphi)^T &= \Psi \begin{pmatrix} \mathbf{P}_D^H & \mathbf{P}_D^H \Psi \end{pmatrix} \\ &= \Psi \begin{pmatrix} \sum_{q=1}^H \mathbf{u}_q & \mathbf{e}_q \end{pmatrix} \begin{pmatrix} \sum_{q=1}^H \mathbf{e}_q & \mathbf{u}_q \end{pmatrix} \Psi^T = \Psi \begin{pmatrix} \sum_{q=1}^H \mathbf{u}_q & \mathbf{u}_q \end{pmatrix} \Psi^T \\ &= \left(\sum_{q=1}^N \sigma_q \begin{pmatrix} \mathbf{v}_q & \mathbf{u}_q \end{pmatrix} \right) \begin{pmatrix} \sum_{q=1}^H \mathbf{u}_q & \mathbf{u}_q \end{pmatrix} \begin{pmatrix} \sum_{q=1}^N \sigma_q \begin{pmatrix} \mathbf{u}_q & \mathbf{v}_q \end{pmatrix} \end{pmatrix} \\ &= \sum_{q=1}^H \sigma_q^2 \begin{pmatrix} \mathbf{v}_q & \mathbf{v}_q \end{pmatrix} = \mathbf{M}_H. \end{aligned}$$

Hence, taking $\varphi = \mathbf{P}_D^H \xi$ minimises the inner product loss $L_k(\varphi, D)$.

9 Proof of Proposition 2

In this section we prove Proposition 2.

Proposition 2. $L_k^\lambda(\Phi, D) = (1 - \lambda) \frac{1}{M} \sum_{m=1}^M L_k(\varphi_m, D) + \lambda L_k(\bar{\Phi}, D)$.

Proof. It suffices to show that

$$L_k(\bar{\Phi}, D) = \frac{1}{M} \sum_{m=1}^M L_k(\varphi_m, D) - \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{V} \left(f\varphi_m(\mathbf{x}_i) - \varphi_m(\mathbf{x}_j) g_{m=1}^M \right).$$

This may be seen by subtracting

$$(1 - \lambda) \frac{1}{M} \sum_{m=1}^M L_k(\varphi_m, D) + \lambda L_k(\bar{\Phi}, D)$$

from

$$\frac{1}{M} \sum_{m=1}^M L_k(\varphi_m, D) - \lambda \frac{1}{N^2} \sum_{i,j=1}^N \mathbb{V} \left(f\varphi_m(\mathbf{x}_i) - \varphi_m(\mathbf{x}_j) g_{m=1}^M \right),$$

Thus, it suffices to show that for any $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{X}^2$,

$$\left(\mathbb{E}[\bar{\Phi}(\mathbf{x}_i), \bar{\Phi}(\mathbf{x}_j)] - k_c(\mathbf{x}_i, \mathbf{x}_j) \right)^2 = \frac{1}{M} \sum_{m=1}^M \left(\mathbb{E}[\varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j)] - k_c(\mathbf{x}_i, \mathbf{x}_j) \right)^2 \mathbb{V} \left(f\varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) g_{m=1}^M \right).$$

Note that by construction $\mathbb{E}[\bar{\Phi}(\mathbf{x}_i), \bar{\Phi}(\mathbf{x}_j)] = (1/M) \sum_{m=1}^M \mathbb{E}[\varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j)]$. Thus, the result follows from the observation that for any sequence $f a_m g_{m=1}^M \in \mathbb{R}$ and $b \in \mathbb{R}$ if we let $\bar{a} = (1/M) \sum_{m=1}^M a_m$ then we have

$$(\bar{a} - b)^2 = \frac{1}{M} \sum_{m=1}^M (a_m - b)^2 - \frac{1}{M} \sum_{m=1}^M (a_m - \bar{a})^2.$$

This decomposition is known as the ambiguity-decomposition in the ensemble literature and was observed by Krogh and Vedelsby in the context of ensembles [15].

10 Proof of Theorem 1 & Propositions 3, 4, 5

In this section we prove Theorem 1 which justifies the procedure given by Algorithm 1 (the MBM algorithm), along with several supporting propositions. Throughout this section $\| \cdot \|$ denotes the Frobenius norm. Recall that we assume the existence of an empirical kernel map $\psi \in \mathcal{H}_k^H$.

Proposition 3. *Given a rank R empirical kernel map ψ , the minimum for $L_k^\lambda(\Phi, D)$ is attained by $\Phi = \sum_{m=1}^M \varphi_m \mathcal{G}_{m=1}^M$ with each φ_m of the form $\varphi_m(\mathbf{x}) = \mathbf{W}_m^T \psi(\mathbf{x})$ for some matrix $\mathbf{W}_m \in \mathbb{R}^{H \times R}$.*

Recall that $\xi : X \rightarrow \mathcal{H}_k$ denote the canonical embedding given by $\xi(\mathbf{x})(y) = k(\mathbf{x}, y)$.

Lemma 1. *Given $f \in \mathcal{H}_k$ there exists $\mathbf{a} \in \mathbb{R}^R$ such that $\tilde{f} = \mathbf{a}^T \psi$ satisfies $\tilde{f}(x_n) = f(x_n)$ for $n = 1, \dots, N$.*

Proof. Write $f = \sum_{q=1}^N \beta_q \xi(\mathbf{x}_q) + g$ where $g \in \mathcal{H}_k$ is orthogonal to $\sum_{n=1}^N \xi(\mathbf{x}_n) \mathcal{G}_{n=1}^N$. By orthogonality, combined with the reproducing property, for each $n = 1, \dots, N$ we have $\langle g, \xi(\mathbf{x}_n) \rangle = 0$. Hence, $f(\mathbf{x}_n) = \sum_{q=1}^N \beta_q \xi(\mathbf{x}_q)(\mathbf{x}_n)$ for $n = 1, \dots, N$. Now let $\mathbf{a} := \sum_{q=1}^N \beta_q \psi(\mathbf{x}_q)$ and define $\tilde{f} := \mathbf{a}^T \psi$. Then, for each $n = 1, \dots, N$ we have

$$\begin{aligned} \tilde{f}(\mathbf{x}_n) &= \mathbf{a}^T \psi(\mathbf{x}_n) = \sum_{q=1}^N \beta_q \psi(\mathbf{x}_q)^T \psi(\mathbf{x}_n) \\ &= \sum_{q=1}^N \beta_q k(\mathbf{x}_q, \mathbf{x}_n) = \sum_{q=1}^N \beta_q \xi(\mathbf{x}_q)(\mathbf{x}_n) = f(\mathbf{x}_n). \end{aligned}$$

Proof (Proof of Proposition 3). Take $\Phi = \sum_{m=1}^M \varphi_m \mathcal{G}_{m=1}^M \in F(H, M)$. By Lemma 1, for each $m = 1, \dots, M$, there exists matrices $\mathbf{W}_m \in \mathbb{R}^{H \times R}$ such that $\tilde{\varphi}_m := \mathbf{W}_m^T \psi_m$ satisfies $\tilde{\varphi}_m(\mathbf{x}_n) = \varphi_m(\mathbf{x}_n)$ for each $n = 1, \dots, N$. It follows that $\tilde{\Phi} = \sum_{m=1}^M \tilde{\varphi}_m \mathcal{G}_{m=1}^M$ satisfies $\tilde{\Phi}(x_n) = \Phi(x_n)$ for all $\mathbf{x}_n \in D$, so $L_k^\lambda(\tilde{\Phi}, D) = L_k^\lambda(\Phi, D)$.

This leads to the objective of minimising

$$C^\lambda(W, \Psi) := \sum_{m=1}^M \left\| \mathbf{F}_m^T \mathbf{F}_m - \Psi^T \Psi \right\|^2 + \lambda \sum_{m=1}^M \left\| \mathbf{F}_m^T \mathbf{F}_m - \frac{1}{M} \sum_{q=1}^M \mathbf{F}_q^T \mathbf{F}_q \right\|^2 + L_k^\lambda(\Phi_{\mathcal{W}}, D),$$

where $\Psi = [\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_N)] \in \mathbb{R}^{R \times N}$, $\mathbf{F}_m = \mathbf{W}_m \Psi$ and $\Phi_{\mathcal{W}} = \sum_{m=1}^M \tilde{\varphi}_m \mathcal{G}_{m=1}^M$. Recall Definition 5.

Definition 5. *Define $RT_r : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{r \times d}$ by*

$$RT_r(\mathbf{M}) = \operatorname{argmin}_{\mathbf{F} \in \mathbb{R}^{r \times d}} \left\{ \left\| \mathbf{F}^T \mathbf{F} - \mathbf{M} \right\|^2 \right\}.$$

Dax has shown that when $\mathbf{M} \in \mathbb{R}^{d \times d}$ is symmetric, $RT_r(\mathbf{M}) \in \mathbb{R}^{r \times d}$ may be computed via the singular value decomposition for any $r \leq D$ [7]. We may take an eigen-decomposition of $\mathbf{M} = \sum_{i=1}^d \gamma_i \mathbf{u}_i \mathbf{u}_i^T$, where $\gamma_i \in \mathbb{R}$ are eigen-values and \mathbf{u}_i are unit eigen-vectors of \mathbf{M} . We can compute the approximate root as follows,

$$RT_r(\mathbf{M}) = \text{diag} \left(\sqrt{\max\{\gamma_1, 0\}}, \dots, \sqrt{\max\{\gamma_r, 0\}} \right) [\mathbf{u}_1, \dots, \mathbf{u}_r]^T. \quad (2)$$

Note that in the special case in which \mathbf{M} is positive semi-definite i.e. $\gamma_i \geq 0$ for all i , this result is the classical Eckart and Young Theorem [9].

We first prove Proposition 5.

Proposition 5. *Suppose that ψ is an empirical kernel map of rank R . Take $\tilde{\Psi} = (RT_R(\Psi\Psi^T))^T \in \mathbb{R}^{R \times R}$. For all $W = \frac{1}{M} \sum_{m=1}^M \mathbf{W}_m \mathbf{g}_{m=1}^M$ with $\mathbf{W}_m \in \mathbb{R}^{H \times R}$ we have $C_\lambda(W, \tilde{\Psi}) = C_\lambda(W, \Psi)$. Moreover, computing $\tilde{\Psi}$ is $O(R^2 \times N)$ in time complexity and $O(R^2)$ in space complexity.*

Proof (Proof of Proposition 5). Take a singular value decomposition of Ψ , $\Psi = \sum_{l=1}^{\rho} \sigma_l \mathbf{u}_l \mathbf{v}_l^T$, let $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_\rho] \in \mathbb{R}^{D \times \rho}$, $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_\rho) \in \mathbb{R}^{\rho \times \rho}$ and $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_\rho] \in \mathbb{R}^N \times \rho$. Hence, $\Psi = \mathbf{U} \Sigma \mathbf{V}^T$. We define $\tilde{\Psi} \in \mathbb{R}^{D \times \rho}$ by $\tilde{\Psi} = \mathbf{U} \Sigma$. Note that $\tilde{\Psi}$ may be computed using the singular value decomposition in $O(R^2 \times N)$ time and $O(R^2)$ space complexity. Note also that since the columns of \mathbf{V} are orthonormal and $\|\cdot\|_F$ is the Frobenius norm, for any matrix $\mathbf{A} \in \mathbb{R}^{\rho \times \rho}$ we have $\|\mathbf{V} \mathbf{A} \mathbf{V}^T\|_F = \|\mathbf{A}\|_F$. Hence, we have

$$\begin{aligned} C_\lambda(W, \tilde{\Psi}) &= (1 - \lambda) \frac{1}{M} \sum_{m=1}^M \left\| \left(\mathbf{W}_m \tilde{\Psi} \right)^T \left(\mathbf{W}_m \tilde{\Psi} \right) - \tilde{\Psi}^T \tilde{\Psi} \right\|^2 \\ &\quad + \lambda \left\| \frac{1}{M} \sum_{m=1}^M \left(\mathbf{W}_m \tilde{\Psi} \right)^T \left(\mathbf{W}_m \tilde{\Psi} \right) - \tilde{\Psi}^T \tilde{\Psi} \right\|^2 \\ &= (1 - \lambda) \frac{1}{M} \sum_{m=1}^M \left\| \Sigma \mathbf{U}^T \mathbf{W}_m^T \mathbf{W}_m \mathbf{U} \Sigma - \Sigma \mathbf{U}^T \mathbf{U} \Sigma \right\|^2 \\ &\quad + \lambda \left\| \Sigma \mathbf{U}^T \left(\frac{1}{M} \sum_{m=1}^M \mathbf{W}_m^T \mathbf{W}_m \right) \mathbf{U} \Sigma - \Sigma \mathbf{U}^T \mathbf{U} \Sigma \right\|^2 \\ &= (1 - \lambda) \frac{1}{M} \sum_{m=1}^M \left\| \mathbf{V} \Sigma \mathbf{U}^T \mathbf{W}_m^T \mathbf{W}_m \mathbf{U} \Sigma \mathbf{V}^T - \mathbf{V} \Sigma \mathbf{U}^T \mathbf{U} \Sigma \mathbf{V}^T \right\|^2 \\ &\quad + \lambda \left\| \mathbf{V} \Sigma \mathbf{U}^T \left(\frac{1}{M} \sum_{m=1}^M \mathbf{W}_m^T \mathbf{W}_m \right) \mathbf{U} \Sigma \mathbf{V}^T - \mathbf{V} \Sigma \mathbf{U}^T \mathbf{U} \Sigma \mathbf{V}^T \right\|^2 = C_\lambda(W, \Psi). \end{aligned}$$

We now prove Proposition 4.

Proposition 4. *Suppose we take $m \geq 1$, $\lambda > 0$, $M \geq 1$, fix \mathbf{W}_q for $q \neq m$, and let*

$$\mathbf{T}_m = \frac{M}{M - \lambda} \Psi^T \left(\mathbf{I}_D - \frac{\lambda}{M} \sum_{q \neq m} \mathbf{W}_q^T \mathbf{W}_q \right) \Psi.$$

Take $\mathbf{F}_m = R T_H(\mathbf{T}_m)$. Setting $\mathbf{W}_m = \mathbf{F}_m \Psi^\dagger$ minimises $C^\lambda(W, \Psi)$ with respect to \mathbf{W}_m , under the constraint that \mathbf{W}_q remains fixed for $q \neq m$, where Ψ^\dagger denotes the pseudo-inverse of Ψ .

Proposition 4 is a special case of the more general Proposition 10 which will be useful in proving Theorem 1.

Suppose we take $m \geq 1$, $\lambda > 0$, $M \geq 1$ and fix \mathbf{W}_q for $q \neq m$ and take $\epsilon > 0$. Fix a matrix $\mathbf{W}_m^{old} \in \mathbb{R}^{H \times D}$ and let $\Theta = (\mathbf{W}_m^{old} \Psi)^T (\mathbf{W}_m^{old} \Psi)$. Take $c = ((1 - \lambda) + \lambda/M + \epsilon)^{-1}$. Define $\mathbf{S}_m = \sum_{q \neq m} \mathbf{F}_q^T \mathbf{F}_q$ where $\mathbf{F}_q = \mathbf{W}_q \Psi$, and $\mathbf{Q} = \Psi^T \Psi$,

$$\mathbf{T} = c \left(\mathbf{Q} - \frac{\lambda}{M} \mathbf{S}_m + \epsilon \Theta \right).$$

Take $\mathbf{F}_m = R T_H(\mathbf{T})$. Then taking $\mathbf{W}_m = \mathbf{F}_m \Psi^\dagger$ minimises

$$C^\lambda(W, \Psi) + \frac{\epsilon}{M} \|(\mathbf{W}_m \Psi)^T (\mathbf{W}_m \Psi) - \Theta\|^2,$$

with respect to \mathbf{W}_m , under the constraint that \mathbf{W}_q remains fixed for $q \neq m$.

Proof (Proof of Proposition 10). We note that with $\mathbf{F}_q = \mathbf{W}_q \Psi$ for all $q = 1, \dots, M$ we have

$$\begin{aligned} C^\lambda(W, \Psi) + \frac{\epsilon}{M} \|(\mathbf{W}_m \Psi)^T (\mathbf{W}_m \Psi) - \Theta\|^2 \\ = (1 - \lambda) \frac{1}{M} \sum_{q=1}^M \|\mathbf{F}_q^T \mathbf{F}_q - \mathbf{Q}\|^2 + \lambda \left\| \frac{1}{M} \sum_{q=1}^M \mathbf{F}_q^T \mathbf{F}_q - \mathbf{Q} \right\|^2 + \frac{\epsilon}{M} \|\mathbf{F}_m^T \mathbf{F}_m - \Theta\|^2. \end{aligned}$$

By removing terms independent of \mathbf{F}_m and multiplying by M we see that this is equivalent to minimising

$$(1 - \lambda) \|\mathbf{F}_m^T \mathbf{F}_m - \mathbf{Q}\|^2 + \frac{\lambda}{M} \|\mathbf{F}_m^T \mathbf{F}_m - (M \mathbf{Q} - \mathbf{S}_m)\|^2 + \epsilon \|\mathbf{F}_m^T \mathbf{F}_m - \Theta\|^2.$$

Moreover, this expression is equal up to a constant term (independent of \mathbf{F}_m) to the following,

$$\left((1 - \lambda) + \frac{\lambda}{M} + \epsilon \right) \|\mathbf{F}_m^T \mathbf{F}_m\|^2 - 2 \left\| \mathbf{F}_m^T \mathbf{F}_m \left(\mathbf{Q} - \frac{\lambda}{M} \mathbf{S}_m + \epsilon \Theta \right) \right\|^2.$$

In addition, by adjusting by another constant and multiplying through by $c = ((1 - \lambda) + \frac{\lambda}{M} + \epsilon)^{-1}$ we see that this is equivalent to minimising $\|\mathbf{F}_m^T \mathbf{F}_m - \mathbf{T}\|^2$

where $\mathbf{T} = c \left(\mathbf{Q} - \frac{\lambda}{M} \mathbf{S}_m + \epsilon \mathbf{\Theta} \right)$. Thus, we minimise $\mathbf{F}_m \in \mathbb{R}^{H \times N}$ by taking $\mathbf{F}_m = R T_H(\mathbf{T})$. Now take $\mathbf{W}_m = \mathbf{F}_m \mathbf{\Psi}^\mathcal{Y}$. It suffices to show that $\mathbf{W}_m \mathbf{\Psi} = \mathbf{F}_m$.

To prove this consider the singular value decomposition $\mathbf{\Psi} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ where $\mathbf{\Psi} = \sum_{l=1}^{\rho} \sigma_l \mathbf{u}_l \mathbf{v}_l^T$, with orthonormal vectors $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_\rho] \in \mathbb{R}^{D \times \rho}$, $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_\rho) \in \mathbb{R}^{\rho \times \rho}$ and orthonormal vectors $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_\rho] \in \mathbb{R}^{N \times \rho}$. Here ρ denotes the rank of $\mathbf{\Psi}$ and for all $l = 1, \dots, \rho$ we have $\sigma_l > 0$. Thus, the pseudo-inverse is given by $\mathbf{\Psi}^\mathcal{Y} = \mathbf{V} \mathbf{\Sigma}^{-1} \mathbf{U}^T$.

Given any vector $\mathbf{z} \in \mathbb{R}^N$ which lies in the orthogonal complement to the span of $\hat{r}\mathbf{v}_1, \dots, \mathbf{v}_\rho$ (equivalently, in the null space of \mathbf{V}^T) we have $\mathbf{\Psi} \mathbf{z} = \mathbf{0}$. Noting that \mathbf{T} may be written

$$\mathbf{T} = c \mathbf{\Psi}^T \left(\mathbf{I}_D - \frac{\lambda}{M} \sum_{q \neq m} \mathbf{W}_q^T \mathbf{W}_q + \epsilon (\mathbf{W}_m^{\text{old}})^T (\mathbf{W}_m^{\text{old}}) \right) \mathbf{\Psi}.$$

Hence for \mathbf{z} in the null space of \mathbf{V}^T we have $\mathbf{T} \mathbf{z} = \mathbf{0}$. By equation 2 this implies $\mathbf{F}_m \mathbf{z} = \mathbf{0}$. Note also that $\mathbf{\Psi}^\mathcal{Y} \mathbf{\Psi} = \mathbf{V} \mathbf{V}^T$, so for \mathbf{z} in the null space of \mathbf{V}^T we have $\mathbf{\Psi}^\mathcal{Y} \mathbf{\Psi} \mathbf{z} = \mathbf{0}$ whereas for \mathbf{w} in the linear span of $\hat{r}\mathbf{v}_1, \dots, \mathbf{v}_\rho$ we have $\mathbf{\Psi}^\mathcal{Y} \mathbf{\Psi} \mathbf{w} = \mathbf{w}$. Hence, if we take $\mathbf{W}_m = \mathbf{F}_m \mathbf{\Psi}^\mathcal{Y}$, then for \mathbf{z} in the null space of \mathbf{V}^T we have

$$\mathbf{F}_m \mathbf{z} = \mathbf{0} = \mathbf{F}_m \mathbf{\Psi}^\mathcal{Y} \mathbf{\Psi} \mathbf{z}.$$

On the other hand, if \mathbf{w} is in the linear span of $\hat{r}\mathbf{v}_1, \dots, \mathbf{v}_\rho$ then we have $(\mathbf{\Psi}^\mathcal{Y} \mathbf{\Psi}) \mathbf{w} = \mathbf{w}$, so

$$\mathbf{F}_m \mathbf{w} = \mathbf{F}_m \mathbf{\Psi}^\mathcal{Y} \mathbf{\Psi} \mathbf{w}.$$

Thus, by taking an orthogonal decomposition we see if we take $\mathbf{W}_m = \mathbf{F}_m \mathbf{\Psi}^\mathcal{Y}$ then indeed

$$\mathbf{W}_m \mathbf{\Psi} = (\mathbf{F}_m \mathbf{\Psi}^\mathcal{Y}) \mathbf{\Psi} = \mathbf{F}_m.$$

This completes the proof of Proposition 10 which implies Proposition 4.

We now combine Propositions 5 and 10 to prove Theorem 1.

First we recall Algorithm 1 and the statement of Theorem 1.

Theorem 1. *Given $E \in \mathbb{N}$, let $\mathcal{F}^E \in F(H, M)$ denote the set obtained by training with Algorithm 1, for E epochs. Then for all $E \in \mathbb{N}$, $L_k^\lambda(\mathcal{F}^{E+1}, D) < L_k^\lambda(\mathcal{F}^E, D)$, unless \mathcal{F}^E is a critical point of $L_k^\lambda(\mathcal{F}, D)$, in which case $L_k^\lambda(\mathcal{F}^{E+1}, D) = L_k^\lambda(\mathcal{F}^E, D)$.*

Proof (Proof of Theorem 1). Firstly, by Proposition 5, the replacement of $\mathbf{\Psi}$ with $(R T_\rho(\mathbf{\Psi} \mathbf{\Psi}^T))^T$ at the start of Algorithm 1, while computationally essential, has no effect upon the loss function. Secondly, by Proposition 10, each update in the inner loop of Algorithm 1, minimises

$$f(\mathbf{W}_m) = N^2 L_k^\lambda(\mathcal{F}, D) + \frac{\epsilon}{M} \left\| (\mathbf{W}_m \mathbf{\Psi})^T (\mathbf{W}_m \mathbf{\Psi}) - \mathbf{\Theta} \right\|^2,$$

Inputs: A data set $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, a rank R empirical kernel map ψ , a number of modules M , a number of dimensions per module H , a diversity parameter λ and $\epsilon > 0$.
 Compute $\Psi = [\psi(\mathbf{x}_1), \dots, \psi(\mathbf{x}_N)]$;
 Update $\Psi = (RT_\rho(\Psi\Psi^T))^T$;
 Randomly initialise $\mathbf{F}_m \in \mathbb{R}^{H \times R}$ for $m = 1, \dots, M$;
 Compute $\mathbf{Q} = \Psi^T\Psi$ and $\mathbf{S} = \sum_{m=1}^M \mathbf{F}_m^T \mathbf{F}_m$;
 Compute $c = ((1 - \lambda) + \lambda/M + \epsilon)^{-1}$;
for $e = 1, \dots, E$ **do**
 for $m = 1, \dots, M$ **do**
 Compute $\mathbf{S}_m = \mathbf{S} - \mathbf{F}_m^T \mathbf{F}_m$;
 Compute $\mathbf{T} = c \cdot (\mathbf{Q} - (\lambda/M) \mathbf{S}_m + \epsilon \cdot \mathbf{F}_m^T \mathbf{F}_m)$;
 Update $\mathbf{F}_m = RT_H(\mathbf{T})$;
 Update $\mathbf{S} = \mathbf{S}_m + \mathbf{F}_m^T \mathbf{F}_m$;
 end
end
 Compute $\mathbf{W}_m = \mathbf{F}_m \Psi^\nu$ for $m = 1, \dots, M$;
Output: $\Phi = \{\mathbf{W}_m \cdot \psi\}_{m=1}^M$.

Algorithm 2: The module-by-module (MBM) algorithm.

where $\Theta = (\mathbf{W}_m^{old} \Psi)^T (\mathbf{W}_m^{old} \Psi)$. In particular, we have $f(\mathbf{W}_m) = f(\mathbf{W}_m^{old})$. Thus, by the definition of Θ we see that replacing \mathbf{W}_m^{old} cannot increase $L_k^\lambda(\Phi, D)$, as $\|(\mathbf{W}_m \Psi)^T (\mathbf{W}_m \Psi) - \Theta\|^2$ is minimised by $\mathbf{W}_m = \mathbf{W}_m^{old}$. Hence, we have $L_k^\lambda(\Phi^{E+1}, D) \leq L_k^\lambda(\Phi^E, D)$ for all $E \geq 1$.

Now suppose $L_k^\lambda(\Phi^{E+1}, D) = L_k^\lambda(\Phi^E, D)$. By Proposition 10, the combined cost $L_k^\lambda(\Phi, D) + \frac{\epsilon}{M} \|(\mathbf{W}_m \Psi)^T (\mathbf{W}_m \Psi) - \Theta\|^2$ is minimised at each stage, so we must have $(\mathbf{W}_m^{E+1} \Psi)^T (\mathbf{W}_m^{E+1} \Psi) = (\mathbf{W}_m^E \Psi)^T (\mathbf{W}_m^E \Psi)$ for every m (otherwise taking $\mathbf{W}_m = \mathbf{W}_m^E = \mathbf{W}_m^{old}$ would lead to a lower value of $L_k^\lambda(\Phi, D) + \frac{\epsilon}{M} \|(\mathbf{W}_m \Psi)^T (\mathbf{W}_m \Psi) - \Theta\|^2$). Hence, we see that $\mathbf{W}_m = \mathbf{W}_m^E$ is a critical point for

$$f(\mathbf{W}_m) = L_k^\lambda(\Phi, D) + \frac{\epsilon}{M} \|(\mathbf{W}_m \Psi)^T (\mathbf{W}_m \Psi) - \Theta\|^2,$$

Since \mathbf{W}_m^E is clearly a critical point for the map $\mathbf{W}_m \mapsto \|(\mathbf{W}_m \Psi)^T (\mathbf{W}_m \Psi) - \Theta\|^2$, it follows that $\mathbf{W}_m = \mathbf{W}_m^E$ is also a critical point for $L_k^\lambda(\Phi, D)$. Since this holds for each $m = 1, \dots, M$ we see that $\mathcal{W}^E = f \mathbf{W}_m^E g_{m=1}^E$ is a critical point for $L_k^\lambda(\Phi, D)$. Equivalently, Φ^E is a critical point for $L_k^\lambda(\Phi, D)$.

11 Proof of Theorem 2

In this section we prove Theorem 2 which connects the modular loss $L_k^\lambda(\Phi, D)$, defined in Section 3, with the drop-module loss $L_{k,p}^{\text{drop}}(\Phi, D)$ defined in Section 5.

Theorem 2. *The drop-module inner product loss at p is equivalent to the modular inner product loss at $\lambda = Mp/(1+p(M-1))$. To be precise, for $\Phi_0 \in F(H, M)$ we have*

$$\left. \frac{\partial L_{k,p}^{\text{drop}}(\Phi, D)}{\partial \Phi} \right|_{\Phi=\Phi_0} = p \left. \frac{\partial L_k^\lambda(\Phi, D)}{\partial \Phi} \right|_{\Phi=(\frac{p}{p/\lambda})\Phi_0}.$$

Proof (Proof of Theorem 2). Recall from Section 5 that

$$L_{k,p}^{\text{drop}}(\Phi, D) = \mathbb{E} \left(L_k(\Phi_\eta, D) : \eta \sim B(p)^M \right),$$

where

$$\begin{aligned} L_k(\Phi_\eta, D) &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N (\langle \Phi_\eta(\mathbf{x}_i), \Phi_\eta(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j))^2 \\ &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \left(\frac{1}{M} \sum_{m=1}^M \eta_m \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j) \right)^2. \end{aligned}$$

Now fix $q \in \{1, \dots, M\}$ and take a pair $\mathbf{x}_i, \mathbf{x}_j \in D$. Then,

$$\begin{aligned} & \frac{\partial}{\partial \varphi_q(\mathbf{x}_i)} \left(\frac{1}{M} \sum_{m=1}^M \eta_m \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j) \right)^2 \\ &= \frac{\partial \langle \varphi_q(\mathbf{x}_i), \varphi_q(\mathbf{x}_j) \rangle}{\partial \varphi_q(\mathbf{x}_i)} \frac{\partial}{\partial \langle \varphi_q(\mathbf{x}_i), \varphi_q(\mathbf{x}_j) \rangle} \left(\frac{1}{M} \sum_{m=1}^M \eta_m \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j) \right)^2 \\ &= \varphi_q(\mathbf{x}_j) \frac{2\eta_q}{M} \left(\frac{1}{M} \sum_{m=1}^M \eta_m \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j) \right). \end{aligned}$$

It follows that

$$\begin{aligned} & \frac{\partial}{\partial \varphi_q(\mathbf{x}_i)} \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M \eta_m \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j) \right)^2 \right] \\ &= \varphi_q(\mathbf{x}_j) \frac{2p}{M} \left(\frac{1}{M} \sum_{m=1}^M p \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle + \frac{1}{M} (1-p) \langle \varphi_q(\mathbf{x}_i), \varphi_q(\mathbf{x}_j) \rangle - k(\mathbf{x}_i, \mathbf{x}_j) \right). \end{aligned}$$

Hence, from the definition of $L_{k,p}^{\text{drop}}(\Phi, D)$, for each $\mathbf{x}_i \in D$ and $q \in \{1, \dots, M\}$ we have

$$\begin{aligned} & \frac{\partial}{\partial \varphi_q(\mathbf{x}_i)} \left(L_{k,p}^{\text{drop}}(\Phi, D) \right) \\ &= \frac{1}{N^2} \varphi_q(\mathbf{x}_i) \frac{2p}{M} \left(\frac{1}{M} \sum_{m=1}^M p \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_i) \rangle_i + \frac{1}{M} (1-p) \langle \varphi_q(\mathbf{x}_i), \varphi_q(\mathbf{x}_i) \rangle_i - k(\mathbf{x}_i, \mathbf{x}_i) \right) \\ & \quad + \frac{1}{N^2} \sum_{j=1}^N \varphi_q(\mathbf{x}_j) \frac{4p}{M} \left(\frac{1}{M} \sum_{m=1}^M p \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle_{ij} + \frac{1}{M} (1-p) \langle \varphi_q(\mathbf{x}_i), \varphi_q(\mathbf{x}_j) \rangle_{ij} - k(\mathbf{x}_i, \mathbf{x}_j) \right). \end{aligned}$$

On the other hand,

$$\begin{aligned} L_k^\lambda(\Phi, D) &= (1-\lambda) \frac{1}{M} \sum_{m=1}^M L_k(\varphi_m, D) + \lambda L_k(\bar{\Phi}, D) \\ &= \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \left((1-\lambda) \frac{1}{M} \sum_{m=1}^M (\langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle_{ij} - k(\mathbf{x}_i, \mathbf{x}_j))^2 \right. \\ & \quad \left. + \lambda \left(\frac{1}{M} \sum_{m=1}^M \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle_{ij} - k(\mathbf{x}_i, \mathbf{x}_j) \right)^2 \right). \end{aligned}$$

Hence for each $\mathbf{x}_i \in D$ and $q \in \{1, \dots, M\}$ we have

$$\begin{aligned} & \frac{\partial}{\partial \varphi_q(\mathbf{x}_i)} \left(L_k^\lambda(\Phi, D) \right) \\ &= \frac{1}{N^2} \varphi_q(\mathbf{x}_i) \frac{2}{M} \left(\frac{1}{M} \sum_{m=1}^M \lambda \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_i) \rangle_i + (1-\lambda) \langle \varphi_q(\mathbf{x}_i), \varphi_q(\mathbf{x}_i) \rangle_i - k(\mathbf{x}_i, \mathbf{x}_i) \right) \\ & \quad + \frac{1}{N^2} \sum_{j=1}^N \varphi_q(\mathbf{x}_j) \frac{4}{M} \left(\frac{1}{M} \sum_{m=1}^M \lambda \langle \varphi_m(\mathbf{x}_i), \varphi_m(\mathbf{x}_j) \rangle_{ij} + (1-\lambda) \langle \varphi_q(\mathbf{x}_i), \varphi_q(\mathbf{x}_j) \rangle_{ij} - k(\mathbf{x}_i, \mathbf{x}_j) \right). \end{aligned}$$

Thus, with $\lambda = Mp/(1+p(M-1))$ and $\Phi_0 = \{\varphi_m^0\}_{m=1}^M$ we have

$$\begin{aligned} & \frac{\partial}{\partial \varphi_q(\mathbf{x}_i)} \left(L_k^\lambda(\Phi, D) \right) \Big|_{\Phi = \left(\frac{p}{1+p(M-1)} \right) \Phi_0} \\ &= \frac{1}{N^2} \varphi_q^0(\mathbf{x}_i) \frac{2}{M} \left(\frac{1}{M} \sum_{m=1}^M p \langle \varphi_m^0(\mathbf{x}_i), \varphi_m^0(\mathbf{x}_i) \rangle_i + \frac{1}{M} (1-p) \langle \varphi_q^0(\mathbf{x}_i), \varphi_q^0(\mathbf{x}_i) \rangle_i - k(\mathbf{x}_i, \mathbf{x}_i) \right) \\ & \quad + \frac{1}{N^2} \sum_{j=1}^N \varphi_q^0(\mathbf{x}_j) \frac{4}{M} \left(\frac{1}{M} \sum_{m=1}^M p \langle \varphi_m^0(\mathbf{x}_i), \varphi_m^0(\mathbf{x}_j) \rangle_{ij} + \frac{1}{M} (1-p) \langle \varphi_q^0(\mathbf{x}_i), \varphi_q^0(\mathbf{x}_j) \rangle_{ij} - k(\mathbf{x}_i, \mathbf{x}_j) \right). \end{aligned}$$

Combining this with the above computation of $\partial \left(L_{k,p}^{\text{drop}}(\Phi, D) \right) / \partial \varphi_q(\mathbf{x}_i)$ gives

$$\frac{\partial \left(L_{k,p}^{\text{drop}}(\Phi, D) \right)}{\partial \varphi_q(\mathbf{x}_i)} \Big|_{\Phi=\Phi_0} = p \frac{\partial \left(L_k^\lambda(\Phi, D) \right)}{\partial \varphi_q(\mathbf{x}_i)} \Big|_{\Phi=(\frac{p}{\lambda}) \Phi_0}.$$

Since this holds for all $\mathbf{x}_i \in D$ and $q \in \{1, \dots, M\}$, this completes the proof of the theorem.

12 Experimental results

12.1 The data sets

We used the following eight data sets in our experiments.

Background The MNIST + Background images data set from Larochelle et al. [17]: 10 classes, 784 features per example, 12000 training examples, 50000 test examples.

Convex The Convex data set from Larochelle et al. [17]: 2 classes, 784 features per example, 8000 training examples, 50000 test examples.

MNIST The original MNIST data set as used in Larochelle et al. [17]: 10 classes, 784 features per example, 12000 training examples, 50000 test examples.

NUS Wide The NUS-Wide-Object data set from Chua et al. [5] with all extracted features concatenated: 31 classes, 639 features per example, 17928 training examples, 12072 test examples.

NVDR The NVDR data set from Wu et al. [28] with all extracted features concatenated: 24 classes, 418 features per example, 6438 training examples, 6429 test examples.

Random The MNIST + Random background data set from Larochelle et al. [17]: 10 classes, 784 features per example, 12000 training examples, 50000 test examples.

Rectangles The Rotated MNIST digits data set from Larochelle et al. [17]: 2 classes, 784 features per example, 1200 training examples, 50000 test examples.

Rotations The Rotated MNIST digits data set from Larochelle et al. [17]: 10 classes, 784 features per example, 12000 training examples, 50000 test examples.

12.2 Test time and the number of dimensions per module (H)

The following figures show the test time for both image retrieval and classification as a function of the number of dimensions per module H . Note that $H \cdot M = 300$ throughout, so M is inversely proportional to H . Test time for classification increases with H as this corresponds to the number of dimensions per module, and the predictions for each module may be made in parallel. The aggregation of predictions corresponds to taking the most common label from a list, which

adds little relative overhead. In contrast, the test time for image retrieval is a non-monotonic function of the number of dimensions per module - it typically falls for very small H , before rising again as H grows. Recall from Section 6.2 that the computation of nearest neighbours for image retrieval consists of two steps. First a collection κ of nearest neighbours $\mathcal{I}_{\kappa,n}^{\varphi_m}(\mathbf{x})$ is computed in each H -dimensional space separately. Second, the sets $\mathcal{I}_{\kappa,n}^{\varphi_m}(\mathbf{x})$ are combined and a subset $\mathcal{I}_{\kappa,n}^{\Phi}(\mathbf{x}) = \bigcup_{m=1}^M \mathcal{I}_{\kappa,n}^{\varphi_m}(\mathbf{x})$ of size κ is extracted so as to minimise the average squared distance $(1/M) \sum_{m=1}^M \|\varphi_m(\mathbf{x}_q) - \varphi_m(\mathbf{x})\|_2^2$. The first step is easily distributable over the M modules. Hence, the time taken is monotonically increasing in H . The second stage is not straightforwardly distributable and the size of the search space $\bigcup_{m=1}^M \mathcal{I}_{\kappa,n}^{\varphi_m}(\mathbf{x})$ will typically grow with M . This adds an additional overhead which accounts for fact that image retrieval test time can typically be seen to fall with H when H is small (ie. when M is large). Note that the size of this additional overhead is bounded by $M \kappa$, where as the time taken by the first stage, consisting of multiple nearest neighbour searches done in parallel, increases with the size of the training set. Hence, the effect of the additional overhead for small values of H is at its most dramatic when the size of the training set is small, such as for the Rectangles data set, where the training set consists of just 1200 images.

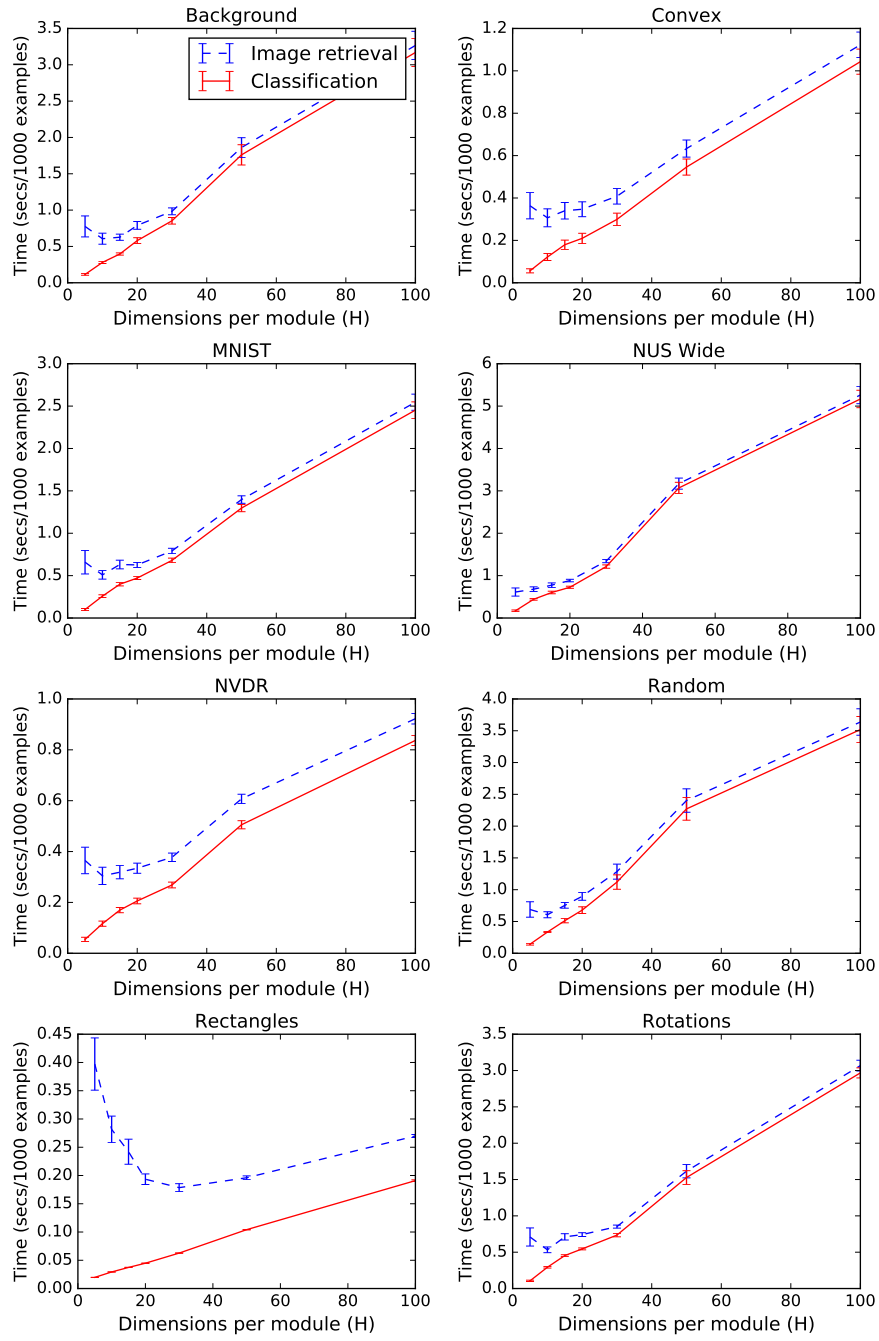


Fig. 6. Classification accuracy as a function of H (See Section 6.2).

12.3 Image retrieval precision and the number of dimensions per module (H)

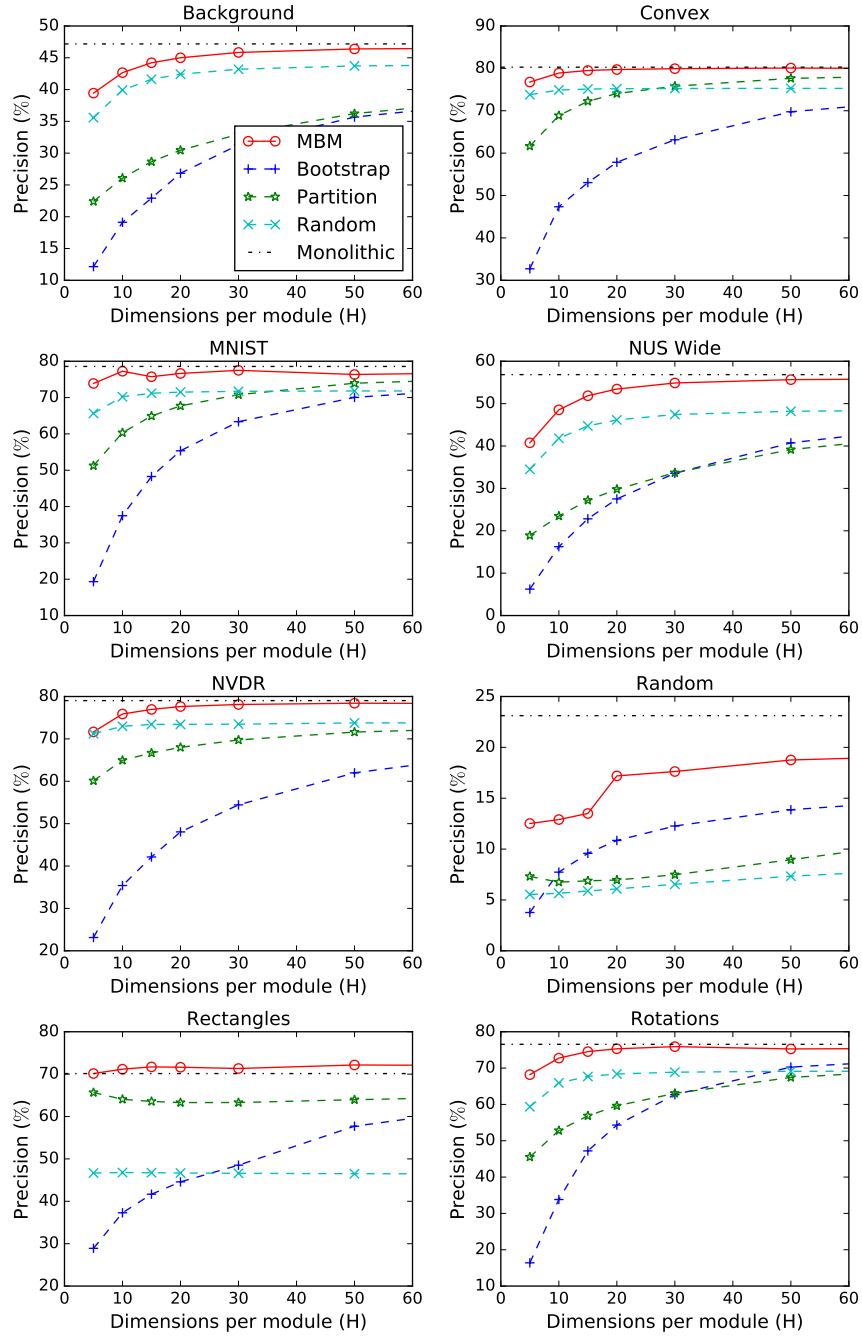


Fig. 7. Precision as a function of H (See Section 6.2).

12.4 Image retrieval precision and the diversity parameter

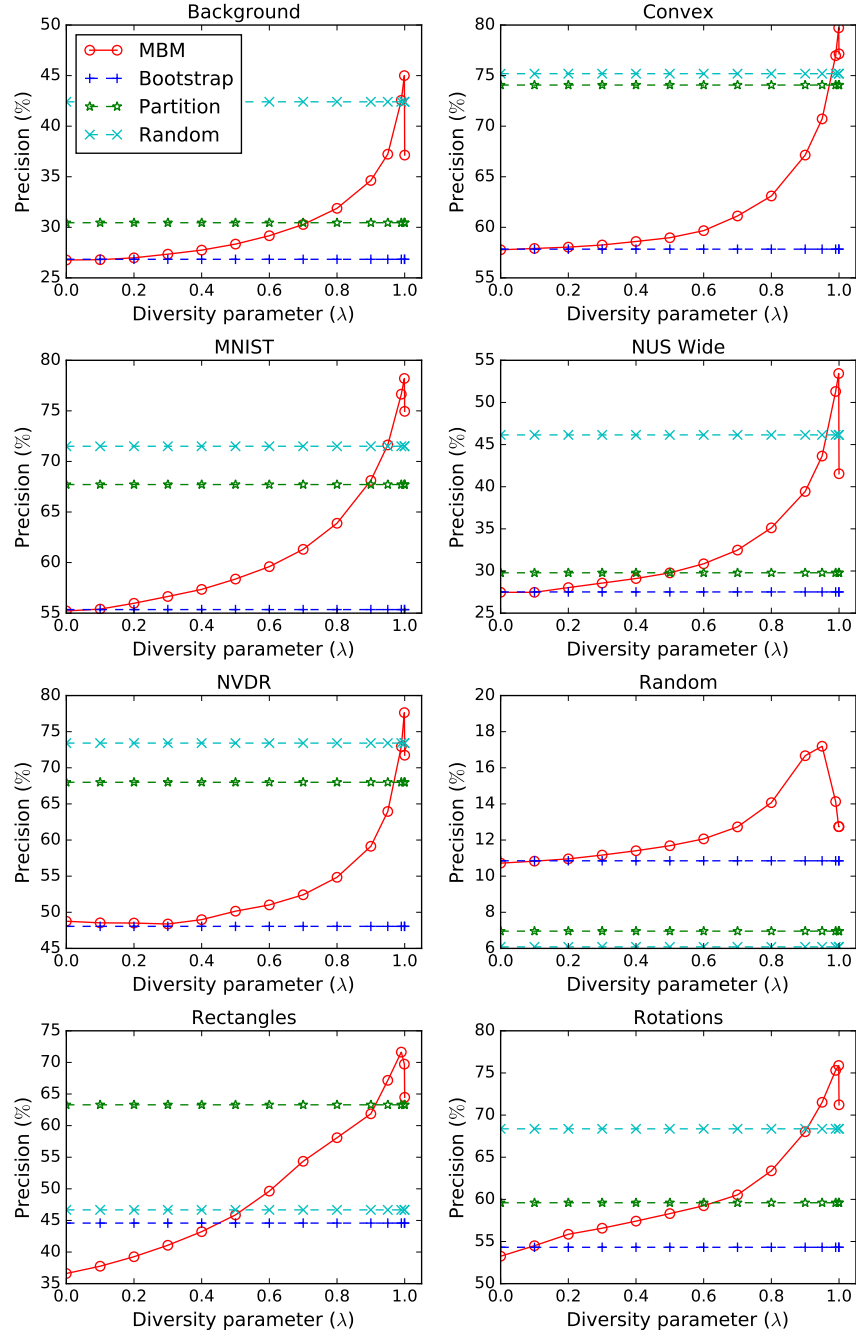


Fig. 8. Precision as a function of the diversity parameter λ (See Section 6.2).

12.5 Classification accuracy and the number of dimensions per module

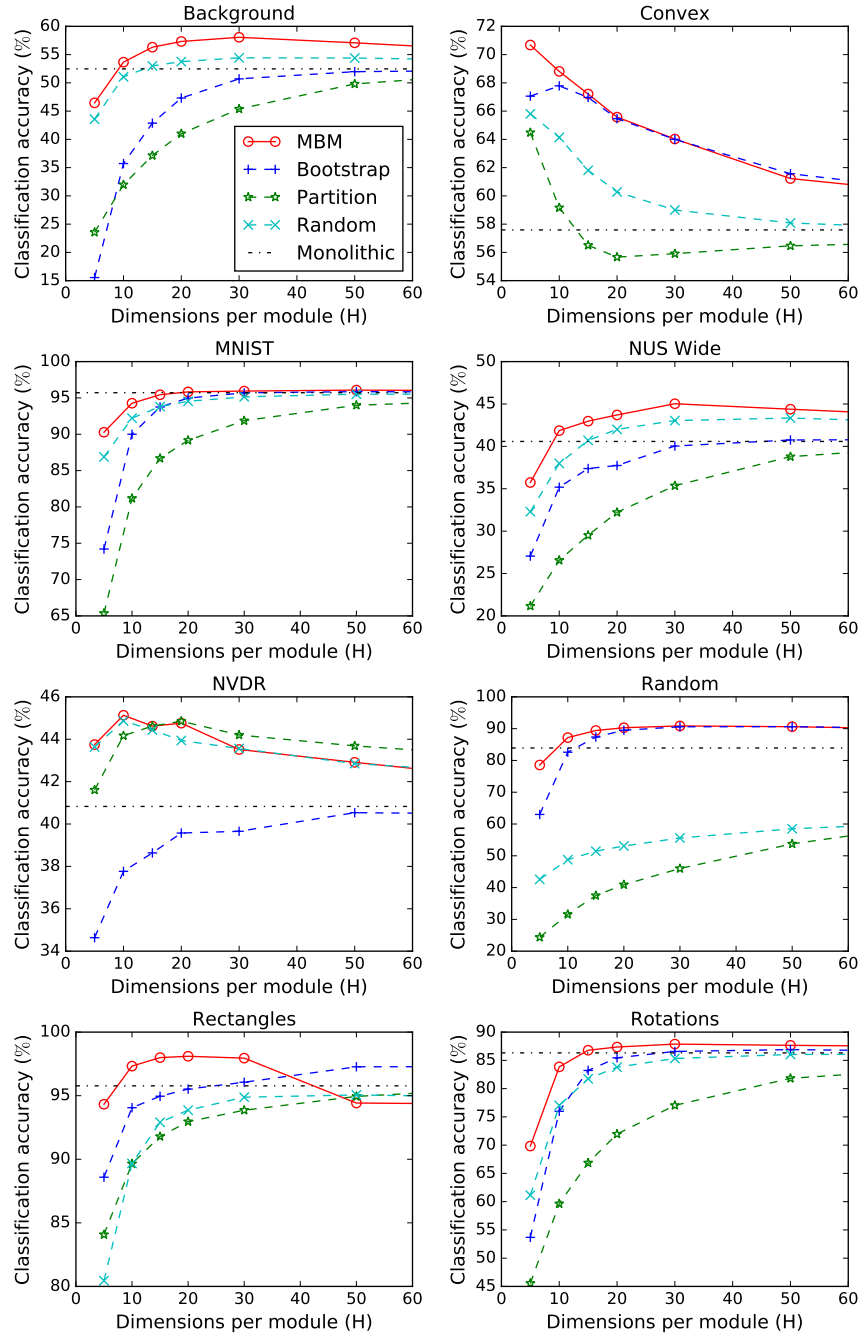


Fig. 9. Classification accuracy as a function of H (See Section 6.2).

12.6 Classification accuracy and the diversity parameter

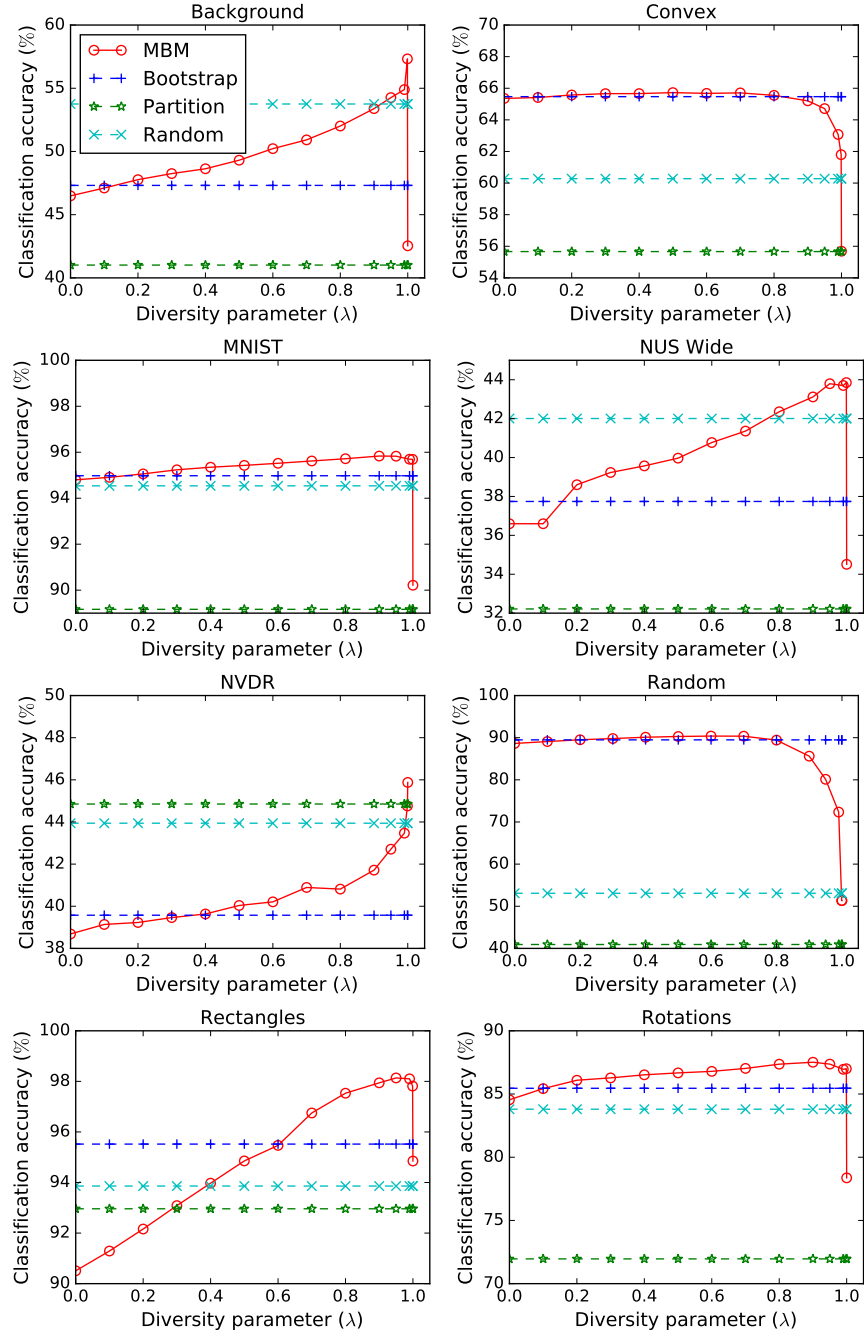


Fig. 10. Classification accuracy as a function of H (See Section 6.2).